

# LOGIC IN COMPUTER SCIENCE

by

Benji MO

Some people are always critical of vague statements.  
I tend rather to be critical of precise statements.  
They are the only ones which can  
correctly be labeled wrong.  
– Raymond Smullyan

August 2020

Supervisor: Professor Hantao Zhang

## TABLE OF CONTENTS

	Page
LIST OF FIGURES . . . . .	viii
CHAPTER	
1 Introduction to Logic . . . . .	1
1.1 Logic is Everywhere . . . . .	2
1.1.1 Statement or Proposition . . . . .	2
1.1.2 A Brief History of Logic . . . . .	4
1.2 Logical Fallacies in Arguments . . . . .	5
1.2.1 Formal Fallacies . . . . .	5
1.2.2 Informal Fallacies . . . . .	7
1.3 A Brief Review of Mathematical Logic . . . . .	13
1.3.1 Set Theory . . . . .	13
1.3.2 Model Theory . . . . .	21
1.3.3 Proof Theory . . . . .	22
1.3.4 Computability Theory . . . . .	25
1.4 Exercise Problems . . . . .	29
2 Propositional Logic . . . . .	32
2.1 Syntax . . . . .	32
2.1.1 Logical Operators . . . . .	32
2.1.2 Formulas . . . . .	34
2.2 Semantics . . . . .	36
2.2.1 Interpretations . . . . .	36
2.2.2 Models, Satisfiability, and Validity . . . . .	38
2.2.3 Equivalence . . . . .	40
2.2.4 Entailment . . . . .	43
2.2.5 Theorem Proving and the SAT Problem . . . . .	45
2.3 Normal Forms . . . . .	46
2.3.1 Negation Normal Form (NNF) . . . . .	48
2.3.2 Conjunctive Normal Form (CNF) . . . . .	49
2.3.3 Disjunctive Normal Form (DNF) . . . . .	51
2.3.4 Full DNF and Full CNF from Truth Table . . . . .	53
2.3.5 Binary Decision Diagram (BDD) . . . . .	54
2.4 Optimization Problems . . . . .	59
2.4.1 Minimum Set of Operators . . . . .	59

2.4.2	Logic Minimization . . . . .	61
2.4.3	Maximum Satisfiability . . . . .	69
2.5	Using Propositional Logic . . . . .	70
2.5.1	Bitwise Operators . . . . .	70
2.5.2	Specify Problems in Propositional Logic . . . . .	72
2.6	Exercise Problems . . . . .	77
3	Proof Procedures for Propositional Logic . . . . .	85
3.1	Semantic Tableau . . . . .	87
3.1.1	Tableau: A Tree Structure for DNF . . . . .	88
3.1.2	$\alpha$ -Rules and $\beta$ -Rules . . . . .	89
3.2	Deductive Systems . . . . .	92
3.2.1	Inference Rules and Proofs . . . . .	92
3.2.2	Hilbert Systems . . . . .	94
3.2.3	Natural Deduction . . . . .	96
3.2.4	Inference Graphs . . . . .	99
3.3	Resolution . . . . .	100
3.3.1	Resolution Rule . . . . .	100
3.3.2	Resolution Strategies . . . . .	103
3.3.3	Preserving Satisfiability . . . . .	105
3.3.4	Completeness of Resolution . . . . .	108
3.3.5	A Resolution-based Decision Procedure . . . . .	110
3.3.6	Clause Deletion Strategies . . . . .	112
3.4	Boolean Constraint Propagation (BCP) . . . . .	115
3.4.1	BCP: a Simplification Procedure . . . . .	115
3.4.2	BCP: a Decision Procedure for Horn Clauses . . . . .	117
3.4.3	Unit Resolution versus Input Resolution . . . . .	118
3.4.4	Head/Tail Literals for <i>BCP</i> . . . . .	120
3.5	Exercise Problems . . . . .	124
4	Propositional Satisfiability . . . . .	128
4.1	The DPLL Algorithm . . . . .	129
4.1.1	Recursive Version of <i>DPLL</i> . . . . .	129
4.1.2	All-SAT and Incremental SAT Solvers . . . . .	132
4.1.3	<i>BCP<sub>w</sub></i> : Implementation of Watch Literals . . . . .	133
4.1.4	Iterative Implementation of <i>DPLL</i> . . . . .	136
4.2	Conflict-Driven Clause Learning (CDCL) . . . . .	138
4.2.1	Generating Clauses from Conflicting Clauses . . . . .	139
4.2.2	DPLL with CDCL . . . . .	140
4.2.3	Unsatisfiable Cores . . . . .	143
4.2.4	Random Restart . . . . .	144
4.2.5	Branching Heuristics for DPLL . . . . .	145
4.3	Use of SAT Solvers . . . . .	147

4.3.1	Specify SAT Instances in DIMACS Format . . . . .	148
4.3.2	Sudoku Puzzle . . . . .	149
4.3.3	Latin Square Problems . . . . .	151
4.3.4	Graph Problems . . . . .	152
4.4	Local Search Methods and MaxSAT . . . . .	154
4.4.1	Local Search Methods for SAT . . . . .	154
4.4.2	2SAT vesus Max2SAT . . . . .	157
4.5	Maximum Satisfiability . . . . .	159
4.5.1	Weight MaxSAT and Hybrid MaxSAT . . . . .	159
4.5.2	The Branch-and-Bound Algorithm . . . . .	161
4.5.3	Simplification Rules and Lower Bound . . . . .	163
4.5.4	Use of Hybrid MaxSAT Solvers . . . . .	166
4.6	Exercise Problems . . . . .	168
5	First Order Logic . . . . .	172
5.1	Syntax of First Order Languages . . . . .	172
5.1.1	Terms and Formulas . . . . .	172
5.1.2	The Quantifiers . . . . .	176
5.1.3	Unsorted and Many-Sorted Logics . . . . .	178
5.2	Semantics . . . . .	180
5.2.1	Interpretation . . . . .	180
5.2.2	Models, Satisfiability, and Validity . . . . .	184
5.2.3	Entailment and Equivalence . . . . .	186
5.3	Proof Methods . . . . .	189
5.3.1	Semantic Tableau . . . . .	189
5.3.2	Natural Deduction . . . . .	192
5.4	Conjunctive Normal Form . . . . .	192
5.4.1	Prenex Normal Form . . . . .	193
5.4.2	Skolemization . . . . .	195
5.4.3	Skolemizing Non-Prenex Formulas . . . . .	198
5.4.4	Clausal Form . . . . .	199
5.4.5	Herbrand Models for CNF . . . . .	200
5.5	Exercise Problems . . . . .	202
6	Unification and Resolution . . . . .	206
6.1	Unification . . . . .	206
6.1.1	Substitutions and Unifiers . . . . .	206
6.1.2	Combining Substitutions . . . . .	208
6.1.3	Rule-Based Unification . . . . .	209
6.1.4	Practically Linear Time Unification . . . . .	212
6.2	Resolution . . . . .	219
6.2.1	The Resolution and Factoring Rules . . . . .	219
6.2.2	A Refutational Proof Procedure . . . . .	221

6.3	Ordered Resolution . . . . .	223
6.3.1	Simplification Orders . . . . .	223
6.3.2	Completeness of Ordered Resolution . . . . .	228
6.4	Prover9: A Resolution Theorem Prover . . . . .	231
6.4.1	Input Formulas to Prover9 . . . . .	231
6.4.2	Inference Rules and Options . . . . .	234
6.4.3	Simplification Orders in Prover9 . . . . .	237
6.5	Exercise Problems . . . . .	238
7	Equational Logic . . . . .	241
7.1	Equality of Terms . . . . .	241
7.1.1	Axioms of Equality . . . . .	242
7.1.2	Semantics of “=” . . . . .	243
7.1.3	Theory of Equations . . . . .	245
7.2	Rewrite Systems . . . . .	247
7.2.1	Rewrite Rules . . . . .	247
7.2.2	Termination of Rewrite Systems . . . . .	248
7.2.3	Confluence of Rewriting . . . . .	249
7.2.4	The Knuth-Bendix Completion Procedure . . . . .	251
7.3	Inductive Theorem Proving . . . . .	257
7.3.1	Inductive Theorems . . . . .	257
7.3.2	Structural Induction . . . . .	258
7.3.3	Induction on Two Variables . . . . .	260
7.3.4	Multi-Sort Algebraic Specifications . . . . .	261
7.4	Resolution with Equality . . . . .	264
7.4.1	Paramodulation . . . . .	264
7.4.2	Simplification Rules . . . . .	266
7.4.3	Equality in Prover9 . . . . .	268
7.5	Mace4: Finite Model Finding in FOL . . . . .	270
7.5.1	Use of Mace4 . . . . .	270
7.5.2	Finite Model Finding by SAT Solvers . . . . .	274
7.6	Exercise Problems . . . . .	276
8	Prolog: Programming in Logic . . . . .	278
8.1	Prolog’s Working Principle . . . . .	278
8.1.1	Horn Clauses in Prolog . . . . .	278
8.1.2	Resolution Proofs in Prolog . . . . .	280
8.1.3	A Goal-Reduction Procedure . . . . .	282
8.2	Prolog’s Data Types . . . . .	287
8.2.1	Atoms, Numbers, and Variables . . . . .	287
8.2.2	Compound Terms and Lists . . . . .	289
8.2.3	Popular Predicates over Lists . . . . .	290
8.2.4	Sorting Algorithms in Prolog . . . . .	293
8.3	Recursion in Prolog . . . . .	294

8.3.1	Program Termination . . . . .	295
8.3.2	Focused Recursion . . . . .	297
8.3.3	Tail Recursions . . . . .	298
8.3.4	A Prolog program for <i>N</i> -Queen Puzzle . . . . .	299
8.4	Beyond Clauses and Logic . . . . .	300
8.4.1	The Cut Operator ! . . . . .	301
8.4.2	Negation as Failure . . . . .	302
8.4.3	Beyond Clauses . . . . .	304
8.5	Exercise Problems . . . . .	305
9	Undecidable Problems in First-Order Logic . . . . .	307
9.1	Turing Machines . . . . .	308
9.1.1	Formal Definition of Turing Machines . . . . .	309
9.1.2	High-level Description of Turing Machines . . . . .	312
9.1.3	Recognizable vs Decidable . . . . .	313
9.2	Decidability of Problems . . . . .	313
9.2.1	Encoding of Decision Problems . . . . .	313
9.2.2	Decidable Problems . . . . .	315
9.2.3	Undecidable Problems . . . . .	317
9.3	Turing Completeness . . . . .	320
9.3.1	Turing Completeness of Prolog . . . . .	321
9.3.2	Turing Completeness of Rewrite Systems . . . . .	323
9.4	Exercise Problems . . . . .	324
10	Hoare Logic . . . . .	325
10.1	Hoare Triples . . . . .	326
10.1.1	Hoare Rules . . . . .	327
10.1.2	Examples of Formal Verification . . . . .	332
10.1.3	Hoare Rule for Function Calls . . . . .	334
10.1.4	Partial and Total Correctness . . . . .	334
10.2	Automated Generation of Assertions . . . . .	335
10.2.1	Verification Conditions . . . . .	337
10.2.2	Proof of Theorem 10.2.4 . . . . .	338
10.2.3	Implementing <i>vc</i> in Prolog . . . . .	342
10.3	Obtaining Good Loop Invariants . . . . .	346
10.3.1	Invariants from Generalizing Postconditons . . . . .	347
10.3.2	Program Synthesis from Invariants . . . . .	349
10.3.3	Choosing A Good Language for Assertions . . . . .	353
10.4	Exercise Problems . . . . .	353

## LIST OF FIGURES

Figure	Page
2.1.1 The formula tree for $(\neg p \wedge q) \rightarrow (p \wedge (q \vee \neg r))$ . . . . .	35
2.3.1 The BDD of INF $ite(x_1, ite(x_2, 1, 0), ite(x_2, ite(x_3, 1, 0), ite(x_3, 0, 1)))$ derived from $\overline{x_1 x_2 x_3} \vee x_1 x_2 \vee x_2 x_3$ . . . . .	57
2.3.2 The first BDD uses $a_1 > b_1 > a_2 > b_2 > a_3 > b_3$ and the second BDD uses $a_1 > a_2 > a_3 > b_1 > b_2 > b_3$ for the same formula $(a_1 \wedge b_1) \vee (a_2 \wedge b_2) \vee (a_3 \wedge b_3)$ . . . . .	57
2.4.1 K-map for $f(x, y) = x\bar{y} + \bar{x}y + \bar{x}\bar{y}$ . . . . .	63
2.4.2 The K-map for obtaining $f(A, B, C) = AB + \bar{A}C$ . . . . .	64
2.4.3 The K-map for $f(P, Q, R, S) = m_0 + m_2 + m_5 + m_7 + m_8 + m_{10} +$ $m_{13} + m_{15}$ . The simplified formula is $f = \overline{QS} + QS$ . . . . .	65
2.4.4 The K-map for $f(A, B, C) = \overline{ABC} + \overline{A\bar{B}C} + \overline{AB\bar{C}} + \overline{A\bar{B}C}$ . The simplified CNF is $f = (A + B + C)(\bar{A} + \bar{B})(\bar{B} + \bar{C})$ . . . . .	66
3.1.1 The tableaux for $A = p \wedge (\neg q \vee \neg p)$ and $B = (p \vee q) \wedge (\neg p \wedge \neg q)$ . .	89
3.2.1 The proof tree of $\neg q$ . . . . .	93
4.1.1 The result of BCP is shown in the box near each node. . . . .	130
4.3.1 A typical example of Sudoku puzzle and its solution. . . . .	149
4.4.1 The implication graph for 2CNF C and the SCC. . . . .	158
4.5.1 Finding a maximal clique in a graph. . . . .	161
6.1.1 The term graph for $t = f(f(f(a, x_1), x_2), x_3)$ and $s = f(x_3, f(x_2, f(x_1, a)))$ . 214	214
8.1.1 $\vee$ -nodes in the $\wedge$ - $\vee$ tree for the query <code>dfs(a, N)</code> . . . . .	284

## CHAPTER 1

### INTRODUCTION TO LOGIC

When you see the word “logic”, what would you think of it? In the *Merriam-Webster* dictionary, the word *logic* is defined as follows:

1. (a) a science that deals with the principles and criteria of validity of inference and demonstration: the science of the formal principles of reasoning  
e.g., a professor of *logic*
- (b) a branch or variety of logic  
e.g., Boolean *logic*
- (c) a branch of semiotics, especially: *syntactics*
- (d) the formal principles of a branch of knowledge  
e.g., the *logic* of grammar.
2. (a) a particular mode of reasoning viewed as valid or faulty.  
e.g., She spent a long time explaining the situation, but he failed to see her *logic*.
- (b) *Relevance, Propriety*  
e.g., Could not understand the *logic* of such an action.
- (c) interrelation or sequence of facts or events when seen as inevitable or predictable  
e.g., By the *logic* of events, anarchy leads to dictatorship.
- (d) the arrangement of circuit elements (as in a computer) needed for computation

In this book, we will take “logic” in the sense of 1.(a): Logic is the science of the formal principles of reasoning. When these principles are introduced, we will use “logic” in the sense of 1.(b) and introduce several branches of logic, such as *propositional logic*, *equational logic*, *first-order logic*, etc. These logics belong to the domain of *mathematical logic* (formerly known as *symbolic logic*).

## 1.1 Logic is Everywhere

Logic has been called the *calculus of computer science*, because logic plays a fundamental role in computer science, starting from the construction of computers, to the computing devices beyond our ability of construction, such as computer architecture (digital gates, hardware verification), software engineering (specification, verification), programming languages (semantics, type theory, abstract data types, object-oriented programming), databases (relational algebra), artificial intelligence (automated theorem proving, knowledge representation), algorithms and theory of computation (complexity, computability), etc.

Logic also plays important roles in other fields, such as mathematics and philosophy. In mathematics, logic includes both the mathematical study of logic and the applications of formal logic to other areas of mathematics. The unifying themes in mathematical logic include the study of the expressive power of formal systems and the deductive power of formal proof systems. Mathematical logic is often divided into the fields of (a) set theory, (b) model theory, (c) proof theory, and (d) computability theory. These areas share basic results on logic.

Symbolic logic in the late nineteenth century and mathematical logic in the twentieth are also called *formal logic*. Topics traditionally treated by logic not being part of formal logic have tended to be termed either *philosophy of logic* or *philosophical logic*. Philosophical logic dealt with arguments in the natural language used by humans.

**Example.** Given the premises that (a) “all men are mortal” and (b) “Socrates is a man”, we may draw the conclusion that (c) “Socrates is mortal”, by the inference rule of *Modus Ponens*.

Philosophical logic is the investigation, critical analysis and intellectual reflection on issues arising in logic and is the branch of studying questions about reference, predication, identity, truth, quantification, existence, entailment, modality, and necessity.

### 1.1.1 Statement or Proposition

Logic comes from natural languages as most sentences in natural languages are statements or propositions. A *proposition* is a statement or assertion that expresses a judgment or opinion. Here statement, assertion, and proposition are synonyms. Here are some examples:

- That apple is delicious.

- Today is either Monday or Tuesday.
- He has good memory and runs fast.

Every statement can be either true or false. True or false are called the *truth values* of a statement. Some sentences in natural languages are not statements, such as commands or exclamatory sentences. For example, “Run fast!” *Coach shouts*. The first part of the sentence, “Run fast!”, is not a statement; the second part, “Coach shouts”, is a statement. In fact, a sentence is a statement if and only if it has a truth value.

In natural languages, we can combine, or relate statements with words such as “not” (negation), “and” (conjunction), “or” (disjunction), “implies” (implication), “if-then”, etc. That is, a statement can be obtained from other statements by these words. In logic, these words are called *logical operators*, or equivalently, *logical connectives*. A statement is *composed* if it can be expressed as a composition of several simpler statements; otherwise it is *simple*. In the above three examples of statements, the first statement is simple, but the other two are composed. That is, “today is either Monday or Tuesday” is the composition of “today is Monday” and “today is Tuesday”, using the logical operator “or”. We often regard “that apple is not delicious” as a composed statement: It is the negation of a simple statement.

Logical connectives are indispensable for expressing the relationship between statements. For example, the following statement is a composition of several simple statements.

(\*) Since I don’t have a good memory, I don’t know that today is either Monday or Tuesday.

To see this clearly, we need to introduce symbols, such as  $p, q, r, s$ , and  $t$ , to denote them:

1.  $p$ : I have a good memory.
2.  $q$ : Today is Monday.
3.  $r$ : Today is Tuesday.
4.  $s$ :  $q$  or  $r$ .
5.  $t$ : I know  $s$ .

To express that the statement (\*) is composed, let us use these symbols for logical operators:  $\neg$  for negation,  $\vee$  for “OR”,  $\wedge$  for “AND”, and  $\rightarrow$  for implication. Then (\*) becomes the following formula:

$$(\neg p) \rightarrow (\neg t)$$

Here each of the symbols  $p, q, r, s,$  and  $t$  is called a *propositional variable*, which denotes a statement, either simple or composed. Naturally, these propositional variables, like statements, can take on only the truth values, i.e., true and false.

Propositional variables are also called *Boolean variables* after their inventor, the nineteenth century mathematician George Boole. Boolean logic includes any logic in which the considered truth values are *true* and *false*. The study of propositional variables with logic operators is called *Propositional logic*, which is the most simple one in the family of Boolean logics. On the other hand, probability logic is not a Boolean logic because probability values are used to represent various degrees of the truth values. Questions regarding degrees of truth or the sense-content of sentences are ignored in Boolean logic. Despite this simplification, or indeed because of it, such a method is scientifically successful. One does not even have to know exactly what the truth values true and false actually are.

### 1.1.2 A Brief History of Logic

Sophisticated theories of logic were developed in many cultures, including China, India, Greece and the Islamic world. Greek methods, particularly Aristotelian logic (or term logic) as found in the *Organon*, found wide application and acceptance in Western science and mathematics for millennia.

Aristotle (384 – 322 BC) was the first logician to attempt a systematic analysis of logical syntax, of noun (or term), and of verb. He demonstrated the principles of reasoning by employing variables to show the underlying logical form of an argument. He sought relations of dependence which characterize necessary inference, and distinguished the validity of these relations, from the truth of the premises. He was the first to deal with the principles of contradiction and excluded middle in a systematic way. Aristotle has had an enormous influence in Western thought. and he developed the theory of the *syllogism*, where three important principles are applied for the first time in history: the use of symbols, a purely formal treatment, and the use of an axiomatic system. Aristotle also developed the theory of fallacies, as a theory of non-formal logic.

Christian and Islamic philosophers such as Boethius (died 524), Ibn Sina (Avicenna, died 1037) and William of Ockham (died 1347) further developed Aristotle’s

logic in the Middle Ages, reaching a high point in the mid-fourteenth century, with Jean Buridan (1301-1358/62 AD).

In the 18th century, attempts to treat the operations of formal logic in a symbolic or algebraic way had been made by philosophical mathematicians, including Leibniz and Lambert, but their works remained isolated and little known.

In the beginning of the 19th century, logic was studied with rhetoric, through the syllogism, and with philosophy. Mathematical logic emerged in the mid-19th century as a field of mathematics independent of the traditional study of logic. The development of the modern “symbolic” or “mathematical” logic during this period by the likes of Boole, Frege, Russell, and Peano is the most significant in the two-thousand-year history of logic, and is arguably one of the most important and remarkable events in human intellectual history.

## 1.2 Logical Fallacies in Arguments

In this section, we will give a brief introduction to logical fallacies, which are the subjects of philosophical logic. This introduction is based on *Wikipedia's* page on this topic. We feel that the knowledge on fallacies will help us to write better reports and papers by avoiding fallacies.

In philosophical logic, an *argument* has similar meaning as a proof in formal logic, as the argument consists of premises and conclusions. A *fallacy* is the use of invalid or otherwise faulty reasoning in the construction of an argument. A fallacious argument may be deceptive by appearing to be better than it really is. Some fallacies are committed intentionally to manipulate or persuade by deception, while others are committed unintentionally due to carelessness or ignorance.

Fallacies are commonly divided into “formal” and “informal”. A formal fallacy can be expressed neatly in a formal logic, such as propositional logic, while an informal fallacy cannot be expressed in a formal logic.

### 1.2.1 Formal Fallacies

In philosophical logic, a formal fallacy is also called *deductive fallacy*, *logical fallacy*, or *non sequitur* (Latin for “it does not follow”). This is a pattern of reasoning rendered invalid by a flaw in its logical structure.

**Example.** Give the premises that (a) my car is some car, and (b) some cars are red, we draw the conclusion that (c) my car is red.

This is a typical example of a conclusion that does not follow logically from

premises or that is based on irrelevant data. Here are some common logical fallacies:

- **Affirming the consequent:** Any argument with the invalid structure of: If A then B. B, therefore A.

**Example.** If I get a B on the test, then I will get the degree. I got the degree, so it follows that I must have received a B. In fact, I got an A.

- **Denying the antecedent:** Any argument with the invalid structure of: If A then B. Not A, therefore not B.

**Example.** If it's a dog then it's a mammal. It's not a dog, so it must not be a mammal. In fact, it's a cat.

- **Affirming a disjunct:** Any argument with the invalid structure of: If A or B. A, therefore not B.

**Example.** I am working or I am at home. I am working, so I must not be at home. In fact, I am working at home.

- **Denying a conjunct:** Any argument with the invalid structure of: It is not the case that both A and B. Not A, therefore B.

**Example.** I cannot be both at work and at home. I am not at work, so I must be at home. In fact, I am at a park.

- **Undistributed middle:** Any argument with the invalid structure of: Every A has B. C has B, so C is A.

**Example.** Every bird has a beak. That creature has a beak, so that creature must be a bird. In fact, the creature is a dinosaur.

### 1.2.2 Informal Fallacies

A formal fallacy occurs when the structure of the argument is incorrect, despite of the truth of the premises. A valid argument always has a correct formal structure and if the premises are true, the conclusion must be true. A valid argument is a formally correct argument that use true premises. When we use false premises, the formal fallacies disappear, but the argument may be regarded as a fallacy as the conclusion is invalid.

As an application of modus ponens, the following example contains no formal fallacies:

**Example.** If you took that course on CD player repair right out of high school, you would be doing well and having a vacation on the moon right now.

Even though there is no logic error in the argument, the conclusion is invalid because the premise is contrary to the fact. With a false premise, you can make any conclusion, so that the composed statement is always true. However, an always true statement has no value in reasoning.

By contrast, an argument with a formal fallacy could still contain all true premises:

**Example.** (a) If someone owns Fort Knox, then he is rich. (b) Bill Gates is rich. Therefore, (c) Bill Gates owns Fort Knox.

Although, (a) and (b) are true statements, (c) does not follow from (a) and (b) because the argument commits the formal fallacy of “affirming the consequent”.

There are numerous kinds of informal fallacies that use an incorrect relation between premises and conclusion. These fallacies can be grouped approximately into four groups:

- fallacies of improper premise;
- fallacies of faulty generalizations;
- fallacies of questionable cause; and
- relevance fallacies.

### 1.2.2.1 Improper Premise

- **Begging the question:** Providing what is essentially the conclusion of the argument as a premise.

**Example.** Everyone wants the new iPhone because it is the hottest new gadget on the market!

- **Circular reasoning:** The reasoner begins with what he or she is trying to end up with.

**Example.** You must obey the law, because it's illegal to break the law.

- **Complex question:** Someone asks a question that presupposes something that has not been proven or accepted by all the people involved.

**Example.** “Who is the King of France?” would commit the complex question fallacy because while it assumes there is a place called France (true), it also assumes France currently has a king (false).

### 1.2.2.2 Faulty Generalizations

The arguments in which the premises are related to the conclusions yet only weakly support the conclusions. A faulty generalization is thus produced.

- **Accident:** an exception to a generalization is ignored.

**Example.** Cutting people with knives is a crime. Surgeons cut people with knives. Therefore, surgeons are criminals.

- **Cherry picking.** Only select favor evidence in order to persuade the audience to accept a position, and evidence that would go against the position is withheld. The stronger the withheld evidence, the more fallacious the argument.

**Example.** My political candidate gives 10% of his income to the needy, and volunteers one day a week at a homeless shelter. Therefore, he is honest and morally straight.

- **Weak analogy:** The analogy is poorly suited.

**Example.** To say humans are immortal is like saying a car can run forever.

- **Hasty generalization:** Basing a broad conclusion on a small sample or the making of a determination without all of the information required to do so.

**Example.** You send a message back home that everyone in this new country is rude because the first person you meet in the airport is rude.

- **Misleading vividness:** involves describing an occurrence in vivid detail, even if it is an exceptional occurrence, to convince someone that it is a problem.

**Example.** We have had four days of temperatures that were 100 degrees or more! Global warming is getting dramatically worse!

### 1.2.2.3 Questionable Cause

This is a general type of error with many variants. Its primary basis is the confusion of association with causation, either by inappropriately deducing (or rejecting) causation or a broader failure to properly investigate the cause of an observed effect.

- **Faulty cause/effect:** A false cause/effect fallacy occurs when one cites to sequential events as evidence that the first caused the second. The argument generally looks like this: “Event A happened. Event B happened after A. Therefore, A caused B.”

**Example:** Every day, I eat cereal for breakfast. One time, I had a muffin instead, and there was a major earthquake in my city. I’ve eaten cereal ever since.

- **Complex cause:** It is assumed that there is one, simple cause of an outcome when in reality it may have been caused by a number of only jointly sufficient causes.

**Example.** President Trump has been in office for a month and gas prices have been skyrocketing.

- **Slippery slope:** In a slippery slope argument, one insists that one evidence will lead to a chain reaction resulting in an undesirable end or ends, without direct evidence that this course of events will happen.

**Example.** Today late for ten minutes, tomorrow late for an hour, and then someday you will simply cease to show up.

- **Furtive fallacy:** Outcomes are asserted to have been caused by the hidden misconduct or wrongdoing of decision makers. “Furtive” means secretive, hidden, conspiratorial, sly, and sneaky. Conspiracy theory is a typical example of furtive fallacy.
- **Gambler’s fallacy:** The incorrect belief that separate, independent events can affect the likelihood of another random event. If a fair coin lands on heads 10 times in a row, the belief that it is “due to land on tails” is incorrect. Parents might believe incorrectly that the next birth is a son, after having three girls.

#### 1.2.2.4 Fallacies of Relevance

- **Appeal to ignorance:** Assuming that a claim is true because it has not been or cannot be proven false or vice versa.

**Example.** “You can’t prove that there isn’t a mirror universe of our own, so there must be one out there somewhere!”

- **Appeal to common sense:** Asserting that your conclusion or facts are just “common sense” when, in fact, they are not.

**Example.** “I cannot imagine how this could be true; therefore, it must be false.”

- **Argument from repetition:** Repeating an argument until nobody cares to discuss it any more.

**Example.** Stupid people are so annoying. They prove their stupidity by saying stupid things.

- **Argument from silence:** Assuming that a claim is true based on the absence of textual or spoken evidence from an authoritative source.

**Example.** Jay: “Dude, where are my car keys?” Silent Bob says nothing. Jay: “I KNEW you took them!”

- **Appeal to motive:** Attacking the arguer instead of the argument, dismissing an idea by questioning the motives of its proposer. The word “motive” can be replaced by character, race, status, etc. This type of fallacies is called “ad hominem”.

**Example.** “This is why a woman shouldn’t do a man’s job!”

- **Appeal to emotion:** An argument is made due to the manipulation of emotions, rather than the use of valid reasoning. The emotion can be fear, flattery, pity, ridicule, etc.

**Example.** “There must be objective rights and wrongs in the universe. If not, how can you possibly say that torturing babies for fun could ever be right?”

- **Appeal to authority:** An argument uses the testimony of an authority in order to warrant their conclusion, but the authority appealed to is not an expert in the field in question. The word *authority* can be replaced by majority, nature, novelty, tradition, wealth, etc.

**Example.** “Einstein said ‘God does not play dice with the universe,’ therefore God must exist.”

- **Straw man:** An argument based on misrepresentation of an opponent’s position, especially to attack a weaker version of it rather than the argument actually presented. Also known as the *weak man* fallacy.

**Example.** Parent: “No dessert until you finish your chicken and vegetables!” Child: “You only love me when I eat.”

- **Two wrongs make a right:** It occurs when it is assumed that if one wrong is committed, another wrong will rectify it.

**Example.** Jimmy stole Tommy’s lunch in the past. Therefore, it is acceptable for Tommy to steal Jimmy’s lunch today.

Some informal fallacies do not belong to the above four groups. For example, the *false dilemma* fallacy and the *middle ground fallacy*.

### *False Dilemma Fallacy*

False Dilemma presents a choice between two mutually exclusive options, implying that there are no other options. One option is clearly worse than the other, making the choice seem obvious. Also known as the *either/or* fallacy, false dilemmas are a type of informal logical fallacy in which a faulty argument is used to persuade an audience to agree. False dilemmas are everywhere.

- Vote for me or live through four more years of higher taxes.
- America: Love it or leave it.
- Either we let every immigrant into our country, or we close the borders for everyone.
- Subscribe to our streaming service or be stuck with cable.
- You can either eat at this restaurant or have a sad TV dinner alone.
- Would you rather invest in your future or enjoy your money now?

### *The Middle Ground Fallacy*

This fallacy assumes that a compromise between two extreme conflicting points is always true. Arguments of this style ignore the possibility that one or both of the extremes could be completely true or false – rendering any form of compromise between the two invalids as well.

**Example.** Mary thinks the best way to improve sales is to redesign the entire company website, but John is firmly against making any changes to the website. Therefore, the best approach is to redesign some portions of the website.

An argument could contain both an informal fallacy and a formal fallacy yet lead to a conclusion that happens to be true. For example, again affirming the consequent, now also from an untrue premise:

If a scientist makes a statement about science, it is correct. It is true that quantum mechanics is deterministic. Therefore, a scientist has made a statement about it.

Having an understanding of these basic logical fallacies can help us more confidently parse the arguments we participate in and witness on a daily basis, separating fact from sharply dressed fiction.

### 1.3 A Brief Review of Mathematical Logic

Mathematical Logic is roughly divided into four areas:

- set theory,
- model theory,
- proof theory, and
- computability theory.

Each area has a distinct focus, although many techniques and results are shared between multiple areas. The borderlines between these areas, and the lines between mathematical logic and other fields of mathematics, are not always sharp. Each area contains rich materials which can be studied in graduate courses of multiple levels. As an introductory book on logic, we will introduce some basic concepts, which are closely related to this book, in these areas. This introduction is indeed less than a tip of the iceberg.

#### 1.3.1 Set Theory

A *set* is a structure, representing an unordered collection of zero or more distinct objects. All elements in a set are unequal (distinct) and unordered. Set theory is a branch of mathematical logic that studies sets and deals with operations between, relations among, and statements about sets. Although any type of object can be collected into a set, set theory is applied most often to objects that are relevant to mathematics. The language of set theory can be used to define nearly all mathematical objects.

##### 1.3.1.1 Basic Concepts and Notations

Let  $U$  be the collection of all objects under consideration.  $U$  is often called the *universal set* or *universe of discourse*. For example,  $U$  can be the set of natural numbers or the population of the world. Any object  $x$  under consideration is a *member* of  $U$ , written as  $x \in U$ . A *set* is any collection of objects from  $U$ . The *empty set* is often denoted by  $\emptyset$  or simply  $\{\}$ . Given two sets  $A$  and  $B$ , the following are the popular set operations:

- $x \in A$ : it is true iff (e.g., if and only if)  $x$  a member of  $A$ .

- $A \cup B$ : the *union* of  $A$  and  $B$ . For any object  $x$ ,  $x \in A \cup B$  iff  $x$  is in  $A$  or  $B$  or both.
- $A \cap B$ : the *intersection* of  $A$  and  $B$ .  $x \in A \cap B$  iff  $x \in A$  and  $x \in B$ .
- $A - B$ : the *set difference* of  $A$  and  $B$ .  $x \in A - B$  iff  $x \in A$  but not  $x \in B$ .
- $\bar{A}$ : the *complement* of  $A$  (with respect to the universal set  $U$ ).  $\bar{A} = U - A$ .  $A - B$  is also called the *complement* of  $B$  with respect to  $A$ :  $A - B = A \cap \bar{B}$ .
- $A \subseteq B$ :  $A$  is a *subset* of  $B$ .  $A \subseteq B$  is true iff every member of  $A$  is a member of  $B$ . Naturally,  $A \subseteq U$  and  $B \subseteq U$  for the universal set  $U$ .
- $A \subset B$ :  $A$  is a *proper subset* of  $B$ .  $A \subset B$  is true iff  $A \subseteq B$  and  $A \neq B$ .
- $A = B$ :  $A = B$  iff  $A \subseteq B$  and  $B \subseteq A$ , that is,  $A$  and  $B$  contain exactly the same elements.
- $A \times B$ : the *Cartesian product* of  $A$  and  $B$ . An ordered pair  $(a, b)$  is in  $A \times B$  iff  $a \in A$  and  $b \in B$ .
- $A^i$ : the *Cartesian product* of  $i$  copies of  $A$ , where  $i$  is a natural number (i.e., non-negative integers). When  $i = 0$ ,  $A^0$  denotes the *empty sequence*;  $A^1 = A$ ;  $A^2 = A \times A$ ,  $A^3 = A \times A \times A$ , etc.
- $|A|$ : the *cardinality* of a set. When  $A$  is finite,  $|A|$  is the number of members in  $A$ .  $|A|$  is called a *cardinal number* if  $A$  is infinite.
- $\mathcal{P}(A)$ : the *power set* of  $A$ .  $\mathcal{P}(A)$  is the set whose members are all of the possible subsets of  $A$ . That is, for any set  $X$ ,  $X \in \mathcal{P}(A)$  iff  $X \subseteq A$ . When  $|A| = n$ ,  $|\mathcal{P}(A)| = 2^n$ .

Elementary set theory can be studied informally and intuitively, and so can be taught in primary schools using Venn diagrams. Many mathematical concepts can be defined precisely using only set theoretic concepts, such as relations and functions.

### 1.3.1.2 Relations and Functions

Given two sets  $A$  and  $B$ , a relation between  $A$  and  $B$  is a subset of  $A \times B$ . In many applications,  $A = B$ . For example, the less-than relation,  $<$ , on  $N$ , is

$$< = \{\langle 0, 1 \rangle, \langle 0, 2 \rangle, \langle 0, 3 \rangle, \langle 1, 2 \rangle, \langle 0, 4 \rangle, \langle 1, 3 \rangle, \dots\}$$

which is a subset of  $N \times N$ . In general, give a relation  $R$  over set  $S$ , i.e.,  $R \subseteq S \times S$ , we assume that  $R(a, b) = \top$  if  $\langle a, b \rangle \in R$  and  $R(a, b) = \perp$  if  $\langle a, b \rangle \notin R$ . Of course, we are used to  $0 < 1 = \top$  and  $1 < 0 = \perp$ , but awkward to write  $< (0, 1) = \top$  and  $< (1, 0) = \perp$ .

A relation  $\leq$  over the set  $S$  is *reflexive* if  $a \leq a$  for any  $a, b \in S$ ;  $\leq$  is *antisymmetry* if  $a \leq b$  and  $b \leq a$  imply  $a = b$ ;  $\leq$  is *transitive* if  $a \leq b$  and  $b \leq c$  imply  $a \leq c$ ;  $\leq$  is *connex* (i.e., *comparable*) if either  $a \leq b$  or  $b \leq a$  is true. We say  $\leq$  is a *total order* if  $\leq$  is antisymmetry, transitive, and connex. We say  $\leq$  is a *partial order* if  $\leq$  is antisymmetry, transitive, and reflexive. Since the connexity implies the reflexivity, a total order is always a partial order but the inverse is not always true (some pair of elements may not comparable in a partial order).

The strict part of a partial order  $\leq$  is  $<$ :  $a < b$  iff  $a \leq b$  and  $a \neq b$ . Given a partial order  $\leq$ , a *minimal* element of  $X \subseteq S$  is an element  $m \in X$  such that if  $x \in X$  and  $x \leq m$ , then  $x = m$ . A partial order  $\leq$  is *well-founded* if there is a minimal element  $m \in X$  for any nonempty  $X \subseteq S$ . For example, the ordinary relation  $<$  on integers is well-founded on the set of natural numbers but not well-founded on the set of integers.

A *function*  $f$  is a relation  $R$  between  $A$  and  $B$  with the property that for every  $a \in A$ , there exists a unique  $b \in B$ , such that  $\langle a, b \rangle$  is in the relation  $R$ . We say  $b$  is the output of  $f$  on the input  $a$ . The function  $f$  is denoted by  $f : A \longrightarrow B$ , where  $A$  is the *domain* of  $f$  and  $B$  is the *range* of  $f$ . The domain refers to the set of possible input values for  $f$ . The range is the set of possible output values of  $f$ . If a function  $g$  takes two values as input, say  $g(a, a') = b$ , we denote  $g$  by  $g : A \times A' \longrightarrow B$ , and the domain of  $g$  is of course  $A \times A'$ . By convention, we write  $g(a, a') = b$  instead of  $g(\langle a, a' \rangle) = b$ , where  $a \in A$ ,  $a' \in A'$ , and  $b \in B$ . This practice can be generalized to functions with more than two inputs.

When both  $A$  and  $B$  are finite, the set  $F = \{f \mid f : A \longrightarrow B\}$  is also finite. Let  $|A| = m$  and  $|B| = n$ , then  $|F| = n^m$ , because there are  $n$  choices for  $f(a)$ , and there are  $m$  choices for  $a \in A$ . For example, if  $B = \{0, 1\}$  and  $A = B \times B$ , then  $|B| = 2$  and  $|A| = 4$ . The number of functions from  $A$  to  $B$  is  $2^4 = 16$ .

A total function  $f : A \longrightarrow B$  is said to be *injective* if for any  $x, y \in A$ ,  $x \neq y$  implies  $f(x) \neq f(y)$ .  $f : A \longrightarrow B$  is said to be *surjective* if for any  $y \in B$ , there exists  $x \in A$  such that  $f(x) = y$ .  $f$  is said to be *bijective* if  $f$  is both injective and surjective. A bijective function is also called a *one-to-one correspondence* between  $A$  and  $B$ . A function  $f$  is bijective iff the inverse of  $f$ , denoted by  $f^{-1}$ , exists, where  $f^{-1}(y) = x$  iff  $f(x) = y$ .

### 1.3.1.3 Countable Sets

People often use  $N$  to denote the set of natural numbers, i.e.,

$$N = \{0, 1, 2, 3, \dots\},$$

$Z$  to denote the set of integers,

$$Z = \{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\},$$

and  $R$  to denote the set of rational numbers. We also use  $Z^+$  to denote the set of all positive integers. These are naturally infinite sets.

For finite sets, the cardinality is a number in  $N$  and we can compare them. How to compare the cardinality of infinite sets? Cantor's idea is to use injective functions: If  $f : A \rightarrow B$  is injective, then  $|A| \leq |B|$ . That is, for every  $x \in A$ , we have a distinct  $f(x) \in B$ , so it cannot be the case that  $|A| > |B|$ .

**Definition 1.3.1.** *A set  $S$  is countable if there exists an injective function  $f : S \rightarrow N$ , where  $N$  is the set of natural numbers. If  $S$  is both countable and infinite,  $S$  is called countably infinite.*

**Example 1.3.2.** Let  $E = \{0, 2, 4, \dots\}$  be the set of all even natural numbers. Let  $f : E \rightarrow N$  with  $f(x) = x$ , then  $f$  is injective, so  $E$  is countable and we may conclude that  $|E| \leq |N|$ . Since  $E$  is infinite,  $E$  is also countably infinite. If we define  $h : N \rightarrow E$  with  $h(x) = 2x$ , then  $h$  is bijective, hence  $|N| \leq |E|$  and we end up with  $|E| = |N|$ . This is not an exception for all countably infinite sets.

**Proposition 1.3.3.** *There is a bijective function between any two countably infinite sets.*

*Proof.* It suffices to show that the proposition holds for the two sets  $A$  and  $N$  (the set of natural numbers). Since  $A$  is countable, there exists an injective function  $f : A \rightarrow N$ . Listing the elements of  $A$  as  $a_0, a_1, a_2, \dots$ , such that  $i < j$  iff  $f(a_i) < f(a_j)$ . Since  $A$  is infinite, define  $g : N \rightarrow A$  as  $g(i) = a_i$ , then  $g$  is a bijection between  $N$  and  $A$ .  $\square$

Today, countable sets form the foundation of *discrete mathematics*, and we first check what sets are countable.

**Proposition 1.3.4.** *The following sets are countable:*

1. Any finite set;
2. Any subset of  $N$ ;

3. The set  $A \times B$ , if both  $A$  and  $B$  are countable;
4. The set  $Z$  of integers;
5. The set of all binary strings of finite length;
6.  $\Sigma^*$ , the set of all strings of finite length built on the symbols from  $\Sigma$ , a finite alphabet.

*Proof.* To show (1) is countable, we just need to define an injective function which map a set of  $n$  elements to the first  $n$  natural numbers.

(2) is countable because if  $S \subseteq N$ , then  $f : S \rightarrow N$ , where  $f(x) = x$ , is an injective function.

(3) is countable when either  $A$  or  $B$  is finite and its proof is left as an exercise. When both  $A$  and  $B$  are countably infinite, without loss of generality, we assume  $A = B = N$ . Let  $g(k) = k(k + 1)/2$ , the sum of the first  $k$  positive integers, and  $f : N \times N \rightarrow N$  be the function:

$$f(i, j) = g(i + j) + j$$

If  $k = f(i, j)$ , then  $i = g(m) + m - k$  and  $j = k - g(m)$ , where  $g(m) \leq k < g(m + 1)$ . That is, the inverse of  $f$  exists and  $f$  is bijective. It is easy to check that  $i + j = m$ . For example, if  $k = 12$ , then  $m = 4$ , and  $i = j = 2$ . Table 1.3.1.3 shows the first few values of  $k = f(i, j)$ .

$i \setminus j$	0	1	2	3	4
0	0	2	5	9	14
1	1	4	8	13	19
2	3	7	12	18	25
3	6	11	17	24	32
4	10	16	23	31	40

Table 1.3.1:  $k = f(i, j) = (i + j)(i + j + 1)/2 + j$ .

To show (4) is countable, we can define a bijective function  $f$  on the integers:  $f(n) =$  if  $n \geq 0$  then  $2n$ ; else  $-2n - 1$ . This function maps 0 to 0; positive integers to even natural numbers, and negative integers to odd natural numbers.

To show (5) is countable, let  $S = \{0, 1\}^*$  be all the binary strings, we may sort  $S$  by length first, then by the decimal value of each string. Thus, the strings will be listed as follows:

$$S = \{\epsilon, 0, 1, 00, 01, 10, 11, 000, 001, \dots\},$$

where  $\epsilon$  denotes the empty string. To find an injection from  $S$  to  $N$ , for any string  $w$  in  $S$ , if  $w = \epsilon$ , let  $h(\epsilon) = 0$ ; and  $h(w)$  will replace every 0 in  $w$  by 2 and leaving 1 intact,

$$h(S) = \{0, 2, 1, 22, 21, 12, 11, 222, 221, \dots\}.$$

Since there is a bijection between  $S$  and  $h(S)$ , and  $h(S) \subset N$ ,  $h : S \rightarrow N$  is injective. So  $S$  is countable.

We can do better by finding a bijection between  $S$  and  $N$ . For any binary string  $s$ , let the length of  $s$  be  $n$ , i.e.,  $|s| = n$ . There are  $1 + 2 + 2^2 + \dots + 2^{n-1} = 2^n - 1$  strings shorter than  $s$ . Let  $v(s)$  be the decimal value of  $s$  (assuming  $v(\epsilon) = 0$ ), then there are  $v(s)$  strings of length  $n$  before  $s$ . So the position of  $s$  in the sorted list is  $2^n + v(s)$ . Define  $f : S \rightarrow N$ :  $f(s) = 2^{|s|} + v(s) - 1$ , then  $f$  must be a bijection.

To show (6) is countable, let  $k = \lceil \log_2(|\Sigma|) \rceil$ . Since  $|\Sigma| \leq 2^k$ , for each symbol  $a \in \Sigma$ , let  $h(a)$  be a distinct binary string of length  $k$ . For each string  $w \in \Sigma^*$ , let  $h(w) = h(a_1) \cdots h(a_n)$  if  $w = a_1 \cdots a_n$ . Then  $h : \Sigma^* \rightarrow \{0, 1\}^*$  is an injective function, so  $\Sigma^*$  is countable.  $\square$

The last three cases in the above proposition are examples of countably infinite sets.

#### 1.3.1.4 Uncountable Sets

In 1874, in his first set theory article, Georg Cantor introduced the term “countable set” and proved that the set of real numbers is uncountable by the famous diagonalization method, thus showing that not all infinite sets are countable. In 1878, he used bijection to define and compare cardinalities, contrasting sets that are countable with those that are uncountable.

**Proposition 1.3.5.** *The following sets are uncountable:*

1.  $B$ : the set of binary strings of infinite length;
2.  $R_1$ : the set of non-negative real numbers less than 1;
3.  $\mathcal{P}(N)$ : the power set of the natural numbers;
4.  $F$ : The set of Boolean functions,  $f : N \rightarrow \{0, 1\}$ , over the natural numbers.

*Proof.* For the first set  $B$ , we will use Cantor’s diagonalization method to show that  $B$  is uncountable by a refutational proof. Assume that  $B$  is countable, then there exists a one-to-one correspondence between  $B$  and  $N$ . For each  $i \in N$ , let the corresponding string in  $B$  be  $s_i$ , and the  $j^{\text{th}}$  symbol in  $s_i$  be  $s_{ij}$ . Now, we construct

$s_0 =$	<b>0</b>	0	0	0	0	0	0	0	0	0	...
$s_1 =$	0	<b>0</b>	0	0	0	0	0	0	1	0	...
$s_2 =$	0	0	<b>0</b>	0	0	1	0	0	0	0	...
$s_3 =$	0	0	0	<b>1</b>	0	0	0	0	0	0	...
$s_7 =$	0	0	0	1	<b>0</b>	1	0	0	0	0	...
$s_6 =$	0	0	0	0	1	<b>0</b>	0	0	0	0	...
$s_4 =$	0	0	0	1	0	1	<b>0</b>	1	0	0	...
$s_5 =$	0	0	1	0	0	0	0	<b>1</b>	0	0	...
$s_8 =$	0	1	0	0	1	0	0	1	<b>0</b>	0	...
$s_9 =$	0	1	0	0	1	1	0	1	0	<b>1</b>	...
					.....						

$x = \mathbf{1} \mathbf{1} \mathbf{1} \mathbf{0} \mathbf{1} \mathbf{1} \mathbf{1} \mathbf{0} \mathbf{1} \mathbf{0} \dots$

Table 1.3.2: Illustration of Cantor's diagonalization method: If  $B$  can be enumerated as  $s_0, s_1, s_2, \dots$ , then  $x = \overline{s_{00}} \overline{s_{11}} \overline{s_{22}} \dots$  is a binary string but cannot be in  $B$ .

a binary string  $x$  of infinite length as follows: the  $j^{\text{th}}$  symbol of  $x$ , i.e.,  $x_j$ , is the complement of  $s_{jj}$ . That is, if  $s_{jj} = 0$ , then  $x_j = 1$ ; if  $s_{jj} = 1$ , then  $x_j = 0$ . Clearly,  $x \in B$ . Suppose  $x$ 's corresponding number in  $N$  is  $k$ , i.e.,  $x = s_k$ . However,  $x$  differs from  $s_k$  on the  $k^{\text{th}}$  symbol, because  $x_k$  is the complement of  $s_{kk}$ . The contradiction comes from the assumption that  $B$  is countable. Table 1.3.1.4 illustrates the idea of Cantor's diagonalization method, where  $\bar{0} = 1$  and  $\bar{1} = 0$ .

Once we knew  $B$  is uncountable, if we have a bijection from  $B$  to  $R_1$ , then  $R_1$  must be uncountable. The bijection  $f : B \rightarrow R_1$  can be simply  $f(s) = 0.s$ , as each number of  $R_1$  can be represented by a binary number, where  $s \in B$ .

To show  $\mathcal{P}(N)$  is uncountable, let  $f : B \rightarrow \mathcal{P}(N)$  be  $f(s) = \{j \mid j^{\text{th}} \text{ symbol of } s \text{ is } 1\}$ . Obviously  $f$  is bijective.

Similarly, to show  $F$  is uncountable, let  $f : B \rightarrow F$  be  $f(s) = f$  if for every  $j$ ,  $f(j) =$  the  $j^{\text{th}}$  symbol of  $s$ . Obviously  $f$  is bijective, too.  $\square$

The result that  $\mathcal{P}(N)$  is uncountable while  $N$  is countable can be generalized: for every set  $S$ , the power set of  $S$ , i.e.,  $\mathcal{P}(S)$ , has a larger cardinality than  $S$  itself. This result implies that the notion of the *set of all sets* is an inconsistent notion. If  $S$  were the set of all sets then  $\mathcal{P}(S)$  would at the same time be bigger than  $S$  and a subset of  $S$ .

**Theorem 1.3.6. (Cantor's Theorem)**  $|A| < |\mathcal{P}(A)|$  for any set  $A$ .

*Proof.* Let us prove it by contradiction. If  $|A| = |\mathcal{P}(A)|$ , there is a bijection  $f : A \rightarrow \mathcal{P}(A)$ . Define  $S = \{a \in A \mid a \notin f(a)\}$ , then  $S \subseteq A$  and  $S \in \mathcal{P}(A)$ . There

must exist  $b \in A$  such that  $f(x) = S$  because  $f$  is surjective. There only two possibilities:  $b \in S$  or  $b \notin S$ .

1. If  $b \in S$ , by definition of  $S$ ,  $b \notin f(b) = S$ .
2. If  $b \notin S$ , by definition of  $S$ ,  $b \in f(b) = S$ .

Both cases lead to a contradiction, so  $f$  cannot exist. It cannot be the case that  $|A| > |P(A)|$  because  $g(a) = \{a\}$  is an injection from  $A$  to  $P(A)$ .  $\square$

Cantor's theorem implies that there is no such thing as the "set of all sets". Suppose  $A$  were the set of all sets. Since every element of  $P(A)$  is a set, so  $P(A) \subseteq A$ . Thus  $|P(A)| \leq |A|$ , a contradiction to Cantor's theorem.

Cantor chose the symbol  $\aleph_0$  for the cardinality of  $n$ ,  $|N|$ , a cardinal number.  $\aleph_0$  is read as aleph-null, after the first letter of the Hebrew alphabet. The cardinality of the reals is often denoted by  $\aleph_1$ , or  $c$  for the continuum of real numbers. Cantor's theorem implies that there are infinitely many infinite cardinal numbers, and that there is no largest cardinal number.

$$\begin{aligned} \aleph_0 &= |N| \\ \aleph_1 &= |P(N)| = 2^{\aleph_0} > \aleph_0 \\ \aleph_2 &= |P(P(N))| = 2^{\aleph_1} > \aleph_1 \\ \aleph_3 &= |P(P(P(N)))| = 2^{\aleph_2} > \aleph_2 \\ &\dots \end{aligned}$$

The famous inconsistent example related to set theory is so-called *Russell's paradox*, discovered by Bertrand Russell in 1901:

**Paradox 1.3.7. (Russell's paradox)** Let  $T = \{S \mid S \notin S\}$ , the set of all sets  $S$  that does not contain itself. What is the case of  $T \in T$ ?

If  $T \in T$  is false, then  $T \in T$  by the condition  $S \notin S$ . If  $T \in T$  is true, then  $T \notin T$  because every member of  $T$  satisfies  $S \notin S$ . This paradox showed that some attempted formalizations of the naïve set theory created by Georg Cantor led to a contradiction. Russell's paradox shows that the concept of "a set contains itself" is invalid. If " $S \in S$ " is always false, then  $T = \{S \mid S \notin S\}$  is the "set of all sets", which does not exist, either.

A layman's version of Russell's paradox is called the *Barber's paradox*.

**Paradox 1.3.8. (Barber's paradox)** In a village, there is only one barber who shaves all those, and those only, who do not shave themselves. The question is, does the barber shave himself?

Answering this question results in a contradiction. The barber cannot shave himself as he only shaves those who do not shave themselves. Conversely, if the barber does not shave himself, then he fits into the group of people who would be shaved by the barber, and thus, as the barber, he must shave himself.

The discovery of paradoxes in informal set theory caused mathematicians to wonder whether mathematics itself is inconsistent, and to look for proofs of consistency. Ernst Zermelo (1904) gave a proof that every set could be well-ordered, using the *axiom of choice*, which drew heated debate and research among mathematicians and the pioneers of set theory. In 1908, Zermelo provided the first set of axioms for set theory. These axioms, together with the additional axiom of replacement proposed by Abraham Fraenkel, are now called Zermelo–Fraenkel set theory (ZF). Besides ZF, many set theories are proposed since then, to rid paradoxes from set theory.

### 1.3.2 Model Theory

Model theory is the study of mathematical structures (e.g. groups, fields, algebras, graphs, logics) in a formal language. Every formal language has its syntax and semantics. Models are a semantic structure associated with syntactic structures in a formal language. Theories are then introduced based on models. Following this approach, every formal logic is defined inside of a formal language.

#### 1.3.2.1 Syntax and Semantics

The syntax of a formal language specifies how various components of the language, such as symbols, words, and sentences, are defined. For example, the language for propositional logic uses only propositional variables and logic operators as its symbols, and well-formed formulas built on these symbols as its sentences. In model theory, a set of sentences in a formal language is one of the components that form a theory. The language for logic often contains the two constant symbols, either 1 and 0, or  $\top$  and  $\perp$ , which are interpreted as true and false, respectively. They are usually considered to be special logical operators which take no arguments, not propositional variables.

The semantics of a language specifies the meaning of various components of the language. For example, if we use symbol  $q$  to stand for the statement “Today is Monday”, and  $r$  for “Today is Tuesday”, then the meaning of  $q$  is “Today is Monday”;  $q$  can be used to denote a thousand different statements, just like a thousand Hamlets in a thousand people’s eyes. On the other hand, a formal meaning of the formula  $q \vee r$  is the truth value of  $q \vee r$ , which can be decided uniquely when

the truth values of  $q$  and  $r$  are given. In model theory, the formal meaning of a sentence is explored: it examines semantical elements (meaning and truth) by means of syntactical elements (formulas and proofs) of a corresponding language.

In model theory, semantics and model are synonyms. A *model* of a theory is an interpretation that satisfies the sentences of that theory. Abstract algebras are often used as models. In a summary definition, dating from 1973,

$$\text{model theory} = \text{abstract algebra} + \text{logic},$$

*Abstract algebra*, also called *universal algebra*, is a broad field of mathematics, concerned with sets of abstract objects associated various operations and properties.

### 1.3.2.2 Boolean Algebra

In abstract algebra, the most relevant algebra related to the logic discussed in this book is *Boolean algebra*, which is approximately equivalent to propositional logic. Many syntactic concepts of Boolean algebra carry over to propositional logic with only minor changes in notation and terminology, while the semantics of propositional logic are defined via Boolean algebras in a way that the tautologies (theorems) of propositional logic correspond to equational theorems of Boolean algebra.

In Boolean algebra, the values of the variables are the truth values true and false, usually denoted 1 and 0 respectively. The main operations of Boolean algebra are the multiplication  $\cdot$  (conjunction), the addition  $+$  (disjunction), and the inverse  $i$  (negation). It is thus a formalism for describing logical operations in the same way that elementary algebra describes numerical operations, such as addition and multiplication. In fact, a *Boolean algebra* is any set with binary operations  $+$  and  $\cdot$  and a unary operation  $i$  thereon satisfying the Boolean laws (equations), which define the properties of the logical operations.

Boolean algebra was introduced by George Boole in his first book *The Mathematical Analysis of Logic* (1847). According to Huntington, the term “Boolean algebra” was first suggested by Sheffer in 1913. Sentences that can be expressed in propositional logic have an equivalent expression in Boolean algebra. Thus, Boolean algebra is sometimes used to denote propositional logic performed in this way. However, Boolean algebra is not sufficient to capture logic formulas using quantifiers, like those from first order logic.

### 1.3.3 Proof Theory

In a similar way to model theory, proof theory is situated in an interdisciplinary area among mathematics, philosophy, and computer science. According to

*Wikipedia*, proof theory is a major branch of mathematical logic that represents proofs as formal mathematical objects, facilitating their analysis by mathematical techniques.

### 1.3.3.1 Axioms and Theorems

In a formal logic, the *axioms* are a set of sentences which are assumed to be true. Typically, the axioms contain by default the definitions of all logical operators. For example, the negation operator  $\neg$  is defined by  $\neg(\top) = \perp$  and  $\neg(\perp) = \top$ .

The *theorems* are the formulas which can be proved to be true from the axioms. For example, for any propositional variable  $p$ ,  $\neg(\neg(p)) = p$  is a theorem. Using the case analysis method, if the truth value of  $p$  is  $\top$ , then  $\neg(\neg(\top)) = \neg(\perp) = \top$ ; if the truth value of  $p$  is  $\perp$ , then  $\neg(\neg(\perp)) = \neg(\top) = \perp$ . It is easy to see that different sets of axioms would lead to different theorems.

There are two important properties concerning a set of axioms:

- **Consistency** The axiom set is *consistent* if every formula in the axiom can be true at the same time.
- **Independence** The axiom set is *independent* if no axiom is a theorem of the other axioms.

### 1.3.3.2 Proof Procedures

Given a set  $\mathbf{A}$  of axioms and a formula  $B$ , we need a procedure  $P$  that can answer the question if  $B$  is a theorem of  $\mathbf{A}$ : If  $P(\mathbf{A}, B)$  returns “yes”, we say  $B$  is proved, if  $P(\mathbf{A}, B)$  returns “no”, we say  $B$  is disproved; if it returns “unknown”, we do not know if  $B$  is a theorem or not.  $P$  is called a *proof procedure*. This procedure can be carried out by hand, or executed on a computer.

There are two important properties concerning a proof procedure  $P$ :

- **Soundness** The proof procedure  $P$  is *sound* iff whenever  $P(\mathbf{A}, B)$  returns true,  $B$  is a theorem of  $\mathbf{A}$ ; whenever  $P(\mathbf{A}, B)$  returns false,  $B$  is not a theorem of  $\mathbf{A}$ .
- **Completeness** The proof procedure  $P$  is *complete* iff for any theorem  $B$  of  $\mathbf{A}$ ,  $P(\mathbf{A}, B)$  returns true.

If a procedure  $P(\mathbf{A}, B)$  is sound and always halts on any input  $\mathbf{A}$  and  $B$  with “yes” or “no”, i.e.,  $P$  is an algorithm, then  $P$  is called a *decision procedure* for the

logic. This definition of decision procedure is narrow; in general, any algorithm which can answer a decision question is called a *decision procedure*.

**Theorem 1.3.9.** *Every decision procedure is sound and complete.*

*Proof.* Given any axiom set  $\mathbf{A}$  and any formula  $B$ , since  $P$  is a decision procedure,  $P(\mathbf{A}, B)$  will return a truth value. If  $B$  is a theorem of  $\mathbf{A}$ ,  $P(\mathbf{A}, B)$  will return true because  $P$  is sound.  $\square$

Some proof procedures may stop with “yes”, “no”, or “unknown”; they cannot be called a decision procedure. A sound and complete proof procedure may loop on a formula that is not a theorem, thus cannot be a decision procedure by definition. From the above theorem, we see that the soundness and halting properties of a proof procedure implies the completeness. A sound and complete proof procedure is also called a *semi-decision procedure* for the logic, because it will halt on any formula which is a theorem. For any logic, we always look for its decision procedure. For example, propositional logic has many different decision procedures. If a decision procedure does not exist, we look for a semi-decision procedure. For some logics, even a semi-decision procedure cannot exist. Clearly, these discussions will only make sense once we have formally defined the logic in question.

### 1.3.3.3 Inference Rules

In structural proof theory, which is a major area of proof theory, inferences rules are used to construct a proof. In logic, an *inference rule* is a pattern which takes formulas as premises, and returns a formula as conclusion. For example, the rule of *modus ponens* (MP) takes two premises, one in the form “If  $p$  then  $q$ ” (i.e.,  $p \rightarrow q$ ) and another in the form  $p$ , and returns the conclusion  $q$ . The rule is *sound* if the premises are true under an interpretation, then so is the conclusion.

An inference system  $S$  consists of a set of axioms  $A$  and a set of inference rules  $R$ . The soundness of  $S$  comes from the soundness of every inference rule in  $R$ . A *proof* of formula  $F_n$  in  $S = (A, R)$  is a sequence of formulas  $F_1, F_2, \dots, F_n$  such that each  $F_i$  is either in  $A$  or can be generated by an inference rule of  $R$ , using the formulas before  $F_i$  in the sequence as the premises.

**Example 1.3.10.** Let  $S = (A, R)$ ,  $A = \{p \rightarrow (q \rightarrow r), p, q\}$ , and  $R = \{MP\}$ . A proof of  $r$  from  $A$  using  $S$  is given below:

1.  $p \rightarrow (q \rightarrow r)$  *axiom*
2.  $p$  *axiom*
3.  $q \rightarrow r$  *MP, 1, 2*
4.  $q$  *axiom*
5.  $r$  *MP, 3, 4*

One obvious advantage of such a proof is that these proofs can be checked either by hand or automatically by computer. Checking formal proofs is usually simple, whereas finding proofs (automated theorem proving) is generally hard. An informal proof in the mathematics literature, by contrast, requires weeks of peer review to be checked, and may still contain errors. Informal proofs of everyday mathematical practice are unlike the proofs of proof theory. They are rather like high-level sketches that would allow an expert to reconstruct a formal proof at least in principle, given enough time and patience. For most mathematicians, writing a fully formal proof is too pedantic and long-winded to be in common use.

In proof theory, proofs are typically presented as recursively-defined data structures such as lists (as shown here), or trees, which are constructed according to the axioms and inference rules of the logical system. As such, proof theory is syntactic in nature, in contrast to model theory, which is semantic in nature.

Every inference system can be converted to a proof procedure if a *fair strategy* is adapted to control the use of inference rules. Fair strategy ensures that if an inference rule is applicable at some point to produce a new formula, this new formula will be derived eventually. If every inference rule is sound, the corresponding proof procedure is sound. Extra effort is needed in general to show that the proof procedure is complete, or halting.

We have defined a formal notion of proof as a sequence of formulae, starting with axioms and continuing with theorems derived from earlier members of the sequence by rules of inference.

Besides structural proof theory, some of the major areas of proof theory include ordinal analysis, provability logic, reverse mathematics, proof mining, automated theorem proving, and proof complexity.

### 1.3.4 Computability Theory

According to *Wikipedia*, *computability theory*, used to be called *recursion theory*, is a branch of mathematical logic, of computer science, and of the theory of computation that originated in the 1930s with the study of computable functions. The field has since expanded to include the study of generalized computability and definability. In these areas, computability theory overlaps with proof theory and set theory.

#### 1.3.4.1 Recursion

The main reason for the name of “recursion theory” is that recursion is used to construct objects and functions.

**Example 1.3.11.** Given a constant symbol  $0$  of type  $T$  and a function symbol  $s : T \rightarrow T$ , the objects of type  $T$  can be recursively constructed as follows:

1.  $0$  is an object of type  $T$ ;
2. If  $n$  is an object of type  $T$ , so is  $s(n)$ ;
3. Nothing else will be an object of  $T$ .

If we let  $s^0(0) = 0$ ,  $s^1(0) = s(0)$ ,  $s^2(0) = s(s(0))$ , etc., then  $T$  can be expressed as

$$T = \{0, s(0), s^2(0), s^3(0), \dots, s^i(0), \dots\}.$$

Obviously, there exists a bijection  $f$  between  $T$  and  $N$ , the set of natural numbers. That is, for any  $i \in N$ ,  $f(i) = s^i(0)$ . By giving  $s$  the meaning of *successor* of a natural number, mathematicians use  $T$  as a definition of  $N$ .

Functions can be recursively defined in a similar way. Let  $pre$ ,  $add$ ,  $sub$ ,  $mul$  be the predecessor, addition, subtraction, and multiplication functions over the set of natural numbers:

$$\begin{aligned} pre(0) &= 0; \\ pre(s(x)) &= x. \\ add(0, y) &= y; \\ add(s(x), y) &= s(add(x, y)). \\ sub(x, 0) &= x; \\ sub(x, s(y)) &= sub(pre(x), y). \\ mul(0, y) &= 0; \\ mul(s(x), y) &= add(mul(x, y), y). \end{aligned}$$

Note that the subtraction defined above is different from the subtraction over the integers: we have  $sub(s^i(0), s^j(0)) = 0$  when  $i \leq j$ . This function holds the property that  $sub(x, y) = 0$  iff  $x \leq y$ .

The above functions are examples of so-called *primitive recursive functions*, which are total and computable functions. The set of *general recursive functions* include all primitive recursive functions, and some functions which are not computable.

#### 1.3.4.2 Backus-Naur Form (BNF)

In computer science, *Backus-Naur form* (BNF) is a notation technique for *context-free grammars*, often used to describe the syntax of languages used in computing, such as computer programming languages, document formats, instruction

sets and communication protocols. John Backus and Peter Naur are the first to use this notation to describe ALGOL, an early programming language. BNF is applied wherever formal descriptions of languages are needed. We will use BNF describe the syntax of logic as a language.

BNF can be extended to define recursively constructed objects in a succinct style. For example, the set of natural numbers is expressed in BNF as follows:

$$\langle N \rangle ::= 0 \mid s(\langle N \rangle)$$

The symbol  $::=$  can be understood as “is defined by”, and “ $\mid$ ” separates alternative parts of a definition. The above formula can be read as follows: the member of  $N$  is either 0, or  $s$  applying to (another) member of  $N$ ; nothing else will be in  $N$ . From this example, we see that to define a set  $N$ , we use any member of  $N$  is denoted by  $\langle N \rangle$  and is defined by the items on the right of  $::=$ . For a recursive definition, we have a basic case (such as 0) and a recursive case (such as  $s(\langle N \rangle)$ ). Items like  $\langle N \rangle$  are called *variables* in a context-free grammar and items like 0,  $s$ , ( $($ , and  $)$  are called *terminals*. A context-free grammar specifies what set of strings of terminals can be derived for each variable. For the above example, the set of strings derived from  $\langle N \rangle$  is

$$\{0, s(0), s^2(0), s^3(0), \dots, s^i(0), \dots\}$$

This set is the same as the objects defined in Example 1.3.11, but BNF gives us a formal and compact definition.

**Example 1.3.12.** *To define the set  $\Sigma^*$  of all strings (of finite length) over the set  $\Sigma$  of symbols, called the alphabet, we may use BNF. For instance, if  $\Sigma = \{a, b, c\}$ , then  $\Sigma^*$  is the set of all binary strings:*

$$\begin{aligned} \langle symbol \rangle & ::= a \mid b \mid c \\ \langle \Sigma^* s \rangle & ::= \epsilon \mid \langle symbol \rangle \langle \Sigma^* \rangle \end{aligned}$$

*That is, the empty string  $\epsilon \in \Sigma^*$ ; if  $w \in \Sigma^*$ , for every symbol  $a \in \Sigma$ ,  $aw \in \Sigma^*$ ; nothing else will be in  $\Sigma^*$ . Thus,  $\Sigma^*$  contains exactly all the strings built on  $\Sigma$ :*

$$\Sigma^* = \{\epsilon, a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, aaa, \dots\}.$$

Since there is a one-to-one corresponding between a string and a sequence, sequences can be also recursively defined in a similar way. In fact, all recursively constructed objects can be defined this way.

### 1.3.4.3 Computable Functions

One of the basic questions addressed by computability theory is the following:

- What does it mean for a function on the natural numbers to be computable?

The answer to the question is that the function is computable if there exists an algorithm which computes the function. Logicians and computer scientists, including Turing and Church, gave different definitions of algorithms based on different computing models. These computing models can be Turing machines, Church's  $\lambda$ -calculus, general recursive functions, Markov algorithms, Chomsky's grammars, Minski's counter machine, or von Neumann model. One point in common is that, by algorithm, we mean that the computing process will halt in a finite number of steps.

In section 1.3.1.4, we saw that the set of functions, which take a natural number and return a Boolean value, is uncountable. However, the set of algorithms is countable, because every algorithm can be recorded by a text file of finite length (which can be regarded as a binary string), and the set of binary strings is countable. In other words, we have only a countable set of algorithms but the functions on the natural numbers are uncountable. As a result, there exist massive number of functions on the natural numbers which are not computable.

One notable non-computable problem is to decide if a Turing machine  $T$  will halt on an input  $x$ . If we encode  $T$  and  $x$  as a single positive integer  $\langle T, x \rangle$ , and define the function  $f(\langle T, w \rangle) = 1$  if  $T$  halts on  $w$  and 0, otherwise, then  $f$  is a non-computable function on the natural numbers. Note that the Turing machine  $T$  can be replaced by any other equivalent computing model and the result still holds. This follows from the *Church–Turing thesis* (also known as computability thesis, the Church–Turing conjecture), which states that a function on the natural numbers can be calculated by an effective method if and only if it is computable by a Turing machine.

The concept of *Turing completeness* says that if a computing model has the ability to simulate a Turing machine. A computing model that is Turing complete is theoretically capable of doing all tasks done by computers; on the other hand, nearly all computers are Turing complete if the limitation of finite memory is ignored. Some logics are also Turing complete as they can also be used to simulate a Turing machine. As a result, some problems such as the halting problem for such logics are undecidable. Computability theory will help us to decide if there exist or not decision procedures for some logic, and we will see some examples in Chapter 9.

One of the main aims of this book is to present algorithms invented for logic and the software tools built up based on these algorithms. That is why we will

invest a disproportionate amount of effort (by the standards of conventional logic textbooks) in the algorithmic aspect of logic. Set theory and model theory help us to understand and specify the problems; proof theory helps us to create algorithms; and computability theory helps us to know what problems have algorithms and what problems do not have algorithms at all.

#### 1.4 Exercise Problems

1. Decide if the following sentences are statements or not. If they are, use propositional symbols and logic connectives to express them.
  - (a) If I win the lottery, I'll be poor.
  - (b) The man asked, "shut the door behind you."
  - (c) Today is the first day of the school.
  - (d) Today is either sunny or cloudy.
2. Answer the questions in the following examples regarding logical fallacies.
  - (a) During their argument, Anthony tells Marie that he hates her and wished that she would get hit by a car. Later that evening, Anthony receives a call from a friend who tells him that Marie is in the hospital because she was struck by a car. Anthony immediately blames himself and reasons that if he hadn't made that comment during their fight, Marie would not have been hit. What logical fallacy has Anthony committed?
  - (b) Louise is running for class president. In her campaign speech she says, "My opponent does not deserve to win. She is a smoker and she cheated on her boyfriend last year." What fallacy has Louise committed?
  - (c) Bill: "I think capital punishment is wrong." Adam: "No it isn't. If it was wrong, it wouldn't be legal." Of which fallacy is this an example?
  - (d) Ricky is watching television when he sees a commercial for foot cream. The commercial announcer says, "This is the best foot cream on the market because no other foot cream makers have been able to prove otherwise!" What fallacy has the announcer committed?
  - (e) Maria has been working at her current job for more than 30 years at the same wage. She desperately wants a raise so she approaches her boss to ask for one. She says, "You are one of the kindest people I know. You are smart and good-looking and I really love your shoes." What type of fallacy is this?

- (f) Jeff's mom is concerned when she finds out that he skipped class one day. She tells him that she is concerned that since he skipped one class, he will start skipping more frequently. Then he will drop out altogether, never graduate or get into college, and end up unemployed and living at home for the rest of his life. What type of fallacy has Jeff's mom committed?
- (g) Dana is trying to raise money for her university's library. In her address to the board of trustees, she says, "We must raise tuition to cover the cost of new books. Otherwise the library will be forced to close." Of what fallacy is this an example?
- (h) Jeff is preparing to create a commercial for a new energy drink. He visits a local high school and surveys students in an English class about their beverage preferences. The majority of the class says they prefer grape flavored drinks, so Jeff tells his superiors that grape is the flavor favored most by high school students. What error in reasoning has Jeff made?
3. Let  $A = \{0, 1\}$ .  $\mathcal{P}(A) = ?$   $\mathcal{P}(\mathcal{P}(A)) = ?$
4. Let  $a, b$  be real numbers and  $a < b$ . Define a function  $f : (0, 1) \rightarrow (a, b)$  by  $f(x) = a + x(b - a)$ , where  $(a, b)$  denotes the interval  $\{x \mid a < x < b\}$ . Show that  $f$  is bijective.
5. Answer the following questions with justification.
- Is the set of all odd integers countable?
  - Is the set of real numbers in the interval  $(0, 0.1)$  countable?
  - Is the set of angles in the interval  $(0^\circ, 90^\circ)$  countable?
  - Is the set of all points in the plane with rational coordinates countable?
  - Is the set of all Java programs countable?
  - Is the set of all words using an English alphabet countable?
  - Is the set of sands on the earth countable?
  - Is the set of atoms in the solar system countable?
6. **Barber's paradox:** In a village, there is only one barber who shaves all those, and those only, who do not shave themselves. The question is, does the barber shave himself? Let the relation  $s(x, y)$  be that " $x$  shaves  $y$ " and  $b$  denote the barber.

- (a) Define  $A$  to be the set of the people shaved by the barber using  $s(x, y)$  and  $b$ .
  - (b) Define  $B$  to be the set of the people who does not shave himself.
  - (c) Use  $A$  and  $B$  to show that the answer to the question “does the barber shave himself?” leads a contradiction.
7. Prove that any subset of a countable set is countable.
  8. Prove that the set of rational numbers is countable.
  9. Prove that the set  $N \times \{a, b, c\}$  is countable, where  $N$  is the set of natural numbers. Your proof should not use the property that countable sets are closed under Cartesian product.
  10. Decide with a proof that the set  $N^k$  is countable or not, where  $N$  is the set of natural numbers and  $k \in N$ .
  11. Decide with a proof that the set of all functions  $f : \{0, 1\} \rightarrow N$  is countable or not.
  12. Prove that if both  $A$  and  $B$  are countable, so is  $A \cup B$ .
  13. Prove that if  $A$  is uncountable and  $B$  is countable, then  $A - B$  is uncountable.
  14. Prove that the set  $\{a, b, c\}^*$  of all strings (of finite length) on  $\{a, b, c\}$  is countable.
  15. If a language is any set of strings in  $\Sigma^*$ , where  $\Sigma \neq \emptyset$ , prove that the set of all languages is uncountable.
  16. Suppose the set of natural numbers is constructed by the constant 0 and the successor function  $s$ . Provide a definition of  $<$  and  $\leq$  over the natural numbers. No functions are allowed to use.
  17. Provide a definition of power function,  $pow(x, y)$ , which computes  $x^y$ , over the natural numbers built up by 0 and the successor function  $s$ . You are allowed to use the functions  $add$  (addition) and  $mul$  (multiplication).
  18. Provide a BNF (Backus–Naur form) for the set of binary trees, where each non-empty node contains a natural number. You may use  $null : Tree$  for the empty tree and  $node : Nat, Tree, Tree \rightarrow Tree$  for a node of the tree.

## CHAPTER 2

### PROPOSITIONAL LOGIC

Following model theory, we will present propositional logic as a language. We first introduce the syntax of this language and then its semantics.

#### 2.1 Syntax

The syntax of propositional logic specifies what symbols and what sentences are used in the logic. Propositional logic uses *propositional variables* and *logical operators* as the only symbols. In the first chapter, we said that propositional variables are symbols for denoting statements and take only true and false as truth values. For convenience, we will use 1 for true and 0 for false. We will treat  $\top$  and  $\perp$  as the two nullary logical operators such that the truth value of  $\top$  is always 1 and the truth value of  $\perp$  is always 0. Note that  $\top$  and  $\perp$  are syntactic symbols while 1 and 0 are semantic symbols.

##### 2.1.1 Logical Operators

Mathematicians use the words “not”, “and”, “or”, etc., for operators that change or combine propositions. The meaning of these logical operators can be specified as a function which takes Boolean values, and returns a Boolean value. These functions are also called *Boolean functions*. For example, if  $p$  is a proposition, then so is  $\neg p$  and the truth value of the proposition  $\neg p$  is determined by the truth value of  $p$  according to the meaning of  $\neg$ :  $\neg(1) = 0$  and  $\neg(0) = 1$ . That is, if the value of  $p$  is 1, then the value of  $\neg(p)$  is 0; if the value of  $p$  is 0, then the value of  $\neg(p)$  is 1.

In general, the definition of a Boolean function is displayed by a table where the output of the Boolean function is given for each possible truth values of the input variables. Such a table is called *truth table*. For example, the meanings of  $\wedge$ ,  $\vee$ , and  $\rightarrow$  are defined by the following truth table: This table has four lines, since there are four pairs of Boolean values for the two variables:

According to this table, the truth value of  $p \wedge q$  is 1 when both  $p$  and  $q$  are 1. The truth value of  $p \vee q$  is 1 when either  $p$  or  $q$ , or both are true. This is not always the intended meaning of “or” in everyday dialog, but this is the standard definition in logic. So if a logician says, “You may have cake, or you may have ice cream,” he means that you could have both. If you want to exclude the possibility of having

$p$	$q$	$p \wedge q$	$p \vee q$	$p \rightarrow q$	$p \oplus q$	$p \leftrightarrow q$
1	1	1	1	1	0	1
1	0	0	1	0	1	0
0	1	0	1	1	1	0
0	0	0	0	1	0	1

Table 2.1.1: Definitions of  $\wedge$ ,  $\vee$ ,  $\rightarrow$ ,  $\oplus$ , and  $\leftrightarrow$ 

both cake and ice cream, you should combine them with the exclusive-or operator,  $\oplus$ , which is also defined in Table 2.1.1. The exclusive disjunction  $\oplus$  corresponds to addition modulo 2 and is therefore given the symbol  $+$  when the boolean values are assumed to be 0 and 1.

In Table 2.1.1, we also give the definition of  $\leftrightarrow$  (if and only if). If  $p$  denotes “Tom wears a blue jean” and  $q$  denotes “Sam wears a blue jean”, then the formula  $p \leftrightarrow q$  asserts that “Tom and Sam always wear a blue jean at the same time”. That is, the value of  $p \leftrightarrow q$  is 1 when  $p$  and  $q$  have the same truth value: either both are 1 or both are 0. Note that the  $p \oplus q$  and  $p \leftrightarrow q$  always have opposite values, that is,  $\neg(p \oplus q)$  and  $p \leftrightarrow q$  always have the same value.

The truth table for implications is also given in Table 2.1.1. This operator is worth remembering, because a large fraction of all mathematical statements are of the if-then form. Now let’s figure out the truth of the following example: If elephants fly, I’ll be on Mars. What is the truth value of this statement? Elephants do not fly, so the truth value of  $p$  is 0 and we fall to the last two lines of Table 2.1.1: In both lines,  $p \rightarrow q$  is 1. In section 1.2.2 on Informal Fallacies, we pointed out that with a false premise, you can make any conclusion, so that the composed statement is always true. However, an always true statement has no value in reasoning, and it does not imply the causal connection between the premise and the conclusion of an implication. In logic, it is important to accept the fact that logical implications ignore causal connections.

Let  $B = \{0, 1\}$ , any logical operator is a function  $f : B^n \rightarrow B$ . There are two nullary operators (when  $n = 0$ ):  $\top = 1$  and  $\perp = 0$ . There are four unary operators (when  $n = 1$ ):  $f_1(x) = 0$ ;  $f_2(x) = 1$ ;  $f_3(x) = x$ ; and  $f_4(x) = \neg x$ . There are 16 binary operators (when  $n = 2$ ), and  $\wedge$ ,  $\vee$ ,  $\rightarrow$ ,  $\oplus$ , and  $\leftrightarrow$  are some of them. We may use operators when  $n > 2$ . For example, let  $ite : B^3 \rightarrow B$  be the if-then-else operator, if we use a truth table for the meaning of  $ite(p, q, r)$ , then the truth table will have eight lines. In general, there are  $2^{2^n}$   $n$ -ary Boolean functions, as the number of functions  $f : A \rightarrow B$  is  $|B|^{|A|}$ , where  $|A| = |B^n| = |B|^n = 2^n$ .

### 2.1.2 Formulas

**Definition 2.1.1.** *The propositional formulas are the set of strings built up by propositional variables and logical operators. Let  $op$  be the binary operators and  $V_p$  be the propositional variables used in the current application, then the formulas for this application can be defined by the following BNF grammar:*

$$\begin{aligned} \langle op \rangle & ::= \wedge \mid \vee \mid \rightarrow \mid \oplus \mid \leftrightarrow \\ \langle V_p \rangle & ::= p \mid q \mid r \mid s \mid t \\ \langle Formulas \rangle & ::= \top \mid \perp \mid \langle V_p \rangle \mid \neg \langle Formulas \rangle \mid (\langle Formulas \rangle \langle op \rangle \langle Formulas \rangle) \end{aligned}$$

We may add/delete operators or propositional variables at will in the above definition. According to the definition, every member of *Formulas* is either the constants  $\top$  or  $\perp$ , a propositional variable of  $V_p$ , or the negation of a well-formula, or a binary operator applying to two formulas. Here are some examples:

$$\top, \quad p, \quad \neg q, \quad (p \wedge \neg q), \quad (r \rightarrow (p \wedge \neg q)), \quad ((p \wedge \neg q) \vee r)$$

Throughout this chapter, we will use  $p, q, \dots$ , for propositional variables, and  $A, B, \dots$ , for formulas.

For every binary operator, we introduce a pair of parentheses, and some of them may be unnecessary and we will use a precedence relation over all the operators to remove them. The typical precedence is given as the following list, where operators with higher precedence go first. Note that we assign  $\oplus$  and  $\leftrightarrow$  the same precedence.

$$\neg, \wedge, \vee, \rightarrow, \{\oplus, \leftrightarrow\}.$$

Thus  $(r \rightarrow (p \wedge \neg q))$  can be simply written as  $r \rightarrow p \wedge \neg q$ , and we will not take it as  $(r \rightarrow p) \wedge \neg q$ , because  $\wedge$  has higher precedence than  $\rightarrow$ . Similarly,  $((p \wedge \neg q) \vee r)$  is written as  $p \wedge \neg q \vee r$  without ambiguity.

When a well-formed expression, like the formulas defined above, is represented by a string (a sequence of symbols), we use parentheses and commas to identify the structures of the expression. Such a structure can be easily identified by a tree, where we do not need parentheses and commas.

**Definition 2.1.2.** *A tree for formula  $A$  is defined as follows:*

1. *If  $A$  is  $p, \top$ , or  $\perp$ , the tree is a single node with  $A$  as its label.*
2. *If  $A = \neg B$ , then the root node is labeled with  $\neg$  and has one branch pointing to the formula tree of  $B$ .*

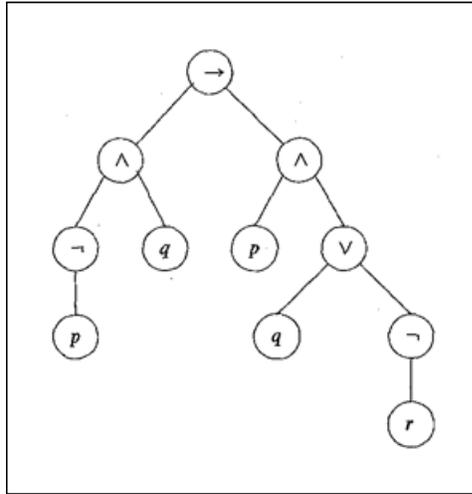


Figure 2.1.1: The formula tree for  $(\neg p \wedge q) \rightarrow (p \wedge (q \vee \neg r))$

3. IF  $A = B \text{ op } C$ , then the root node is labeled with  $\text{op}$  and has two branches pointing to the formula trees of  $B$  and  $C$ , respectively.

For example, the formula tree of  $(\neg p \wedge q) \rightarrow (p \wedge (q \vee \neg r))$  is displayed in Figure 2.1.1. The tree representation of a formula does not need parentheses.

A formula and its formula tree are the two representations of the same object. We are free to choose one of the representations for the convenience of discussion.

**Definition 2.1.3.** Given a formula  $A$ , we say  $B$  is a subformula  $A$  if either  $B = A$  or  $A = (A_1 \text{ op } A_2)$ , where  $\text{op}$  is a binary operator, or  $A = \neg A_1$  and  $B$  is a subformula of  $A_1$  or  $A_2$ . If  $B$  is a subformula of  $A$  and  $B \neq A$ , then  $B$  is a proper subformula of  $A$ .

For example,  $\neg p \vee \neg(q \vee r)$  contains  $p$ ,  $q$ ,  $r$ ,  $\neg p$ ,  $q \vee r$ , and  $\neg(q \vee r)$  as proper subformulas. Intuitively, a subformula is just a formula represented by any subtree in the formula tree.

**Definition 2.1.4.** Given a subformula  $B$  of  $A$  and another formula  $C$ , we use  $A[B \leftarrow C]$  to denote the formula obtained by substituting all occurrences of  $B$  by  $C$  in  $A$ . If  $B$  is a propositional variable, then  $A[B \leftarrow C]$  is called an instance of  $A$ .

For example, if  $A$  is  $\neg p \vee \neg(q \vee r)$ ,  $B$  is  $\neg(q \vee r)$ , and  $C$  is  $\neg q \wedge \neg r$ , then  $A[B \leftarrow C]$  is  $\neg p \vee (\neg q \wedge \neg r)$ .  $A[p \leftarrow p \wedge q] = \neg(p \wedge q) \vee \neg(q \vee r)$  is an instance of  $A$ .

## 2.2 Semantics

In the previous section, we talked about Boolean values, truth tables, and meanings of logical operators, which are, strictly speaking, semantic concepts of the logic. Now, let us talk about meaning of propositional variables and formulas.

The informal meaning of a propositional variable is the statement that symbol represents. In different applications, the same symbol can represent different statements, thus has different meanings. However, each statement can have only two truth values, i.e., true, or false. We will use these truth values of propositional variables, extending to the formulas which contain them, as the semantics of propositional logic.

### 2.2.1 Interpretations

A propositional variable may be assigned true or false. This truth value assignment is considered to be the *semantics* of the variable. For a formula, an assignment of truth values to every propositional variable is said to be an *interpretation* of the formula. If  $A$  is a formula built on the propositional variables  $V_p$ , then an interpretation of  $A$  is the function  $\sigma : V_p \rightarrow \{1, 0\}$ . It is easy to check that there are  $2^{|V_p|}$  distinct interpretations.

Suppose  $\sigma$  is an interpretation over  $V_p = \{p, q\}$  such that  $\sigma(p) = 1$  and  $\sigma(q) = 0$ . We may write  $\sigma = \{p \mapsto 1, q \mapsto 0\}$ . We may use  $(1, 0)$  for the same  $\sigma$ , assuming  $V_p$  is a list  $(p, q)$  and  $\sigma(p, q) = (1, 0)$ .

An alternative representation of an interpretation  $\sigma$  is to use a subset of  $V_p$ . Given  $\sigma$ , let  $X_\sigma = \{x \in V_p \mid \sigma(x) = 1\}$ . Then there is bijection between  $\sigma$  and  $X_\sigma$ . So we may use  $X_\sigma$  to represent  $\sigma$ . For example, if  $V_p = \{p, q\}$ , then the four interpretations over  $V_p$  is the power set of  $V_p$ :

$$2^{V_p} = \{\emptyset, \{p\}, \{q\}, \{p, q\}\}$$

Another alternative representation of  $\sigma$  is to use a set of literals. A *literal* is either a propositional variable or the negation of a variable. Let

$$Y_\sigma = X_\sigma \cup \{\neg y \in V_p, y \notin X_\sigma\}$$

then  $Y_\sigma$  is a set of literals such that every variable of  $V_p$  appears in  $Y_\sigma$  exactly once. If  $V_p = \{p, q\}$  and  $X_\sigma = \emptyset$ , then  $Y_\sigma = \{\neg p, \neg q\}$ ; if  $X_\sigma = \{p\}$ ,  $Y_\sigma = \{p, \neg q\}$ , etc. That is, by adding the negations of the missing variables in  $X_\sigma$ , we obtain such a representation  $Y_\sigma$  for  $\sigma$ .

Given  $\sigma$ , we can check the truth value of the formula under  $\sigma$ , denoted by  $\sigma(A)$ . This notation means we treat  $\sigma$  as a function from the formulas to  $\{0, 1\}$ , as

the unique extension of  $\sigma : V_p \rightarrow \{0, 1\}$ . On the other hand, the notation  $A\sigma$  is used to denote the formula obtained by substituting each propositional variable  $p$  by  $\sigma(p)$ .

**Example 2.2.1.** Given the formula  $A = p \wedge \neg q \vee r$ , and an interpretation  $\sigma = \{p \mapsto 1, q \mapsto 1, r \mapsto 0\}$ , we represent  $\sigma$  by  $\sigma = \{p, q, \neg r\}$ . Replace  $p$  by 1,  $q$  by 1, and  $r$  by 0 in  $A$ , we have  $A\sigma = 1 \wedge \neg 1 \vee 0$ . Applying the definitions of  $\neg$ ,  $\wedge$ , and  $\vee$ , we obtain the truth value 0. In this case, we say the formula is *evaluated* to 0 under  $\sigma$ , denoted by  $\sigma(A) = 0$ . Given another interpretation  $\sigma' = \{p, q, r\}$ , the same formula will be evaluated to 1, i.e.,  $\sigma'(A) = 1$ .

Recall that we use  $A[B \leftarrow C]$  to denote an instance of  $A$  where every occurrence of  $B$  is replaced by  $C$ . An interpretation  $\sigma$  is a special case of the substitution of formulas where  $B$  is a variable and  $C$  is either 1 or 0. For example, if  $\sigma = \{p, \neg q\}$ , then  $A\sigma = A[p \leftarrow 1][q \leftarrow 0]$ . Strictly speaking,  $A\sigma$  is not a propositional formula: it is the meaning of  $A$  under  $\sigma$ . To obtain  $\sigma(A)$ , we may use the algorithm *eval*.

**Algorithm 2.2.2.** The algorithm *eval* will take a formula  $A$  and an interpretation  $\sigma$ , and returns a Boolean value.

```

proc eval( $A, \sigma$ )
  if  $A = \top$  return 1;
  if  $A = \perp$  return 0;
  if  $A \in V_p$  return  $\sigma(A)$ ;
  if  $A$  is  $\neg B$  return  $\neg \text{eval}(B, \sigma)$ ;
  else  $A$  is  $(B \text{ op } C)$  return  $\text{eval}(B, \sigma) \text{ op } \text{eval}(C, \sigma)$ ;

```

**Example 2.2.3.** Let  $\sigma = \{p, \neg q\}$ . Then the execution of  $\text{eval}(p \rightarrow p \wedge \neg q, \sigma)$  will return 1:

```

eval( $p \rightarrow p \wedge \neg q, \sigma$ ) calls
  eval( $p, \sigma$ ), which returns 1; and
  eval( $p \wedge \neg q, \sigma$ ), which calls
    eval( $p, \sigma$ ), which returns 1; and
    eval( $\neg q, \sigma$ ), which calls
      eval( $q, \sigma$ ), which returns 0; and
      returns  $\neg 0$ , i.e., 1;
    returns  $1 \wedge 1$ , i.e., 1;
  returns  $1 \rightarrow 1$ , i.e., 1.

```

What *eval* does is to travel the formula tree of  $A$  bottom-up, if the node is labeled with a variable, use  $\sigma$  to get the truth value; otherwise, compute the *truth value of this node* under  $\sigma$  using the operator at that node with the truth values from its children nodes. The process of running *eval* is exactly what we do when constructing a truth table for  $p \rightarrow p \wedge \neg q$ . The truth values under  $p$  and  $q$  in the truth table give us all the interpretations  $\sigma$  and the truth values of  $A = \neg q, p \wedge \neg q$ , or  $p \rightarrow p \wedge \neg q$  are the values of  $eval(A, \sigma)$ .

$p$	$q$	$\neg q$	$p \wedge \neg q$	$p \rightarrow p \wedge \neg q$
0	0	1	0	1
0	1	0	0	1
1	0	1	1	1
1	1	0	0	0

**Theorem 2.2.4.** *Algorithm  $eval(A, \sigma)$  returns 1 iff  $\sigma(A) = 1$ , and runs in  $O(|A|)$  time, where  $|A|$  denotes the number of symbols, excluding parentheses, in  $A$ .*

$|A|$  is also the number of nodes in the formula tree of  $A$ . The proof is left as an exercise.

**Corollary 2.2.5.** *For any formulas  $A$  and  $B$ , and any interpretation  $\sigma$ , the following equations hold.*

$$\begin{aligned}\sigma(A \vee B) &= \sigma(A) \vee \sigma(B) \\ \sigma(A \wedge B) &= \sigma(A) \wedge \sigma(B) \\ \sigma(\neg A) &= \neg\sigma(A)\end{aligned}$$

The truth value of a formula in propositional logic reflects the two foundational principles of Boolean logics: the *principle of bi-valence*, which allows only two truth values, and the *principle of extensibility* that the truth value of a general formula depends only on the truth values of its parts, not on their informal meaning.

### 2.2.2 Models, Satisfiability, and Validity

Interpretations play a very important role in propositional logic and introduce many important concepts.

**Definition 2.2.6.** Given a formula  $A$  and an interpretation  $\sigma$ ,  $\sigma$  is said to be a model of  $A$  if  $eval(A, \sigma) = 1$ . If  $A$  has a model,  $A$  is said to be *satisfiable*.

Given a set  $V_p$  of propositional variables, we will use  $IV_p$  to denote the set of all interpretation over  $V_p$ . For example, if  $V_p = \{p, q\}$ , then  $IV_p = \{\{p \mapsto 1, q \mapsto 1\}, \{p \mapsto 1, q \mapsto 0\}, \{p \mapsto 0, q \mapsto 1\}, \{p \mapsto 0, q \mapsto 0\}\}$ , or abbreviated as

$IV_p = \{(1, 1), (1, 0), (0, 1), (0, 0)\}$  if  $p$  and  $q$  are understood from the context. In general,  $|IV_p| = 2^{|V_p|}$ , as each variable has two values. We may use *eval* to look for a model of  $A$  by examining every interpretation in  $IV_p$ .

**Definition 2.2.7.** *Given a formula  $A$ , let  $\mathcal{M}(A)$  be the set of all models of  $A$ . If  $\mathcal{M}(A) = \emptyset$ ,  $A$  is said to be unsatisfiable; If  $\mathcal{M}(A) = IV_p$ , i.e., every interpretation is a model of  $A$ , then  $A$  is said to be valid, or tautology. We write  $\models A$  if  $A$  is valid.*

If there is a truth table for a formula  $A$ ,  $\mathcal{M}(A)$  can be easily obtained from the truth table, as the truth values under the variables of  $A$  are all the interpretations and the models of  $A$  are those rows in the table where the truth value of  $A$  is 1.

**Example 2.2.8.** Let  $V_p = \{p, q\}$ . If  $A$  is  $\neg p \rightarrow \neg q$ , then  $\mathcal{M}(A) = \{(1, 1), (1, 0), (0, 0)\}$ . If  $B$  is  $p \vee (\neg q \rightarrow \neg p)$ , then  $\mathcal{M}(B) = IV_p$ . From  $\mathcal{M}(A)$  and  $\mathcal{M}(B)$ , we know that both  $A$  and  $B$  are satisfiable, and  $B$  is also valid, i.e.,  $\models B$  is true.

A valid formula is one which is always true in every interpretation, no matter what truth values its variables may take, that is, every interpretation is its model. The simplest example is  $\top$ , or  $p \vee \neg p$ . There are many formulas that we want to know if they are valid, and this is done by so-called *theorem proving*, either by hand or automatically.

We may think about valid formulas as capturing fundamental logical truths. For example, the transitivity property of implication states that if one statement implies a second one, and the second one implies a third, then the first implies the third. In logic, the transitivity can be expressed as the formula

$$(p \rightarrow q) \wedge (q \rightarrow r) \rightarrow (p \rightarrow r)$$

The validity of the above formula confirms the truth of this property of implication. There are many properties like this for other logical operators, such as the commutativity and associativity of  $\wedge$  and  $\vee$ , and they all can be stated by a tautology.

**Theorem 2.2.9. (Substitution)** *Let  $A$  be a tautology, so is any instance of  $A$ . That is, if  $p$  is a propositional variable in  $A$  and  $B$  is any formula, then  $A[p \leftarrow B]$  is a tautology.*

*Proof.* For any interpretation  $\sigma$  of  $A[p \leftarrow B]$ , let  $\sigma(p) = \sigma(B)$ . Applying  $\sigma$  to the formula trees of  $A$  and  $A[p \leftarrow B]$ , then the truth values of all the nodes of  $A$  must be identical to those of the corresponding nodes of  $A[p \leftarrow B]$ . Since  $A$  is a tautology, the root node of  $A$  must be 1 under  $\sigma$ , so the root node of  $A[p \leftarrow B]$  must the same truth value, i.e., 1. Since  $\sigma$  is arbitrary,  $A[p \leftarrow B]$  must be a tautology.  $\square$

The significance of the above theorem is that if we have a tautology for one variable, the tautology holds when the variable is substituted by any formula. For example, from  $p \vee \neg p$ , we have  $A \vee \neg A$  for any  $A$ . On the other hand, when we try to prove a tautology involving symbols  $A, B, \dots$ , we may treat each of these symbols as a propositional variable. For example, to prove  $(A \wedge B) \leftrightarrow (B \wedge A)$ , we may prove instead  $(p \wedge q) \leftrightarrow (q \wedge p)$ .

An unsatisfiable formula is one which does not have any model, that is, no interpretation is its model. The simplest example is  $0$ , or  $p \wedge \neg p$ . Validity and unsatisfiability are dual concepts, as stated by the following proposition.

**Proposition 2.2.10.** *A formula  $A$  is valid iff  $\neg A$  is unsatisfiable.*

From the view of theorem proving, proving the validity of a formula is as hard (or as easy) as proving its unsatisfiability.

The following theorem shows the relationship between logic and set theory.

**Theorem 2.2.11.** *For any formulas  $A$  and  $B$ ,*

$$\begin{aligned} (a) \quad \mathcal{M}(A \vee B) &= \mathcal{M}(A) \cup \mathcal{M}(B); \\ (b) \quad \mathcal{M}(A \wedge B) &= \mathcal{M}(A) \cap \mathcal{M}(B); \\ (c) \quad \mathcal{M}(\neg A) &= IV_p - \mathcal{M}(A). \end{aligned}$$

*Proof.* To show  $\mathcal{M}(A \vee B) = \mathcal{M}(A) \cup \mathcal{M}(B)$ , we need to show that, for any interpretation  $\sigma$ ,  $\sigma \in \mathcal{M}(A \vee B)$  iff  $\sigma \in \mathcal{M}(A) \cup \mathcal{M}(B)$ . By definition,  $\sigma \in \mathcal{M}(A \vee B)$  iff  $\sigma(A \vee B) = 1$ . By Corollary 2.2.5,  $\sigma(A \vee B) = \sigma(A) \vee \sigma(B)$ . So  $\sigma(A \vee B) = 1$  iff  $\sigma(A) = 1$  or  $\sigma(B) = 1$ .  $\sigma(A) = 1$  or  $\sigma(B) = 1$  iff  $\sigma(A) \vee \sigma(B) = 1$ .  $\sigma(A) \vee \sigma(B) = 1$  iff  $\sigma \in \mathcal{M}(A) \cup \mathcal{M}(B)$ .

The proof of the other two is left as exercise. □

**Example 2.2.12.** Let  $V_p = \{p, q\}$ .  $\mathcal{M}(p) = \{(1, 1), (1, 0)\}$ ,  $\mathcal{M}(q) = \{(1, 1), (0, 1)\}$ ;  $\mathcal{M}(\neg p) = IV_p - \mathcal{M}(p) = \{(0, 1), (0, 0)\}$ , and  $\mathcal{M}(\neg p \vee q) = \mathcal{M}(\neg p) \cup \mathcal{M}(q) = \{(1, 1), (0, 1), (0, 0)\}$ .

### 2.2.3 Equivalence

In natural language, one statement can be stated in different forms with the same meaning. For example, “if I won the lottery, I must be rich now.” This meaning can be also expressed “I am not rich, because I didn’t win the lottery.” Introducing  $p$  for “I won the lottery” and  $q$  for “I’m rich”, the first statement is  $p \rightarrow q$  and the second statement is  $\neg q \rightarrow \neg p$ . It happens that  $\mathcal{M}(p \rightarrow q) = \mathcal{M}(\neg q \rightarrow \neg p)$ ,

as Table 2.2.1 shows that both formulas have the same set of models. Each line of the truth table defines an interpretation; if a formula is true in the same line, that interpretation is a model. The table also reveals that  $\mathcal{M}(p \rightarrow q) = \mathcal{M}(\neg p \vee q)$ . In logic, these formulas are equivalent.

$p$	$q$	$\neg p$	$\neg q$	$p \rightarrow q$	$\neg p \vee q$	$\neg q \rightarrow \neg p$	$q \rightarrow p$
1	1	0	0	1	1	1	1
1	0	0	1	0	0	0	1
0	1	1	0	1	1	1	0
0	0	1	1	1	1	1	1

Table 2.2.1: The equivalence of the implication  $p \rightarrow q$  and its contrapositive  $\neg q \rightarrow \neg p$ .

The formula  $\neg q \rightarrow \neg p$  is called the *contrapositive* of the implication  $p \rightarrow q$ . The truth table shows that an implication and its contrapositive are equivalent: they are just different ways of saying the same thing. In contrast, the *converse* of  $p \rightarrow q$  is the formula  $q \rightarrow p$ , which has a different set of models (as shown by the second and third lines in Table 2.2.1).

**Definition 2.2.13.** *Given two formulas  $A$  and  $B$ ,  $A$  and  $B$  are said to be logically equivalent if  $\mathcal{M}(A) = \mathcal{M}(B)$ , denoted by  $A \equiv B$ .*

$A \equiv B$  means that, for every interpretation  $\sigma$ ,  $\sigma(A) = \sigma(B)$ , so  $\sigma \in \mathcal{M}(A)$  iff  $\sigma \in \mathcal{M}(B)$ . The relation  $\equiv$  over formulas is an equivalence relation as  $\equiv$  is reflexive, symmetric, and transitive. The relation  $\equiv$  is also a congruence relation as  $A \equiv C$  and  $B \equiv D$  imply that  $\neg A \equiv \neg C$  and  $A \circ B \equiv C \circ D$  for any binary operator  $\circ$ . This property allows us to obtain an equivalent formula by replacing a part of the formula by an equivalent one. The relation  $\equiv$  plays a very important role in logic as it is used to simplify formulas, to convert formulas into equivalent normal forms, or to provide alternative definitions for logical operators.

In arithmetic one writes simply  $s = t$  to express that the terms  $s$ ,  $t$  represent the same function. For example,  $(x + y)^2 = x^2 + 2xy + y^2$  expresses the equality of values of the left-and right-hand terms for all  $x, y$ . In propositional logic, however, we use the equality sign like  $A = B$  only for the syntactic identity of the formulas  $A$  and  $B$ . Therefore, the equivalence of formulas must be denoted differently, such as  $A \equiv B$ .

**Theorem 2.2.14.** *For any formulas  $A$ ,  $B$ , and  $C$ , where  $B$  is a subformula of  $A$ , and  $B \equiv C$ , then  $A \equiv A[B \leftarrow C]$ .*

*Proof.* For any interpretation  $\sigma$ ,  $\sigma(B) = \sigma(C)$ , since  $B \equiv C$ . Apply  $\sigma$  to the formula trees of  $A$  and  $A[B \leftarrow C]$  and compare the truth values of all corresponding nodes of the two trees, ignoring the proper subtrees of  $B$  and  $C$ . Since  $\sigma(B) = \sigma(C)$ , they must have the same truth values, that is,  $\sigma(A) = \sigma(A[B \leftarrow C])$ .  $\square$

The equivalence relation is widely used to simplify formulas and has real practical importance in computer science. Formula simplification in software can make a program easier to read and understand. Simplified programs may also run faster, since they require fewer operators. In hardware, simplifying formulas can decrease the number of logic gates on a chip because digital circuits can be described by logical formulas. Minimizing the logical formulas corresponds to reducing the number of gates in the circuit. The payoff of gate minimization is potentially enormous: a chip with fewer gates is smaller, consumes less power, has a lower defect rate, and is cheaper to manufacture.

Suppose a formula  $A$  contains  $k$  propositional variable, then  $A$  can be reviewed as one of Boolean function  $f : \{1, 0\}^k \rightarrow \{1, 0\}$ . For example,  $\neg p \vee q$  contains two variables and can be regarded as a Boolean function  $f(p, q)$ . The truth table (Table 2.2.1) reveals that  $f(p, q)$ , i.e.,  $\neg p \vee q$ , always has the same truth value as  $p \rightarrow q$ , so  $f$  and  $\rightarrow$  are the same function. As a result, we may use  $\neg p \vee q$  as the definition of  $\rightarrow$ . As another example, the if-then-else function,  $ite : \{1, 0\}^3 \rightarrow \{1, 0\}$ , can be defined by  $ite(1, B, C) = B$  and  $ite(0, B, C) = C$ , instead of using a truth table of eight lines.

The following proposition lists many equivalent pairs of the formulas.

**Proposition 2.2.15.** *Assume  $A$ ,  $B$ , and  $C$  are any formulas.*

$$\begin{array}{ll}
A \vee A \equiv A; & A \wedge A \equiv A; \\
A \vee B \equiv B \vee A; & A \wedge B \equiv B \wedge A; \\
(A \vee B) \vee C \equiv A \vee (B \vee C); & (A \wedge B) \wedge C \equiv A \wedge (B \wedge C); \\
A \vee \perp \equiv A; & A \wedge \top \equiv A; \\
A \vee \neg A \equiv \top; & A \wedge \neg A \equiv \perp; \\
A \vee \top \equiv \top; & A \wedge \perp \equiv \perp; \\
A \vee (B \wedge C) \equiv (A \vee B) \wedge (A \vee C); & A \wedge (B \vee C) \equiv (A \wedge B) \vee (A \wedge C); \\
\neg \top \equiv \perp; \quad \neg \perp \equiv \top; & \neg \neg A \equiv A; \\
\neg(A \vee B) \equiv \neg A \wedge \neg B; & \neg(A \wedge B) \equiv \neg A \vee \neg B; \\
A \rightarrow B \equiv \neg A \vee B; & A \rightarrow B \equiv \neg B \rightarrow \neg A; \\
A \oplus B \equiv (A \vee B) \wedge (\neg A \vee \neg B); & A \leftrightarrow B \equiv (A \vee \neg B) \wedge (\neg A \vee B).
\end{array}$$

One way to show that two formulas are equivalent if we can prove the validity of a formula, as given by the following proposition:

**Proposition 2.2.16.** *For any formulas  $A$  and  $B$ ,  $A \equiv B$  iff  $\models A \leftrightarrow B$ ;  $\models A$  iff  $A \equiv \top$ .*

*Proof.*  $A \equiv B$  means, for any interpretation  $\sigma$ ,  $\sigma(A) = \sigma(B)$ , which means  $\sigma(A \leftrightarrow B) = 1$ , i.e., any  $\sigma$  is a model of  $A \leftrightarrow B$ . The proof of the second part is omitted.  $\square$

The above proposition shows the relationship between  $\equiv$ , which is a semantic notation, and  $\leftrightarrow$ , which is a syntactical symbol. Moreover, Theorem 2.2.9 allows us to prove the equivalences using the truth table method, treating  $A$ ,  $B$ , and  $C$  as propositional variables. For example, from  $(p \rightarrow q) \equiv (\neg p \vee q)$ , we know  $(p \rightarrow q) \leftrightarrow (\neg p \vee q)$  is tautology; from  $(p \rightarrow q) \leftrightarrow (\neg p \vee q)$ , we know  $(A \rightarrow B) \leftrightarrow (\neg A \vee B)$  is a tautology for any formulas  $A$  and  $B$ , thus  $(A \rightarrow B) \equiv (\neg A \vee B)$ .

#### 2.2.4 Entailment

In natural languages, given a set of premises, we like to know what conclusions can be drawn from the premises. If we represent these premises by a single formula, it will be the conjunction of the premises. Since  $\wedge$  is commutative and associative, a conjunction of formulas can be conveniently written as a set  $S$  of formulas. That is, we assume that  $A_1 \wedge A_2 \wedge \cdots \wedge A_n$  is equivalent to  $S = \{A_1, A_2, \cdots, A_n\}$ . For example,  $\{A, B\} \equiv A \wedge B$ .

To catch the relation between the premises and the conclusion in logic, we have the notion of *entailment*.

**Definition 2.2.17.** *Given two formulas  $A$  and  $B$ , we say  $A$  entails  $B$ , or  $B$  is a logical consequence of  $A$ , denoted by  $A \models B$ , if  $\mathcal{M}(A) \subseteq \mathcal{M}(B)$ .*

Thus,  $\models A$  and  $\top \models A$  have the same meaning, i.e.,  $A$  is valid.

The above definition allows many irrelevant formulas to be logical consequences of  $A$ , including all tautologies and logically equivalent formulas. Despite of this irrelevant relationship between  $A$  and  $B$ , the concept of entailment is indispensable in logic. For instance, an inference rule is a pattern which takes formulas as premises, and returns a formula as conclusion. To check the soundness of the inference rule, we let the premises be represented by formula  $A$  and the conclusion represented by  $B$ , and check if  $A \models B$ , i.e.,  $B$  is a logical consequence of  $A$ . If the number of variables is small, we may use the truth table method to check that if every model of  $A$  remains to be a model for  $B$ .

**Example 2.2.18.** The modus ponens rule says that given the premises  $p \rightarrow q$  and  $p$ , then we draw the conclusion  $q$ . Let  $A$  be  $(p \rightarrow q) \wedge p$ , then  $\mathcal{M}(A) = \{(1, 1)\}$ , and  $\mathcal{M}(q) = \{(1, 1), (0, 1)\}$ . Since  $\mathcal{M}(A) \subseteq \mathcal{M}(q)$ , so  $q$  is a logical consequence of  $A$ , or  $A \models B$ .

**Proposition 2.2.19.** *The relation  $\models$  is transitive, that is, if  $A \models B$  and  $B \models C$ , then  $A \models C$ .*

The proof is left as an exercise.

**Definition 2.2.20.** *Given a formula  $A$ , the set  $\{B \mid A \models B\}$  is called the theory of  $A$  and is denoted by  $\mathcal{T}(A)$ . Every member of  $\mathcal{T}(A)$  is called a theorem of  $A$ .*

In the above definition,  $A$  can be a set of formulas to represent the conjunction of its members.

**Proposition 2.2.21.** *The following three statements are equivalent: (a)  $A \models B$ ; (b)  $\mathcal{M}(A) \subseteq \mathcal{M}(B)$ ; and (c)  $\mathcal{T}(B) \subseteq \mathcal{T}(A)$ .*

*Proof.* By definition,  $A \models B$  iff  $\mathcal{M}(A) \subseteq \mathcal{M}(B)$ . Also by definition,  $B \models C$  for any  $C \in \mathcal{T}(B)$ . If  $A \models B$ , then  $A \models C$  because of the transitivity of  $\models$ , so  $C \in \mathcal{T}(A)$ . If  $\mathcal{T}(B) \subseteq \mathcal{T}(A)$ , since  $B \in \mathcal{T}(B)$ , so  $B \in \mathcal{T}(A)$ , thus  $A \models B$ .  $\square$

**Corollary 2.2.22.** *For any formula  $A$ ,  $\mathcal{T}(\top) \subseteq \mathcal{T}(A) \subseteq \mathcal{T}(\perp)$ .*

*Proof.* This holds because  $\perp \models A$  and  $A \models \top$ . Note that  $\mathcal{T}(\top)$  contains every tautologies and  $\mathcal{T}(\perp)$  contains every formula.  $\square$

**Corollary 2.2.23.** *For any formulas  $A$  and  $B$ , the following statements are equivalent: (a)  $A \equiv B$ ; (b)  $\mathcal{M}(A) = \mathcal{M}(B)$ ; and (c)  $\mathcal{T}(A) = \mathcal{T}(B)$ .*

*Proof.* Since  $A \equiv B$  iff  $A \models B$  and  $B \models A$ , this corollary holds by applying Proposition 2.2.21 twice.  $\square$

Intuitively, we may regard a formula  $A$  as a constraint to mark some interpretations as “no-model” and  $\mathcal{M}(A)$  records those interpretations which are models of  $A$ . The stronger the constraint, the less the set of models. The strongest constraint is  $\perp$  (no interpretations left) and the weakest constraint is  $\top$  (no constraints). On the other hand,  $\mathcal{T}(A)$  collects all constraints which do not mark any model in  $\mathcal{M}(A)$  as “no-model”.

**Proposition 2.2.24.** *Let  $A$  and  $B$  be sets of formulas. If  $A \subseteq B \subseteq \mathcal{T}(A)$ , then  $\mathcal{M}(A) = \mathcal{M}(B)$ .*

*Proof.* Let  $C = B - A$ . Since  $B \subseteq \mathcal{T}(A)$ , so  $C \subseteq \mathcal{T}(A)$ , or  $A \models C$ , which implies  $\mathcal{M}(A) \subseteq \mathcal{M}(C)$ . Thus  $\mathcal{M}(B) = \mathcal{M}(A \cup C) = \mathcal{M}(A \wedge C) = \mathcal{M}(A) \cap \mathcal{M}(C) = \mathcal{M}(A)$ , because  $\mathcal{M}(A) \subseteq \mathcal{M}(C)$ .  $\square$

**Corollary 2.2.25.** *If  $C \in \mathcal{T}(A)$ , then  $\mathcal{T}(A) = \mathcal{T}(A \wedge C)$ .*

*Proof.* Let  $B = A \cup \{C\}$ , then  $A \subseteq B \subseteq \mathcal{T}(A)$ , so  $\mathcal{M}(A) = \mathcal{M}(B) = \mathcal{M}(A \wedge C)$ . Hence  $\mathcal{T}(A) = \mathcal{T}(A \wedge C)$  by Corollary 2.2.23.  $\square$

**Corollary 2.2.26.** *If  $\perp \in \mathcal{T}(A)$ , then  $A$  is unsatisfiable.*

*Proof.* Let  $B = A \cup \{\perp\}$ , then  $A \subseteq B \subseteq \mathcal{T}(A)$ , so  $\mathcal{M}(A) = \mathcal{M}(B) = \mathcal{M}(A \wedge \perp) = \mathcal{M}(\perp) = \emptyset$ .  $\square$

One task of theorem proving is to show that, given a set  $A$  of axioms,  $B$  is a theorem of  $A$ , that is,  $B \in \mathcal{T}(A)$ . If we use the operator  $\rightarrow$ , to show that  $B$  is a logical consequence of  $A$  by showing  $(A \rightarrow B) \in \mathcal{T}(\top)$ , as given by the following proposition.

**Proposition 2.2.27.** *For any formulas  $A$  and  $B$ ,  $A \models B$  iff  $\models A \rightarrow B$ .*

*Proof.*  $A \models B$  means, for any interpretation  $\sigma$ , if  $\sigma(A) = \top$ , then  $\sigma(B)$  must be  $\top$ ; if  $\sigma(A) = \perp$ , we do not care about the value of  $\sigma(B)$ . The two cases can be combined as  $\sigma(A) \rightarrow \sigma(B) = 1$ , or  $\sigma(A \rightarrow B) = 1$ , which means  $A \rightarrow B$  is valid, because  $\sigma$  is any interpretation.  $\square$

The above proposition shows the relationship between  $\models$  (a semantic notation) and  $\rightarrow$  (a syntactical symbol). The closure of the entailment relation under substitution generalizes the fact that from  $p \wedge q \models p$ , all entailment of the form  $A \wedge B \models A$  arise from substituting  $A$  for  $p$  and  $B$  for  $q$ , just like the closure of the equivalence relation under substitution.

### 2.2.5 Theorem Proving and the SAT Problem

The concepts of validity, unsatisfiability, or entailment give rise difficult problems of computation. These three problems have the same degree of difficulty and can be represented by theorem proving, which asks to check if a formula  $A \in \mathcal{T}(\top)$  ( $A$  is valid), or  $\perp \in \mathcal{T}(A)$  ( $A$  is unsatisfiable), or  $A$  is in  $\mathcal{T}(B)$  ( $B$  entails  $A$ ). Theorem proving by truth table runs out of steam pretty quickly: a formula with  $n$  variables has a truth table of  $2^n$  interpretations, so the effort grows exponentially with the number of variables. For a formula with just 30 variables, that is already over a billion interpretations to check!

The general problem of deciding whether a formula is  $A$  satisfiable is called the SAT problem. One approach to SAT is to construct a truth table and check whether or not a  $\top$  ever appears, but as with testing validity, this approach quickly

bogs down for formulas with many variables because truth tables grow exponentially with the number of variables.

The good news is that SAT belongs to the class of  $NP$  problems, where the solutions of these problems can be efficiently checked. For example, given an interpretation  $\sigma$ , we can check in linear time (in terms of the size of  $A$ ) if  $\sigma$  is a model of  $A$ , thus showing  $A$  is satisfiable. This can be done using  $eval(A, \sigma)$ . If a problem can be solved in polynomial time (with respect to the size of inputs), it belongs to the class  $P$ . SAT is not known to be in  $P$ , and theorem proving is not known to be in  $NP$ .

The bad news is that SAT is known to be the hardest problem in  $NP$ . That is, if there exists a polynomial-time algorithm to solve SAT, then every problem in  $NP$  can be solved in polynomial-time. These problems include many other important problems involving scheduling, routing, resource allocation, and circuit verification across multiple disciplines including programming, algebra, finance, and political theory. This is the famous result of Cook and Levin who proved that SAT is the first NP-complete problem. Does this mean that we cannot find an algorithm which takes  $O(n^{100})$  time, where  $n$  is the size of  $A$ , to decide  $A$  is satisfiable? No one knows. This is the most famous open problem in computer science:  $P = NP$ ? It is also one of the seven Millennium Problems: the Clay Institute will award you \$1,000,000 if you solve the  $P = NP$  problem. If  $P = NP$ , then SAT will be in  $P$ ; theorem proving will be also in  $P$ . If  $P \neq NP$ , then theorem proving won't be in  $NP$ : it belongs the complement of the NP-complete problems. An awful lot hangs on the answer and the general consensus is  $P \neq NP$ .

In the last twenty years, there have been exciting progress on SAT-solvers for practical applications like digital circuit verification. These programs find satisfying assignments with amazing efficiency even for formulas with thousands of variables. In the next two chapters, we will study the methods for theorem proving and SAT-solvers.

### 2.3 Normal Forms

By “normal form”, we mean every formula can be converted into an equivalent formula with certain syntactical restriction on the use and position of Boolean operators. If all equivalent formulas has a unique normal form, that normal form is also called “canonical form”.

Since it is expensive to use truth table to prove the equivalence relation when the number of variables is big, an alternative approach is to use algebra to prove equivalence. A lot of different operators may appear in a propositional formula, so a useful first step is to get rid of all but three:  $\wedge, \vee$  and  $\neg$ . This is easy because

each of the operators is equivalent to a simple formula using only these three. For example,  $A \rightarrow B$  is equivalent to  $\neg A \vee B$ .

Using the equivalence relations discussed in the previous section, any propositional formula can be proved equivalent to a canonical form. What has this got to do with equivalence? That's easy: to prove that two formulas are equivalent, convert them both to canonical forms over the set of variables that appear in any formula. Now if two formulas are equivalent to the same canonical form, then the formula are certainly equivalent. Conversely, if two formulas are equivalent, they will have the same canonical form. We can also use canonical form to show a formula is valid or unsatisfiable: it is valid if its canonical form is  $\top$ ; it is unsatisfiable if its canonical form is  $\perp$ .

In this section, we present four normal forms:

- negation normal form (NNF);
- conjunctive normal form (CNF);
- disjunctive normal form (DNF); and
- ITE normal form (INF).

Of the four normal forms, there exist canonical forms for the last three. However, the canonical forms derived from DNF or CNF are essentially copies of truth tables, thus not effective enough for showing the equivalence relation. Only the last one, INF, may have a compact size so that equivalence of two formulas can be conveniently established in this canonical form.

We will define a normal form by giving a set of transformation rules, that is a set of pairs of formulas, one is the left hand side and the other is the right hand side of the rule. Most importantly, the left side of each rule is logically equivalent to its right side. Each rule deals with a formula which is not in normal form and the application of the rule on the formula will result in a new formula closer to normal form. To convert any formula into a formula in normal form, we pick a rule, matching the left side of the rule to a subformula and replacing the subformula by the right side of the rule.

A general procedure for obtaining normal forms goes as follows:

1. If the formula is already in normal form then stop with success.
2. Otherwise it contains a subformula violating the normal form criteria.
3. Pick such a subformula and find the left hand side of a rule that matches.

4. Replace the subformula by the right hand side of the rule and continue.

### 2.3.1 Negation Normal Form (NNF)

A formula  $A$  is said to be a *literal* if  $A$  is a variable (positive literal) or  $A$  is  $\neg p$  (negative literal), where  $p$  is a variable.

**Definition 2.3.1.** A propositional formula  $A$  is a negation normal form (NNF) if the negation symbol  $\neg$  appears in literals (as subformulas of  $A$ ).

In a NNF, if we replace negative literals by a variable, there will be no  $\neg$  in the resulting formula.

**Proposition 2.3.2.** Every propositional formula can be transformed into an equivalent NNF.

The proof is done by the following algorithm, which uses only equivalence relations in the transformation.

Equivalence rules for obtaining NNF

**Algorithm 2.3.3.** The algorithm to convert a formula into NNF takes the following three groups of rules:

1. Use the following equivalences to remove  $\rightarrow$ ,  $\oplus$ , and  $\leftrightarrow$  from the formula.

$$\begin{aligned} A \rightarrow B &\equiv \neg A \vee B; \\ A \oplus B &\equiv (A \vee B) \wedge (\neg A \vee \neg B); \\ A \leftrightarrow B &\equiv (A \vee \neg B) \wedge (\neg A \vee B). \end{aligned}$$

2. Use the de Morgan laws to push  $\neg$  down to variables:

$$\begin{aligned} \neg\neg A &\equiv A; \\ \neg(A \vee B) &\equiv \neg A \wedge \neg B; \\ \neg(A \wedge B) &\equiv \neg A \vee \neg B. \end{aligned}$$

3. Use  $A \vee \perp \equiv A$ ,  $A \vee \top \equiv \top$ ,  $A \wedge \perp \equiv \perp$ , and  $A \wedge \top \equiv A$ , to simplify the formula.

For example,  $\top$  and  $\neg p \vee q$  are NNF.  $\neg(\neg p \vee q)$  is not NNF; we may transform it into  $p \wedge \neg q$ , which is NNF. If the input formula contains only  $\neg$ ,  $\vee$ ,  $\wedge$ , and  $\rightarrow$ , it takes only linear time to obtain NNF by working recursively top-down. If  $\leftrightarrow$  or  $\oplus$  are present, the computing time and the size of the output formula may be exponential in terms of the original size, because the arguments of  $\leftrightarrow$  and  $\oplus$  are duplicated during the transformation.

**Definition 2.3.4.** Given a NNF  $A$ , if we replace every occurrence of  $\wedge$  by  $\vee$ ,  $\vee$  by  $\wedge$ , positive literal by its negation, and negative literal  $\neg p$  by its counterpart  $p$ , we obtain a formula  $B$ , which is called the dual of  $A$ , denoted by  $B = \text{dual}(A)$ .

For example, if  $A$  is  $p \wedge (q \vee (\neg p \wedge \neg r))$ , then  $\text{dual}(A) = \neg p \vee (\neg q \wedge (p \vee r))$ . Both  $A$  and  $\text{dual}(A)$  are NNF.

**Proposition 2.3.5.** If  $A$  is in NNF, then  $\text{dual}(A)$  is also in NNF and  $\text{dual}(A) \equiv \neg A$ .

The proposition can be easily proved by DeMorgan's law.

### 2.3.2 Conjunctive Normal Form (CNF)

**Definition 2.3.6.** In propositional logic, a clause is either  $\perp$  (the empty clause), a literal (unit clause) or a disjunction of literals.

A clause is also called *maxterm* in the community of circuit designs. Since  $\vee$  is commutative and associative, we may represent a clause by a list of literals. We may use  $A \vee A \equiv A$  to remove duplicate literals in a clause.

**Definition 2.3.7.** A formula is said to be a conjunctive normal form (CNF) if it is a conjunction of clauses.

A CNF is also called a *product of sums* (POS) expression. The idea of conjunctive normal form is that each type of connective appears at a distinct height in the formula. The negation signs,  $\neg$ , are lowest, with only propositions as their arguments. The disjunction signs,  $\vee$ , are in the middle, with only literals as their arguments. And the conjunction signs,  $\wedge$ , are at the top, taking the disjunctions as their arguments. Of course, some operators may be missing in CNF. For example,  $p \wedge q$  is a CNF, which contains two unit clauses:  $p$  and  $q$ .

Equivalence rules for obtaining CNF

- stage 1: The equivalence rules for obtaining NNF.
- stage 2:  $A \vee (B \wedge C) \equiv (A \vee B) \wedge (A \vee C)$ , where the arguments of  $\wedge$  and  $\vee$  can be switched.

**Proposition 2.3.8.** Every propositional formula can be transformed into an equivalent CNF.

To show the proposition is true, we first convert the formula into NNF, and then apply the distribution law to move  $\vee$  under  $\wedge$ . This process will terminate because we have only a finite number of  $\vee$  and  $\wedge$ .

**Example 2.3.9.** To convert  $p_1 \leftrightarrow (p_2 \leftrightarrow p_3)$  into CNF, we first use  $A \leftrightarrow B \equiv (\neg A \vee B) \wedge (A \vee \neg B)$  to remove the second  $\leftrightarrow$ :

$$p_1 \leftrightarrow ((\neg p_2 \vee p_3) \wedge (p_2 \vee \neg p_3)),$$

and then remove the first  $\leftrightarrow$ :

$$(\neg p_1 \vee ((\neg p_2 \vee p_3) \wedge (p_2 \vee \neg p_3))) \wedge (p_1 \vee \neg((\neg p_2 \vee p_3) \wedge (p_2 \vee \neg p_3))).$$

Its NNF is

$$(\neg p_1 \vee ((\neg p_2 \vee p_3) \wedge (p_2 \vee \neg p_3))) \wedge (p_1 \vee ((p_2 \wedge \neg p_3) \vee (\neg p_2 \wedge p_3))).$$

We then obtain the following non-tautology clauses:

$$\neg p_1 \vee \neg p_2 \vee p_3, \quad \neg p_1 \vee p_2 \vee \neg p_3, \quad p_1 \vee p_2 \vee p_3, \quad p_1 \vee \neg p_2 \vee \neg p_3.$$

Since  $\vee$  is commutative and associative and  $X \vee X \equiv X$ , a clause can be represented by a set of literals. In this case, we use  $|$  for  $\vee$  and  $\bar{p}$  for  $\neg p$  in a clause. For example, the above clauses will be displayed as

$$(\bar{p}_1 | \bar{p}_2 | p_3), \quad (\bar{p}_1 | p_2 | \bar{p}_3), \quad (p_1 | p_2 | p_3), \quad (p_1 | \bar{p}_2 | \bar{p}_3).$$

**Definition 2.3.10.** *If a clause includes every variable exactly once, the clause is said to be full. A CNF is called a full CNF if every clause of the CNF is full.*

If a clause  $C$  misses a variable  $r$ , we may replace  $C$  by two clauses:  $\{(C | r), (C | \bar{r})\}$ . Repeating this process for every missing variable, we will obtain an equivalent full CNF. Transforming a set of clauses into an equivalent full CNF is sound because for any formula  $A$  and  $B$ , the following relation is true:

$$A \equiv (A \vee B) \wedge (A \vee \neg B),$$

which can be proved by truth table.

Equivalence rules for obtaining full CNF

- stage 1: The equivalence rules for NNF.
- stage 2:  $A \vee (B \wedge C) \equiv (A \vee B) \wedge (A \vee C)$ , where the arguments of  $\wedge$  and  $\vee$  can be switched.
- stage 3:  $A \equiv (A \vee p) \wedge (A \vee \neg p)$ , if  $A$  is a clause and  $p$  does not appear in  $A$ .

**Proposition 2.3.11.** *Every propositional formula can be transformed into an equivalent and unique full CNF (up to the commutativity and associativity of  $\wedge$  and  $\vee$ ). That is, full CNF is a canonical form.*

Given a CNF  $A$ , each clause can be regarded as a constrain on what interpretations cannot be a model of  $A$ . For example, if  $A$  contains the unit clause  $(\neg p)$ , then all the interpretations where  $p \mapsto 1$  cannot be a model of  $A$  and these are half of all interpretations. If  $A$  contains a binary clause  $(p \mid q)$ , then all interpretations where  $p \mapsto 0, q \mapsto 0$  (there are a quarter of such interpretations) cannot be a model of  $A$ . If  $A$  has a clause  $C$  containing every variable exactly once, then this clause removes only interpretation, i.e.,  $\neg C$ , as a model of  $A$ . If  $A$  has  $m$  models over  $V_p$ , i.e.,  $|\mathcal{M}(A)| = m$ , then a full CNF of  $A$  contains  $2^{|V_p|} - m$  clauses, as each clause removes exactly one interpretation as a model of  $A$ . That is exactly the idea for a proof of the above proposition.

The concept of full CNF is useful for theoretic proofs. Unfortunately, the size of a full CNF is exponential in terms of number of variables, too big for any practical application. CNFs are useful to specify many problems in practice, because practical problems are often specified by a set of constraints. Each constraint can be specified by a set of clauses, and the conjunction of all clauses from each constraint gives us a complete specification of the problem.

### 2.3.3 Disjunctive Normal Form (DNF)

If we regard  $\wedge$  as multiplication, and  $\vee$  as addition, then a CNF is a “product of sums. On the other hand, a “sum of products” is a disjunctive normal form (DNF). A formula  $A$  is said to be a *minterm* if  $A$  is either  $\top$ , a literal, or a conjunction of literals. Product is also called *product* or *product term* in the community of circuit designs.

Like clauses, we remove duplicates in a minterm and represent them by a set of literals.

**Definition 2.3.12.** *A formula  $A$  is a disjunctive normal form (DNF) if  $A$  is a disjunction of products.*

A DNF is also called a *sum of products* (SOP) expression. A formula like  $p \wedge \neg q$  or  $p \vee \neg q$  can be both a CNF and a DNF.

Equivalence rules for obtaining DNF

- stage 1: The equivalence rules for NNF.
- stage 2:  $A \wedge (B \vee C) \equiv (A \wedge B) \vee (A \wedge C)$ , where the arguments of  $\wedge$  and  $\vee$  can be switched.

**Proposition 2.3.13.** *Every propositional formula can be transformed into an equivalent DNF.*

To show the proposition is true, we first convert the formula into NNF in the first stage, and then apply the distribution law in the second stage to move  $\wedge$  under  $\vee$ .

Once a formula  $A$  is transformed into DNF, it is very easy to check if  $A$  is satisfiable: If DNF contains a product in which each variable appears at most once, then  $A$  is satisfiable, because we can simply assign 1 to each literal of this product and the product becomes true. Since it is a hard problem to decide if a formula is satisfiable, it is also hard to obtain DNF, due to its large size.

Like CNF, DNF is not a canonical form. For example, the following equivalence is true:

$$p \wedge \neg q \vee \neg p \wedge \neg r \vee q \wedge r \equiv \neg p \wedge q \vee p \wedge r \vee \neg q \wedge \neg r$$

However, we cannot convert the both sides to the same formula.

**Definition 2.3.14.** *A product is full if it contains every variable exactly once. A DNF is full if every product in the DNF is full.*

A full minterm in a DNF  $A$  is equivalent to a model of  $A$ . For example, if  $V_p = \{p, q, r\}$ , then the product  $p \wedge q \wedge \neg r$  specifies the model  $\{p \mapsto 1, q \mapsto 1, r \mapsto 0\}$ . Obviously, A full DNF  $A$  has  $m$  full products iff  $A$  has  $m$  models. Like full CNF, every propositional formula is equivalent to a unique full DNF, that is, full DNF is a canonical form.

Equivalence Rules for obtaining full DNF

- stage 1: The equivalence rules for NNF.
- stage 2:  $A \wedge (B \vee C) \equiv (A \wedge B) \vee (A \wedge C)$ , where the arguments of  $\wedge$  and  $\vee$  can be switched.

- stage 3:  $A \equiv (A \wedge p) \vee (A \wedge \neg p)$ , if  $A$  is a product and  $p$  does not appear in  $A$ .

**Proposition 2.3.15.** *Every propositional formula can be transformed into an equivalent and unique full DNF (up to the commutativity and associativity of  $\wedge$  and  $\vee$ ).*

#### 2.3.4 Full DNF and Full CNF from Truth Table

A truth table is often used to define a Boolean function. In the truth table, the truth value of the function is given for every interpretation of the function. Each interpretation corresponds to a full minterm and the function can be defined as a sum (disjunction) of all the minterms for which the function is true in the corresponding interpretation.

**Example 2.3.16.** Suppose we want to define a Boolean function  $f : B^3 \rightarrow B$ , where  $B = \{0, 1\}$  and the truth value of  $f(a, b, c)$  is given in the following truth table. Following the convention of EDA, we will use  $+$  for  $\vee$ ,  $\cdot$  for  $\wedge$ , and  $\bar{A}$  for  $\neg A$ .

$a$	$b$	$c$	$f$	minterms	clauses
0	0	0	0	$m_0 = \bar{a}\bar{b}\bar{c}$	$M_0 = a + b + c$
0	0	1	1	$m_1 = \bar{a}\bar{b}c$	$M_1 = a + b + \bar{c}$
0	1	0	0	$m_2 = \bar{a}b\bar{c}$	$M_2 = a + \bar{b} + c$
0	1	1	1	$m_3 = \bar{a}bc$	$M_3 = a + \bar{b} + \bar{c}$
1	0	0	0	$m_4 = a\bar{b}\bar{c}$	$M_4 = \bar{a} + b + c$
1	0	1	1	$m_5 = a\bar{b}c$	$M_5 = \bar{a} + b + \bar{c}$
1	1	0	0	$m_6 = ab\bar{c}$	$M_6 = \bar{a} + \bar{b} + c$
1	1	1	0	$m_7 = abc$	$M_7 = \bar{a} + \bar{b} + \bar{c}$

For each interpretation  $\sigma$  of  $(a, b, c)$ , we define a minterm  $m_i$ , where  $i$  is the decimal value of  $\sigma$ , when  $\sigma$  is viewed as a binary number. For each  $m_i$ , we also define a clause  $M_i$  and the relation between  $m_i$  and  $M_i$  is that  $\neg m_i \equiv M_i$ . For the function  $f$  defined in the truth table,  $f = m_1 + m_3 + m_5$ , or  $f = \bar{a}\bar{b}c + \bar{a}bc + a\bar{b}c$ .

In the above example, for each model of  $f$ , we create a minterm (product term) in its full DNF: If the variable is true in the model, the positive literal of this variable is in the product; if a variable is false, the negative literal of the variable is in the product.  $f$  is then defined by a full DNF consisting of these minterms. In other words, every Boolean function can be defined by a DNF. We can also define the same function by a full CNF. This can be processed as follows: At first, we define the negation of  $f$  by a full DNF from the truth table:

$$\bar{f} = m_0 + m_2 + m_4 + m_6 + m_7$$

Apply the negation on the both sides of the above equation, we have

$$f = M_0 \cdot M_2 \cdot M_4 \cdot M_6 \cdot M_7$$

where  $\cdot$  stands for  $\wedge$  and  $\overline{m_i} \equiv M_i$ . In other words,  $f$  can be defined by a full CNF. Note that  $M_0 \cdot M_2 \cdot M_4 \cdot M_6 \cdot M_7$  is the dual of  $m_0 + m_2 + m_4 + m_6 + m_7$ . In general, the dual of DNF is CNF and the dual of CNF is DNF.

Once we know the principle, we can construct a full CNF directly from the truth table of a formula. That is, for each non-model interpretation in the truth table, we create one full clause: If a variable is true in the interpretation, the negative literal of this variable is in the clause; if a variable is false, the positive literal of the variable is in the clause.

**Example 2.3.17.** The truth table of  $(p \vee q) \rightarrow \neg r$  has eight lines and three interpretations of  $\langle p, q, r \rangle$  are non-models:  $\langle 1, 1, 1 \rangle$ ,  $\langle 1, 0, 1 \rangle$ , and  $\langle 0, 1, 1 \rangle$ . From these three interpretations, we obtain three full clauses:

$$(\neg p \vee \neg q \vee \neg r) \wedge (\neg p \vee q \vee \neg r) \wedge (p \vee \neg q \vee \neg r)$$

### 2.3.5 Binary Decision Diagram (BDD)

Let  $ite : \{1, 0\}^3 \rightarrow \{1, 0\}$  be the "if-then-else" operator such that, for any formula  $A$  and  $B$ ,  $ite(1, A, B) \equiv A$  and  $ite(0, A, B) \equiv B$ . In fact,  $ite$ , 1 and 0 can be used to represent any logical operator. For example,  $\neg y \equiv ite(y, 0, 1)$ . The following table shows that every binary logical operator can be represented by  $ite$ . Following the convention on BDD, we will use 1 for  $\top$  and 0 for  $\perp$  in every propositional formula.

**Definition 2.3.18.** A formula is said to be an *ite normal form (INF)* if the formula uses logical operators no more than  $ite$ , 1, and 0, and, for each occurrence of  $ite(C, A, B)$  in the formula,  $C$  is a propositional variable, and  $A$  and  $B$  are INFs.

By definition, 1 and 0 are INFs, but  $p$  is not;  $ite(p, 1, 0)$  is.

**Proposition 2.3.19.** Every propositional formula can be transformed into an equivalent INF.

For example, let  $A$  be  $\neg p \vee q \wedge \neg r$ , its equivalent INF is  $ite(p, ite(q, ite(r, 0, 1), 0), 1)$ . To prove the above proposition, we provide the algorithm below to do the transformation.

output	formula	ite formula
0000	0	0
0001	$x \wedge y$	$ite(x, y, 0)$
0010	$x \wedge \neg y$ ( $x > y$ )	$ite(x, \neg y, 0)$
0011	$x$	$ite(x, 1, 0)$
0100	$\neg x \wedge y$ ( $x < y$ )	$ite(x, 0, y)$
0101	$y$	$ite(y, 1, 0)$
0110	$x \oplus y$	$ite(x, \neg y, y)$
0111	$x \vee y$	$ite(x, 1, y)$
1000	$x \downarrow y$	$ite(x, 0, \neg y)$
1001	$x \leftrightarrow y$	$ite(x, y, \neg y)$
1010	$\neg y$	$ite(y, 0, 1)$
1011	$x \vee \neg y$ ( $x \geq y$ )	$ite(x, 1, \neg y)$
1100	$\neg x$	$ite(x, 0, 1)$
1101	$\neg x \vee y$ ( $x \leq y$ )	$ite(x, y, 1)$
1110	$x \uparrow y$	$ite(x, \neg y, 1)$
1111	1	1

Table 2.3.1: Binary logical operator by *ite*. The first column shows the output of the function on the input pairs  $(x, y) = (0, 0), (0, 1), (1, 0),$  and  $(1, 1)$ . The second column shows its conventional form and the last column gives the equivalent of ite form. In column 3,  $y$  stands for  $ite(y, 1, 0)$  and  $\neg y$  for  $ite(y, 0, 1)$ .

**Algorithm 2.3.20.** The algorithm to convert a formula into INF is a recursive procedure *convertINF*:

```

proc convertINF( $A$ )
  if  $A = 1$  return 1;
  if  $A = 0$  return 0;
  else if  $A$  contains  $p \in V_p$ 
    return ite( $p$ , convertINF( $A[p \leftarrow 1]$ ), convertINF( $A[p \leftarrow 0]$ ));
  else return simplify( $A$ ).

```

Note that  $A[p \leftarrow 1]$  (or  $A[p \leftarrow 0]$ ) stands for the formula resulting from replacing every occurrence of  $p$  by 1 (or 0) in  $A$ . The subroutine *simplify*( $A$ ) will return the truth value of formula  $A$  which has no propositional variables. The formula returned by *convertINF* contains only *ite*, 0, and 1 as the logical operators. All the propositional variables appear as the first argument of *ite*.

**Example 2.3.21.** let  $A$  be  $\neg p \vee q \wedge \neg r$ , then  $A[p \leftarrow 1] = \neg 1 \vee q \wedge \neg r \equiv q \wedge \neg r$ ;  $A[p \leftarrow 0] = \neg 0 \vee q \wedge \neg r \equiv 1$ . So  $\text{convertINF}(A) = \text{ite}(p, \text{convertINF}(q \wedge \neg r), 1)$ .  $\text{convertINF}(q \wedge \neg r) = \text{ite}(q, \text{convertINF}(\neg r), 0)$ , and  $\text{convertINF}(\neg r) = \text{ite}(r, 0, 1)$ . Combining all, we have  $\text{convertINF}(A) = \text{ite}(p, \text{ite}(q, \text{ite}(r, 0, 1), 0), 1)$ .

A *binary decision diagram* (BDD) is a directed acyclic graph (DAG)  $G$ , which represents an INF by maximally sharing common subformulas.  $G$  has exactly two leaf nodes (no outgoing edges) with labels 1 and 0, respectively, and a set of internal nodes (with out-going edges); each internal node represents INF  $\text{ite}(p, A, B)$ , that is, the node is labeled with  $p$  and has the two out-going edges to the nodes representing  $A$  and  $B$ , respectively.

Given an order on the propositional variables, if we use this order to choose each variable in the procedure *convertINF*, the resulting INF is an *ordered BDD* (OBDD). For the above example, the order  $p > q > r$  is used. If we use  $r > q > p$ , then for  $A = \neg p \vee q \wedge \neg r$ ,  $\text{convertINF}(A) = \text{ite}(r, \text{ite}(p, 0, 1), \text{ite}(q, 1, \text{ite}(p, 1, 0)))$ . The choice of orders has a big impact on the sizes of OBDD, as one order gives you the linear size (in terms of number of variables); another order may give you an exponential size. The example in Figure 2.3.2 shows the two OBDDs for  $(a_1 \wedge b_1) \vee (a_2 \wedge b_2) \vee (a_3 \wedge b_3)$  using two orders. If we generalize this example to  $(a_1 \wedge b_1) \vee (a_2 \wedge b_2) \vee \dots \vee (a_n \wedge b_n)$ , the size of the first OBDD will be  $O(n)$ , while the second OBDD will have an exponential size.

Since  $\text{ite}(C, A, A) \equiv A$ , we may use this equivalence repeatedly to simplify an INF, until no more simplification can be done. The resulting INF is called *reduced*. The corresponding BDD is called *reduced* if, additionally, identical subformulas of an INF are represented by a unique node.

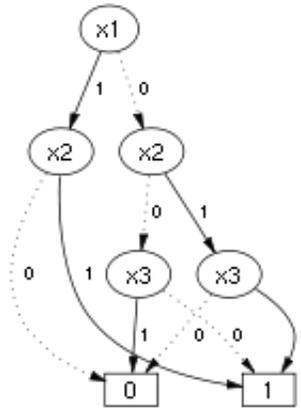


Figure 2.3.1: The BDD of  $\text{INF } \text{ite}(x_1, \text{ite}(x_2, 1, 0), \text{ite}(x_2, \text{ite}(x_3, 1, 0), \text{ite}(x_3, 0, 1)))$  derived from  $\overline{x_1 x_2 x_3} \vee x_1 x_2 \vee x_2 x_3$ .

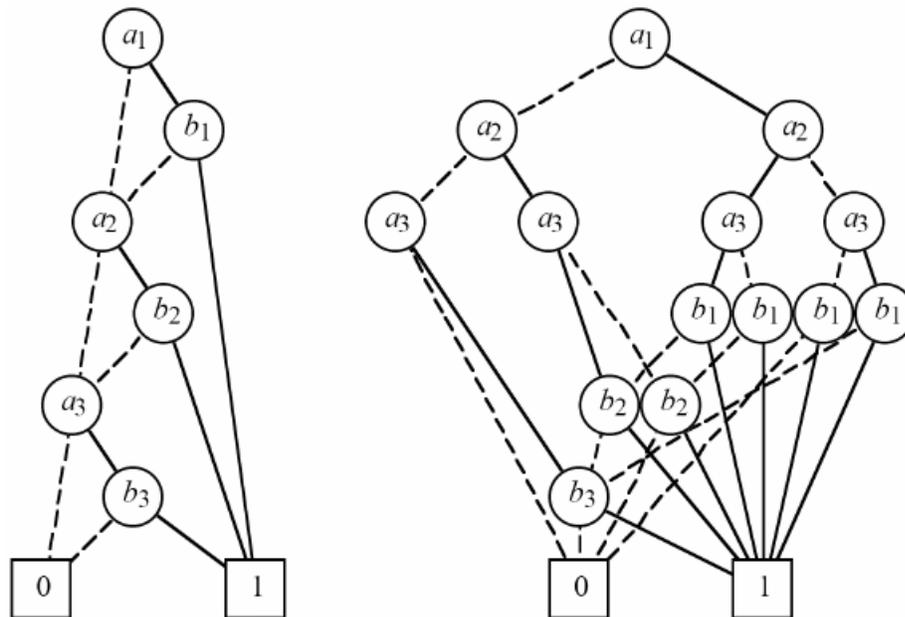


Figure 2.3.2: The first BDD uses  $a_1 > b_1 > a_2 > b_2 > a_3 > b_3$  and the second BDD uses  $a_1 > a_2 > a_3 > b_1 > b_2 > b_3$  for the same formula  $(a_1 \wedge b_1) \vee (a_2 \wedge b_2) \vee (a_3 \wedge b_3)$ .

**Proposition 2.3.22.** *All logical equivalent INFs, using the same order on propositional variables, have a unique reduced OBDD (ROBDD).*

*Proof.* Since a ROBDD is a graph representing an INF  $A$ , we may do an induction proof on the number of variables in  $A$ . If  $A$  contains no variables, then 1 or 0 are represented by the unique nodes. Let  $p_n > p_{n-1} > \dots > p_2 > p_1$  be the order on the  $n$  variables in  $A$ . Then  $A \equiv ite(p_n, B, C)$  for INFs  $B, C$ , and  $B \not\equiv C$ . If  $n = 1$ , then  $\{B, C\} = \{0, 1\}$ , and  $A$  is represented by a unique node. If  $n > 1$ , because  $B$  and  $C$  are INFs without variable  $p_n$ , by induction hypothesis,  $B$  and  $C$  are represented by the unique ROBDDs, respectively. Thus,  $A = ite(p_n, B, C)$  has a unique ROBDD representation.  $\square$

It is easy to apply  $ite(p, A, A) \equiv A$  as a reduction rule on INF. To implement the sharing of common subformulas, we may use a hash table which contains the triple  $\langle v, x, y \rangle$  for each node  $ite(v, x, y)$  in a ROBDD, where  $v$  is a propositional variable and both  $x$  and  $y$ ,  $x \neq y$ , are nodes in the ROBDD. The hash table provides two operations:  $lookupHashTable(v, x, y)$  checks if  $\langle v, x, y \rangle$  exists in the hash table: if yes, the node is returned; otherwise, *null* is returned. The other operation is  $saveCreateNode(v, x, y)$ , which creates a new node for the triple  $\langle v, x, y \rangle$  and saves the node in the hash table. The pseudo-code of the algorithm for creating ROBDD is given below.

**Algorithm 2.3.23.** The algorithm *ROBDD* takes a formula  $A$  and returns a ROBDD node for  $A$ .  $vars(A)$  returns the set of propositional variables of  $A$ ,  $topVariable(S)$  returns the maximal variable of  $S$  by the given variable order.  $simplify(A)$  returns the truth value of a formula without variables.

```

proc ROBDD( $A$ )
  if  $vars(A) = \emptyset$  return  $simplify(A)$ ;
   $v := topVariable(vars(A))$ ;
   $x := ROBDD(A[v \leftarrow 1])$ ;
   $y := ROBDD(A[v \leftarrow 0])$ ;
  if  $(x = y)$  return  $x$ ; // reduction
   $p := lookupHashTable(v, x, y)$ ;
  if  $(p \neq null)$  return  $p$ ; // sharing
  return  $saveCreateNode(v, x, y)$ ; // a new node is created

```

The input  $A$  to *ROBDD* can contain the nodes of a ROBDD, so that we can perform various logical operations on ROBDDs.

**Example 2.3.24.** Assume  $a > b > c$ , and let  $A$  be  $ite(F, G, H)$ , where  $F = ite(a, 1, b)$ ,  $G = ite(a, c, 0)$ , and  $H = ite(b, 1, d)$ . Then  $ROBDD(A)$  will call  $ROBDD(A[a \leftarrow 1])$  and  $ROBDD(A[a \leftarrow 0])$ , which returns  $c$  and  $J = ite(b, 0, H)$ , respectively. Finally,  $ROBDD(A)$  returns  $ite(a, c, J)$ .

ROBDDs have used for presenting Boolean functions, symbolic simulation of combinational circuits, equivalence checking and verification of Boolean functions, and finding and counting models of propositional formulas.

## 2.4 Optimization Problems

There exist many optimization problems in propositional logic. For example, one of them is to find a variable order for a formula, so that the resulting ROBDD is minimal. In this section, we introduce three optimization problems.

### 2.4.1 Minimum Set of Operators

To represent the sentences like  $A$  provided  $B$ , which represent the converse implication, we do not need a formula like  $A \Leftarrow B$ , because we may use simply  $B \rightarrow A$ . This and similar reasons explain why only a few of the sixteen binary Boolean functions require notation. In CNF or DNF, we need only three Boolean operators:  $\neg$ ,  $\vee$ , and  $\wedge$ ; all other Boolean operators can be defined in terms of them.

**Definition 2.4.1.** *A set  $S$  of Boolean operators is said to be sufficient if every other Boolean operators can be defined using only the operators in  $S$ .  $S$  is said to be minimally sufficient if  $S$  is sufficient and no proper subset of  $S$  is a sufficient set of operators.*

Obviously,  $\{\neg, \vee, \wedge\}$  is a sufficient set of operators. For example,  $A \oplus B \equiv (A \vee \neg B) \wedge (\neg A \vee B)$  and  $ite(A, B, C) \equiv (A \wedge B) \vee (\neg A \wedge C)$ . Is the set  $\{\neg, \vee, \wedge\}$  minimally sufficient? The answer is no, because  $A \vee B \equiv \neg(\neg A \wedge \neg B)$ . Thus,  $\{\neg, \wedge\}$  is a sufficient set of operators.  $\{\neg, \vee\}$  is another sufficient set, because  $A \wedge B \equiv \neg(\neg A \wedge \neg B)$ . To show that  $\{\neg, \wedge\}$  is minimally sufficient, we need to show that neither  $\{\neg\}$  nor  $\{\wedge\}$  is sufficient. Since  $\neg$  takes only one argument and cannot be used alone to define  $\wedge$ ,  $\{\neg\}$  cannot be sufficient. We will see later that  $\{\wedge\}$  is not sufficient.

From IND (ITE normal form), we might think  $\{ite\}$  is minimally sufficient. That is not true, because we also use 1 and 0 (stand for  $\top$  and  $\perp$ ) which are nullary operators. For the set  $\{\neg, \wedge\}$ ,  $\perp \equiv p \wedge \neg p$  for any propositional variable  $p$ .

Do we have a minimally sufficient set of Boolean operators which contains a single operator? The answer is yes and there are two binary operators as candidates. They are  $\uparrow$  (nand, the negation of and) and  $\downarrow$  (nor, the negation of or):

$$A \uparrow B \equiv \neg(A \wedge B) \qquad A \downarrow B \equiv \neg(A \vee B)$$

**Proposition 2.4.2.** *Both  $\{\uparrow\}$  and  $\{\downarrow\}$  are minimally sufficient.*

*Proof.* To show  $\{\uparrow\}$  is sufficient, we just need to define  $\neg$  and  $\wedge$  using  $\uparrow$ :  $\neg(A) \equiv A \uparrow A$  and  $A \wedge B \equiv \neg(A \uparrow B)$ . The proof that  $\{\downarrow\}$  is sufficient is left as an exercise.  $\square$

In fact, among all the sixteen binary operators, only  $\uparrow$  and  $\downarrow$  have such property.

**Theorem 2.4.3.** *For any binary operator which can be used to define both  $\neg$  and  $\wedge$ , that operator must be either  $\uparrow$  or  $\downarrow$ .*

*Proof.* Let  $o$  be any binary operator. If  $\neg$  can be defined by  $o$ , then  $\neg A \equiv t(o, A)$ , where  $t(o, A)$  is the formula represented by a binary tree  $t$  whose leaf nodes are labeled with  $A$  and whose internal nodes are labeled with  $o$ . Furthermore, we assume that  $t$  is such a binary tree with the least number of nodes. If we replace  $A$  by 1, then  $t(o, 1)$  should have the truth value 0, because  $\neg 1 = t(o, 1) = 0$ . If any internal node of  $t(o, 1)$  has the truth value 1, we may replace that node by a leaf node 1; if any internal node other than the root of  $t$  has truth value 0, we may use the subtree at this internal node as  $t$ . In both cases, this is a violation to the assumption that  $t$  is such a tree with the least of nodes. As a result,  $t$  should have only one internal node, i.e.,  $1 \ o \ 1 = \neg 1 = 0$ . Similarly, we should have  $0 \ o \ 0 = \neg 0 = 1$ .

There are four cases regarding the truth values of  $1 \ o \ 0$  and  $0 \ o \ 1$ .

- case 1:  $1 \ o \ 0 = 0$  and  $0 \ o \ 1 = 1$ . Then  $A \ o \ B = \neg A$  and we cannot use this  $o$  to define  $A \wedge B$ , as the output from  $o$  on  $A$  and  $B$  will be one of  $A$ ,  $\neg A$ ,  $B$ , or  $\neg B$ .
- case 2:  $1 \ o \ 0 = 1$  and  $0 \ o \ 1 = 0$ . Then  $A \ o \ B = \neg B$ , and we still cannot use this  $o$  to define  $A \wedge B$  for the same reason as in Case 1.
- case 3:  $1 \ o \ 0 = 0 \ o \ 1 = 1$ .  $o$  is  $\uparrow$ .
- case 4:  $1 \ o \ 0 = 0 \ o \ 1 = 0$ .  $o$  is  $\downarrow$ .

$p$	$q$	$p \circ_1 q$	$p \circ_2 q$	$p \circ_3 q$	$p \circ_4 q$
1	1	0	0	0	0
0	1	1	0	1	0
1	0	0	1	1	0
0	0	1	1	1	1

Table 2.4.1:  $p \circ_1 q \equiv \neg p$  (the negation of projection on the first argument);  $p \circ_2 q \equiv \neg q$  (the negation of projection on the second argument);  $p \circ_3 q \equiv p \uparrow q$  (NAND); and  $p \circ_4 q \equiv p \downarrow q$  (NOR).

The last two cases give us the definitions of  $\uparrow$  and  $\downarrow$ , as shown in Table 2.4.1.  $\square$

**Corollary 2.4.4.** *Both  $\{\neg, \wedge\}$  and  $\{\neg, \vee\}$  are minimally sufficient.*

In Boolean algebra, the property of any sufficient set of operators can be specified using equations with variables and the operator. To specify the property of  $\uparrow$ , we may use three equations:  $0 \uparrow Y = 1$ ,  $X \uparrow 0 = 1$ , and  $1 \uparrow 1 = 0$ . For  $\uparrow$ , what is the minimal number of equations needed? And what is the minimal number of symbols in the equations? These questions have been answered firmly: only one equation is needed and only three variables (for a total of either occurrences) and six uses of  $\uparrow$  are needed. We state below the result without proof.

**Proposition 2.4.5.** *The Boolean algebra can be generated by the following equation:*

$$((x \uparrow y) \uparrow z) \uparrow (x \uparrow ((x \uparrow z) \uparrow x)) = z$$

The above equation contains three variables and 14 symbols (excluding parentheses). For the set  $\{\neg, \vee\}$ , a single axiom also exists to specify the Boolean algebra:

**Proposition 2.4.6.** *The Boolean algebra can be generated by the following equation:*

$$\neg(\neg(\neg(x \vee y) \vee z) \vee \neg(x \vee \neg(\neg z \vee \neg(z \vee u)))) = z.$$

The above equation contains four variables and 21 symbols (excluding parentheses).

## 2.4.2 Logic Minimization

Logic minimization, also called *logic optimization*, is a part of circuit design process and its goal is to obtain the smallest combinational circuit that is represented by a Boolean formula. Logic minimization seeks to find an equivalent representation

of the specified logic circuit under one or more specified constraints. Generally the circuit is constrained to minimum chip area meeting a pre-specified delay. Decreasing the number of gates will reduce the power consumption of the circuit. Choosing gates with less transistors will reduce the circuit area. Decreasing the number of nested levels of gates will reduce the delay of the circuit. Logic minimization will reduce substantially the cost of circuits and improve its quality.

#### 2.4.2.1 Karnaugh Maps for Minimizing DNF and CNF

In the early days of electronic design automation (EDA), a Boolean function is often defined by a truth table and we can derive full DNF (sum of products) or full CNF (product of sums) directly from the truth table, as shown in Example 2.3.16. The obtained DNF or CNF need to be minimized, in order to use the least number of gates to implement this function.

Following the convention of EDA, we will use  $+$  for  $\vee$ ,  $\cdot$  for  $\wedge$ , and  $\bar{A}$  for  $\neg A$ .

**Example 2.4.7.** Let  $f : \{0, 1\}^2 \rightarrow \{0, 1\}$  be defined by  $f(A, B) = \overline{AB} + A\bar{B} + \bar{A}B$ . We need three AND gates and one OR gate of 3-input. However,  $f(A, B) \equiv \bar{A} + \bar{B}$ . Using  $\bar{A} + \bar{B}$ , we need only one OR gate.

The equivalence relations used in the above simplification process are  $AB + A\bar{B} \equiv A$  and  $A + \bar{A}B \equiv A + B$  (or  $\bar{A} + AB \equiv \bar{A} + B$ ). For the above example,  $\overline{AB} + A\bar{B} \equiv \bar{B}$  and  $\bar{B} + \bar{A}B \equiv \bar{B} + \bar{A}$ .

We want to find an equivalent circuit of minimum size possible, and this is a difficult computation problem. There are many tools such as *Karnaugh maps* available to achieve this goal.

The Karnaugh map (K-map) is a popular method of simplifying DNF or CNF formulas, originally proposed by Maurice Karnaugh in 1953. The Karnaugh map reduces the need for extensive calculations by displaying the output of the function graphically and taking advantage of humans' pattern-recognition capability. It is suitable for Boolean functions with 2 to 4 variables. When the number of variables is greater than 4, it is better to use an automated tool.

The Karnaugh map uses a two-dimensional grid where each cell represents an interpretation (or equivalently, a full minterm) and the cell contains the truth value of the function for that interpretation. The position of each cell contains all the information of an interpretation and the cells are arranged such that adjacent cells differ by exactly one truth value in the interpretation. Adjacent 1s in the Karnaugh map represent opportunities to simplify the formula. The minterms for the final formula are found by encircling groups of 1s in the map. Minterm groups must be rectangular and must have an area that is a power of two (i.e., 1, 2, 4, 8, ...).

A	B	F
0	0	1
0	1	1
1	0	1
1	1	0

	A		
	B		
		0	1
	0	1	1
	1	1	0

Figure 2.4.1: K-map for  $f(x, y) = x\bar{y} + \bar{x}y + \bar{x}\bar{y}$ .

Minterm rectangles should be as large as possible without containing any 0s (it may contain don't-care cells). Groups may overlap in order to make each one larger. A least number of groups that cover all 1s will represent the minterms of a DNF of the function. These minterms can be used to write a minimal DNF representing the required function and thus implemented by the least number of AND gates feeding into an OR gate. The map can also be used to obtain a minimal CNF that is implemented by OR gates feeding into an AND gate.

Let  $f(x, y) = x\bar{y} + \bar{x}y + \bar{x}\bar{y}$ . There are four minterms of  $x$  and  $y$ , the K-map has four cells, that is, four minterms corresponding to the four interpretations on  $x$  and  $y$ . As the output of  $f$ , three cells have value 1 and one cell (i.e.,  $xy$ ) has value 0. Initially, each cell belongs to a distinct group of one member (Figure 2.4.1).

**The merge operation on K-map:** Merge two adjacent and disjoint groups of the same size and of the same truth value into one larger group.

This operation creates groups of size 2, 4, 8, etc., and cells are allowed to appear in different groups at the same time. For example, cells  $x\bar{y}$  and  $\bar{x}\bar{y}$  in Figure 2.4.1 are adjacent, they are merged into one group  $\{x\bar{y}, \bar{x}\bar{y}\}$ , which can be represented by  $\bar{y}$  as a shorthand of the group. Similarly, cells  $\bar{x}y$  and  $\bar{x}\bar{y}$  can be merged into  $\{\bar{x}y, \bar{x}\bar{y}\}$ , which can be represented by  $\bar{x}$ . Note that  $\bar{x}\bar{y}$  is used twice in the two merge operations. No more merge operations can be performed, and the final result is  $f = \bar{x} + \bar{y}$  ( $f$  is the NAND function). This result can be proved logically as follows:

$$\begin{aligned}
 f &= x\bar{y} + \bar{x}y + \bar{x}\bar{y} \\
 &= (x\bar{y} + \bar{x}\bar{y}) + (\bar{x}y + \bar{x}\bar{y}) \\
 &= (x + \bar{x})\bar{y} + \bar{x}(y + \bar{y}) \\
 &= \bar{y} + \bar{x}
 \end{aligned}$$

A	B	C	f
0	0	0	0
0	0	1	1
0	1	0	0
0	1	1	1
1	0	0	0
1	0	1	0
1	1	0	1
1	1	1	1

		AB			
		00	01	11	10
C	0	0	2	1	4
	1	1	3	1	5

Figure 2.4.2: The K-map for obtaining  $f(A, B, C) = AB + \overline{AC}$ .

**Example 2.4.8.** Consider the function  $f$  defined by the truth table in Figure 2.4.2. The corresponding K-map has four cells with truth value 1 as  $f$  has four models shown in the truth table. There are three possible merge operations: (i)  $\overline{A}BC$  and  $\overline{A}B\overline{C}$  merge into  $\overline{A}C$ ; (ii)  $ABC$  and  $AB\overline{C}$  merge into  $AB$ ; and (iii)  $\overline{A}BC$  and  $ABC$  merge into  $BC$ . Since all the cells with value 1 are in the groups represented by  $\overline{A}C$  and  $AB$ ,  $BC$  is redundant. So the simplified function is  $f = \overline{A}C + AB$ , instead of  $f = \overline{A}C + AB + BC$ . On the other hand, if  $g = AB + AC + BC$ , then none of the three minterms in  $g$  is redundant, because if you delete one, you will miss some 1s, even though the full minterm  $ABC$  is covered by all the three minterms.

From the above examples, it is clear that after each merge operation, two minterms are merged into one with one less variable. That is, the larger the group, the less the number of variables in the minterm to represent the group. This is possible because adjacent cells represent interpretations which differ by exactly one truth value.

Because K-maps are two-dimensional, some adjacent relations of cells are not shown on K-maps. For example, in Figure 2.4.2,  $m_0 = \overline{A}B\overline{C}$  and  $m_4 = A\overline{B}\overline{C}$  are adjacent,  $m_1 = \overline{A}BC$  and  $m_5 = A\overline{B}C$  are adjacent. Strictly speaking, K-maps are *toroidally* connected, which means that rectangular groups can wrap across the edges. Cells on the extreme right are actually “adjacent” to those on the far left, in the sense that the corresponding interpretations only differ by one truth value. We need a 3-dimensional graph to completely show the adjacent relation of 3 variables; a 4-dimensional graph to completely show the adjacent relation of 4 variables; etc.

Based on the idea of K-maps, an automated tool will create a graph of  $2^n$  nodes for  $n$  variables, each node represents one interpretation, and two nodes have an edge iff their interpretations differ by one truth value. The merge operation is then implemented on this graph to generate minterms for a least size DNF. The same graph is also used for generating Gray code (by finding a Hamiltonian path

PQ \ RS	00	01	11	10
00	1	0	0	1
01	0	1	1	0
11	0	1	1	0
10	1	0	0	1

Figure 2.4.3: The K-map for  $f(P, Q, R, S) = m_0 + m_2 + m_5 + m_7 + m_8 + m_{10} + m_{13} + m_{15}$ . The simplified formula is  $f = \overline{Q}\overline{S} + QS$ .

in the graph).

**Example 2.4.9.** The K-map for  $f(P, Q, R, S) = m_0 + m_2 + m_5 + m_7 + m_8 + m_{10} + m_{13} + m_{15}$ , is given in Figure 2.4.3, where  $m_0 = \overline{P}\overline{Q}\overline{R}\overline{S}$ ,  $m_2 = \overline{P}\overline{Q}R\overline{S}$ ,  $m_5 = \overline{P}Q\overline{R}S$ ,  $m_7 = \overline{P}QRS$ ,  $m_8 = P\overline{Q}\overline{R}\overline{S}$ ,  $m_{10} = P\overline{Q}RS$ ,  $m_{13} = PQ\overline{R}\overline{S}$ , and  $m_{15} = PQRS$  ( $i$  in  $m_i$  is the decimal value of  $PQRS$  in the corresponding interpretation). At first,  $m_0$  and  $m_2$  are adjacent and can be merged into  $\overline{P}\overline{Q}\overline{S}$ ;  $m_8$  and  $m_{10}$  are also adjacent and can be merged into  $P\overline{Q}\overline{S}$ . Then  $\overline{P}\overline{Q}\overline{S}$  and  $P\overline{Q}\overline{S}$  can be further merged into  $\overline{Q}\overline{S}$ . Similarly,  $m_5$  and  $m_7$  merge into  $\overline{P}QS$ ;  $m_{13}$  and  $m_{15}$  merge into  $PQS$ . Then  $\overline{P}QS$  and  $PQS$  merge into  $QS$ . The final result is  $f = \overline{Q}\overline{S} + QS$ . It would also have been possible to derive this simplification by carefully applying the equivalence relations, but the time it takes to do that grows exponentially with the number of minterms.

Note that the formula obtained from a K-map is not unique in general. For example,  $g(P, Q, R, S) = P\overline{Q} + R\overline{S} + \overline{P}Q\overline{S}$  or  $g(P, Q, R, S) = P\overline{Q} + R\overline{S} + \overline{Q}R\overline{S}$  are possible outputs for  $g$  from K-maps. This means the outputs of K-maps are not canonical.

K-maps can be used to generate simplified CNF for a function. In Example 2.3.16, we have shown how to obtain a full CNF from a full DNF. Using the same idea, if we simplify the full DNF first, the CNF obtained by negating the simplified DNF will a simplified one.

**Example 2.4.10.** Let  $f(A, B, C) = \overline{A}BC + \overline{A}B\overline{C} + A\overline{B}C + A\overline{B}\overline{C}$  and its K-map

		BC			
		00	01	11	10
A	0	0	1	0	1
	1	1	1	0	0
		0	1	3	2
		4	5	7	6

Figure 2.4.4: The K-map for  $f(A, B, C) = \overline{A}BC + \overline{A}B\overline{C} + A\overline{B}C + A\overline{B}\overline{C}$ . The simplified CNF is  $f = (A + B + C)(\overline{A} + \overline{B})(\overline{B} + \overline{C})$ .

is given in Figure 2.4.4. The negation of  $f$  is

$$\overline{f} = \overline{A}BC + \overline{A}B\overline{C} + A\overline{B}C + ABC$$

Since  $\overline{A}BC$  and  $ABC$  can merge into  $BC$  (shown in the figure),  $A\overline{B}C$  and  $A\overline{B}\overline{C}$  can merge into  $AB$ , the simplified  $\overline{f}$  is  $\overline{f} = \overline{A}BC + AB + BC$ . The simplified CNF for  $f$  is thus  $f = (A + B + C)(\overline{A} + \overline{B})(\overline{B} + \overline{C})$ .

Karnaugh maps can also be used to handle “don’t care” values and can simplify logic expressions in software design. Boolean conditions, as used for example in conditional statements, can get very complicated, which makes the code difficult to read and to maintain. Once minimized, canonical sum-of-products and product-of-sums expressions can be implemented directly using AND and OR logic operators.

#### 2.4.2.2 Other Issues in Logic Minimalization

Simplified Boolean functions obtained from K-maps are DNFs or CNFs. To implement them in a circuit using CMOS (complementary metal–oxide–semiconductor), we need further optimization. One of the most important characteristics of CMOS devices is low static power consumption. CMOS devices draw significant power only momentarily during switching between on and off states. Consequently, CMOS devices do not produce as much waste heat as other forms of logic, like transistor–transistor logic (TTL), which normally has some standing current even when not changing state. These characteristics allow CMOS to integrate a high density of logic functions on a chip. It was primarily for this reason that CMOS became the most widely used technology to be implemented in VLSI (very large scale integration) chips.

Using CMOS, the NAND/NOR/NOT gates are cheaper to implement than the AND/OR gates. For a basic unbuffered NAND/NOR/NOT gate, we need only

two transistors per input. It works when that gate is one among many others, driving a few similar gates. Many practical discrete gates, however, will be buffered and they will need additional transistors to improve the output strength and prevent the output strength from depending on active inputs. Typically, AND/OR gates will be implemented using the NAND/NOR/NOT gates and need more transistors than NAND/NOR gates.

**Example 2.4.11.** If we directly implement  $f(A, B, C) = \overline{A}C + \overline{B}$  in CMOS, we need two inverters (for negation), one OR gate (as the negation of  $\downarrow$ ) and one AND gate (as the negation of  $\uparrow$ ). That is, we need four inverters (each uses two transistors), one NOR gate and one NAND gate (each uses four transistors), for a total of 16 transistors. We may convert  $f(A, B, C) = \overline{A}C + \overline{B}$  into the simplified and equivalent formula using only  $\uparrow$  and  $\neg$ :  $f(A, B, C) = (\overline{A} \uparrow C) \uparrow B$ , which uses a total of 10 transistors (two for the inverter; four each for the two NAND gates).

To convert DNF or CNF into a formula using only NAND/NOR/NOT, we may use the algorithm *NNF2NOR*.

**Algorithm 2.4.12.** The algorithm *NNF2NOR* takes a formula  $A$  in NNF as input and returns a formula  $B$  which is equivalent to  $A$  and uses only NOR/NAND/NOT operators. It calls a mutually recursive algorithm *neg2NOR*, which takes a formula  $A$  in NNF as input and returns a formula  $B$  which is equivalent to  $\neg A$  and uses only NOR/NAND/NOT operators.

```

proc NNF2NOR( $A$ )
  if  $A \in V_p$  or  $A = \neg B$  return  $A$ ;
  if  $A = B \vee C$  return neg2NOR( $B$ )  $\uparrow$  neg2NOR( $C$ );
  if  $A = B \wedge C$  return neg2NOR( $B$ )  $\downarrow$  neg2NOR( $C$ );

proc neg2NOR( $A$ )
  if  $A \in V_p$  return  $\neg A$ ;
  if  $A = \neg B$  return  $B$ ;
  if  $A = B \vee C$  return NNF2NOR( $B$ )  $\downarrow$  NNF2NOR( $C$ );
  if  $A = B \wedge C$  return NNF2NOR( $B$ )  $\uparrow$  NNF2NOR( $C$ );

```

For input  $A = (\overline{p} \vee q) \vee r$ , *NNF2NOR*( $A$ ) returns  $(\overline{p} \downarrow q) \uparrow \overline{r}$ , which needs 12 transistors to implement. However, this result is not optimal, because another equivalent formula, i.e.,  $p \uparrow (q \downarrow r)$ , needs eight transistors to implement in CMOS. Relatively speaking, finding optimal formulas from DNF/CNF is an easy task. The

trick is to get rid of as many negation in literals as possible, by switching between  $\uparrow$  and  $\downarrow$ , using the relation

$$\begin{aligned}(\neg A \uparrow B) \downarrow \neg C &\equiv A \downarrow (B \uparrow C); \\ (\neg A \downarrow B) \uparrow \neg C &\equiv A \uparrow (B \downarrow C).\end{aligned}$$

If a 3-input NAND gate is available, the same  $A$  can be represented by  $\uparrow(p, \bar{q}, \bar{r})$ , which needs 10 transistors to implement. Note that  $\uparrow(A, B, C)$ ,  $(A \uparrow B) \uparrow C$ , and  $A \uparrow (B \uparrow C)$  are all logically different.

DNFs (sum of products) and CNFs (product of sums) are called *2-level* formulas. A *3-level formula* is a product of DNFs or a sum of CNFs. A *4-level formula* is a product of 3-level formulas or a sum of 3-level formulas, and so on. Finding equivalent formulas of high levels can possibly reduce the size of formulas.

**Example 2.4.13.** Let  $f(a, b, c, d, f, g) = adf + aef + bdf + bef + cdf + cef + g$ .  $f$  is a DNF and we need 6 AND gates of 3-input and 3 OR gates of 3-input (to implement the 7 input OR) for a total of 9 gates. If we allow 3-level formulas, then  $f(a, b, c, d, f, g) \equiv (a + b + c)(d + e)f + g$ , which needs one OR gate of 3 input, 2 OR gates of 3 input, one AND gate of 3 input, for a total of 4 gates.

Besides using higher level formulas, CMOS technology allows to use other types of gates besides OR, AND, and inverters (for negation). For example, if we want to implement a function  $f(a, b) = a\bar{b} + \bar{a}b$ , since  $f(a, b) \equiv a \oplus b$ , one XOR gate will suffice. Many circuit design systems provide a complex of libraries with a large number of gates. For example, the EE271 library allows the use of the following gates for combinational circuits:

inv, nand2, nand3, nand4, nor2, nor3, xor, xnor, ao12, ao122, oai2, oai22, xnor

where “xnor” is the negation of “xor”, i.e., the logical equal ( $\leftrightarrow$ ), “ao1” stands for “AND-OR-Invert”, “oai” stands for “OR-AND-Invert”, and they define the following functions:

$$\begin{aligned}\text{aoi2:} & \quad F = \overline{(A \wedge B) \vee C} \\ \text{aoi22:} & \quad F = \overline{(A \wedge B) \vee (C \wedge D)} \\ \text{oai2:} & \quad F = \overline{(A \vee B) \wedge C} \\ \text{oai22:} & \quad F = \overline{(A \vee B) \wedge (C \vee D)}\end{aligned}$$

AOI and OAI gates can be readily implemented in CMOS circuitry. AOI gates are particularly advantaged in that the total number of transistors is less than

if the AND, NOT, and OR functions were implemented separately. This results in increased speed, reduced power, smaller area, and potentially lower fabrication cost. For example, an AOI2 gate can be constructed with 6 transistors in CMOS compared to 10 transistors using a 2-input NAND gate (4 transistors), an inverter (2 transistors), and a 2-input NOR gate (4 transistors), because  $F = \overline{(A \wedge B)} \vee C \equiv \overline{A \uparrow B} \downarrow C$ .

Finding an equivalent and minimal circuit with high-level formulas of a rich set of gates becomes even more challenging and is an on-going research area in EDA.

### 2.4.3 Maximum Satisfiability

When a CNF is unsatisfiable, we are still interested in finding an interpretation which makes most clauses to be true. This is so-called the *maximum satisfiability problem*, which is the problem of determining the maximum number of clauses, of a given Boolean formula in CNF, that can be made true by an interpretation (an assignment of truth values to the variables of the formula). It is a generalization of the Boolean satisfiability problem, which asks whether there exists a truth assignment that makes all clauses true.

**Example 2.4.14.** Let  $A$  be a set of unsatisfiable clauses  $\{(\bar{p}), (p \vee \bar{q}), (p \vee q), (p \vee \bar{r}), (p \vee r)\}$ . If we assign 0 to  $p$ , two clauses will be false no matter what values are assigned to  $q$  and  $r$ . On the other hand, if  $p$  is assigned 1, only clause  $(\bar{p})$  will be false. Therefore, if  $A$  is given as an instance of the MAX-SAT problem, the solution to the problem is 4 (four clauses are true in one interpretation).

More generally, one can define a weighted version of MAX-SAT as follows: given a conjunctive normal form formula with non-negative weights assigned to each clause, find truth values for its variables that maximize the combined weight of the satisfied clauses. The MAX-SAT problem is an instance of weighted MAX-SAT where all weights are 1.

Find the final exam schedule for a big university is a challenging job. There are hard constraints such as no students can take two exams at the same, and no classrooms can be hold two exams at the same time. There are also software constraints such that one student is preferred to take at most two exams in one day, or two exams are separated by two hours. To specify such a problem as an instance of weighted MAX-SAT is easy: We just give a very large weight to hard constraints and a small weight to soft constraints.

Various methods have developed to solve MAX-SAT or weighted MAX-SAT problems. These methods are related to solving SAT problems and we will present

them later in the book.

## 2.5 Using Propositional Logic

Propositions and logical operators arise all the time in computer programs. All programming languages such as C, C++, and Java use symbols like “&&”, “||”, and “!” in place of  $\wedge$ ,  $\vee$ , and  $\neg$ . Consider the Boolean expression appearing in an if-statement

```
if ( x > 0 || (x <= 0 && y > 100 ) ) ...
```

If we let  $p$  denote  $x > 0$  and  $q$  be  $y > 100$ , then the Boolean expression is abstracted as  $p \vee \neg p \wedge q$ . Since  $(p \vee \neg p \wedge q) \equiv (p \vee q)$ , the original code can be rewritten as follows:

```
if ( x > 0 || y > 100 ) ...
```

In other words, the knowledge about the logic may help us to write neater codes.

Consider another piece of the code:

```
if ( x >= 0 && A[x] == 0 ) ...
```

Let  $p$  be  $x >= 0$  and  $q$  be  $A[x] == 0$ , then the Boolean expression is  $p \wedge q$ . Since  $p \wedge q \equiv q \wedge p$ , it is natural to think that  $( x >= 0 \ \&\& \ A[x] == 0 )$  can be replaced by  $( A[x] == 0 \ \&\& \ x >= 0 )$  and the program will be the same. Unfortunately, this is not the case: the execution of  $( A[x] == 0 \ \&\& \ x >= 0 )$  will be aborted when  $x < 0$ . This example shows that “&&” in programming languages are not commutative in general. Neither is “||” as  $( x < 0 \ || \ A[x] == 0 )$  is different from  $( A[x] == 0 \ || \ x < 0 )$ . These abnormalities of the logic not only appear in array indices, but also when assignments are allowed in Boolean expressions. For example,  $x > 0 \ || \ (x = c) > 10$  is different from  $(x = c) > 10 \ || \ x > 0$  in C, where  $x = c$  is an assignment command. As a programmer we need to be careful about these abnormalities.

### 2.5.1 Bitwise Operators

In programming languages like C, C++, and Java, an integer is represented by a list of 16-, 32-, or 64-bits (for C, a char is regarded as an 8-bit integer or a byte). A bitwise operation operates on one or two integers at the level of their individual bits and several of them are logical operations on bits. It is a fast and simple action, directly supported by the processor, and should be used whenever possible.

$(a)_2$	=	$(0, 1, 0, 0, 0, 0, 1, 1)$
$(b)_2$	=	$(0, 0, 0, 1, 0, 0, 1, 0)$
<b>name</b>	<b>symbol</b>	<b>example</b>
NOT	$\sim$	$(\sim a)_2 = (1, 0, 1, 1, 1, 1, 0, 0)$
OR	$ $	$(a   b)_2 = (0, 1, 0, 1, 0, 0, 1, 1)$
AND	$\&$	$(a \& b)_2 = (0, 0, 0, 0, 0, 0, 1, 0)$
XOR	$\wedge$	$(a \wedge b)_2 = (0, 1, 0, 1, 0, 0, 0, 1)$
shift right	$\gg$	$(a \gg 2)_2 = (0, 0, 0, 1, 0, 0, 0, 0)$
shift left	$\ll$	$(b \ll 2)_2 = (0, 1, 0, 0, 1, 0, 0, 0)$

Table 2.5.1: Examples of bitwise operators.

For illustration purpose, we will use 8-bit integers as examples. In general, for any 8-bit integer  $x$ , the binary representation of  $x$  is a 8-bit vector  $(x)_2 = (b_7, b_6, b_5, b_4, b_3, b_2, b_1, b_0)$ , where  $b_i \in \{0, 1\}$  and  $x = \sum_{i=0}^7 b_i * 2^i$ . For example, for  $a = 67$ , its bit vector is  $(a)_2 = (0, 1, 0, 0, 0, 0, 1, 1)$ , because  $67 = 64 + 2 + 1 = 2^6 + 2^1 + 2^0$ , and for  $b = 18$ ,  $(b)_2 = (0, 0, 0, 1, 0, 0, 1, 0)$ , since  $18 = 16 + 2 = 2^4 + 2^1$ . In fact,  $(a)_2$  and  $(b)_2$  are internal representation of  $a$  and  $b$  in computer.

The following bitwise operators are provided, treating each bit as a Boolean value and performing on them simultaneously. We will use the bit vectors of  $a = 67$ , and  $b = 18$  as examples.

If necessary, other bitwise operations can be implemented using the above four operations. For example a bitwise implication on  $a$  and  $b$  can be implemented as  $\sim a | b$ .

The bit shifts are also considered bitwise operations, because they treat an integer as a bit vector rather than as a numerical value. In these operations the bits are moved, or shifted, to the left or right. Registers in a computer processor have a fixed width, so some bits will be “shifted out” of the register at one end, while the same number of bits are “shifted in” from the other end. Typically, the “shifted out” bits are discarded and the “shifted in” bits are all 0’s. When the leading bit is the sign of an integer, the left shift operation becomes more complicated and its discussion is out of the scope of this book. In the table, we see that the shift operators take two integers:  $x \ll y$  will shift  $(x)_2$  left by  $y$  bits;  $x \gg y$  will shift  $(x)_2$  right by  $y$  bits.

The value of the bit vector  $x \ll y$  is equal to  $x * 2^y$ , provided that this multiplication does not cause an overflow. E.g., for  $b = 18$ ,  $b \ll 2 = 18 * 4 = 72$ . Similarly, the value of  $x \gg y$  is equal to the integer division  $x/2^y$ . E.g., for  $a = 67$ ,

$a \gg 2 = \lfloor 67/4 \rfloor = 16$ . Shifting provides an efficient way of multiplication and division when the second operand is a power of two. Typically, bitwise operations are substantially faster than division, several times faster than multiplication.

Suppose you are the manager for a social club of 64 members. For each gathering, you need to record the attendance. If each member has a unique member id  $\in \{0, 1, \dots, 63\}$ , you may use a set of member ids to record the attendance. Instead of a list of integers, you may use a single 64-bit integer, say  $s$ , to store the set. That is, if a member id is  $i$ , then the  $i^{\text{th}}$  bit of  $s$  is 1. For this purpose, we need a way to set a bit in  $s$  to be 1. Suppose  $s = 0$  initially, we may use

$$s = s \mid (1 \ll i)$$

to set the  $i^{\text{th}}$  bit of  $s$  to be 1. To set the value of the  $i^{\text{th}}$  bit of  $s$  to be 0, use

$$s = s \& \sim(1 \ll i).$$

To get the value of the  $i^{\text{th}}$  bit of  $s$ , use  $s \& (1 \ll i)$ . For another gathering, we use another integer  $t$ . Now it is easy to compute the union or intersection of these sets:  $s \mid t$  is the union and  $s \& t$  is the intersection. For example, if you want to check that nobody attended both gatherings, you can check if  $s \& t = 0$ , which is more efficient than using lists for sets.

### 2.5.2 Specify Problems in Propositional Logic

Propositional logic is a very simple language but rich enough to specify most decision problems in computer science, because these problems belong to the class of NP problems and SAT is NP-complete. In fact, propositional logic is a formal language much more rigorous than natural languages. Since natural languages are ambiguous, expressing a sentence in logic helps us to understand the exact meaning of the sentence. For example, the two sentences

1. Students and seniors pay half price,
2. Students or seniors pay half price

evidently have the same meaning.

How to explain this? Let “being student” and “being senior” be abbreviated by  $p$  and  $q$ , and “pay half price” by  $r$ . Let  $A$  be  $(p \rightarrow r) \wedge (q \rightarrow r)$ ,  $B$  be  $(p \vee q) \rightarrow r$ , then  $A$  and  $B$  express somewhat more precisely the factual content of sentences 1 and 2, respectively. It is easy to show that the formulas  $A$  and  $B$  are logically equivalent. The everyday-language statements of  $A$  and  $B$  obscure the structural

difference of  $A$  and  $B$  through an apparently synonymous use of the words “and” and “or”.

We will show to use specify various problems in propositional logic through examples.

**Example 2.5.1.** The so-called  $n$ -queen example is stated as follows: Given a chessboard of size  $n \times n$ , where  $n > 0$ , how to place  $n$  queens on the board such that no two queens attack each other? For  $n = 2, 3$ , there is no solution; for other  $n$ , there exists a solution. To specify this problem in logic, we need  $n^2$  propositional variables. Let the variables be  $p_{i,j}$ , where  $1 \leq i, j \leq n$ , and the meaning of  $p_{i,j}$  is that  $p_{i,j}$  is true iff there is a queen at row  $i$  and column  $j$ . The following properties ensure that there are  $n$  queens on the board such that no two queens attack each other:

1. There is at least one queen in each row.

$$\bigwedge_{1 \leq i \leq n} (p_{i,1} \vee p_{i,2} \vee \cdots \vee p_{i,n})$$

2. No more than one queen in each row.

$$\bigwedge_{1 \leq i \leq n} (\bigwedge_{1 \leq j < k \leq n} (\neg p_{i,j} \vee \neg p_{i,k}))$$

3. No more than one queen in each column.

$$\bigwedge_{1 \leq i \leq n} (\bigwedge_{1 \leq j < k \leq n} (\neg p_{j,i} \vee \neg p_{k,i}))$$

4. No more than one queen on each diagonal.

$$\bigwedge_{1 \leq i < k \leq n} (\bigwedge_{1 \leq j, l \leq n, k-i=|l-j|} (\neg p_{i,j} \vee \neg p_{k,l}))$$

The formula following each property expresses the property formally, and can be easily converted into clauses. When  $n = 2, 3$ , the formulas are unsatisfiable. For  $n \geq 4$ , the formulas are satisfiable and each model gives us a solution. For example, when  $n = 4$ , we have two solutions:  $\{p_{1,2}, p_{2,4}, p_{3,1}, p_{4,3}\}$  and  $\{p_{1,3}, p_{2,1}, p_{3,4}, p_{4,2}\}$ . If we assign the four variables in a solution to be true and the other variables false, we obtain a model. Note that we ignore the property that “there is at least one queen in each row,” because it is redundant.

The next example is a logical puzzle.

**Example 2.5.2.** Three people are going to the beach, each using a different mode of transportation (car, motorcycle, and boat) in a different color (blue, orange, green). Who’s using what? The following clues are given:

1. Abel loves orange, but she hates travel on water.
2. Bob did not use the green vehicle.
3. Carol drove the car.

The first step to solve a puzzle using propositional logic is to make sure what properties are assumed and what conditions are provided by clues. For this puzzle, the first assumed property is that “everyone uses a unique transportation tool”.

To specify this property in logic, there are several ways of doing it. One way is to use a set of propositional variables to specify the relation of “who uses what transportation”: Let  $p_{x,y}$  denote a set of nine variables, where  $x \in \{a, b, c\}$  ( $a$  for Abel,  $b$  for Bob, and  $c$  for Carol), and  $y \in \{c, m, b\}$  ( $c$  for car,  $m$  for motorcycle, and  $b$  for boat). Now the property that “everyone uses a unique transportation tool” can be specified formally as the following clauses:

- |   |                                      |
|---|--------------------------------------|
| (1) $(p_{a,c} \mid p_{a,m} \mid p_{a,b})$           | Abel uses $c$ , $m$ , or $b$ ;       |
| (2) $(p_{b,c} \mid p_{b,m} \mid p_{b,b})$           | Bob uses $c$ , $m$ , or $b$ ;        |
| (3) $(p_{c,c} \mid p_{c,m} \mid p_{c,b})$           | Carol uses $c$ , $m$ , or $b$ ;      |
| (4) $(\overline{p_{a,c}} \mid \overline{p_{b,c}})$  | Abel and Bob cannot both use $c$ ;   |
| (5) $(\overline{p_{a,m}} \mid \overline{p_{b,m}})$  | Abel and Bob cannot both use $m$ ;   |
| (6) $(\overline{p_{a,b}} \mid \overline{p_{b,b}})$  | Abel and Bob cannot both use $b$ ;   |
| (7) $(\overline{p_{a,c}} \mid \overline{p_{c,c}})$  | Abel and Carol cannot both use $c$ ; |
| (8) $(\overline{p_{a,m}} \mid \overline{p_{c,m}})$  | Abel and Carol cannot both use $m$ ; |
| (9) $(\overline{p_{a,b}} \mid \overline{p_{c,b}})$  | Abel and Carol cannot both use $b$ ; |
| (10) $(\overline{p_{b,c}} \mid \overline{p_{c,c}})$ | Bob and Carol cannot both use $c$ ;  |
| (11) $(\overline{p_{b,m}} \mid \overline{p_{c,m}})$ | Bob and Carol cannot both use $m$ ;  |
| (12) $(\overline{p_{b,b}} \mid \overline{p_{c,b}})$ | Bob and Carol cannot both use $b$ .  |

It appears quite cumbersome that 12 clauses are needed to specify one property. Note that the goal of using logic is to find a solution automatically. For computers, 12 million clauses are not too many.

From the clues of the puzzle, we have the following two unit clauses:

- |                             |                             |
|-----------------------------|-----------------------------|
| (13) $(p_{c,c})$            | Carol drove the car;        |
| (14) $(\overline{p_{a,b}})$ | Abel hates travel on water. |

These two unit clauses forces  $p_{c,c}$  to be assigned 1 and  $p_{a,b}$  to be 0. Then we can decide the truth values of the rest seven variables:

$p_{a,c} \mapsto 0$	from	$p_{c,c} \mapsto 1, (7);$
$p_{b,c} \mapsto 0$	from	$p_{c,c} \mapsto 1, (11);$
$p_{a,m} \mapsto 1$	from	$p_{a,c} \mapsto 0, p_{a,b} \mapsto 0, (1);$
$p_{b,m} \mapsto 0$	from	$p_{a,m} \mapsto 1, (5);$
$p_{c,m} \mapsto 0$	from	$p_{a,m} \mapsto 1, (8);$
$p_{b,b} \mapsto 1$	from	$p_{b,c} \mapsto 0, p_{b,m} \mapsto 0, (2);$
$p_{c,b} \mapsto 0$	from	$p_{b,b} \mapsto 1, (12).$

The second assumed property is that “every transportation tool has a unique color”. If we read the puzzle carefully, the clues involve the relation that “who will use what color of the transportation”. So it is better to use  $q_{x,z}$  to represent this relation that person  $x$  uses color  $y$ , where  $x \in \{a, b, c\}$  and  $z \in \{b, o, g\}$  ( $b$  for blue,  $o$  for orange, and  $g$  for green). We may specify the property that “everyone uses a unique color” in a similar way as we specify “everyone uses a unique transportation tool”. The clues give us the following unit clauses:  $(q_{a,o})$ , and  $(\overline{q_{b,g}})$ . From these two unit clauses, we can derive the truth values of 12  $q_{x,y}$  variables.

The next puzzle looks quite different from the previous one, but its specification is similar.

**Example 2.5.3.** Four people sit on a bench of four seats for a group photo: two Americans (A), a Briton (B), and a Canadian (C). They asked that (i) two Americans do not want to sit next to each other; (ii) the Briton likes to sit next to the Canadian. Suppose we use propositional variables  $X_y$ , where  $X \in \{A, B, C\}$  and  $1 \leq y \leq 4$ , with the meaning that “ $X_y$  is true iff the people of nationality  $X$  sits at seat  $y$  of the bench”. How do you specify the problem in CNF over  $X_y$  such that the models of the CNF matches exactly all the sitting solutions of the four gentlemen?

The CNF will specify the following conditions:

1. Every seat takes at least one person:  $(A_1 \mid B_1 \mid C_1), (A_2 \mid B_2 \mid C_2), (A_3 \mid B_3 \mid C_3), (A_4 \mid B_4 \mid C_4)$ . The truth of these clauses implies that there are at least four people using these seats.
2. Each seat can take at most one person:  $(\neg A_1 \mid \neg B_1), (\neg A_1 \mid \neg C_1), (\neg B_1 \mid \neg C_1), (\neg A_2 \mid \neg B_2), (\neg A_2 \mid \neg C_2), (\neg B_2 \mid \neg C_2), (\neg A_3 \mid \neg B_3), (\neg A_3 \mid \neg C_3), (\neg B_3 \mid \neg C_3), (\neg A_4 \mid \neg B_4), (\neg A_4 \mid \neg C_4), (\neg B_4 \mid \neg C_4)$ . Combining with the first condition, this condition implies that there are exactly four persons sitting on the bench.

3. At most two Americans on the bench:  $(\neg A_1 \mid \neg A_2 \mid \neg A_3)$ ,  $(\neg A_1 \mid \neg A_2 \mid \neg A_4)$ ,  $(\neg A_1 \mid \neg A_3 \mid \neg A_4)$ ,  $(\neg A_2 \mid \neg A_3 \mid \neg A_4)$ . These clauses say that for every three seats, one of them must not be taken by an American.
4. At most one Briton on the bench:  $(\neg B_1 \mid \neg B_2)$ ,  $(\neg B_1 \mid \neg B_3)$ ,  $(\neg B_1 \mid \neg B_4)$ ,  $(\neg B_2 \mid \neg B_3)$ ,  $(\neg B_2 \mid \neg B_4)$ ,  $(\neg B_3 \mid \neg B_4)$ . These clauses say that for every two seats, one of them must not be taken by the Canadian.
5. At most one Canadian on the bench:  $(\neg C_1 \mid \neg C_2)$ ,  $(\neg C_1 \mid \neg C_3)$ ,  $(\neg C_1 \mid \neg C_4)$ ,  $(\neg C_2 \mid \neg C_3)$ ,  $(\neg C_2 \mid \neg C_4)$ ,  $(\neg C_3 \mid \neg C_4)$ .
6. The Americans do not want to sit next to each other:  $(\neg A_1 \mid \neg A_2)$ ,  $(\neg A_2 \mid \neg A_3)$ ,  $(\neg A_3 \mid \neg A_4)$ .
7. The Briton likes to sit next to a Canadian:  $(\neg C_1 \mid B_2)$ ,  $(\neg C_2 \mid B_1 \mid B_3)$ ,  $(\neg C_3 \mid B_2 \mid B_4)$ ,  $(\neg C_4 \mid B_3)$ ,  $(\neg B_1 \mid C_2)$ ,  $(\neg B_2 \mid C_1 \mid C_3)$ ,  $(\neg B_3 \mid C_2 \mid C_4)$ ,  $(\neg B_4 \mid C_3)$ . If we read  $(\neg C_2 \mid B_1 \mid B_3)$  as  $C_2 \rightarrow B_1 \vee B_3$ , then it means that if the Canadian takes the second seat, then the Briton must sit at either seat 1 or seat 3.

Some conditions such as there are at least two Americans, and at least one Briton and one Canadian, can also be easily specified by clauses. However, these clauses are logical consequence of the other clauses, as conditions 1 and 2 ensure that there are four persons sitting on the bench. If there are at most one Briton and one Canadian, then condition 2 is also redundant, and we may remove them safely. Note that the input clauses serve as a set of axioms. This set is minimal if no clauses are logical consequence of other clauses. Of course, without a through search, we cannot identify all the clauses that are logical consequences of the other clauses.

The next problem is called *Knights and Knaves*, which was coined by Raymond Smullyan in his 1978 work *What Is the Name of This Book?* The problem is actually a collection of similar logic puzzles where some characters can only answer questions truthfully, and others only falsely. The puzzles are set on a fictional island where all inhabitants are either *knights*, who always tell the truth, or *knaves*, who always lie. The puzzles involve a visitor to the island who meets small groups of inhabitants. Usually the aim is for the visitor to deduce the inhabitants' type from their statements, but some puzzles of this type ask for other facts to be deduced. The puzzles may also be to determine a yes-no question which the visitor can ask in order to discover a particular piece of information.

**Example 2.5.4.** You meet two inhabitants: Zoey and Mel. Zoey tells you that “Mel is a knave”. Mel says, “Neither Zoey nor I are knaves.” Can you determine who is a knight and who is a knave?

Let  $p$  stands for “Zoey is knight”, with the understanding that if  $p$  is true, then Zoey is a knight; if  $p$  is false, then Zoey is a knave. Similarly, define  $q$  as “Mel is knight”. Using a function *says*, what they said can be expressed as

$$(1) \text{ says}(Zoey, \neg q); \quad (2) \text{ says}(Mel, p \wedge q),$$

together with the two axioms:

$$\text{says}(knight, X) \equiv X; \quad \text{says}(knave, X) \equiv \neg X;$$

Now using a truth table to check all the truth values of  $p$  and  $q$ , only when  $p \mapsto 1$  and  $q \mapsto 0$ , both (1) and (2) are true. That is, we replace Zoey by knight and Mel by knave in (1) and (2), to obtain  $\neg q$  and  $\neg(p \wedge q)$ , respectively. Then  $\{p \mapsto 1, q \mapsto 0\}$  is a model of  $\neg q \wedge \neg(p \wedge q)$ . Thus the solution is “Zoey is a knight and Mel is a knave.”

**Example 2.5.5.** You meet two inhabitants: Peggy and Zippy. Peggy tells you that “of Zippy and I, exactly one is a knight”. Zippy tells you that only a knave would say that Peggy is a knave. Can you determine who is a knight and who is a knave?

Let  $p$  be “Peggy is knight” and  $q$  be “Zippy is knight”. Then what they said can be expressed as

$$(1) \text{ says}(Peggy, p \oplus q); \quad (2) \text{ says}(Zippy, \text{says}(knave, \neg p)).$$

The only solution is  $p \mapsto 0$  and  $q \mapsto 0$ , i.e., both Peggy and Zippy are knaves.

## 2.6 Exercise Problems

1. How many logical operators which take five arguments? That is, how many functions of type  $f : B^5 \rightarrow B$ , where  $B = \{0, 1\}$ ?
2. Suppose  $V_p = \{p, q, r\}$ . Then  $\mathcal{M}((p \rightarrow q) \wedge r) = ?$  and  $\mathcal{M}((p \rightarrow q) \wedge \neg q) = ?$
3. Prove Theorem 2.2.11, for any formulas  $A$  and  $B$ ,  $\mathcal{M}(A \wedge B) = \mathcal{M}(A) \cap \mathcal{M}(B)$  and  $\mathcal{M}(\neg A) = IV_p - \mathcal{M}(A)$ .
4. Provide the truth table for defining  $ite : B^3 \rightarrow B$ , the if-then-else operator.

5. Provide a BNF grammar which defines the well-formed formulas without unnecessary parentheses, assuming the precedence relation from the highest to the lowest:  $\neg$ ,  $\wedge$ ,  $\vee$ ,  $\rightarrow$ ,  $\{\oplus, \leftrightarrow\}$ .

6. Prove by the truth table method that the following formulas are valid.

- (a)  $(A \wedge (A \rightarrow B)) \rightarrow B$
- (b)  $(\neg B \rightarrow \neg A) \rightarrow (A \rightarrow B)$
- (c)  $(A \oplus B) \leftrightarrow ((A \vee B) \wedge (\neg A \vee \neg B))$
- (d)  $(A \leftrightarrow B) \leftrightarrow ((A \vee \neg B) \wedge (\neg A \vee B))$

7. Prove by the truth table method that the following entailments hold.

- (a)  $(A \wedge (A \rightarrow B)) \models B$
- (b)  $(A \rightarrow (B \rightarrow C)) \models (A \rightarrow B) \rightarrow (A \rightarrow C)$
- (c)  $(\neg B \rightarrow \neg A) \models (A \rightarrow B)$
- (d)  $((A \vee B) \wedge (\neg A \vee C)) \models (B \vee C)$

8. Prove by the truth table method that the following logical equivalences hold.

- (a)  $\neg(A \vee B) \equiv \neg A \wedge \neg B$
- (b)  $(A \vee B) \wedge (A \vee \neg B) \equiv A$
- (c)  $(A \leftrightarrow B) \equiv ((A \wedge B) \vee (\neg A \wedge \neg B))$
- (d)  $(A \oplus B) \equiv ((A \wedge \neg B) \vee (\neg A \wedge B))$
- (e)  $(\overline{A} \downarrow B) \uparrow \overline{C} \equiv A \uparrow (B \downarrow C)$
- (f)  $(\overline{A} \uparrow B) \downarrow \overline{C} \equiv A \downarrow (B \uparrow C)$

9. Prove the equivalence of the following formulas by using other equivalence relations:

- (a)  $A \oplus B \equiv \neg(A \leftrightarrow B) \equiv \neg A \leftrightarrow B \equiv A \leftrightarrow \neg B$
- (b)  $(A \rightarrow B) \wedge (A \rightarrow C) \equiv \neg A \vee (A \wedge B \wedge C)$
- (c)  $A \rightarrow (B \rightarrow (C \vee D)) \equiv (A \rightarrow C) \vee (B \rightarrow D)$ .

10. Prove Theorem 2.2.4. **Hint:** Induction on the structure of  $A$ .

11. Prove that for any formulas  $A$ ,  $B$ , and  $C$ , if  $A \models B$ , then  $A \wedge C \models B \wedge C$ .

12. Define all the sixteen binary Boolean functions by formulas using only  $\top$ ,  $\perp$ ,  $\neg$ , and  $\wedge$ .

13. Prove that the **nor** operator  $\downarrow$ , where  $x \downarrow y = \neg(x \vee y)$ , can serve as a minimally sufficient set of Boolean operators.

14. Convert the following formulas into equivalent CNFs:

- (a)  $p \rightarrow (q \wedge r)$ ;
- (b)  $(p \rightarrow q) \rightarrow r$ ;
- (c)  $\neg(\neg p \vee q) \vee (r \rightarrow \neg s)$ ;
- (d)  $p \vee (\neg q \wedge (r \rightarrow \neg p))$ ;
- (e)  $\neg(((p \rightarrow q) \rightarrow p) \rightarrow q)$ ;
- (f)  $(p \rightarrow q) \leftrightarrow (p \rightarrow r)$ ;

15. Convert the following formulas into equivalent DNFs:

- (a)  $p \rightarrow (q \wedge r)$ ;
- (b)  $(p \rightarrow q) \rightarrow r$ ;
- (c)  $\neg(\neg p \vee q) \vee (r \rightarrow \neg s)$ ;
- (d)  $p \vee (\neg q \wedge (r \rightarrow \neg p))$ ;
- (e)  $\neg(((p \rightarrow q) \rightarrow p) \rightarrow q)$ ;
- (f)  $(p \rightarrow q) \leftrightarrow (p \rightarrow r)$ ;

16. Construct full CNF for the following formulas, where  $\bar{a}$  is  $\neg a$ ,  $+$  is  $\vee$  and  $\wedge$  is omitted as production.

- (a)  $F = \bar{a}\bar{b}c + \bar{a}b\bar{c} + a\bar{b}\bar{c} + \bar{a}\bar{b}\bar{c}$ ;
- (b)  $G = abc + \bar{a}bc + a\bar{b}c + ab\bar{c}$ ;
- (c)  $H = a \oplus b \oplus c$ ;
- (d)  $I = (a \wedge b) \rightarrow (b \rightarrow c)$ .

17. Construct full DNF for the following formulas, where  $\bar{a}$  is  $\neg a$ ,  $+$  is  $\vee$  and  $\wedge$  is omitted as production.

- (a)  $F = (\bar{a} + \bar{b} + c)(\bar{a} + b + \bar{c})(a + \bar{b} + \bar{c})(\bar{a} + \bar{b} + \bar{c})$ ;
- (b)  $G = (a + b + c)(\bar{a} + b + c)(a + \bar{b} + c)(a + b + \bar{c})$ ;
- (c)  $H = a \oplus b \oplus c$ ;
- (d)  $I = (a \wedge b) \rightarrow (b \rightarrow c)$ .

18. Assume  $a > b > c$ , construct ROBDDs for the following formulas, where  $'$  is  $\neg$ ,  $+$  is  $\vee$  and  $\wedge$  is omitted as production.

- (a)  $F = abc + \bar{a}\bar{b}c + a\bar{b}\bar{c}$ ;  
 (b)  $G = \bar{a}bc + abc + \bar{a}\bar{b}c$ ;  
 (c)  $H = a \oplus b \oplus c$ ;  
 (d)  $I = (a \wedge b) \rightarrow (b \rightarrow c)$ .
19. Let  $\mathcal{F}$  be the set of all formulas built on the set of  $n$  propositional variables. Prove that there are exactly  $2^{2^n}$  classes of equivalent formulas among  $\mathcal{F}$ .
20. (a) Prove that  $\{\rightarrow, \neg\}$  is a minimally sufficient set of Boolean operators.  
 (b) Prove that  $\{\rightarrow, \perp\}$  is a minimally sufficient set of Boolean operators.
21. Prove that  $\{\wedge, \vee, \rightarrow, \top\}$  is not a sufficient set of Boolean operators.
22. Prove that (a)  $\{\wedge, \vee, \top, \perp\}$  is not a sufficient set of Boolean operators; (b) if we add any function  $o$  which cannot be expressed in terms of  $\wedge, \vee, \top$ , and  $\perp$ , then  $\{\wedge, \vee, \top, \perp, o\}$  is sufficient.
23. Identify the pairs of equivalent formulas from the following candidates:  
 (a)  $p \uparrow (q \uparrow r)$ , (b)  $(p \uparrow q) \uparrow r$ , (c)  $\uparrow(p, q, r)$ , (d)  $p \downarrow (q \downarrow r)$ , (e)  $(p \downarrow q) \downarrow r$ ,  
 (f)  $\downarrow(p, q, r)$ , (g)  $p \uparrow (q \downarrow r)$ , (h)  $(p \uparrow q) \downarrow r$ , (i)  $p \downarrow (q \uparrow r)$ , (j)  $(p \downarrow q) \uparrow r$
24. Try to find an equivalent DNF of the minimal size for the following DNF formulas using K-maps:  
 (a)  $f(A, B, C) = \bar{A}\bar{B}\bar{C} + \bar{A}B + AB\bar{C} + AC$ ;  
 (b)  $f(A, B, C) = \bar{A}B + B\bar{C} + BC + A\bar{B}\bar{C}$ ;  
 (c)  $f(A, B, C, D) = \bar{A}\bar{B}CD + \bar{A}BC + B\bar{C}\bar{D} + B\bar{C}D$ ;  
 (d)  $f(A, B, C, D) = m_0 + m_1 + m_5 + m_7 + m_8 + m_{10} + m_{14} + m_{15}$ , where  $m_i$  is a full minterm of  $A, B, C$  and  $D$ , and  $i$  is the decimal value of the binary string  $ABCD$ .
25. Using the K-map method and the procedure  $NNF2NOR(A)$ , to find an equivalent formula, which uses only NAND/NOR/NOT operators, for each of the formulas in the previous question.
26. Try to find an equivalent CNF of the minimal size for the following Boolean functions using K-maps:  
 (a)  $f(A, B, C) = \bar{A}\bar{B} + \bar{A}B + AB\bar{C} + AC$ ;

- (b)  $f(A, B, C) = \overline{A}BC + AB\overline{C} + A\overline{B}C + \overline{A}\overline{B}\overline{C}$ ;
- (c)  $f(A, B, C, D) = m_0 + m_1 + m_3 + m_4 + m_5 + m_7 + m_{12} + m_{13} + m_{15}$ ;
- (d)  $f(A, B, C, D) = m_0 + m_1 + m_2 + m_4 + m_5 + m_6 + m_8 + m_9 + m_{12} + m_{13} + m_{14}$ .
27. Try to find a minimal and equivalent formula, which uses only NAND/NOR/NOT operators, for each of the formulas in the previous question.
28. A function  $f : B^n \rightarrow B$ , where  $B = \{0, 1\}$ , is called *linear* if  $f(x_1, x_2, \dots, x_n) = a_0 + a_1x_1 + \dots + a_nx_n$  for suitable coefficients  $a_0, \dots, a_n \in B$ . Here  $+$  denotes the addition modulo 2, and the not written multiplication is the multiplication over integers.
- (a) Show that the above representation of a linear function  $f$  is unique.
- (b) Determine the number of  $n$ -ary linear functions.
29. Suppose  $x$  is a 32-bit positive integer. Provide a java program for each of the following programming tasks, using as much bit-wise operations as possible:
- (a) Check whether  $x$  is even or odd.
- (b) Get the position of the highest bit of 1 in  $x$ .
- (c) Get the position of the lowest bit of 1 in  $x$ .
- (d) Count trailing zeros in  $x$  as a binary number.
- (e) Shift  $x$  right by 3 bits and put the last three bits of  $x$  at the beginning of the result from the shifting. The whole process is called “rotate  $x$  right” by 3 bits.
- (f) Convert a string of decimal digits into a binary number without using multiplication by 10.
30. Suppose a formula  $A$  contains 15 variables, each variable is represented by an integer from 0 to 14. We use the last 15 bits of an integer variable  $x$  to represent an interpretation of  $A$ , such that  $i^{th}$  bit is 1 iff we assign 1 to variable  $i$ . Modify the algorithm  $eval(A, \sigma)$ , where  $\sigma$  is replaced by  $x$  using bit-wise operations. Design an algorithm which lists all the interpretations of  $A$  and apply each interpretation to  $A$  using  $eval(A, x)$ . The pseudo-code of the algorithm should be like a java program.
31. The meaning of the following propositional variables is given below:

- $t$ : “taxes are increased”,
- $e$ : “expenditures rise”,
- $d$ : “the debt ceiling is raised”,
- $c$ : “the cost of collecting taxes rises”,
- $g$ : “the government borrows more money”,
- $i$ : “interest rates increase”.

(a) express the following statements as propositional formulas; (b) convert the formulas into CNF.

Either taxes are increased or if expenditures rise then the debt ceiling is raised. If taxes are increased, then the cost of collecting taxes rises. If a rise in expenditures implies that the government borrows more money, then if the debt ceiling is raised, then interest rates increase. If taxes are not increased and the cost of collecting taxes does not increase then if the debt ceiling is raised, then the government borrows more money. The cost of collecting taxes does not increase. Either interest rates do not increase or the government does not borrow more money.

32. Specify the following puzzle in propositional logic and convert the specification into CNF. Find the model of your CNF and use this model to construct a solution of the puzzle.

Four kids, Abel, Bob, Carol, and David, are eating lunch on a hot summer day. Each has a big glass of water, a sandwich, and a different type of fruit (apple, banana, orange, and grapes). Which fruit did each child have? The given clues are:

- (a) Abel and Bob have to peel their fruit before eating.
- (b) Carol doesn't like grapes.
- (c) Abel has a napkin to wipe the juice from his fingers.

33. Specify the following puzzle in propositional logic and convert the specification into CNF. Find the model of your CNF and use this model to construct a solution of the puzzle.

Four electric cars are parked at a charge station and their mileages are 10K, 20K, 30K, and 40K, respectively. Their charge costs are \$15, \$30, \$45, and \$60, respectively. Figure out how the charge cost for each car from the following hints:

- (a) The German car has 30K miles.
  - (b) The Japanese car has 20K miles and it charged \$15.
  - (c) The French car has less miles than the Italian car.
  - (d) The vehicle that charged \$30 has 10K more mileages than the Italian car.
  - (e) The car that charged \$45 has 20K less mileage than the German car.
34. Specify the following puzzle in propositional logic and convert the specification into CNF. Find the model of your CNF and use this model to construct a solution of the puzzle.
- Four boys are waiting for Thanksgiving dinner and each of them has a unique favorite food. Their names are Larry, Nick, Philip, and Tom. Their ages are all different and fall in the set of  $\{8, 9, 10, 11\}$ . Find out which food they are expecting to eat and how old they are.
- (a) Larry is looking forward to eating turkey.
  - (b) The boy who likes pumpkin pie is one year younger than Philip.
  - (c) Tom is younger than the boy that loves mashed potato.
  - (d) The boy who likes ham is two years older than Philip.
35. Six people sit on a long bench for a group photo. Two of them are Americans (A), two are Canadians (C), one Italian (I), and one Spaniard (S). Suppose we use propositional variables  $X_y$ , where  $x \in \{A, C, I, S\}$  and  $1 \leq y \leq 6$ , and  $X_y$  is true iff the people of nationality  $X$  sits at seat  $y$  of the bench. Please specify the problem formally in a propositional formula using  $X_y$ , and convert the formula into CNF such that each model of the CNF gives us one solution for all the six people to sit on the bench, enforcing the rules that (a) the people from the same country will not sit beside one another; (b) the Italian must be surrounded by Americans; and (c) each American sits next to a Canadian. Please list all the models of your CNF.
36. A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie. You meet two inhabitants: Sally and Zippy. Sally claims, "I and Zippy are not the same." Zippy says, "Of I and Sally, exactly one is a knight." Can you determine who is a knight and who is a knave?

37. A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie. You meet two inhabitants: Homer and Bozo. Homer tells you, "At least one of the following is true: that I am a knight or that Bozo is a knight." Bozo claims, "Homer could say that I am a knave." Can you determine who is a knight and who is a knave?
38. A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie. You meet two inhabitants: Bart and Ted. Bart claims, "I and Ted are both knights or both knaves." Ted tells you, "Bart would tell you that I am a knave." Can you determine who is a knight and who is a knave?

### CHAPTER 3

## PROOF PROCEDURES FOR PROPOSITIONAL LOGIC

In Chapter 1, we introduced the concept of “proof procedure”: Given a set  $A$  of axioms (i.e., the formulas assumed to be true) and a formula  $B$ , a proof procedure  $P(A, B)$  will answer the question whether  $B$  is a theorem of  $A$ : If  $P(A, B)$  returns true, we say  $B$  is proved, if  $P(A, B)$  returns false, we say  $B$  is disproved. For propositional logic, the axiom set  $A$  is regarded to be equivalent to the conjunction of formulas and a theorem is any formula  $B$  such that  $A \models B$ , i.e.,  $B \in \mathcal{T}(A)$ . That is, a proof procedure is also called a *theorem prover*. The same procedure  $P$  can be used to show  $B$  is valid (i.e.,  $B$  is a tautology) by calling  $P(\top, B)$ , or that  $C$  is unsatisfiable by calling  $P(C, \perp)$ . That is, proving the entailment relation is no harder than proving validity or proving unsatisfiability, as the following theorem shows. The three problems are equivalent and belong to the class of co-NP-complete problems. In practice, a proof procedure is designed to solve one of the three problems as its direct goal. A proof procedure is called

- *theorem prover*  $P(X, Y)$  if it is to show that formula  $Y$  is a theorem (an entailment) of  $X$  (assuming a set of formulas is equivalent to the conjunction of all its members);
- *tautology prover*  $T(A)$  if it is to show formula  $A$  is valid (a tautology);
- *refutation prover*  $R(B)$  if it is to show formula  $B$  is unsatisfiable.

If we have one of the above three procedures, i.e.,  $P(X, Y)$ ,  $T(A)$ , and  $R(B)$ , we may implement the other two as follows.

- If  $P(X, Y)$  is available, implement  $T(A)$  as  $P(\top, A)$ , and  $R(B)$  as  $P(B, \perp)$ .
- If  $T(A)$  is available, implement  $P(X, Y)$  as  $T(X \rightarrow Y)$ , and  $R(B)$  as  $T(\neg B)$ .
- If  $R(B)$  is available, implement  $P(X, Y)$  as  $R(X \wedge \neg Y)$ , and  $T(A)$  as  $R(\neg A)$ .

The correctness of the above implementations is based on the following theorem: Given formulas  $A$  and  $B$ , the following three statements are equivalent:

1.  $A \models B$ ;

2.  $A \rightarrow B$  is valid;
3.  $A \wedge \neg B$  is unsatisfiable.

Thus, one prover can do all the three, either directly or indirectly. In Chapter 2, we have discussed several concepts and some of them can be used to construct proof procedures as illustrated below.

- Truth table: It can be used to construct a theorem prover, a tautology prover, and a refutation prover.
- CNF: It can be used to construct a tautology prover, as a CNF is valid iff every clause is valid and it is very easy to check if a clause is valid.
- DNF: It can be used to construct a refutation prover, as a DNF is unsatisfiable iff every minterm is unsatisfiable and it is very easy to check if a minterm is unsatisfiable.
- ROBDD or INF: It can be used to construct a tautology prover or a refutation prover, because ROBDD is a canonical form.
- Algebraic substitution: It can be used to construct a tautology prover or a refutation prover, as substituting “equal by equal” preserves the equivalence relation. If we arrive at  $\top$  or  $\perp$ , we can claim the original formula is valid or unsatisfiable.

In terms of its working principle, a proof procedure may have one of three styles: enumeration, reduction, and deduction.

- *enumeration-style*: enumerate all involved interpretations;
- *reduction-style*: use the equivalence relations to transform a set of formulas into a simpler set;
- *deduction-style*: use inference rules to deduce new formulas.

By definition,  $A \models B$  iff  $\mathcal{M}(A) \subseteq \mathcal{M}(B)$ , that is, every model of  $A$  is a model of  $B$ . We may use a truth table to enumerate all the interpretations involving  $A$  and  $B$ , to see if every model of  $A$  is a model of  $B$ . A proof procedure based truth tables is an enumeration-style proof procedure which uses exhaustive search. This type of proof procedures are easy to describe but highly inefficient. For example, if  $A$  contains 30 variables, we need to construct a truth table of  $2^{30}$  rows.

Reduction-style proof procedures use the equivalence relations to transform a set of formulas into desired, simplified formulas, such as reduced ordered binary decision diagrams (ROBDD). Since ROBDD is canonical, we decide that a formula is valid if its ROBDD is  $\top$  (or 1), or unsatisfiable, if its ROBDD is  $\perp$  (or 0). We may also use disjunctive normal forms (DNF) to show a formula is unsatisfiable: Its simplified DNF is  $\perp$ . In the following, we will introduce the semantic tableau method, which is a reduction-style proof procedure and works in the same way as obtaining DNF from the original formula.

Deduction-style proof procedures use inference rules to generate new formulas, until the desired formula is generated. An inference rule typically takes one or two formulas as input and generates one or two new formulas as output. Some of these rules come from the equivalence relations between the input and the output formulas. In general, we require that the output of an inference rule is entailed by the input. There are many deduction-style proof procedures and in this chapter we will present several of them. The goal of this presentation is not for introducing efficient proof procedures, but for introducing a formalism such that other logics can use. This comment applies to the semantics tableau method, too.

Many proof procedures are not clearly classified as enumeration, reduction, or deduction styles. Yet, typically their working principles can be understood as performing enumeration, or reduction, or deduction, or all actions implicitly. For example, converting a formula into DNF can be a reduction, however, converting a formula into full DNF can be an enumeration, as every term in a full DNF corresponds to one model of the formula. In practice, reduction is often used in most deduction-style proof procedures for efficiency.

In Chapter 1, we say that a proof procedure is decision procedure if the procedure terminates for every input, that is, the procedure is an algorithm. Fortunately for propositional logic, every proof procedure introduced in this chapter is a decision procedure.

### 3.1 Semantic Tableau

The semantic tableau method, also called the “truth tree” method, is a proof procedure not just for propositional logic, but also for many other logics, including first-order logic, temporal and modal logics. A tableau for a formula is a tree structure, each node of the tree is associated with a list of formulas derived from the original formula. The tableau method constructs this tree and uses it to prove/refute the satisfiability of the original formula. The tableau method can also determine the satisfiability of finite sets of formulas of various logics. It is the most popular

proof procedure for modal logics.

### 3.1.1 Tableau: A Tree Structure for DNF

For propositional logic, the semantic tableau checks whether a formula is satisfiable or not, by “breaking” complex formulas into smaller ones as we do for transforming the formula into DNF (disjunctive normal form). The method works on a tree whose nodes are labeled with a list of formulas. At each step, we may extend a leaf node by adding one or two successors to this node.

- Each node of the tree is associated with a list of formulas and the comma “,” in the list is a synonym to  $\wedge$ ;
- A node is called *closed* if it contains a complementary pairs of literals (or formulas), or  $\perp$ , or  $\neg\top$ .
- A node is called a *model* node if it is a set of consistent literals, which is equivalent to a minterm in DNF.
- If a node is neither closed nor a model node, then it is *expandable* by the transformation rules for DNF, which apply only to the top logical operation of a formula in the list. The expansion stops when the tree has no expandable nodes.
- The tree links indicate the equivalence relation among formulas. If the result of the transformation rule is a disjunction, two successors of the current node are created and the formulas of the current node is equivalent to the disjunction of those formulas in the two successors. This type of the transformation rules is called  $\beta$ -rules; the other rules are called  $\alpha$ -rules, which generate only one child and the formulas of the current node is equivalent to those of the child node.

To prove the unsatisfiability of a formula, the method starts by generating the tree and check that every leaf node is closed, that is, no model node and no expandable nodes as leaves. In this case, we say the tree is *closed*.

**Example 3.1.1.** For the formulas  $A = p \wedge (\neg q \vee \neg p)$  and  $B = (p \vee q) \wedge (\neg p \wedge \neg q)$ , their tableaux are shown in Figure 3.1.1. For  $A$ , the tableau has one model node (marked by  $\odot$ ) and one closed node (marked by  $\times$ ). For  $B$ , it has only two closed nodes. Thus  $A$  is satisfiable as it has a model node; and  $B$  is unsatisfiable as its tableau is closed.

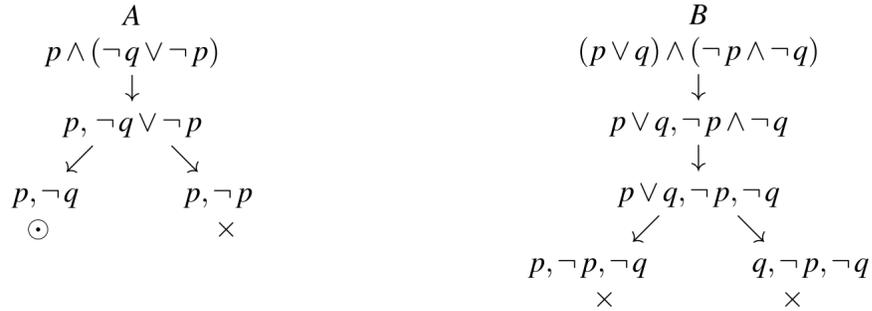


Figure 3.1.1: The tableaux for  $A = p \wedge (\neg q \vee \neg p)$  and  $B = (p \vee q) \wedge (\neg p \wedge \neg q)$

The tableaux can be also drawn linearly, using proper labels for the parent-child relation: The root node has the empty label. If a node has label  $x$ , its first child has label  $x1$  and its second child has label  $x2$ , and so on. For example, we may display the tableaux of  $A$  and  $B$  in Figure 3.1.1 as follows:

$p \wedge (\neg q \vee \neg p)$	$(p \vee q) \wedge (\neg p \wedge \neg q)$
1 : $p, (\neg q \vee \neg p)$	1 : $(p \vee q), (\neg p \wedge \neg q)$
11 : $p, \neg q \quad \odot$	11 : $(p \vee q), \neg p, \neg q$
12 : $p, \neg p \quad \times$	111 : $p, \neg p, \neg q \quad \times$
	112 : $q, \neg p, \neg q \quad \times$
(a) Proof of $A$	(b) Proof of $B$

The labels of nodes can uniquely identify the positions (or address) of these nodes in a tableau.

### 3.1.2 $\alpha$ -Rules and $\beta$ -Rules

Here is the listing of  $\alpha$ -rules and  $\beta$ -rules, which apply to the top logical operators in a formula.

$\alpha$	$\alpha_1, \alpha_2$
$A \wedge B$	$A, B$
$\neg(A \vee B)$	$\neg A, \neg B$
$\neg(A \rightarrow B)$	$A, \neg B$
$\neg(A \oplus B)$	$(A \vee \neg B), (\neg A \vee B)$
$A \leftrightarrow B$	$(A \vee \neg B), (\neg A \vee B)$
$\neg(A \uparrow B)$	$A, B$
$A \downarrow B$	$\neg A, \neg B$
$\neg\neg A$	$A$

$\beta$	$\beta_1$	$\beta_2$
$\neg(A \wedge B)$	$\neg A$	$\neg B$
$A \vee B$	$A$	$B$
$A \rightarrow B$	$\neg A$	$B$
$A \oplus B$	$A, \neg B$	$\neg A, B$
$\neg(A \leftrightarrow B)$	$A, \neg B$	$\neg A, B$
$A \uparrow B$	$\neg A$	$\neg B$
$\neg(A \downarrow B)$	$A$	$B$

If we start with a formula in NNF (negation normal form), all the  $\alpha$ -rules and  $\beta$ -rules with the leading symbol  $\neg$  are not needed.

For every  $\alpha$ -rule, we transform formula  $\alpha$  into  $\alpha_1, \alpha_2$ ; for every  $\beta$ -rule, we transform formula  $\beta$  into  $\beta_1$  and  $\beta_2$ . It is easy to check that  $\alpha \equiv \alpha_1 \wedge \alpha_2$  and  $\beta \equiv \beta_1 \vee \beta_2$ . Those rules are what we used for transforming a formula into DNF. Since a tableau maintains the property that the formulas of the current node is equivalent to the disjunction of the formulas in its children nodes (the formulas in the same node are considered in conjunction), a tableau is a tree-like representation of a formula that is a disjunction of conjunctions, i.e., a DNF.

**Example 3.1.2.** To show that Frege's formula,  $(p \rightarrow (q \rightarrow r)) \rightarrow (p \rightarrow q) \rightarrow (p \rightarrow r)$ , is valid, we show that the semantic tableau of its negation is closed.

	$\neg((p \rightarrow (q \rightarrow r)) \rightarrow (p \rightarrow q) \rightarrow (p \rightarrow r))$	$\alpha \neg \rightarrow$
1 :	$(p \rightarrow (q \rightarrow r)), \neg(p \rightarrow q) \rightarrow (p \rightarrow r)$	$\alpha \neg \rightarrow$
11 :	$(p \rightarrow (q \rightarrow r)), (p \rightarrow q), \neg(p \rightarrow r)$	$\alpha \neg \rightarrow$
111 :	$(p \rightarrow (q \rightarrow r)), (p \rightarrow q), p, \neg r$	$\beta \rightarrow$
1111 :	$\neg p, (p \rightarrow q), p, \neg r$	$\times$
1112 :	$(q \rightarrow r), (p \rightarrow q), p, \neg r$	$\beta \rightarrow$
11121 :	$(q \rightarrow r), \neg p, p, \neg r$	$\times$
11122 :	$(q \rightarrow r), q, p, \neg r$	$\beta \rightarrow$
111221 :	$\neg q, q, p, \neg r$	$\times$
111222 :	$r, q, p, \neg r$	$\times$

In Chapter 1, we pointed out that there are two important properties concerning a proof procedure: soundness and completeness. For a refutation prover like semantic tableau, the soundness means that if the procedure says the formula is unsatisfiable, then the formula must be unsatisfiable. The completeness means that if the formula is unsatisfiable, then the procedure should be able to give us this result.

**Theorem 3.1.3.** *The semantic tableau is a decision procedure for propositional logic.*

*Proof.* Suppose the list of formulas at one node is  $\alpha, \beta$ , and  $\gamma$  and the formula represented by this node is  $A = \alpha \wedge \beta \wedge \gamma$ . If  $\alpha$  is transformed by a  $\alpha$ -rule to  $\alpha_1$  and  $\alpha_2$ , then from  $\alpha \equiv \alpha_1 \wedge \alpha_2$ ,  $A \equiv \alpha_1 \wedge \alpha_2 \wedge \beta \wedge \gamma$ . If  $\beta$  is transformed by a  $\beta$ -rule to  $\beta_1$  and  $\beta_2$ , then from  $\beta \equiv \beta_1 \vee \beta_2$ ,  $A \equiv (\alpha \wedge \beta_1 \wedge \gamma) \vee (\alpha \wedge \beta_2 \wedge \gamma)$ . In other words, the logical equivalence is maintained for every parent-child link in the semantic tableau. By a simple induction on the structure of the tree, the formula at the root of the tree is equivalent to the disjunction of all the formulas represented at the leaf nodes.

That is, the procedure modifies the tableau in such a way that the disjunction of the formulas represented by the leaf nodes of the resulting tableau is equivalent to the original one. One of these conjunctions may contain a pair of complementary literals, in which case that conjunction is equivalent to false. If all conjunctions are equivalent to false, the original formula is unsatisfiable. This shows the soundness of the procedure.

If a formula is unsatisfiable, then every term in its DNF must be equivalent to false and the corresponding node in the tableau is closed. Thus we will have a closed tableau. This shows the completeness of the procedure.

Since neither  $\alpha$ -rules nor  $\beta$ -rules can be used forever (each rule reduces either one occurrence of a binary operator or two negations), the procedure must be terminating.  $\square$

Of course, semantic tableau can be used to show the satisfiability of a formula. If there is a model node in the tableau, we may assign 1 to every literal in this node so that the formula represented by this node is true under this interpretation. The original formula will be true under this interpretation because it is equivalent to the disjunction of the formulas of the leaf nodes. Thus, this interpretation is a model of the original formula. This shows that the procedure can be used to find a model when the formula is satisfiable.

Semantic tableau specifies what rules to be used but it does not specify which rule should be used first or which node should be expanded first. In general,  $\alpha$ -rules are better be applied before  $\beta$ -rules, so that the number of nodes in a tree can be reduced. In practice, the order of formulas in a list can be changed and the order of  $\beta_1$  and  $\beta_2$  can be switched. This gives the user the freedom to construct trees of different shapes, or to find a model node quicker. The procedure may terminate once a model node is found, if the goal is to show that the original formula is satisfiable.

Semantic tableaux are much more expressive and easy to use than truth-tables, though that is not the reason for their introduction. Sometimes, a tableau uses more space than a truth table. For example, an unsatisfiable full CNF of  $n$  variables will generate a tableau of more than  $n!$  nodes, which is large than a truth table of  $2^n$  lines. The beauty of semantic tableaux lies on the simplicity of presenting a proof procedure using a set of  $\alpha$  and  $\beta$  rules. We will see in later chapters how new rules are added into this procedure so that a proof procedure for other logics can be obtained.

### 3.2 Deductive Systems

In logic, a deductive system  $\mathcal{S}$  consists of a set of inference rules, where each rule takes premises as input and returns a conclusion (or conclusions) as output. Popular inference rules in propositional logic include modus ponens (MP), modus tollens (MT), and contraposition (CP), which can be displayed, respectively, as follows:

$$\frac{A \rightarrow B \quad A}{B} \text{ (MP)} \quad \frac{A \rightarrow B \quad \neg B}{\neg A} \text{ (MT)} \quad \frac{A \rightarrow B}{\neg B \rightarrow \neg A} \text{ (CP)}$$

#### 3.2.1 Inference Rules and Proofs

In general, an inference rule can be specified as

$$\frac{P_1 \quad P_2 \quad \cdots \quad P_k}{C}$$

where  $P_1, P_2, \dots, P_k$  are the premises and  $C$  is the conclusion. We may also write it as

$$P_1, P_2, \dots, P_k \vdash C$$

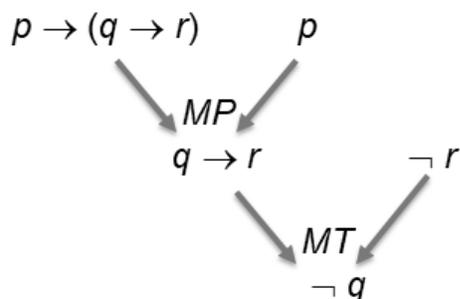
If the premises are empty ( $k = 0$ ), we say this inference rule is an *axiom rule*.

An inference rule is *sound* if its premises entail the conclusion, that is,

$$P_1, P_2, \dots, P_k \models C$$

for every inference rule  $\{P_1, P_2, \dots, P_k\} \vdash C$  (or equivalently,  $(P_1 \wedge P_2 \wedge \cdots \wedge P_k) \rightarrow C$  is a tautology). For example, MP (modus ponens) is sound because  $\{A \rightarrow B, A\} \models B$ ; modus tollens is sound because  $\{A \rightarrow B, \neg B\} \models \neg A$ .

Given a set  $A$  of axioms, which are the formulas assumed to be true, a *proof* of formula  $B$  in  $\mathcal{S}$  is a sequence of formulas  $F_1, F_2, \dots, F_n$  such that  $F_n = B$  and each  $F_i$  is either a formula in  $A$  or can be generated by an inference rule of  $\mathcal{S}$ , using the formulas before  $F_i$  in the sequence as the premises. If such a proof exists, we denote it by  $A \vdash_{\mathcal{S}} B$ ; the subscript  $\mathcal{S}$  can be dropped if  $\mathcal{S}$  is understood from the context. The proof procedure  $P(A, B)$  based on  $\mathcal{S}$  is simply to show  $A \vdash_{\mathcal{S}} B$ . Obviously, this is a deduction-style proof procedure.

Figure 3.2.1: The proof tree of  $\neg q$ 

**Example 3.2.1.** Let  $A = \{p \rightarrow (q \rightarrow r), p, \neg r\}$ , and  $\mathcal{S} = \{MP, CT\}$ . A proof of  $A \vdash_{\mathcal{S}} \neg q$  is given below:

1.  $p \rightarrow (q \rightarrow r)$  assumed
2.  $p$  assumed
3.  $q \rightarrow r$  MP, 1, 2
4.  $\neg r$  assumed
5.  $\neg q$  MT, 3, 4

The formulas in a proof can be rearranged so that all axioms appear in the beginning of the sequence. For example, the above proof can be rewritten as the following.

1.  $p \rightarrow (q \rightarrow r)$  assumed
2.  $p$  assumed
3.  $\neg r$  assumed
4.  $q \rightarrow r$  MP, 1, 2
5.  $\neg q$  MP, 3, 4

Every proof can be represented by a directed graph where each node is a formula in the proof and if a formula  $B$  is derived from  $A_1, A_2, \dots, A_k$  by an inference rule, then there are edges  $(A_i, B)$ ,  $1 \leq i \leq k$ , in the graph. This proof graph must be acyclic as we assume that no formula can be removed from a proof without destroying the proof. Every topological sort of the nodes of this graph should give us a linear proof. Since every acyclic direct graph (DAG) can be converted into a tree (by duplicating some nodes) and still preserve the parent-child relation, we may present a proof by a tree as well. The proof tree for the above proof example is given in Figure 3.2.1. It is clear that all assumed formulas (axioms) do not have the incoming edges. The final formula in the proof does not have outgoing edges. All intermediate formulas have both incoming and outgoing edges.

The soundness of an inference system comes from the soundness of all the inference rules. In general, we require all inference rules be sound to preserves truth of all the derived formulas. For a sound inference system  $\mathcal{S}$ , every derived formula  $B \in \mathcal{T}(A)$ . If  $A$  is a set of tautologies, then  $B$  is a tautology, too.

**Theorem 3.2.2.** *If an inference system  $\mathcal{S}$  is sound and  $A \vdash_{\mathcal{S}} B$ , then  $A \models B$ .*

*Proof.*  $A \vdash_{\mathcal{S}} B$  means there is a proof of  $B$  in  $\mathcal{S}$ . Suppose the proof is  $F_1, F_2, \dots, F_n = B$ . By induction on  $n$ , we show that  $A \models F_i$  for all  $i = 1, \dots, n$ . If  $F_i$  is in  $A$ , then  $A \models F_i$ . If  $F_i$  is derived from an inference rule  $P_1, P_2, \dots, P_m \vdash C$ , using  $F_{j_1}, F_{j_2}, \dots, F_{j_m}$  as premises, by induction hypotheses,  $A \models F_{j_k}$ , because  $j_k < i$ , for all  $1 \leq k \leq m$ . That is,  $A \models \bigwedge_{k=1}^m F_{j_k}$ .

Since the inference rule is sound,  $\bigwedge_{k=1}^m P_k \models C$ . Applying the substitution theorem,  $\bigwedge_{k=1}^m F_{j_k} \models F_i$ . Because  $\models$  is transitive,  $A \models F_i$ .  $\square$

Different inference systems are obtained by changing the axioms or the inference rules. In propositional logic, all these systems are equivalent in the sense that they are sound and complete.

### 3.2.2 Hilbert Systems

*Hilbert systems*, sometimes called Hilbert calculus, or Hilbert-style inference systems, are a type of deduction system attributed to Gottlob Frege and David Hilbert. Variants of Hilbert systems for propositional logic exist and these deductive systems are extended to first-order logic and other logics.

**Definition 3.2.3.** *The Hilbert system  $\mathcal{H}$  contains three axiom rules plus MP (modus ponens):*

$$\begin{aligned} H_1 &: \vdash (A \rightarrow (B \rightarrow A)), \\ H_2 &: \vdash (A \rightarrow (B \rightarrow C)) \rightarrow ((A \rightarrow B) \rightarrow (A \rightarrow C)), \\ H_3 &: \vdash (\neg B \rightarrow \neg A) \rightarrow (A \rightarrow B), \\ \text{MP} &: A, A \rightarrow B \vdash B. \end{aligned}$$

The axiom rule  $H_3$  is not needed if  $\rightarrow$  is the only operator in all formulas. Other logical operators can be defined using either  $\rightarrow$  and  $\neg$ . For example,  $A \wedge B \equiv \neg(A \rightarrow \neg B)$  and  $A \vee B \equiv \neg A \vee B$ .  $\neg A$  can be replaced by  $A \rightarrow \perp$  in  $H_3$ . It will turn out that all tautologies are derivable from the three formulas  $H_1$ ,  $H_2$ , and  $H_3$ .

**Example 3.2.4.** Assume  $A = \emptyset$ , to show that  $A \vdash_{\mathcal{H}} (p \rightarrow p)$ , we have the following

proof in  $\mathcal{H}$ :

- |  |                                       |
|--|---------------------------------------|
| 1. $p \rightarrow ((p \rightarrow p) \rightarrow p)$   | $H_1, A = p, B = (p \rightarrow p)$   |
| 2. $(p \rightarrow ((p \rightarrow p) \rightarrow p)) \rightarrow ((p \rightarrow (p \rightarrow p)) \rightarrow (p \rightarrow p))$ | $H_2, A = C = p, B = p \rightarrow p$ |
| 3. $(p \rightarrow (p \rightarrow p)) \rightarrow (p \rightarrow p)$   | MP from 1, 2                          |
| 4. $(p \rightarrow (p \rightarrow p))$   | $H_1, A = p, B = p$                   |
| 5. $p \rightarrow p$   | MP on 3, 4, $A = p, B = p$ .          |

**Proposition 3.2.5.**  $\mathcal{H}$  is sound: If  $A \vdash_{\mathcal{H}} B$ , then  $A \models B$ .

*Proof.* It is easy to check that every inference rule of  $\mathcal{H}$  is sound. That is, for every axiom rule  $H_i$ ,  $i = 1, 2, 3$ , the result of the rule is a tautology. For MP,  $(p \wedge (p \rightarrow q)) \rightarrow q$  is a tautology. By Theorem 3.2.2, the proposition holds.  $\square$

Hilbert systems originally assumed that  $A = \emptyset$  (or  $A = \top$ ), thus every derived formula is a tautology in a proof. This restriction is relaxed today so that we may take a set of formulas as axioms. The truth of the formulas becomes conditional on the truth of the axioms.

**Example 3.2.6.** Assume  $A = \{p \rightarrow q, q \rightarrow r\}$ , to show that  $A \vdash_{\mathcal{H}} (p \rightarrow r)$ , we have the following proof in  $\mathcal{H}$ :

- |  |              |
|--|--------------|
| 1. $p \rightarrow q$   | assumed      |
| 2. $q \rightarrow r$   | assumed      |
| 3. $((q \rightarrow r) \rightarrow (p \rightarrow (q \rightarrow r)))$                               | $H_1$        |
| 4. $p \rightarrow (q \rightarrow r)$   | MP from 2, 3 |
| 5. $(p \rightarrow (q \rightarrow r)) \rightarrow ((p \rightarrow q) \rightarrow (p \rightarrow r))$ | $H_2$        |
| 6. $(p \rightarrow q) \rightarrow (p \rightarrow r)$   | MP from 4, 5 |
| 7. $p \rightarrow r$   | MP from 1, 6 |

Since  $\mathcal{H}$  is sound, from  $A \vdash_{\mathcal{H}} (p \rightarrow r)$ , we have  $(p \rightarrow q) \wedge (q \rightarrow r) \models_{\mathcal{H}} (p \rightarrow r)$ , or  $((p \rightarrow q) \wedge (q \rightarrow r)) \rightarrow (p \rightarrow r)$  is valid.

Hilbert system is one of the earliest inference systems and has great impact to the creation of many inference systems. However, theorem proving in this system has always been considered a challenge because of its high complexity. Since Hilbert systems can hardly be a practical tool for problem solving, we state without proof the completeness result of Hilbert system.

**Theorem 3.2.7.** *The Hilbert system is sound and complete, that is,  $A \models B$  iff  $A \vdash_{\mathcal{H}} B$ .*

### 3.2.3 Natural Deduction

Both of semantic tableaux and Hilbert systems for constructing proofs have their disadvantages. Hilbert's systems are difficult to construct proofs in; their main uses are metalogical, the small number of rules making it easier to prove results about logic. The tableau method on the other hand is easy to use mechanically but, because of the form of the connective rules and the fact that a tableau starts from the negation of the formula to be proved, the proofs that result are not a natural sequence of easily justifiable steps. Likewise, very few proofs in mathematics are from axioms directly. Mathematicians in practice usually reason in a more flexible way.

Natural deduction is a deductive system in which logical reasoning is expressed by inference rules closely related to the "natural" way of reasoning. This contrasts with Hilbert systems, which instead use axiom rules as much as possible to express the logical laws of deductive reasoning. Natural deduction in its modern form was independently proposed by the German mathematician Gerhard Gentzen in 1934. In Gentzen's natural deductive system, a formula is represented by a set  $A$  of formulas,  $A = \{A_1, A_2, \dots, A_n\}$ , where the comma in  $A$  are understood as  $\vee$ . To avoid the confusion, we will write  $A$  as

$$(A_1 \mid A_2 \mid \dots \mid A_n).$$

$A$  is called a *sequent* by Gentzen. If  $n = 0$ , i.e.,  $A = ()$ , which is equivalent to  $\perp$  or  $(\perp)$ . Since we assume that  $A$  is a set, duplicated formulas are removed and the order of formulas in  $A$  is irrelevant.

**Definition 3.2.8.** *A natural deductive system  $\mathcal{G}$  consists of the following three inference rules, plus the following logical rules which guide the introduction and elimination of logical operators in the sequents.*

name	inference rule
axiom	$\vdash (A \mid \neg A \mid \alpha)$
cut	$(A \mid \alpha), (\neg A \mid \beta) \vdash (\alpha \mid \beta)$
thinning	$(\alpha) \vdash (A \mid \alpha)$

<i>op</i>	<b>introduction</b>	<b>elimination</b>
$\neg$	$(A \mid \alpha) \vdash (\neg\neg A \mid \alpha)$	$(\neg\neg A \mid \alpha) \vdash (A \mid \alpha)$
$\vee$	$(A \mid B \mid \alpha) \vdash (A \vee B \mid \alpha)$	$(A \vee B \mid \alpha) \vdash (A \mid B \mid \alpha)$
$\wedge$	$(A \mid \alpha), (B \mid \alpha) \vdash (A \wedge B \mid \alpha)$	(a) $(A \wedge B \mid \alpha) \vdash (A \mid \alpha)$ (b) $(A \wedge B \mid \alpha) \vdash (B \mid \alpha)$
$\rightarrow$	$(\neg A \mid B \mid \alpha) \vdash (A \rightarrow B \mid \alpha)$	$(A \rightarrow B \mid \alpha) \vdash (\neg A \mid B \mid \alpha)$
$\neg\vee$	$(\neg A \mid \alpha), (\neg B \mid \alpha) \vdash (\neg(A \vee B) \mid \alpha)$	(a) $(\neg(A \vee B) \mid \alpha) \vdash (\neg A \mid \alpha)$ (b) $(\neg(A \vee B) \mid \alpha) \vdash (\neg B \mid \alpha)$
$\neg\wedge$	$(\neg A \mid \neg B \mid \alpha) \vdash (\neg(A \wedge B) \mid \alpha)$	$(\neg(A \wedge B) \mid \alpha) \vdash (\neg A \mid \neg B \mid \alpha)$
$\neg\rightarrow$	$(A \mid \alpha), (\neg B \mid \alpha) \vdash (\neg(A \rightarrow B) \mid \alpha)$	(a) $(\neg(A \rightarrow B) \mid \alpha) \vdash (A \mid \alpha)$ (b) $(\neg(A \rightarrow B) \mid \alpha) \vdash (\neg B \mid \alpha)$

where  $\alpha$  and  $\beta$ , possibly empty, denote the rest formulas in a sequent.

In natural deduction, a sequent is deduced from a set of axioms by applying inference rules repeatedly. Each sequent is inferred from other sequents on earlier lines in a proof according to inference rules, giving a better approximation to the style of proofs used by mathematicians than Hilbert systems.

**Example 3.2.9.** The following is a proof of  $\neg(p \wedge q) \rightarrow (\neg p \vee \neg q)$  in  $\mathcal{G}$ :

1.  $(p \mid \neg p \mid \neg q)$  axiom,  $A = p, \alpha = \neg q$
2.  $(q \mid \neg q \mid \neg p)$  axiom,  $A = q, \alpha = \neg p$
3.  $(p \wedge q \mid \neg p \mid \neg q)$   $\wedge_I$  from 1, 2,  $A = p, B = q, \alpha = \neg p \mid \neg q$
4.  $(\neg\neg(p \wedge q) \mid \neg p \mid \neg q)$   $\neg_I$  from 3,  $A = (p \wedge q), \alpha = (\neg p \vee \neg q)$
5.  $(\neg\neg(p \wedge q) \mid (\neg p \vee \neg q))$   $\vee_I$  from 4,  $A = \neg p, B = \neg q, \alpha = \neg\neg(p \wedge q)$
6.  $(\neg(p \wedge q) \rightarrow (\neg p \vee \neg q))$   $\rightarrow_I$  from 5,  $A = \neg(p \wedge q), B = (\neg p \vee \neg q), \alpha = \perp$

**Proposition 3.2.10.** *The deductive system  $\mathcal{G}$  is sound, i.e., if  $A \vdash_{\mathcal{G}} B$ , then  $A \models B$ .*

*Proof.* It is easy to check that every inference rule of  $\mathcal{G}$  is sound. In particular, the rules of operator introduction preserve the logical equivalence. In the rules of operator elimination, some sequents generate two results. For example,  $(A \wedge B \mid \alpha)$  infers both (a)  $(A \mid \alpha)$  and (b)  $(B \mid \alpha)$ . If we bind the two results together by  $\wedge$ , these rules preserve the logical equivalence, too. On the other hand, the cut rule and the thinning rule do not preserve the logical equivalence; only the entailment relation holds. That is, for the cut rule,  $(A \mid \alpha), (\neg A \mid \beta) \vdash (\alpha \mid \beta)$ , we have

$$(A \vee \alpha) \wedge (\neg A \vee \beta) \models \alpha \vee \beta,$$

and for the thinning rule,  $(\alpha) \vdash (A \mid \alpha)$ , we have  $\alpha \models A \vee \alpha$ . □

Note that the proof in the example above uses only the axiom rule (twice) and the introduction rules for logical operators. This is the case when we are looking for a proof of a tautology. Recall how we would prove that  $\neg(p \wedge q) \rightarrow (\neg p \vee \neg q)$  is a tautology by semantic tableau. Since semantic tableau is a refutation prover, we need to show that the negation of a tautology has a closed tableau.

<i>a.</i>	:	$\neg(\neg(p \wedge q) \rightarrow (\neg p \vee \neg q))$	$\alpha$	$\neg \rightarrow$
<i>b.</i>	:	$\neg(p \wedge q), \neg(\neg p \vee \neg q)$	$\alpha$	$\neg \vee$
<i>c.</i>	:	$\neg(p \wedge q), \neg\neg p, \neg\neg q$	$\alpha$	$\neg$
<i>d.</i>	:	$\neg(p \wedge q), p, \neg\neg q$	$\alpha$	$\neg$
<i>e.</i>	:	$\neg(p \wedge q), p, q$	$\beta$	$\neg \wedge$
<i>f.</i>	1:	$\neg p, p, q$	$\times$	
<i>g.</i>	2:	$\neg q, p, q$	$\times$	

Comparing the proof of  $\mathcal{G}$  and the tableau above, we can see a clear correspondence: Closed nodes (f) and (g) correspond to (1) and (2), which are axioms; (e) to (3), (c) and (d) to (4), (b) to (5), and finally (a) to (6), respectively. That is, the formula represented by each sequent in a proof of  $\mathcal{G}$  can always find its negation in the corresponding node of the semantic tableau.

We show below that  $\mathcal{G}$  can be used as a refutation prover by the same example.

**Example 3.2.11.** Let  $A = \{\neg(\neg(p \wedge q) \rightarrow (\neg p \vee \neg q))\}$ . A proof of  $A \vdash \perp$  in  $\mathcal{G}$  is given below.

1.	( $\neg(\neg(p \wedge q) \rightarrow (\neg p \vee \neg q))$ )	assumed
2.	( $\neg(p \wedge q)$ )	$\neg \rightarrow_E$ (a) from 1, $A = \neg(p \wedge q)$
3.	( $\neg(\neg p \vee \neg q)$ )	$\neg \rightarrow_E$ (b) from 1, $B = (\neg p \vee \neg q)$
4.	( $\neg p \mid \neg q$ )	$\neg \wedge_E$ from 2, $A = p, B = q$
5.	( $\neg\neg p$ )	$\neg \vee_E$ (a) from 3, $A = \neg p$
6.	( $\neg\neg q$ )	$\neg \vee_E$ (b) from 3, $B = \neg q$
7.	( $p$ )	$\neg_E$ from 5, $A = p$
8.	( $q$ )	$\neg_E$ from 6, $A = q$
9.	( $\neg q$ )	cut from 4, 7, $A = p, \alpha = \perp, \beta = \neg q$
10.	( $\quad$ )	cut from 8, 9, $A = q, \alpha = \perp, \beta = \perp$

The above is indeed a proof of  $A \vdash_{\mathcal{G}} \perp$ , where  $\perp$  is used for the empty sequent. In this proof, the used inference rules are the cut rule and the rules for eliminating logical operators. If we apply repeatedly the rules of operator elimination to a formula, we will obtain a list of sequents, each of which represents a clause. In other words, these rules transform the formula into CNF, then the cut rule works on them. The cut rule for clauses has another name, i.e., *resolution*, which is the topic of the next section.

Indeed,  $\mathcal{G}$  contain more rules than necessary: Using only the axiom rule and the rules for operator introduction, it can imitate the dual of a semantic tableau upside down. Using only the cut rule and the rules for operator elimination, it can do what a resolution prover can do. Both semantic tableau and resolution prover are decision procedures for propositional logic. The implication is that natural deduction  $\mathcal{G}$  is a decision procedure, too. We give the main result without a proof.

**Theorem 3.2.12.** *The natural deduction  $\mathcal{G}$  is sound and complete, that is,  $A \models B$  iff  $A \vdash_{\mathcal{G}} B$ .*

### 3.2.4 Inference Graphs

We have seen two examples of deductive systems, Hilbert system  $\mathcal{H}$  and Gentzen's natural deduction  $\mathcal{G}$ . To show the completeness of these deductive systems, we need strategies on how to use the inference rules in these systems and the completeness depends on these strategies. To facilitate the discussion, let's introduce a few concepts.

**Definition 3.2.13.** *Given an inference system  $\mathcal{S}$ , the inference graph of  $\mathcal{S}$  over the axiom set  $A$  is defined as a directed graph  $G = (V, E)$ , where each vertex of  $V$  is a set  $\mathbf{x}$  of formulas, where  $A \in V$ ,  $\mathbf{x} \subseteq \mathcal{T}(A)$ , and  $(\mathbf{x}, \mathbf{y}) \in E$  iff  $\mathbf{y} = \mathbf{x} \cup \{C\}$ , where  $C$  is the conclusion of an inference rule  $\mathbf{r} \in \mathcal{S}$  using some formulas in  $\mathbf{x}$  as the premises of  $\mathbf{r}$ .*

The graph  $G = (V, E)$  defines the search space for  $A \vdash_{\mathcal{S}} B$ , which becomes a search problem: To find a proof of  $A \vdash_{\mathcal{S}} B$  is the same as to find a path from  $A$  to a node of  $V$  containing  $B$ . That is, if all the axioms of  $A$  in a proof appear in the beginning of the proof sequence, then this proof presents a path in the inference graph. Suppose the proof is a sequence of formulas  $F_1, F_2, \dots, F_n$  such that the first  $k$  formulas are from  $A$ , i.e.,  $F_i \in A$  for  $1 \leq i \leq k$  and  $F_i \notin A$  for  $i > k$ . Then this proof represents a directed path in  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{n-k}$  in the inference graph such that  $\mathbf{x}_0 = A$ ,  $F_i \in A$  for  $1 \leq i \leq k$ , and  $F_j \in \mathbf{x}_{j-k}$  for  $k < j \leq n$ . That is, a proof is a succinct way of presenting a directed path in the inference graph  $G$ .

**Example 3.2.14.** In Example 3.2.1, the second proof is copied here:

- |    |                                   |          |
|----|-----------------------------------|----------|
| 1. | $p \rightarrow (q \rightarrow r)$ | assumed  |
| 2. | $p$                               | assumed  |
| 3. | $\neg r$                          | assumed  |
| 4. | $q \rightarrow r$                 | MP, 1, 2 |
| 5. | $\neg r \rightarrow \neg q$       | CT, 4    |
| 6. | $\neg q$                          | MP, 3, 5 |

This proof corresponds to the path  $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4$  in the inference graph, where  $\mathbf{x}_0 = \{1., 2., 3.\}$ ,  $\mathbf{x}_1 = \mathbf{x}_0 \cup \{4.\}$ ,  $\mathbf{x}_2 = \mathbf{x}_1 \cup \{5.\}$ ,  $\mathbf{x}_3 = \mathbf{x}_2 \cup \{5.\}$ , and  $\mathbf{x}_4 = \mathbf{x}_3 \cup \{6.\}$ .

Given a set of inference rules, a proof procedure can be easily constructed if we use a *fair strategy* to search the inference graph. For example, starting from  $A$  in the inference graph, a breadth-first search is a fair strategy, but a depth-first search is not if the graph is infinite. Fair strategy requires that if an inference rule is applicable at some point, this application will happen over time. Any fair search strategies, heuristic search, best-first search,  $A^*$ , etc., can be used for a proof.

If every inference rule is sound, the corresponding proof procedure will be sound. Extra requirement like fair strategy is needed in general to show that the proof procedure is complete or halting. If the inference graph is finite, the termination of the proof procedure is guaranteed, but it does not imply the completeness of the procedure.

### 3.3 Resolution

In the previous section, we introduced two deductive systems, Hilbert systems and natural deduction. These two systems are exemplary for the development of many deductive systems, but are not practical for theorem proving in propositional logic. One of the reasons is the use of axiom rules which provide infinite possibilities because one rule may have infinite instances. Another reason is that, unlike reduction rules, the inference rules may be applied without termination. For automated theorem proving, we need a deductive system which contains very few rules as computers can do simple things fast. Resolution  $\mathcal{R}$  is such a system, which contains a single inference rule, called *resolution*, which is a special case of the cut rule when all sequents are clauses.

#### 3.3.1 Resolution Rule

To use the resolution rule, the formula  $A$  is first transformed into CNF, which is a conjunction of clauses (disjunction of literals). We will use  $A = \{C_1, C_2, \dots, C_m\}$  to denote the conjunction of the clauses; each clause  $C_i$  is denoted by

$$(l_1 \mid l_2 \mid \dots \mid l_k),$$

where each  $l_j$  is a literal (propositional variable or its negation). If  $k = 0$ ,  $C_i$  is the *empty clause*, denoted by  $\perp$ ; if  $k = 1$ ,  $C_i$  is called a *unit clause*; if  $k = 2$ ,  $C_i$  is called a *binary clause*. A clause is said to be *positive* if it is not empty and contains only positive literals; it is *negative* if it is not empty and contains only negative

literals. A clause is a tautology if it contains a pair of complementary literals, i.e., it contains both  $p$  and  $\neg p$  for some variable  $p$ .

The total number of clauses on  $n$  propositional variables is bound by  $3^n$ , because each variable has one of three cases in a non-tautology clause: missing, as a positive literal, or as a negative literal.

Suppose clause  $C_1$  is  $p \mid \alpha$  and  $C_2$  is  $\neg p \mid \beta$ , where  $\alpha$  and  $\beta$  are the rest literals in  $C_1$  and  $C_2$ , respectively, resolution may produce a new clause from  $C_1$  and  $C_2$ . Formally, the resolution rule is defined by the following schema, where  $A$  is a propositional variable.

$$\frac{(A \mid \alpha) \quad (\neg A \mid \beta)}{(\alpha \mid \beta)}$$

Alternatively, we may write the resolution rule as

$$(A \mid \alpha), (\neg A \mid \beta) \vdash (\alpha \mid \beta).$$

The clause  $(\alpha \mid \beta)$  produced by the resolution rule is called *resolvent* by resolving off  $A$  from  $C_1$  and  $C_2$ .  $C_1$  and  $C_2$  are the *parents* of the resolvent. This resolution rule is also called *binary resolution*, as it involves two clauses as premises.

The resolution rule for propositional logic can be traced back to Davis and Putnam (1960). Resolution is extended to first-order logic in 1965 by John Alan Robinson, using the unification algorithm. If we use  $\rightarrow$  instead of  $\mid$  in the resolution rule, we obtain

$$\frac{(\neg\alpha \rightarrow A) \quad (A \rightarrow \beta)}{(\neg\alpha \rightarrow \beta)}$$

which states the transitivity of  $\rightarrow$ . If  $\alpha$  is empty, then we obtain MP:

$$\frac{A \quad A \rightarrow \beta}{\beta}$$

**Example 3.3.1.** Let  $C = \{(p \mid q), (p \mid \neg q), (\neg p \mid q \mid r), (\neg p \mid \neg q), (\neg r)\}$ . A proof of  $A \vdash_{\mathcal{R}} \perp$  is given below:

1.  $(p \mid q)$  assumed
2.  $(p \mid \neg q)$  assumed
3.  $(\neg p \mid q \mid r)$  assumed
4.  $(\neg p \mid \neg q)$  assumed
5.  $(\neg r)$  assumed
6.  $(q \mid r)$   $\mathcal{R}$  from 1, 3,  $A = p, \alpha = q, \beta = (q \mid r)$
7.  $(\neg q)$   $\mathcal{R}$  from 2, 4,  $A = p, \alpha = \neg q, \beta = \neg q$
8.  $(r)$   $\mathcal{R}$  from 6, 7,  $A = q, \alpha = r, \beta = \perp$
9.  $()$   $\mathcal{R}$  from 8, 5,  $A = r, \alpha = \perp, \beta = \perp$

Note that we assume that duplicated literals are removed and the order of literals is irrelevant in a clause. Thus, the resolvent generated from clauses 1 and 3 is  $(q \mid q \mid r)$ , but we keep it as  $(q \mid r)$  by removing a copy of  $q$ . From clauses 1 and 4, we may obtain a resolvent  $(q \mid \neg q)$  on  $p$ , or another resolvent on  $(p \mid \neg p)$  on  $q$ . These two resolvents are tautologies and are not useful in the search for  $\perp$ . You cannot resolve on two variables at the same time to get  $(\ )$  directly from clauses 1 and 4.

**Definition 3.3.2.** *A resolution proof  $P$  from the set of input clauses  $C$  is a refutation proof of  $C \vdash_{\mathcal{R}} \perp$  and the length of  $P$  is the number of resolvents in  $P$ , denoted by  $|P|$ . The set of the input clauses appearing in  $P$  is called the core of  $P$ , denoted by  $\text{core}(P)$ .*

The length of the resolution proof in Example 3.3.1 is 4, which is the length of the list minus the number of assumed clauses in the proof. The core of this proof is the entire input  $C$ .

**Proposition 3.3.3.** *The resolution  $\mathcal{R}$  is sound, i.e., if  $C \vdash_{\mathcal{R}} B$ , then  $C \models B$ .*

*Proof.* Using truth table, it is easy to check that  $((A \vee \alpha) \wedge (\neg A \vee \beta)) \rightarrow (\alpha \vee \beta)$  is a tautology. Thus, the resolution rule is sound. The proposition holds following Theorem 3.2.2.  $\square$

**Corollary 3.3.4.** *If  $C \vdash_{\mathcal{R}} \perp$ , then  $C$  is unsatisfiable.*

*Proof.* Since  $C \vdash_{\mathcal{R}} \perp$  means  $C \models \perp$ , so  $\mathcal{M}(C) = \mathcal{M}(C \wedge \perp) = \mathcal{M}(\perp) = \emptyset$  (see Corollary 2.2.26).  $\square$

This corollary allows us to design a refutation prover: To show  $A$  is valid, we convert  $\neg A$  into an equivalent set  $C$  of clauses and show that  $C \vdash_{\mathcal{R}} \perp$ , where  $C$  is called the *input clauses*. Since the total number of clauses is finite, the inference graph of  $\mathcal{R}$  is finite. It means that the proof procedure will terminate but it does not mean the graph is small. Before we establish the completeness of  $\mathcal{R}$ , let us consider how to build an effective resolution prover.

The successful implementation of a resolution prover needs the integration of search strategies that reduce the search space by pruning unnecessary deduction paths. Some strategies remove redundant clauses as soon as they appear in a derivation. Some strategies avoid generating redundant clauses in the first place. Some strategies sacrifice the completeness for efficiency.

### 3.3.2 Resolution Strategies

Here are some well-known resolution strategies which add restriction to the use of the resolution rule.

- **unit resolution:** One of the two parents is a unit clause.
- **input resolution:** One of the two parents is an input clause (given as the axioms).
- **ordered resolution:** Given an order on propositional variables, the resolved variable must be maximal in both parent clauses.
- **positive resolution:** One of the two parents is positive.
- **negative resolution:** One of the two parents is negative.
- **set-of-support:** The input clauses are partitioned into two sets,  $S$  and  $T$ , where  $S$  is called the set of support, and a resolution is not allowed if both parents are in  $T$ .
- **linear resolution:** The latest resolvent is used as a parent for the next resolution (no restrictions on the first resolution).

A resolution proof is called  $X$  *resolution proof* if every resolution in the proof is an  $X$  resolution, where  $X$  is “unit”, “input”, “ordered”, “positive”, “negative”, “set-of-support”, or “linear”.

**Example 3.3.5.** Let  $C = \{1. (p \mid q), 2. (\neg p \mid r), 3. (p \mid \neg q \mid r), 4. (\neg r)\}$ . In the following proofs, the input clauses are omitted and the parents of each resolvent are shown as a pair of numbers following the resolvent.

(a)	(b)	(c)
5. $(\neg p)$ (2, 4)	5. $(p \mid r)$ (1, 3)	5. $(q \mid r)$ (1, 2)
6. $(p \mid \neg q)$ (3, 4)	6. $(r)$ (2, 5)	6. $(\neg q \mid r)$ (2, 3)
7. $(\neg q)$ (5, 6)	7. $(\ )$ (4, 6)	7. $(r)$ (5, 6)
8. $(q)$ (1, 5)		8. $(\ )$ (4, 7)
9. $(\ )$ (7, 8)		

Note that (a) is a unit and negative resolution proof; (b) is an input, positive and linear resolution proof; (c) is an ordered resolution proof, assuming  $p > q > r$ .

The length of (a) is 5; the length of (b) is 3; and the length of (c) is 4. It is easy to check that (a) is not an input resolution proof, because clause 7 is not obtained by input resolution. (b) is not a unit resolution proof. (c) is neither a unit nor an input resolution; it is neither positive nor negative resolution proof.

Efficiency is an important consideration in automated reasoning and one may sometimes be willing to trade completeness for speed. In unit resolution, one of the parent clauses is always a literal; in input resolution, one of the parent clauses is always selected from the original set. Albeit efficient, neither strategy is complete. For example,  $C = \{(p \mid q), (p \mid \neg q), (\neg p \mid q), (\neg p \mid \neg q)\}$  is unsatisfiable, though no unit resolutions are available. We may use input resolutions to obtain  $p$  and  $\neg p$ , but cannot obtain  $( )$  because  $p$  and  $\neg p$  are not input clauses.

Ordered resolution impose a total ordering on the propositional variables and treats clauses not as sets of literals but a sorted list of literals. Ordered resolution is extremely efficient and complete.

Set-of-support resolution is one of the most powerful strategies employed by (Wos, Carson, Robinson 1965), and its completeness depends on the choice of set of support. To prove  $A \models B$  by resolution, we convert  $A \wedge \neg B$  into clauses. The clauses obtained from  $\neg B$  serve in general as a set of support  $S$  and the clauses derived from  $A$  are in  $T$ . Set-of-support resolution dictates that the resolved clauses are not both from  $T$ . The motivation behind set-of-support is that since  $A$  is usually satisfiable it might be wise not to resolve two clauses from against each other.

Linear resolution always resolves a clause against the most recently derived resolvent. This gives the deduction a simple “linear” structure and it appears that we can have a straightforward implementation. Because linear resolution imposes the restriction on time, we need to add a time stamp to every clause, assuming that all the time stamps of input clauses are 0. Now a clause is represented by  $\langle C, t \rangle$ , where  $t$  is the time stamp of  $C$ , the resolution rule needs to be modified as follows:

$$\frac{\langle (A \mid \alpha), i \rangle \quad \langle (\neg A \mid \beta), j \rangle}{\langle (\alpha \mid \beta), \max(i, j) + 1 \rangle} \text{ if } i \text{ or } j \text{ is maximal.}$$

Now consider the inference graph  $G = (V, E)$  for linear resolution:  $(\mathbf{x}, \mathbf{y}) \in E$  if  $\langle (A \mid \alpha), i \rangle, \langle (\neg A \mid \beta), j \rangle \in \mathbf{x}$ ,  $i$  is maximal in  $\mathbf{x}$ , and  $\mathbf{y} = \mathbf{x} \cup \langle (\alpha \mid \beta), i + 1 \rangle$ . If a clause with the same time stamp appears in two different nodes of  $V$ , it means that we have different sequences (of the same length) to derive this clause by linear resolution. Since the sequences are different, one cannot be replaced by the other. For example, in  $\mathbf{x}$ , we may obtain  $\langle C_3, i + 1 \rangle$  from  $\langle C_1, i \rangle$  and  $\langle C_2, j \rangle$ , but cannot do the same in  $\mathbf{x}'$ , even if  $\langle C_1, i \rangle \in \mathbf{x}'$ , because  $\langle C_2, j \rangle$  may not be in  $\mathbf{x}'$  for any  $j$ . That means that the inference graph has to be a tree structure. The completeness of linear resolution ensures that there exists a path from the root (the input clauses) to a solution node (containing the empty clause) in  $G = (V, E)$ . If we do not use the depth-first search in  $G$ , multiple copies of the same clause may have to be kept in the search. This leaves us the depth-first search as our best option and

we have to inherit all advantages and disadvantages of the depth-first search, such as generating the same resolvent many times, or wasting many futile steps before finding a solution node.

It must be noted that some strategies improve certain aspects of the deduction process at the expense of others. For instance, a strategy may reduce the size of the proof search space at the expense of increasing, say, the length of the shortest proofs.

### 3.3.3 Preserving Satisfiability

The resolution rule does not preserve the logical equivalence, that is, if  $D$  is a resolvent of  $C_1$  and  $C_2$ , then  $C_1 \wedge C_2 \models D$ , or  $C_1 \wedge C_2 \equiv C_1 \wedge C_2 \wedge D$ , but not  $C_1 \wedge C_2 \equiv D$ . For example,  $q$  is a resolvent of  $p$  and  $\neg p \mid q$ ,  $\mathcal{M}(q)$  has two models, and  $\mathcal{M}(p \wedge (\neg p \mid q))$  has only one model. Thus,  $q \not\equiv p \wedge (\neg p \mid q)$ . However, both of them are satisfiable.

**Definition 3.3.6.** *Given two formulas  $A$  and  $B$ ,  $A$  and  $B$  are said to be equally satisfiable, if whenever  $A$  is satisfiable, so is  $B$ , and vice versa. We denote this relation by  $A \approx B$ .*

Obviously,  $A \approx B$  means  $\mathcal{M}(A) = \emptyset$  iff  $\mathcal{M}(B) = \emptyset$  and is weaker than the logical equivalence, which requires  $\mathcal{M}(A) = \mathcal{M}(B)$ .

**Proposition 3.3.7.** *Let  $S = C \cup \{(\alpha \mid \beta)\}$  and  $S' = C \cup \{(\alpha \mid x), (\neg x \mid \beta)\}$  be two sets of clauses, where  $x$  is a variable not appearing in  $S$ . Then  $S \approx S'$ .*

*Proof.* If  $S$  is satisfiable and  $\sigma$  is a model of  $S$ , then  $\sigma(C) = 1$  and  $\sigma(\alpha \mid \beta) = 1$ . If  $\sigma(\alpha) = \perp$ , we define  $\sigma(x) = 1$ ; if  $\sigma(\beta) = 0$ , we define  $\sigma(x) = 0$ . Thus both  $\sigma(\alpha \mid x) = 1$  and  $\sigma(\neg x \mid \beta) = 1$ . So  $\sigma$  is a model of  $S'$ .

On the other hand, if  $\sigma$  is a model of  $S'$ , then  $\sigma$  is also a model of  $S$  without modification, because  $(\alpha \mid \beta)$  is a resolvent of  $(\alpha \mid x)$  and  $(\neg x \mid \beta)$ .  $\square$

In the above proposition, the condition that  $x$  does not appear in  $C$  is necessary. Without it, for example, if  $C = \{(\neg p)\}$ ,  $\alpha$  is empty,  $\beta = q$  and  $x = p$ , then  $S = \{(\neg p), (q)\}$  and  $S' = \{(\neg p), (p), (\neg p \mid q)\}$ . Then  $S$  is satisfiable but  $S'$  is not.

Using the above proposition, we may break a long clause into a set of shorter clauses by introducing some new variables. For example, the clause  $(l_1 \mid l_2 \mid l_3 \mid l_4 \mid l_5)$  can be transformed into three clauses:  $(l_1 \mid l_2 \mid x)$ ,  $(\neg x \mid l_3 \mid y)$ , and  $(\neg y \mid l_4 \mid l_5)$ . This transformation does not preserve the logical equivalence, but preserves the satisfiability.

Another usage of the above proposition is to use resolution to remove  $x$  from a clause set. If we exhaust all resolutions on  $x$ , keep all the resolvents and remove all the clauses containing  $x$ , then the resulting set of clauses will preserve the satisfiability, as stated in Lemma 3.3.11.

**Proposition 3.3.8.** *Let  $B$  be a subformula of  $A$  and  $x$  is a new variable. Then  $A \approx (A[B \leftarrow x] \wedge (x \leftrightarrow B))$ .*

*Proof.* If  $A$  is satisfiable and  $\sigma$  is a model of  $A$ , then we define  $\sigma(x) = \sigma(B)$ . When applying  $\sigma$  to the formula trees of  $A$  and  $A[B \leftarrow x]$ , all the common nodes have the same truth value, so  $\sigma(A) = \sigma(A[B \leftarrow x])$ . Since  $\sigma(x) = \sigma(B)$ , so  $\sigma(x \leftrightarrow B) = 1$ , thus  $\sigma$  is a model of  $A[B \leftarrow x] \wedge (x \leftrightarrow B)$ .

On the other hand, if  $\sigma$  is  $A[B \leftarrow x] \wedge (x \leftrightarrow B)$ ,  $\sigma$  is also a model of  $A$ .  $\square$

The above proposition plays an important role when we convert a formula into CNF and keep the size of CNF as a linear function of the size of the original formula. When we use the distribution law

$$(A \wedge B) \vee C \equiv (A \vee C) \wedge (B \vee C)$$

to transform a formula into CNF,  $C$  becomes duplicated in the result. Because of this, there exists examples that the resulting CNF has an exponential size of the original size. Using the above proposition and introducing new variables, we may keep the size of CNF as four times of the original size, and every clause contains at most three literals. The CNF obtained this way is not logically equivalent to the original formula, but it preserves the satisfiability of the original formula: The CNF is satisfiable iff the original formula is satisfiable. This way of transforming a formula into CNF is called *Tseitin transformation*.

**Example 3.3.9.** Let  $A$  be  $(p \wedge (q \vee \neg p)) \vee \neg(\neg q \wedge (p \vee q))$ , we introduce a new variable for every proper subformula whose node is a binary operator, and convert each of them into at most three clauses.

$$\begin{array}{ll} x_1 \leftrightarrow (q \vee \neg p) & (\overline{x_1} \mid q \mid \overline{p}), (\overline{q} \mid x_1), (p \mid x_1), \\ x_2 \leftrightarrow (p \wedge x_1) & (\overline{x_2} \mid p), (\overline{x_2} \mid x_1), (\overline{p} \mid \overline{x_1} \mid x_2), \\ x_3 \leftrightarrow (p \vee q) & (\overline{x_3} \mid p \mid q), (\overline{p} \mid x_3), (\overline{q} \mid x_3), \\ x_4 \leftrightarrow (\neg q \wedge x_3) & (\overline{x_4} \mid \overline{q}), (\overline{x_4} \mid q), (q \mid \overline{x_3} \mid x_4), \\ x_2 \vee \neg x_4 & (x_2 \mid \overline{x_4}). \end{array}$$

**Example 3.3.10.** In Example 2.3.9, we have converted  $p_1 \leftrightarrow (p_2 \leftrightarrow p_3)$  into four clauses:

$$(\overline{p_1} \mid \overline{p_2} \mid p_3), \quad (\overline{p_1} \mid p_2 \mid \overline{p_3}), \quad (p_1 \mid p_2 \mid p_3), \quad (p_1 \mid \overline{p_2} \mid \overline{p_3}).$$

If we like to convert  $p_1 \leftrightarrow (p_2 \leftrightarrow (p_3 \leftrightarrow p_4))$  into clauses, we may replace  $p_3$  by  $(p_3 \leftrightarrow p_4)$  in the above clauses and then convert each of them into clauses. That is, from  $(\overline{p_1} \mid \overline{p_2} \mid (p_3 \leftrightarrow p_4))$ ,  $(\overline{p_1} \mid p_2 \mid (\overline{p_3 \leftrightarrow p_4}))$ ,  $(p_1 \mid p_2 \mid (p_3 \leftrightarrow p_4))$ , and  $(p_1 \mid \overline{p_2} \mid (\overline{p_3 \leftrightarrow p_4}))$ , we obtain eight clauses, twice as before:

$$\begin{array}{cccc} (\overline{p_1} \mid \overline{p_2} \mid \overline{p_3} \mid p_4), & (\overline{p_1} \mid p_2 \mid p_3 \mid p_4), & (\overline{p_1} \mid \overline{p_2} \mid p_3 \mid \overline{p_4}), & (\overline{p_1} \mid p_2 \mid \overline{p_3} \mid \overline{p_4}), \\ (p_1 \mid p_2 \mid \overline{p_3} \mid p_4), & (p_1 \mid \overline{p_2} \mid p_3 \mid p_4), & (p_1 \mid p_2 \mid p_3 \mid \overline{p_4}), & (p_1 \mid \overline{p_2} \mid \overline{p_3} \mid \overline{p_4}). \end{array}$$

If we replace  $p_4$  by  $(p_4 \leftrightarrow p_5)$ , we will get a set of 16 clauses of length 5. In general,  $p_1 \leftrightarrow (p_2 \leftrightarrow (p_3 \leftrightarrow (\dots(p_{n-1} \leftrightarrow p_n)\dots)))$  will produce  $2^{n-1}$  clauses of length  $n$ . Instead, if we introduce  $n - 3$  new variables for this formula of  $n$  variables, we obtain only  $4(n - 1)$  clauses of length 3. For example, for  $n = 5$ , we introduce two new variables,  $q_1$  and  $q_2$ :

$$(p_1 \leftrightarrow (p_2 \leftrightarrow q_1)) \wedge (q_1 \leftrightarrow (p_3 \leftrightarrow q_2)) \wedge (q_2 \leftrightarrow (p_4 \leftrightarrow q_5)),$$

which will generate  $3 \times 4 = 12$  clauses of length 3.

**Lemma 3.3.11.** *Let  $C$  be a set of clauses and  $x$  be a variable in  $C$ . Let  $D \subset C$  such that each clause of  $D$  contains neither  $x$  nor  $\overline{x}$ . Let  $X = \{(\alpha \mid \beta) \mid (x \mid \alpha), (\overline{x} \mid \beta) \in C\}$ , then  $C \approx D \cup X$ .*

*Proof.* If  $C$  is satisfiable and  $\sigma$  is a model of  $C$ , then  $\sigma(C) = 1$  and  $\sigma(D) = 1$ . If  $\sigma(X) = 0$ , then there exists a clause  $\alpha \mid \beta \in X$  such that  $\sigma(\alpha) = \sigma(\beta) = 0$ . It means either  $\sigma(x \mid \alpha) = 0$  or  $\sigma(\overline{x} \mid \beta) = 0$ , that is a contradiction to the fact that  $\sigma(C) = 1$ , because both  $(x \mid \alpha)$  and  $(\overline{x} \mid \beta)$  are in  $C$ . So  $\sigma(X) = 1$ . It means  $\sigma$  is a model of  $D \cup X$ .

On the other hand, if  $D \cup X$  is satisfiable and  $\sigma$  is its model, we may define  $\sigma(x) = 1$  if  $\sigma(\alpha) = 0$  for some  $x \mid \alpha \in C$ ; in this case, we must have  $\sigma(\beta) = 1$  for all  $(\overline{x} \mid \beta) \in C$ ; otherwise, if  $\sigma(\beta) = 0$  for some  $(\overline{x} \mid \beta) \in C$ , then  $\sigma(\alpha \mid \beta) = 0$ , that is a contradiction to the fact that  $\sigma(X) = 1$ . Similarly, if  $\sigma(\alpha) = 1$  for all  $x \mid \alpha \in C$ , we define  $\sigma(x) = 0$ . Then  $\sigma$  becomes a model of  $C$ .  $\square$

This lemma plays an important role in the completeness proof of resolution: The set  $X$  can be obtained from  $C$  by resolution. From  $C$  to  $D \cup X$ , we have eliminated one variable from  $C$ . Replace  $C$  by  $D \cup X$  and we repeat the same process and remove another variable, until no more variable in  $C$ : either  $C = \top$  or  $C = \perp$ . From the final result, we can tell if the original set of clauses is satisfiable or not.

### 3.3.4 Completeness of Resolution

**Theorem 3.3.12.** *Ordered resolution is complete with any total ordering on the variables.*

*Proof.* Suppose there are  $n$  variables in the input clauses  $S_n$  with the order  $x_n > x_{n-1} > \dots > x_2 > x_1$ . For  $i$  from  $n$  to 1, we apply Lemma 3.3.6 with  $C = S_i$  and  $x = x_i$ , and obtain  $S_{i-1} = D \cup X$ . Every resolvent in  $X$  is obtained by ordered resolution on  $x_i$ , which is the maximal variable in  $S_i$ . Finally, we obtain  $S_0$  which contains no variables. Since  $S_n \approx S_{n-1} \approx \dots \approx S_1 \approx S_0$ , if  $S_n$  is unsatisfiable, then  $S_0$  must contain  $\perp$ ; if  $S_n$  is satisfiable, we arrive at  $S_0 = \emptyset$ , i.e.,  $S_0 \equiv \top$ . In this case, we may use the idea in the proof of Lemma 3.3.6 to assign a truth value for each variable, from  $x_1$  to  $x_n$ , to obtain a model of  $S_n$ .  $\square$

**Example 3.3.13.** Let  $S = \{c_1. (\bar{a}|\bar{b}|\bar{c}), c_2. (a|\bar{e}), c_3. (b|\bar{c}), c_4. (\bar{b}|d), c_5. (c|\bar{d}|\bar{e}), c_6. (\bar{d}|e)\}$  and  $a > b > c > d > e$ . By ordered resolution, we get  $c_7. (\bar{b}|\bar{c}|\bar{e})$  from  $c_1$  and  $c_2$ ;  $c_8. (\bar{c}|d)$  from  $c_3$  and  $c_4$ ;  $c_9. (\bar{c}|\bar{e})$  from  $c_3$  and  $c_7$ ; and  $c_{10}. (\bar{d}|\bar{e})$  from  $c_5$  and  $c_9$ . No other clauses can be generated by ordered resolution, that is,  $c_1$ – $c_{10}$  are saturated by ordered resolution. Since the empty clause is not generated, we may construct a model from  $c_1$ – $c_{10}$  by assigning truth values to the variables from the least to the greatest. Let  $S_x$  be the set of clauses from  $c_1$ – $c_{10}$  such that  $x$  is the maximal variable.

- $S_e = \emptyset$  and we can assign either 0 or 1 to  $e$ , say  $\sigma_e = \{e\}$ .
- $S_d = \{c_6, c_{10}\}$ . From  $c_{10}$ ,  $\bar{d}$  has to be 1, so  $\sigma_d = \sigma_e \cup \{d = 0\} = \{e, \bar{d}\}$ .
- $S_c = \{c_5, c_8, c_9\}$ . From  $c_8$  or  $c_9$ ,  $\bar{c}$  has to be 1, so  $\sigma_c = \sigma_d \cup \{c = 0\} = \{e, \bar{d}, \bar{c}\}$ .
- $S_b = \{c_3, c_7\}$ . Both  $c_3$  and  $c_7$  are satisfied by  $\sigma_c$ , we may assign 0 or 1 to  $b$ , say  $b = 1$ , so  $\sigma_b = \sigma_c \cup \{b\} = \{e, \bar{d}, \bar{c}, b\}$ .
- $S_a = \{c_1, c_2\}$ . From  $c_2$ ,  $a$  has to be 1, so  $\sigma_a = \sigma_b \cup \{a\} = \{e, \bar{d}, \bar{c}, b, a\}$ .

The above example illustrates that if a set of clauses is saturated by ordered resolution and the empty clause is not there, we may construct a model from these clauses. In other words, we may use ordered resolution to construct a decision procedure for deciding if a formula  $A$  is satisfiable or not.

**Algorithm 3.3.14.** The algorithm *orderedResolution* will take a set  $C$  of clauses and an order  $V$  of variables,  $V = (p_1 < p_2 < \dots < p_n)$ . It returns true iff  $C$  is satisfiable. It uses the procedure *sort(A)*, which places the maximal literal of clause  $A$  as the first one in the list.

```

proc orderedResolution( $C, V$ )
   $C := \{\text{sort}(A) \mid A \in C\}$ ;
  for  $i := n$  downto 1 do
     $X := \emptyset$ ;
    for  $(p_i \mid \alpha) \in C$  do //  $p_i$  is maximal in  $(p_i \mid \alpha)$ 
      for  $(\bar{p}_i \mid \beta) \in C$  do //  $p_i$  is maximal in  $(\bar{p}_i \mid \beta)$ 
         $A := (\alpha \mid \beta)$ ;
        if  $A = ()$  return false;
         $X := X \cup \{\text{sort}(A)\}$ ;
     $C := C \cup X$ ;
  return true;

```

Before returning true at the last line, this algorithm can call *createModel* to construct a model of  $C$  when  $C$  is *saturated* by ordered resolution. That is,  $C$  contains the results of every possible ordered resolution among the clauses of  $C$ .

```

proc createModel( $C$ )
   $\sigma := \emptyset$ ; // the empty interpretation
  for  $i := 1$  to  $n$  do
     $v := 0$ ;
    for  $(p_i \mid \alpha) \in C$  do //  $p_i$  is maximal
      if  $\text{eval}(\alpha, \sigma) = \perp$  do
         $v := 1$ ;
     $\sigma := \sigma \cup \{p_i \mapsto v\}$ ;
  return  $\sigma$ ;

```

The resolution algorithm is regarded as a deduction-style proof procedure when the empty clause is its goal. If finding a model is the goal of the resolution algorithm, it is a *saturation-style* proof procedure, because the set of clauses is *saturated* by ordered resolution, that is, all the possible results of the inference rules are generated, before we can build a model.

**Corollary 3.3.15.** *Resolution is a sound and complete decision procedure for propositional logic.*

*Proof.* We have seen that resolution is sound. If  $C$  is unsatisfiable, we may obtain an ordered resolution proof from  $C$ . This proof is also a resolution proof. Thus  $C \models \perp$  implies  $C \vdash_{\mathcal{R}} \perp$ . If  $C$  is satisfiable, we may construct a model for  $C$  as we do for ordered resolution.  $\square$

**Theorem 3.3.16.** *Positive resolution is complete for refutation.*

*Proof.* Suppose  $C$  is unsatisfiable and there is a resolution proof  $C \vdash_{\mathcal{R}} \perp$ . We prove by induction on the length of the proof on how to convert this proof into a positive resolution proof.

Let  $P$  be a resolution proof of  $C \vdash_{\mathcal{R}} \perp$  of minimal length. If  $|P| = 1$ , then the only resolution uses a positive unit clause to obtain  $\perp$  and it is a positive resolution. If  $|P| > 1$ , suppose the last resolution in  $P$  resolves on variable  $q$ . Then  $C \vdash_{\mathcal{R}} q$  and  $C \vdash_{\mathcal{R}} \bar{q}$ . Let  $P_1$  and  $P_2$  be the proofs of  $C \vdash_{\mathcal{R}} q$  and  $C \vdash_{\mathcal{R}} \bar{q}$  presented in  $P$ , respectively.

If we remove all the occurrences of  $q$  from  $P_1$ , we obtain a resolution proof  $P'_1$  and  $core(P'_1)$  contains clauses from  $core(P)$  with  $q$  being removed. Since  $|P'_1| < |P|$ , by induction hypothesis, we have a positive resolution proof of  $core(P'_1) \vdash_R \perp$ . If we add  $q$  back to  $core(P'_1)$ , we obtain a positive resolution proof  $P''_1$  of  $core(P) \vdash_{\mathcal{R}} q$ .

We then remove all the occurrences of  $\bar{q}$  from  $P_2$ , and obtain a resolution proof  $P'_2$  and  $core(P'_2)$  contains clauses from  $core(P)$  with  $\bar{q}$  being removed. By induction hypothesis, there exists a positive resolution proof  $P''_2$  of  $core(P'_2) \vdash_{\mathcal{R}} \perp$ . For each clause  $\alpha \in core(P'_2)$ , if  $\alpha$  is obtained from  $\alpha \mid \bar{q} \in core(P)$ , we add  $\alpha$  as the resolvent of  $\alpha \mid \neq q$  and  $q$  at the beginning of  $P''_2$ . Finally, we append  $P''_1$ , which generates  $q$  by positive resolution, and  $P''_2$  into one proof, the result is a positive resolution proof of  $C \vdash_{\mathcal{R}} \perp$ .  $\square$

The following theorem can be proved using a similar approach.

**Theorem 3.3.17.** *Negative resolution is complete for refutation.*

Completeness of resolution strategies are important to ensure that if the empty clause is not found, the input clauses are unsatisfiable. In this case, one strategy may generate the empty clause quicker than another strategy. When the input clauses are satisfiable, we have to compute all possible resolutions to ensure that no new resolvents before claiming the input clauses are satisfiable. In this case, we say the set of clauses, old and new, is saturated by resolution. Restricted resolutions like ordered resolution, positive resolution, and set-of-support resolution, work much better than unrestricted resolution, when the saturation is needed to show the satisfiability of a set of clauses.

### 3.3.5 A Resolution-based Decision Procedure

If we like to employ heuristics in a decision procedure based on resolution, the following procedure fits the purpose.

**Algorithm 3.3.18.** The algorithm *resolution* will take a set  $C$  of clauses and returns true iff  $C$  is satisfiable. It uses the procedure *resolvable*( $A, B$ ), which decides if clauses  $A$  and  $B$  are allowed to do restricted resolution, and if yes, *resolve*( $A, B$ ) will their resolvent. It uses *pickClause*( $C$ ) to pick out a clause according to a heuristic.

```

proc resolution( $C$ )
   $G := C$ ; //  $G$ : given clauses
   $K = \emptyset$ ; //  $K$ : kept clauses
  while  $G \neq \emptyset$  do
     $A := \text{pickClause}(G)$ ;
     $G := G - \{A\}$ ;
    for  $B \in K$  if resolvable( $A, B$ ) do
       $res := \text{resolve}(A, B)$ ;
      if  $res = ()$  return false;
      if  $res \notin (G \cup K)$ 
         $G := G \cup \{res\}$ ;
       $K := K \cup \{A\}$ ;
  return true;

```

In the procedure, we move each clause from  $G$  to  $K$ , and compute the resolution between this clause and all the clauses in  $K$ . When  $G$  is empty, resolution between any two clauses are done.

With the exception of linear resolution, the above procedure can be used to implement all the resolution strategies introduced in this section. The restriction will be implemented in *resolvable*. For set-of-resolution, a better implementation is that  $G$  is initialized with the set of support and  $K$  is initialized with the rest of the input clauses.

The termination of this procedure is guaranteed because only a finite number of clauses exist. When the procedure returns false, the answer is correct because resolution is sound. When it returns true, the correctness of the answer depends on the completeness of the resolution strategy implemented in *resolvable*.

In *pickClause*, we implement various heuristics for selecting a clause in  $G$  to do resolutions with clauses in  $K$ . Popular heuristics include preferring shorter clauses or older clauses. Preferring shorter clauses because shorter clauses are stronger as constraints on the acceptance of interpretations as models. An empty clause removes all interpretations as models; a unit clause removes half of the interpretations as models; a binary clause removes a quarter of the interpretations as models; and so

on. Preferring older clauses alone would give us the breadth-first strategy. If we wish to find a shorter proof, we may mix this preference with other heuristics.

Despite the result that complete resolution strategies can serve as a decision procedure for positional logic, and some strategies can even build a model when the input clauses are satisfiable, resolution provers are in general not efficient for propositional logic, because the number of possible clauses is exponential in terms of the number of variables. The computer may quickly run of memory before finding a solution for real application problems.

### 3.3.6 Clause Deletion Strategies

There are several strategies which can safely remove redundant clauses during the resolution. They are tautology deletion, subsumption deletion, and pure literal deletion.

A clause is a tautology if it contains a complementary pair of literals of the same variable. For example,  $(p \mid \bar{p})$ . A tautology clause is true in every interpretation and is not needed in any resolution proof.

**Proposition 3.3.19.** *Given any set  $S$  of clauses, let  $S$  be partitioned into  $S' \cup T$ , where  $T$  is the set of tautology clauses in  $S$ . Then  $S \equiv S'$ .*

*Proof.* For any interpretation  $\sigma$ ,  $\sigma(S) = \sigma(S') \wedge \sigma(T) = \sigma(S')$  as  $\sigma(T) = 1$ .  $\square$

**Definition 3.3.20.** *A clause  $A$  subsumes another clause  $B$  if  $B = A \mid \alpha$  for some  $\alpha$ . That is, every literal of  $B$  appears in  $A$ .*

For example,  $(p \mid q)$  subsumes  $(p \mid q \mid r)$ . Intuitively, resolution tries to delete literals by resolution one by one to obtain an empty clause. If we use  $(p \mid q)$  instead of  $(p \mid q \mid r)$ , we avoid the job of removing  $r$ .

**Proposition 3.3.21.** *If  $A$  subsumes  $B$ , then  $S \cup \{A, B\} \equiv S \cup \{A\}$ .*

*Proof.* Since  $A$  subsumes  $B$ ,  $\mathcal{M}(A) \subseteq \mathcal{M}(A) \cup \mathcal{M}(\alpha) = \mathcal{M}(B)$ .  $\mathcal{M}(S \cup \{A, B\}) = \mathcal{M}(S) \cap (\mathcal{M}(A) \cap \mathcal{M}(B)) = \mathcal{M}(S) \cap \mathcal{M}(A) = \mathcal{M}(S \cup \{A\})$ .  $\square$

The above proposition shows that if  $A$  subsumes  $B$ , then  $B$  can be dropped as a constraint in the presence of  $A$ . The following proposition shows that  $B$  is not needed in a resolution proof.

**Proposition 3.3.22.** *If  $A$  and  $B$  are two clauses appearing in the resolution proof of  $S \vdash_{\mathcal{R}} \perp$  such that  $A$  subsumes  $B$ , then  $B$  is no longer needed once  $A$  is present.*

*Proof.* If  $B$  is used in a resolution proof after the presence of  $A$ , we may replace  $B$  by  $A$  to obtain a proof of shorter or equal length, by checking all the resolutions which remove every literal from  $B$ : Suppose the first resolution is between  $B$  and  $\gamma$  and the resolvent is  $B'$ . There are two cases to consider: (1) the resolved variable is not in  $A$ . This resolution is unnecessary, because  $A$  subsumes  $B'$  and we replace  $B'$  by  $A$ . (2) the resolved variable is in  $A$ , then let the resolvent of  $A$  and  $\gamma$  be  $A'$ . It is easy to check that  $A'$  subsumes  $B'$ , and we replace  $B'$  by  $A'$ . We replace the descendants of  $B$  by the corresponding descendants of  $A$ , until the last resolution. The modified resolution proof will have shorter or equal length.  $\square$

**Proposition 3.3.23.** *If there exists a resolution proof using a tautology clause, then there exists a resolution proof without the tautology clause.*

*Proof.* Suppose the tautology clause appearing in the proof is  $(p \mid \bar{p} \mid \alpha)$ . We need another clause  $(\bar{p} \mid \beta)$  (or  $(p \mid \beta)$ ) to resolve off  $p$ . The resolvent of  $(p \mid \bar{p} \mid \alpha)$  and  $(\bar{p} \mid \beta)$  is  $(\bar{p} \mid \alpha \mid \beta)$ , which is subsumed by  $(\bar{p} \mid \beta)$ . By Proposition 3.3.22,  $(\bar{p} \mid \alpha \mid \beta)$  is not needed in the resolution proof. Since every resolvent of a tautology is not needed, so we need any tautology in a resolution proof.  $\square$

**Definition 3.3.24.** *A literal  $l$  is said to be pure in a set  $C$  of clauses if its complement does not appear in  $C$ .*

**Proposition 3.3.25.** *If  $l$  is pure in  $C$ , let  $D$  be the set of clauses obtained from  $C$  by removing all clauses containing  $l$  from  $C$ , then  $C \approx D$ .*

*Proof.* If  $\sigma$  is a model of  $C$ , then  $\sigma$  is also a model of  $D$ . If  $\sigma$  is a model of  $D$ , define  $\sigma(l) = 1$ , then  $\sigma$  is a model of  $C$ .  $\square$

For example, if  $C = \{1. (p \mid q), 2. (\bar{p} \mid r), 3. (p \mid \bar{q} \mid r), 4. (\bar{p} \mid q)\}$ , then  $r$  is pure. Remove 2. and 3. from  $C$ ,  $q$  becomes pure. Remove 1. and 4., we obtain an empty set of clauses, which is equivalent to 1. By the proposition, it means  $C$  is satisfiable. A model can be obtained by setting  $q = r = 1$  in the model.

The above propositions allow to use the following clause deletion strategies without worrying about missing a proof:

- **tautology deletion:** Tautology clauses are discarded.
- **subsumption deletion:** Subsumed clauses are discarded.
- **pure deletion:** Clauses containing pure literals are discarded.

These deletion strategies can be integrated into the algorithm *resolution* as follows.

**Algorithm 3.3.26.** The algorithm *resolution* will take a set  $S$  of clauses and returns true iff  $S$  is satisfiable. In *preprocessing*( $S$ ), we may implement pure-literal check or subsumption to simplify the input clauses as a preprocessing step. *subsumedBy*( $A, S$ ) checks if clause  $A$  is subsumed by a clause in  $S$ .

```

proc resolution( $S$ )
1    $G := preprocessing(S)$ ; //  $G$ : given clauses
2    $K := \emptyset$ ; //  $K$ : kept clauses
3   while  $G \neq \emptyset$  do
4        $A := pickClause(G)$ ;
5        $G := G - \{A\}$ ;
6        $N := \{ \}$ ; // new clauses from  $A$  and  $K$ 
7       for  $B \in K$  if resolvable( $A, B$ ) do
8            $res := resolve(A, B)$ ;
9           if  $res = ()$  return false;
10          if subsumedBy( $res, G \cup K$ ) continue;
11           $N := N \cup \{res\}$ ;
12           $K = K \cup \{A\}$ ;
13          for  $A \in G$  if subsumedBy( $A, N$ ) do
14               $G := G - \{A\}$ ;
15          for  $B \in K$  if subsumedBy( $B, N$ ) do
16               $K := K - \{B\}$ ;
17           $G := G \cup N$ ;
18  return true; //  $K$  is saturated by resolution

```

A life cycle of a clause  $A$  goes as follows in the algorithm *resolution*.

1. If  $A$  is an input clause, it passes the preprocessing check, such as pure-literal check and subsumption check, and goes into  $G$  (line 1).
2. If  $A$  is picked out of  $G$  (lines 4,5), it will pair with every clause  $B$  in  $K$  to try resolution on  $A$  and  $B$ . Once this is done,  $A$  will go to  $K$  as a kept clause (line 12).
3.  $A$  may be subsumed by a new clause when it is in  $G$  (line 13) or in  $K$  (line 15), and will be thrown away.
4. If  $A$  is a new clause (line 8), it has to pass the subsumption check (line 10) and parks in  $N$  (line 11).

5. All new clauses in  $N$  will be used to check if they can subsume some clauses in  $G$  (line 13) and (line 15), then merge into  $G$  (line 17), and start a life cycle the same way as an input clause.
6. If the algorithm stops without finding the empty clause (line 9), the set  $G$  will be empty (line 3), and all the resolutions among the clauses in  $K$  are done. That is,  $K$  is saturated by resolution.

In this algorithm, unit clauses are given high priority to do unit resolution and subsumptions. Dealing with unit clauses can be separated from the main loop of the algorithm as we will see in the next section.

### 3.4 Boolean Constraint Propagation (BCP)

A unit clause  $(A)$  subsumes every clause like  $(A \mid \alpha)$ , where  $A$  appears as a literal. Unit resolution between  $(A)$  and  $(\bar{A} \mid \alpha)$ , where  $\bar{A} = p$  if  $A$  is  $\bar{p}$ , generates the resolvent  $(\alpha)$ , which will subsume  $(\bar{A} \mid \alpha)$ . The rule for replacing  $(\bar{A} \mid \alpha)$  by  $(\alpha)$  at the presence of  $(A)$  is called *unit deletion*. That is, when subsumption deletion is used in resolution, a unit clause  $(A)$  allows us to remove all the occurrences of the variable  $p$  in  $(A)$ , with the unit clause itself as the only exception. New unit clauses may be generated by unit resolution and we can continue to simplify the clauses by the new unit clause and so on. This process is traditionally called *Boolean constraint propagation* (BCP), or *unit propagation*.

#### 3.4.1 BCP: a Simplification Procedure

Clause deletion strategies allow us to remove unnecessary clauses during resolution. BCP also allows us to simplify clauses by removing unnecessary clauses and shortening clauses. In the following, we will describe BCP as a procedure which takes a set  $C$  of clauses as input and returns a pair  $(U, S)$ , where  $U$  is a set of unit clauses and  $S$  is a set non-unit clauses such that  $U$  and  $S$  share no variables and  $C \equiv U \cup S$ . Typically,  $S$  contains a subset of  $C$  or clauses from  $C$  by removing some literals.

**Algorithm 3.4.1.** The algorithm *BCP* will take a set  $C$  of clauses and apply unit resolution and subsumption deletion repeatedly, until no more new clauses can be generated by unit resolution. It will return  $\perp$  if an empty clause is found; otherwise, it returns a simplified set of clauses equivalent to  $C$ .

**proc** *BCP*( $C$ )

```

 $S := C$ ; //  $S$ : simplified clauses
 $U := \emptyset$ ; //  $U$ : unit clauses
while  $S$  has a unit clause  $(A)$  do
     $S := S - \{(A)\}$ ;
     $U := U \cup \{A\}$ ;
    for  $(A \mid \alpha) \in S$  do
         $S := S - \{(A \mid \alpha)\}$ ;
    for  $(\bar{A} \mid \alpha) \in S$  do
        if  $\alpha = ()$  return  $\perp$ ;
         $S := S - \{(\bar{A} \mid \alpha)\} \cup \{(\alpha)\}$ ;
return  $(U, S)$ ;

```

**Example 3.4.2.** Let  $C = \{1 : (x_2 \mid x_5), 2 : (\bar{x}_1 \mid \bar{x}_4), 3 : (\bar{x}_2 \mid x_4), 4 : (x_1 \mid x_2 \mid \bar{x}_3), 5 : (\bar{x}_5)\}$ . Then  $BCP(C)$  will return

$$(\{\bar{x}_5, x_2, x_4, \bar{x}_1\}, \emptyset)$$

because we first add  $\bar{x}_5$  into  $U$ ; clause 1. becomes  $(x_2)$ ,  $x_2$  is then added into  $U$ ; clause 4. is deleted because  $x_2 \in U$ ; clause 3. becomes  $(x_4)$  and  $x_4$  is added into  $U$ ; clause 2. becomes  $(\bar{x}_1)$  and  $\bar{x}_1$  is added into  $U$ . The set  $U$  allows us to construct two models of  $C$ :

$$\{x_1 \mapsto 0, x_2 \mapsto 1, x_3 \mapsto v, x_4 \mapsto 1, x_5 \mapsto 0\}$$

where  $v \in \{0, 1\}$  as  $x_3$  can take either value.

**Proposition 3.4.3.** (a) *If  $BCP(C)$  returns  $\perp$ , then  $C$  is unsatisfiable.* (b) *If  $BCP(C)$  returns  $(U, S)$ , then  $C \equiv U \cup S$ , where  $U$  is a set of unit clauses and  $S$  is a set of non-unit clauses, and  $U$  and  $S$  share no variables.*

*Proof.* In  $BCP(C)$ ,  $S$  is initialized with  $C$ .  $S$  is updated through three ways: (1) unit clause  $A$  is removed from  $S$  into  $U$ ; (2) clause  $(A \mid \alpha)$  is removed from  $S$ ; and (3) clause  $(\bar{A} \mid \alpha)$  is replaced by  $(\alpha)$ .  $(\alpha)$  is the resolvent of  $(A)$  and  $(\bar{A} \mid \alpha)$ , so  $C \models \alpha$ .

(a): If  $\alpha = ()$  in (3),  $\perp$  will be returned and  $C \equiv \perp$  because  $C \models \perp$ .

(b): Initially,  $C \equiv U \cup S$ . Since  $A$  subsumes  $A \mid \alpha$  and  $\alpha$  subsumes  $(\bar{A} \mid \alpha)$ , by Proposition 3.3.21, it justifies that both  $A \mid \alpha$  and  $(\bar{A} \mid \alpha)$  can be safely removed, as it maintains  $C \equiv U \cup S$  for each update of  $U$  and  $S$ . When  $A$  is moved from  $S$  to  $U$ , either  $A$  or  $\bar{A}$  is removed from  $S$ , so no variables of  $A$  appear in  $S$ . The procedure stops only when  $S$  does not have any unit clauses.  $\square$

The procedure  $BCP$  can serve as a powerful simplification procedure when unit clauses are present in the input.  $BCP$  plays an important role in deciding if a formula is satisfiable, which is the topic of the next chapter. We will show later that  $BCP(C)$  can be implemented in time  $O(n)$ , where  $n$  is the number of literals in  $C$ .

### 3.4.2 $BCP$ : a Decision Procedure for Horn Clauses

A *Horn clause* is a clause with at most one positive literal. Horn clauses are named for the logician Alfred Horn, who first pointed out their significance in 1951. A Horn clause is called *definite clause* if it has exactly one positive literal. Thus Horn clauses can be divided into definite clauses and negative clauses. Definite clauses can be further divided into *fact* (a positive unit clause) and *rule* (non-unit definite clauses). Horn clauses have important applications in logic programming, formal specification, and model theory, as they have very efficient decision procedures, using unit resolution or input resolution. That is, unit and input resolutions are incomplete in general, but they are complete for Horn clauses.

**Proposition 3.4.4.**  *$BCP$  is a decision procedure for Horn clauses.*

*Proof.* Suppose  $H$  is a set of Horn clauses and  $\perp \notin H$ . It suffices to show that procedure  $BCP(H)$  returns  $\perp$  iff  $H \equiv \perp$ . If  $BCP(H)$  returns  $\perp$ , by Proposition 3.4.3,  $H \equiv \perp$ . If  $BCP(H)$  returns  $(U, S)$ , by Proposition 3.4.3,  $H \equiv U \cup S$ ,  $U$  is a set of unit clauses,  $S$  is a set of non-unit clauses, and  $U$  and  $S$  share no variables. In this case, we create an assignment  $\sigma$  in which every literal in  $U$  is true, and every variable in  $S$  is false. Note that each variable in  $U$  appears only once;  $U$  and  $S$  do not share any variable. Thus this assignment is consistent and is a model of  $U \cup S$ , because every unit clause in  $U$  is true under  $\sigma$ ; for each every clause in  $S$ , since it is a non-unit Horn clause and must have a negative literal, which is true under  $\sigma$ .  $\square$

**Theorem 3.4.5.** *Unit resolution is sound and complete for Horn clauses.*

*Proof.* Unit resolution is sound because resolution is sound. The only used inference rule in  $BCP$  is unit resolution. Every clause generated by  $BCP$  can be generated by unit resolution. Subsumption only reduces some unnecessary clauses. If the input clauses are unsatisfiable,  $BCP$  will generate the empty clause, which can also be generated by unit resolution. If the input clauses are satisfiable, unit resolution will terminate with a set of clauses saturated by unit resolution, without generating the empty clause, because unit resolution is sound.  $\square$

### 3.4.3 Unit Resolution versus Input Resolution

In the following, we show that input resolution and unit resolution are equivalent, in the sense that if there exists a unit resolution proof, then there exists an input resolution proof and vice versa.

**Theorem 3.4.6.** *For any set  $C$  of clauses, if there exists an input resolution proof from  $C$ , then there exists a unit resolution proof from  $C$ .*

*Proof.* We will prove this theorem by induction on the length of proofs. Let  $P$  be an input resolution proof from  $C$  of minimal length. The last resolvent is  $\perp$  and both its parents are unit clauses and one of them must be an input clause, say  $A$ . For simplicity, assume all the assumed clauses appear in the beginning of  $P$  and  $A$  is the first in  $P$ . If  $|P| = 1$ , then  $P$  must be a unit resolution proof and we are done. If  $|P| > 1$ , then remove  $A$  and the last clause from  $P$ , and remove all occurrences of  $\bar{A}$  from  $P$ , where  $\bar{A} = \bar{p}$  if  $A = p$  and  $\bar{A} = p$  if  $A = \bar{p}$ . The result is an input resolution proof  $P'$  from  $C'$ , where  $C'$  is obtained from  $C$  by removing all occurrences of  $\bar{A}$  in  $C$ . That is,  $(\bar{A} \mid \alpha) \in C$  iff  $(\alpha) \in C'$ . Because  $|P'| < |P|$ , by induction hypothesis, there exists a unit resolution  $Q'$  proof from  $C'$ . If  $(\alpha) \in C'$  comes from  $(\bar{A} \mid \alpha) \in C$ , then change the reason of  $(\alpha) \in Q'$  from “assumed” to “resolution from  $A$  and  $(\bar{A} \mid \alpha)$ ”. Let the modified  $Q'$  be  $Q$ , then  $Q$  is a unit resolution from  $C$ .  $\square$

**Example 3.4.7.** We illustrate the proof of this theorem by Example 3.3.5, where

$$C = \{1. (p \mid q), 2. (\bar{p} \mid r), 3. (p \mid \bar{q} \mid r), 4. (\bar{r})\}.$$

Let  $P$  be (b) of Example 3.3.5. Then  $A = \bar{r}$ ,  $C' = \{1. (p \mid q), 2'. (\bar{p}), 3'. (p \mid \bar{q})\}$  and two input resolutions in  $P'$  are 5. :  $p$  (1, 3') and 6.  $\perp$  (2', 5). A unit resolution  $Q'$  from  $C'$  will contain three unit resolutions:

$$5. (q) (1, 2') \quad 6. (\bar{q}) (2', 3') \quad 7. () (5, 6).$$

Adding two unit resolutions before  $Q'$ , we obtain  $Q$ :

$$2' : (\bar{p}) (2, 4), \quad 3' : (p \mid \bar{q}) (3, 4) \quad 5. (q) (1, 2') \quad 6. (\bar{q}) (2', 3') \quad 7. () (5, 6),$$

which is a unit resolution proof from  $C$ .

**Theorem 3.4.8.** *For any set  $C$  of clauses, if there exists a unit resolution proof from  $C$ , then there exists an input resolution proof from  $C$ .*

The proof of this theorem is left as exercise.

The above two theorems establish the equivalence of input resolution and unit resolution: A unit resolution proof exists iff an input resolution proof exists. Since unit resolution is complete for Horn clauses, we have the following result.

**Corollary 3.4.9.** *Input resolution is sound and complete for Horn clauses.*

**Theorem 3.4.10.** *For any set  $H$  of Horn clauses, let  $H = S \cup T$ , where  $S$  is a set of negative clauses and  $T = H - S$ . Then input resolution is sound and complete for  $H$ , when  $S$  is used as the set of support.*

*Proof.* Resolution is not possible among the clauses in  $S$  as resolution needs a pair of complementary literals from the two parent clauses. Set-of-support resolution does not allow resolution among the clauses in  $T$ . Thus resolution is possible only when one clause is from  $S$  and one from  $T$ , and the resolvent is a negative clause, because the only positive literal from  $T$  is resolved off. For each resolution between  $S$  and  $T$ , we add the resolvent into  $S$ , until either the empty clause  $\perp$  is found or no new resolvents can be generated. If  $\perp$  is generated, then  $H \equiv \perp$  because resolution is sound. If no new resolvents can be generated, we create an assignment  $\sigma$  as follows: call  $BCP$  on  $T$ , since  $T$  is satisfiable, let  $(U, S) = BCP(T)$ . For any unit clause  $p$  in  $U$ , let  $\sigma(p) = 1$ ; for any other unassigned variable  $q$  in  $S \cup S$ , let  $\sigma(q) = \perp$ . We claim that  $\sigma$  is a model of  $H$ . At first,  $\sigma(U \cup S) = 1$ , because  $U$  and  $S$  do not share any variables, and the assignment will make every clause in  $U \cup S$  to be true. Let  $S$  be the set of negative clauses from  $H$  plus all the resolvents between  $S$  and  $T$ .

If  $\sigma(S) = \perp$ , we look for a clause  $C \in S$ , such that  $\sigma(C) = \perp$  and every variable of  $C$  appears as a unit clause in  $T$ . Since  $C$  is a negative clause, let  $C = (\overline{p_1} \mid \overline{p_2} \mid \cdots \mid \overline{p_k})$ . Since  $\sigma(C) = \perp$ ,  $p_i \in U$  and  $\sigma(p_i) = 1$  for  $1 \leq j \leq k$ . If  $p_k \notin T$ ,  $p_k$  is obtained by unit resolution in  $BCP(T)$ , say from  $(\overline{q_1} \mid \overline{q_2} \mid p_k) \in T$  and  $q_1, q_2 \in U$ , then this unit resolution can be avoided by an input resolution between  $(\overline{q_1} \mid \overline{q_2} \mid p_k) \in T$  and  $C \in S$  to generate

$$C' = (\overline{p_1} \mid \overline{p_2} \mid \cdots \mid \overline{p_{k-1}} \mid \overline{q_1} \mid \overline{q_2}) \in S.$$

Obviously,  $\sigma(C') = 0$ . Let  $C$  be  $C'$  and we check again if every variable of  $C$  appears as a unit clause in  $T$ . Finally, we will find a clause  $C \in S$ , such that every variable of  $C$  appears as a unit clause in  $T$ . In this case, input resolution between  $C$  and  $\{(p_1), (p_2), \dots, (p_k)\}$  will generate the empty clause  $(\ )$ , a contradiction coming from the assumption that  $\sigma(S) = 0$ . So  $\sigma(S) = 1$ .  $\square$

Note that input resolution with negative clauses as the set of support on Horn clauses is also negative resolution. If the set of support initially contains a single clause, this input resolution is also a linear resolution. If the positive literal is maximal in each clause, then this resolution is also ordered resolution.

**Corollary 3.4.11.** *If  $H$  is a set of Horn clauses which contains a single negative clause, then the negative, input and linear resolution is complete for  $H$ .*

*Proof.* If  $C$  is the only negative clause, let  $S = \{C\}$  and  $T = H - S$ . Using  $S$  as the set of support, we will get an input resolution proof. Since all negative clauses between  $S$  and  $T$  are from  $C$ , the proof can be arranged as a linear resolution proof.  $\square$

#### 3.4.4 Head/Tail Literals for BCP

For a unit clause ( $A$ ) to be true, the truth value of literal  $A$  needs to be 1. We can use a partial interpretation  $\sigma$  to remember which literals appear in unit clauses.

**Definition 3.4.12.** A partial interpretation is a set  $\sigma$  of literals, where no variables appear in  $\sigma$  more than once, with the understanding that if  $A \in \sigma$ , then  $A$  is true and  $\bar{A}$  is false in  $\sigma$ . When  $\sigma$  contains every variable exactly once, it is called a full interpretation

It is easy to see that a full interpretation  $\sigma$  is equivalent to the interpretation  $\sigma'$ : for every variable  $p$ ,  $\sigma'(p) = 1$  if  $p \in \sigma$  and  $\sigma'(p) = 0$  if  $\bar{p} \in \sigma$ . From the above definition, it is convenient to consider a set of unit clauses as a partial interpretation. For example, if  $U = \{p, \neg q, r\}$ , it represents the interpretation  $\sigma = \{p \mapsto 1, q \mapsto 0, r \mapsto 1\}$ . From now on, we do not distinguish a full interpretation from the old definition of interpretation.

We can evaluate the truth value of a clause or formula in a partial interpretation as we do in a full interpretation.

**Definition 3.4.13.** Given a partial interpretation  $\sigma$  and a clause  $c$ ,

- $c$  is satisfied if one or more of its literals is true in  $\sigma$ .
- $c$  is conflicting if all the literals of  $c$  are false in  $\sigma$ .
- $c$  is unit if all but one of its literals are false and  $c$  is satisfied.
- $c$  is unresolved if  $c$  is not one of the above three cases.

When  $c$  becomes unit, we have to assign 1 to the unassigned literal, say  $A$ , in  $c$ , so that  $c$  becomes true. That is,  $c$  is the *reason* for literal  $A$  being true, and we record this as  $reason(A) = c$ . We also say the literal  $A$  is *implied* by  $c$ .

**Example 3.4.14.** Given the clauses  $\{c_1 : (\bar{x}_1), c_2 : (x_3), c_3 : (x_1 \mid x_2 \mid \bar{x}_3)\}$ ,  $reason(\bar{x}_1) = c_1$ ,  $reason(x_3) = c_2$ , and  $reason(x_2) = c_3$ .  $x_2$  is implied to be true by  $c_3$  because  $c_3$  is equivalent to  $(\bar{x}_1 \wedge x_3) \rightarrow x_2$ .

In *BCP*, we start with an empty partial interpretation  $\sigma$ . If we have found a unit clause  $(A)$ , we extend  $\sigma$  by adding  $A$  into  $\sigma$  so that the unit clause becomes satisfied. To save time, we do not remove clauses containing  $A$  from the clause set. We do not remove  $\bar{A}$  from any clause to create a new clause. In stead, we just remember  $\bar{A}$  is now false in  $\sigma$ .

We need to find out which clauses become unit clauses when some of literals become false. A clause  $c$  of  $k$  literals becomes unit when  $k - 1$  literals of  $c$  are false. If we know two literals of  $c$  are not false, we are certain that  $c$  is not unit. Let us designate these two literals as *head* and *tail* literals of  $c$ , which are two distinct literals of  $c$ , we have the following property.

**Proposition 3.4.15.** *A clause  $c$  of length longer than one is neither unit nor conflicting in a partial interpretation  $\sigma$  iff either one of the head/tail literals of  $c$  is true or neither of them is false in  $\sigma$ .*

*Proof.* Suppose  $c$  contains  $k$  literals and  $k > 1$ . If  $c$  is satisfied, then it contains a true literal and we may use it as one of the head/tail literals. If  $c$  is unit or conflicting, there are  $k - 1$  or all literals of  $c$  are false, and no literal in  $c$  is true and there do not exist two literals which are not false for the head/tail literals. If  $c$  is unsolved, then there must exist two unassigned literals for the head/tail literals.  $\square$

Initially, let the first and last literals be the head/tail literals. When a head literal becomes false, we scan right to find in  $c$  the next literal which is not false and make it the new head literal; similarly, when a tail literal becomes false, we scan left to find the next tail literal. If we cannot find new head or tail literal, we must have found a unit clause or a conflicting clause.

To implement the above idea, we use the following data structures for clauses: For each clause  $c$  containing more than one literal:

- $lits(c)$ : the array of literals storing the literals of  $c$ ;
- $head(c)$ : the index of the first non-false literal of  $c$  in  $lits(c)$ ;
- $tail(c)$ : the index of the last non-false literal of  $c$  in  $lits(c)$ ;

Let  $|c|$  denote the number of literals in  $c$ , then the valid indices for  $lits(c)$  are  $\{0, 1, \dots, |c| - 1\}$ ,  $head(c)$  is initialized with 0, and  $tail(c)$  is with  $|c| - 1$ .

The data structures associated with each literal  $A$  are:

- $val(A)$ : the truth value of  $A$  under partial interpretation,  $val(A) \in \{1, 0, \times\}$ ,  $\times$  (“unassigned”) is the initial value;

- $cls(A)$ : list of clauses  $c$  such that  $A$  is in  $c$  and pointed by either  $head(c)$  or  $tail(c)$  in  $lits(c)$ . We use  $insert(c, cls(A))$  and  $remove(c, cls(A))$  to insert and remove clause  $c$  into/from  $cls(A)$ , respectively.

Assuming global variables are used for clauses and literals, the procedure  $BCP$  can be implemented as follows:

```

proc  $BCP(C)$ 
  // initialization
1:   $U := \emptyset$ ; //  $U$ : stack of unit clauses
2:  for  $c \in C$  do
3:    if  $c = ()$  return  $\perp$ ;
4:    if  $(|c| > 1)$ 
5:       $lits(c) := \text{makeArray}(c)$ ; // create an array of literals
6:       $head(c) := 0$ ;  $insert(c, cls(lits(c)[0]))$ ;
7:       $tail(c) := |c| - 1$ ;  $insert(c, cls(lits(c)[|c| - 1]))$ ;
8:    else  $push(c, U)$ ;
9:  for each literal  $A$ ,  $val(A) := \times$ ; //  $\times =$  "unassigned"
  // major work
10:  $res := BCPht(U)$ ;
  // finishing
11: if  $(res \neq \text{"SAT"})$  return  $\perp$ ;
12:  $U := \{A \mid val(A) = 1\}$ ;
13:  $S := \{clean(c) \mid val(lits(c)[head(c)]) = \times, val(lits(c)[tail(c)]) = \times\}$ ;
14: return  $(U, S)$ ;

```

Note that the procedure  $clean(c)$  at line 13 will remove from  $c$  those literals which are false (under  $val$ ) and return  $\top$  if one of the literals in  $c$  is true. The procedure  $BCPht$  will do the major work of  $BCP$ .

**Algorithm 3.4.16.**  $BCPht(U)$  assumes the input clauses  $C = U \cup S$ , where  $U$  is a stack of unit clauses and  $S$  is a set of non-unit clauses stored using the head/tail data structure, and can be accessed through  $cls(A)$ .  $BCPht$  returns a conflict clause if an empty clause is found during unit propagation; returns "SAT" if no empty clause is found.

```

proc  $BCPht(U)$ 
1  while  $U \neq \emptyset$  do
2     $A := pop(U)$ ;

```

```

3      if  $val(A) = 0$  return  $reason(\overline{A})$ ; // an empty clause is found;
4      else if ( $val(A) = \times$ ) //  $\times$  is “unassigned
5           $val(A) := 1; val(\overline{A}) := 0$ ;
6      for  $c \in cls(\overline{A})$  do //  $\overline{A}$  is either head or tail of  $c$ 
7          if ( $\overline{A} = lits(c)[head(c)]$ )
8               $e_1 := head(c); e_2 := tail(c); step := 1$ ; // scan from left to right
9          else
10              $e_1 := tail(c); e_2 := head(c); step := -1$ ; // scan from right to left
11         while true do
12              $x := e_1 + step$ ; //  $x$  takes values from  $e_1 + step$  to  $e_2$ 
13             if  $x = e_2$  break; // exits the inner while loop
14              $B := lits(c)[x]$ ;
15             if ( $val(B) \neq 0$ )
16                  $remove(c, cls(\overline{A})); insert(c, cls(B))$ ;
17                 if ( $step = 1$ )  $head(c) := x$ ; else  $tail(c) := x$ ;
18                 break; // exit the inner while loop
19             if ( $x = e_2$ ) // no new head or tail
20                  $A := lits(c)[e_2]$ ; // check if  $c$  is unit or conflicting
21                 if  $val(A) = 0$  return  $c$ ; //  $c$  is conflicting;
22                 else if ( $val(A) = \times$ ) //  $c$  is unit
23                      $push(A, U); reason(A) := c; val(A) := 1; val(\overline{A}) := 0$ ;
24     return “SAT”; // no empty clauses are found;

```

**Example 3.4.17.** For the clauses in Example 3.4.2, we have the following data structures after the initialization (the head and tail indices are given after the literal list in each clause):

$$\begin{aligned}
 C &= \{c_1 : \langle [x_2, x_5], 0, 1 \rangle, c_2 : \langle [\overline{x}_1, \overline{x}_4], 0, 1 \rangle, c_3 : \langle [\overline{x}_2, x_4], 0, 1 \rangle, c_4 : \langle [x_1, x_2, \overline{x}_3], 0, 2 \rangle\} \\
 U &= \{\overline{x}_5\}; \text{ the rest clauses are stored as follows :} \\
 cls(x_1) &= \{c_4\}, cls(\overline{x}_1) = \{c_2\}, cls(x_2) = \{c_1\}, cls(\overline{x}_2) = \{c_3\}, \\
 cls(x_3) &= \emptyset, cls(\overline{x}_3) = \{c_4\}, cls(x_4) = \{c_3\}, cls(\overline{x}_4) = \{c_2\}, \\
 cls(x_5) &= \{c_1\}, cls(\overline{x}_5) = \emptyset,
 \end{aligned}$$

and  $val(A) = \times$  (unassigned) for all literal  $A$ . Note that each clause appears twice in the collection of  $cls(A)$ .

Inside of  $BCPht(U)$ , when  $\overline{x}_5$  is popped off  $U$ , we assign  $val(\overline{x}_5) = 1$ ,  $val(x_5) = 0$  (line 5) and check  $c_1 \in cls(x_5)$  (line 6).  $c_1$  generates unit clause  $x_2$  and is added into  $U$  (line 23) with the assignments  $val(x_2) = 1$ ,  $val(\overline{x}_2) = 0$ . When  $x_2$  is popped off  $U$ , we check  $c_3 \in cls(\overline{x}_2)$ .  $c_3$  generates unit clause  $x_4$ , we push  $x_4$  into  $U$ , and

assign  $val(x_4) = 1$ ,  $val(\overline{x_4}) = 0$ . When  $x_4$  is popped from  $U$ , we check  $c_2 \in cls(\overline{x_4})$ .  $c_2$  generates unit clause  $\overline{x_1}$ , we push  $\overline{x_1}$  into  $U$  and assign  $val(\overline{x_1}) = 1$ ,  $val(x_1) = 0$ . When  $\overline{x_1}$  is popped off  $U$ , we check  $c_4 \in cls(x_1)$ .  $c_4$  is true because  $val(x_2) = 1$ . Now the head index of  $c_4$  is updated to 1, and we add  $c_4 = \langle [x_1, x_2, \overline{x_3}], 1, 2 \rangle$  into  $cls(x_2)$ . Since no empty clause is found,  $BCPht$  returns “SAT”. Finally,  $BCP$  returns  $(U, S)$ :

$$(U, S) = (\{\overline{x_5}, x_2, x_4, \overline{x_1}\}, \emptyset).$$

Note that the two assignments  $val(a) := 1; val(\overline{a}) := 0$  at line 23 of  $BCPht$  can be removed as the same work will be done at line 5. However, having them here often improves the performance of  $BCPht$ .  $BCPht$  can be improved slightly by checking  $(val(lits(c)[e_2]) = 1)$  after line 10: if it is true, then skip clause  $c$  as  $c$  is satisfied.

In this implementation of  $BCP$ , once all the clauses  $C$  are read in, each occurrence of literals of non-unit clauses will be visited at most once in  $BCPht$ .

**Proposition 3.4.18.**  $BCP(C)$  runs in  $O(n)$  time, where  $n$  is the sum of the lengths of all clauses in  $C$ .

*Proof.* Obviously, the initialization of  $BCP$  takes  $O(n)$ . For  $BCPht$ , we access a non-unit clause  $c$  when its head (or tail) literal becomes false (line 6). We skip the head (tail) literal (line 11) and check the next literal  $b$  (line 13), if it is false, we continue (line 14); otherwise, we remove  $c$  from  $cls(\overline{A})$  into  $cls(B)$  (line 16) and update the head (tail) index (line 17). If we cannot find a new head (tail) literal (line 19), we have found a new unit clause (line 20). None of the literals of  $c$  is visited more than once by  $BCPht$ . In particular, we never access the clauses from  $cls(A)$  when literal  $A$  is true.  $\square$

Combining the above proposition with Theorem 3.4.4, we have the following result:

**Theorem 3.4.19.**  $BCP$  is a linear time decision procedure for Horn clauses.

### 3.5 Exercise Problems

1. Prove formally that the following statements are equivalent: For any two formulas  $A$  and  $B$ , (a)  $A \models B$ ; (b)  $A \rightarrow B$  is valid; (c)  $A \wedge \neg B$  is unsatisfiable.
2. Use the semantic tableaux method to decide if the following formulas are valid or not.

- (a)  $(A \leftrightarrow (A \wedge B)) \rightarrow (A \rightarrow B)$
- (b)  $(B \leftrightarrow (A \vee B)) \rightarrow (A \rightarrow B)$
- (c)  $(A \oplus B) \rightarrow ((A \vee B) \wedge (\neg A \vee \neg B))$
- (d)  $(A \leftrightarrow B) \rightarrow ((A \vee \neg B) \wedge (\neg A \vee B))$

3. Use the semantic tableaux method to decide if the following entailment relations are true or not.

- (a)  $(A \wedge B) \rightarrow C \models (A \rightarrow C) \vee (B \rightarrow C)$   
 (b)  $(A \wedge B) \rightarrow C \models A \rightarrow C$   
 (c)  $A \rightarrow (B \rightarrow C) \models (A \rightarrow B) \rightarrow (A \rightarrow C)$   
 (d)  $(A \vee B) \wedge (\neg A \vee C) \models B \vee C$

4. Prove by the semantic tableau method that the following statement, “either the debt ceiling isn’t raised or expenditures don’t rise” ( $\neg d \vee \neg e$ ), follows logically from the statements in Problem 31, Chapter 2.
5. Prove that if all the axioms are tautologies and all the inference rules are sound, then all the derived formulas from the axioms and the inference rules are tautologies.
6. Specify the introduction and elimination rules for the operators  $\oplus$  and  $\leftrightarrow$  in natural deductive system.
7. Estimate the upper bound of the size of the inference graph for resolution, if the input clauses contain  $n$  variables.
8. Use the Tseitin encoding to convert the following formulas into 3CNF (CNF in which every clause has at most three literals):

- (a)  $(A \leftrightarrow (A \wedge B)) \rightarrow (A \rightarrow B)$   
 (b)  $(A \oplus B) \rightarrow (A \vee B)$   
 (c)  $((A \wedge B) \rightarrow C) \rightarrow (A \rightarrow C)$   
 (d)  $((A \vee B) \wedge (\neg A \vee C)) \rightarrow (B \vee C)$

9. Use the resolution method to decide if the following formulas are valid or not. If the formula is valid, provide a resolution proof; if it is not valid, provide an interpretation which falsifies the formula.

- (a)  $(A \leftrightarrow (A \wedge B)) \rightarrow (A \rightarrow B)$   
 (b)  $(B \leftrightarrow (A \vee B)) \rightarrow (A \rightarrow B)$   
 (c)  $(A \oplus B) \rightarrow ((A \vee B) \wedge (\neg A \vee \neg B))$   
 (d)  $((\overline{A} \downarrow B) \uparrow \overline{C}) \rightarrow (A \uparrow (B \downarrow C))$

10. Use the resolution method to decide if the following entailment relations are true or not. If the relation holds, provide a resolution proof; if it doesn’t,

provide an interpretation such that the premises are true, but the consequence of the entailment is false.

- (a)  $(A \rightarrow C) \wedge (B \rightarrow D) \models (A \vee B) \rightarrow (C \vee D)$
- (b)  $(A \wedge B) \rightarrow C \models (A \rightarrow C)$
- (c)  $A \rightarrow (B \rightarrow C) \models (A \rightarrow B) \rightarrow (A \rightarrow C)$
- (d)  $(A \vee B) \wedge (\neg A \vee C) \models B \vee C$

11. Apply clause deletion strategies to each of the following sets of clauses, and obtain an equally satisfiable simpler set.

$$\begin{aligned} S_1 &= \{(p \mid q \mid r), (\bar{p} \mid q \mid \bar{r} \mid s), (p \mid \bar{q} \mid s), (\bar{q} \mid r), (q \mid \bar{r})\} \\ S_2 &= \{(p \mid q \mid r \mid s), (\bar{p} \mid q \mid \bar{r} \mid \bar{s}), (p \mid q \mid s), (q \mid \bar{r}), (q \mid s \mid \bar{r})\} \\ S_3 &= \{(p \mid q \mid s), (\bar{p} \mid q \mid \bar{r} \mid \bar{s}), (p \mid r \mid s), (\bar{p} \mid \bar{r}), (\bar{p} \mid q \mid \bar{r})\} \end{aligned}$$

12. Given the clause set  $S$ , where  $S = \{1.(t \mid \bar{e} \mid d), 2.(\bar{t} \mid c), 3.(e \mid \bar{d} \mid i), 4.(\bar{g} \mid \bar{d} \mid i), 5.(t \mid c \mid \bar{d} \mid g), 6.(\bar{c}), 7.(\bar{i} \mid \bar{g}), 8.(c \mid d), 9.(c \mid e)\}$ . Find the following resolution proofs for  $S$ : (a) unit resolution; (b) input resolution; (c) positive resolution; and (d) ordered resolution using the order  $t > i > g > e > d > c$ .
13. Apply various resolution strategies to prove that the following clause set  $S$  entails  $\neg s$ : unit resolution, input resolution, negative resolution, positive resolution, and ordered resolution using the order of  $p > q > r > s$ .

$$S = \{1.(p \mid q \mid \bar{r}), 2.(p \mid \bar{q} \mid s), 3.(\bar{p} \mid q \mid \bar{r}), 4.(\bar{p} \mid \bar{q} \mid r), 5.(\bar{p} \mid q \mid r), 6.(\bar{q} \mid \bar{r}), 7.(r \mid \bar{s})\}$$

Please show the proofs of these resolution strategies.

14. Decide by ordered resolution if the clause set  $S$  is satisfiable or not, where  $S = \{1.(a \mid c \mid \bar{d} \mid f), 2.(a \mid d \mid \bar{e}), 3.(\bar{a} \mid c), 4.(b \mid \bar{d} \mid e), 5.(b \mid d \mid \bar{f}), 6.(\bar{b} \mid \bar{f}), 7.(c \mid \bar{e}), 8.(c \mid e), 9.(\bar{c} \mid \bar{d})\}$ , and the order is  $a > b > c > d > e > f$ . If the clause set is satisfiable, please display all the resolvents from ordered resolution and construct a model from these clauses. If the clause is unsatisfiable, please display an ordered resolution tree of the proof.
15. Prove by resolution that  $S \models \neg f$ , where  $S$  is the set of clauses given in the previous problem. Please display the resolution tree of this proof.
16. If head/tail literals are used for implementing BCP, what literals of the clause set  $S$  in the previous problem will be assigned a truth value during unit resolution before an empty clause is found? in what order?

17. A set  $C$  of clauses is said to be a *2CNF* if every clause contains at most two literals. The *2SAT* problem is to decide if a 2CNF formula is satisfiable or not. Show that *orderedResolution* takes  $O(n^3)$  time and  $O(n^2)$  space to solve 2SAT, where  $n$  is the number of propositional variables in 2CNF.
18. Prove Theorem 3.4.8: For any set  $C$  of clauses, if there exists a unit resolution proof from  $C$ , then there exists an input resolution proof from  $C$ .

## CHAPTER 4

### PROPOSITIONAL SATISFIABILITY

Propositional satisfiability (SAT) is the problem of deciding if a propositional formula is satisfiable, i.e., the formula has a model. Software programs for solving SAT are called SAT solvers. SAT has long enjoyed a special status in computer science. On the theoretical side, it is the first NP-complete problem ever discovered. NP-complete problems are notorious for being hard to solve; in particular, in the worst case, the computation time of any known solution for a problem in this class increases exponentially with the input size. On the practical side, since every NP problem can transform into SAT, an efficient SAT solver can be used to solve many practical problems. SAT found several important applications in the design and verification of hardware and software systems, and in many areas of artificial intelligence and operation research. Thus, there is strong motivation to develop practically useful SAT solvers. However, the NP-completeness is cause for pessimism, since it is unlikely that we will be able to scale the solutions to large practical instances.

In theory, the decision procedures presented in the previous chapter, such as semantic tableau or resolution, can be used as SAT solvers. However, they served largely as academic exercises with little hope of seeing practical use. Resolution as a refutation prover searches for an empty clause. Semantic tableau imitates the transformation of a formula into DNF. None of these decision procedures aims at searching for a model.

Fortunately, in the last two decades or so, several research developments have enabled us to tackle instances with thousands of variables and millions of clauses. SAT researchers introduced techniques such as conflict-driven clause learning, novel branching heuristics, and efficient unit propagation. These techniques form the basis of many modern SAT solvers. Using these ideas, contemporary SAT solvers can often handle practical very large SAT instances.

Annual competitions of SAT solvers promote the continuing advances of SAT solvers. SAT solvers are the engines of modern model checking tools. Contemporary automated verification techniques such as bounded model checking, proof-based abstraction, interpolation-based model checking, are all based on SAT solvers and

their extensions. These tools are used to check the correctness of hardware designs.

#### 4.1 The DPLL Algorithm

In 1960 David and Putnam proposed resolution for propositional logic. For realistic problems, the number of clauses generated by resolution grows quickly. To avoid this explosion, Davis, Logemann, and Loveland suggested to replace the resolution rule with a case split: Pick a variable  $p$  and consider the two cases  $p \mapsto 1$  and  $p \mapsto 0$ . This modified algorithm is commonly referred to as the *DPLL* algorithm, where DPLL stands for Davis, Putnam, Logemann, and Loveland.

*DPLL* can be regarded as a search procedure: The search space is all the interpretations, partial or full. A partial interpretation can be represented by a set of literals, where each literal in the set is assigned to be true and each literal outside of the set is assumed unassigned. *DPLL* starts with an empty interpretation and tries to extend the interpretation by adding either  $p$  or  $\bar{p}$  through the case split. *DPLL* stops when every variable is assigned without contradiction, and the partial interpretation it maintains becomes a model.

##### 4.1.1 Recursive Version of *DPLL*

In the following, we give a recursive version of *DPLL*. The first call to *DPLL* with the input clauses  $C$  and  $\sigma = \{\}$ . It is one of the oldest backtracking procedures. *DPLL* uses the algorithm *BCP* heavily for unit propagation. From the previous chapter, we know that *BCP* implements unit resolution and subsumption deletion, and can be done in linear time. Given a set  $C$  of clauses,  $BCP(C)$  returns  $\perp$  if an empty clause is found by unit resolution; otherwise,  $BCP(C)$  returns  $(U, S)$ , such that  $C \equiv (U \cup X)$ ,  $U$  is a set of unit clauses,  $S$  does not contain any unit clauses,  $U$  and  $S$  share no variables.

**Algorithm 4.1.1.** *DPLL*( $C, \sigma$ ) takes as input a set  $C$  of clauses and a partial interpretation  $\sigma$ ,  $C$  and  $\sigma$  do not share any variables. It returns a model of  $C$  if  $C$  is satisfiable; return  $\perp$  if  $C$  is unsatisfiable. Procedure *pickLiteral*( $S$ ) will pick a literal appearing in  $S$ .

```

proc DPLL( $C, \sigma$ )
1    $res := BCP(C)$ ;
2   if ( $res = \perp$ ) return  $\perp$ ; // no models for  $C$ 
3   // It must be the case  $res = (U, S)$ 
4    $\sigma := \sigma \cup U$ ;
5   if ( $S = \emptyset$ ) return  $\sigma$ ; // a model is found
6    $A := pickLiteral(S)$ ;

```

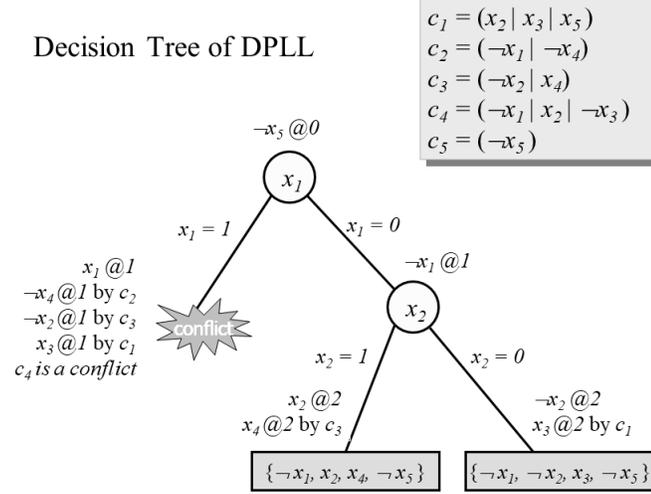


Figure 4.1.1: The result of BCP is shown in the box near each node.

```

7   res := DPLL(S ∪ {(A)}, σ);
8   if (res ≠ ⊥) return res; // a model is found
9   return DPLL(S ∪ {(Ā)}, σ); // return either ⊥ or a model

```

In the above procedure, we assume again  $\bar{A} = p$  if  $A$  is the literal  $\bar{p}$ . If  $BCP(C)$  does not return  $\perp$ , it simplifies  $C$  to  $U \cup S$ , where  $U$  is a set of literals,  $S$  is a set of non-unit clauses, and  $U$  and  $S$  do not share any variables. By Proposition 3.4.3,  $C \equiv U \cup S$ . We identify two types of literals in  $U$ .

**Definition 4.1.2.** In DPLL, a literal  $A \in U$  is said to be a decision literal if  $A$  is the one returned by  $\text{pickLiteral}(S)$  at line 6 of DPLL; otherwise,  $A$  is said to be an implied literal.

Since  $U$  is a set of literals without complementary literals,  $U$  can be regarded as a partial interpretation  $\sigma$  such that  $\sigma(U) = 1$ . Since  $U$  is satisfiable and  $C \equiv U \cup S$ , we conclude that  $C \approx S$ . Furthermore, checking the satisfiability of  $S$  is reduced recursively to testing  $S \wedge A$  and  $S \wedge \neg A$  separately for a literal  $A$  appearing in  $S$ . This transformation, applied recursively, yields a complete decision procedure.

**Example 4.1.3.** Let  $C = \{1 : (x_2 \mid x_3 \mid x_5), 2 : (\bar{x}_1 \mid \bar{x}_4), 3 : (\bar{x}_2 \mid x_4), 4 : (\bar{x}_1 \mid x_2 \mid \bar{x}_3), 5 : (\bar{x}_5)\}$ . If  $\text{pickLiteral}$  picks a variable in the order of  $x_1, x_2, \dots, x_5$ , then

$DPLL(C, \{\})$  calls

$BCP(C)$ , which returns  $(\{\bar{x}_5\}, C_1 = \{1' : (x_2 \mid x_3), 2., 3., 4.\})$ ;  
 $x_1$  is picked as the decision literal of level 1;

$DPLL(C_1 \cup \{(x_1)\}, \{\overline{x_5}\})$ , which calls  
 $BCP(C_1 \cup \{(x_1)\})$ , which returns  $\perp$ ;  
 and returns  $\perp$   
 $DPLL(C_1 \cup \{(\overline{x_1})\}, \{\overline{x_5}\})$ , which calls  
 $BCP(C_1 \cup \{(\overline{x_1})\})$ , which returns  $(\{\overline{x_1}\}, C_2 = \{1', 3.\})$ ;  
 $x_2$  is picked as the decision literal of level 2;  
 $DPLL(C_2 \cup \{(x_2)\}, \{\overline{x_1}, \overline{x_5}\})$ , which calls  
 $BCP(C_2 \cup \{(x_2)\})$ , which returns  $(\emptyset, \{\overline{x_1}, x_2, x_4, \overline{x_5}\})$   
 and returns  $\{\overline{x_1}, x_2, x_4, \overline{x_5}\}$   
 and returns  $\{\overline{x_1}, x_2, x_4, \overline{x_5}\}$   
 and returns  $\{\overline{x_1}, x_2, x_4, \overline{x_5}\}$

Thus a model of  $C$  can be constructed from  $\{\overline{x_1}, x_2, x_4, \overline{x_5}\}$  by adding  $x_3 \mapsto v$  for any value  $v$ .

In the above example, the height of the recursion tree (also called *decision tree*) is two;  $x_1$  and  $x_2$  are decision literals;  $x_4$  and  $\overline{x_5}$  are implied literals. We have called  $BCP$  four times: one call at the root (level 0); two calls at the nodes of level one; and one call at the node of level two. If we want to search for all models, we may choose to try the branch with  $x_2 \mapsto 0$ , as shown in Figure 4.1.1.

**Definition 4.1.4.** *The level of an assigned variable is the depth of a node of the recursion tree of  $DPLL(C, \{\})$  at which the variable is assigned a truth value. We write  $x@l$  if  $x \mapsto 1$  at level  $l$ , and  $\overline{x}@l$  if  $x \mapsto 0$  at level  $l$ .*

In the above example, for the model returned by  $DPLL$ , we may write it as

$$\{\overline{x_1}@1, x_2@2, x_4@2, \overline{x_5}@0\}$$

**Theorem 4.1.5.**  $DPLL(C, \{\})$  is a decision procedure for the satisfiability of  $C$ .

*Proof.* Assuming  $\perp \notin C$ . To prove the theorem, we just need to show that  $DPLL$  is sound and terminating.  $DPLL(C, \sigma)$  will terminate because for each recursive call, the number of variables in  $C$  is decreased by at least one. In  $DPLL(C, \sigma)$ , since  $\sigma$  and  $C$  share no variables,  $(C \cup \sigma) \approx C$ . So in the soundness proof of  $DPLL$ , we ignore the presence of  $\sigma$ .

We will prove by induction on the structure of the recursion tree of  $DPLL$ . As a base case, if  $DPLL(C, \sigma)$  returns  $\perp$ , it means that  $\perp$  is generated by unit resolution inside  $BCP(C)$ . By Proposition 3.4.3,  $C \equiv \perp$ . If  $BCP(C, \sigma)$  returns  $(U, S)$ , by Proposition 3.4.3,  $C \equiv U \cup S$ . As another base case, if  $S = \emptyset$ , then  $C \equiv U$  and  $U$

is the model for the unit clause set  $U$  if  $U$  is regarded as an interpretation, so  $C$  is satisfiable. If  $S \neq \emptyset$ , then since

$$S \equiv (S \wedge A) \vee (S \wedge \neg A)$$

for any formula  $A$ , if either  $DPLL(S \cup \{p\}, \sigma)$  or  $DPLL(S \cup \{\bar{p}\}, \sigma)$  returns a model, by induction hypotheses, it means  $(S \wedge p) \vee (S \wedge \neg p)$  is satisfiable, so is  $S$ . If both  $DPLL(S \cup \{p\}, \sigma)$  and  $DPLL(S \cup \{\bar{p}\}, \sigma)$  return  $\perp$ , by induction hypotheses, it means  $(S \wedge p) \vee (S \wedge \neg p)$  is unsatisfiable, so is  $S$ .  $\square$

#### 4.1.2 All-SAT and Incremental SAT Solvers

Given a satisfiable set of clauses,  $DPLL$  returns a single satisfying assignment. Some applications, however, require us to enumerate all satisfiable assignments of a formula and this problem is called *All-SAT*. It is easy to modify  $DPLL$  so that all the models can be found. That is, we replace line 5 of  $DPLL$  as follows:

```
5:   if ( $S = \emptyset$ ) save( $\sigma$ ); return  $\perp$ ; // save model and continue
```

When the modified  $DPLL$  finishes, all the models are saved. Note that these models are partial models. One partial model represents  $2^k$  full models, where  $k$  is the number of unassigned variables in the partial model.

During the execution of  $DPLL$ , the list of decision literals (labeled with 1) and their complements (labeled with 2) decides the current position in the recursion tree of  $DPLL$ . This list of literals, called *guiding path*, can be very useful for designing an incremental SAT solver: The solver takes the guiding path as input, uses the literals in the guiding path to skip the branches already explored and explores the new branches in the recursion tree for a new model.

**Example 4.1.6.** In Example 4.1.3, a model is found with the decision literals  $\bar{x}_1$  and  $x_2$ . The guiding path for this model is  $[\bar{x}_1 : 2, x_2 : 1]$ , where the label 2 after  $\bar{x}_1$  means this is the second branch of  $x_1$ . The label 1 after  $x_2$  means this is the first branch of  $x_2$ . If the next model is needed, it must be in the subtree with the initial guiding path  $[\bar{x}_1 : 2, x_2 : 2]$ . Using  $[\bar{x}_1 : 2, x_2 : 2]$ , the SAT solver avoids the search space it has traveled before.

Since modern SAT solvers are often designed as a black box, it is cumbersome to use them as an incremental solver for enumerating each model. We can force the SAT solver to find a new model by subsequently blocking all models previously found. A model can be viewed as a conjunction of true literals, one for every

variable. The negation of this conjunction can be represented as a clause, called *blocking clause*. If this clause is added into the set of input clauses, the model will be blocked by this clause and cannot be generated by the SAT solver. Since a blocking clause contains every variable, it is unlikely to become unit and it is expensive to keep all the blocking clauses. Therefore, it is desirable to reduce the size of blocking clauses, i.e., to construct a smaller clause which still blocks the model. If we know the guiding path of this model, the negation of the conjunction of all the decision literals in the guiding path will block the model.

Many applications of SAT solvers require solving a sequence of instances by adding more clauses. Incremental SAT solvers with the guiding path can support this application without any modification, because the search space skipped by the guiding path of the first model will contain no models.

In many application scenarios, it is beneficial to be able to make several SAT checks on the same input clauses under different forced partial interpretation, called “*assumption*”. For instance, people may be interested in questions like “Is the formula  $F$  satisfiable under the assumption  $x \mapsto 1$  and  $y \mapsto 0$ ?” In such applications, the input formula is read in only once, the user implements an iterative loop that calls the same solver instantiation under different sets of assumptions. The calls can be adaptive, i.e., assumptions of future SAT solver calls can depend on the results of the previous solver calls. The SAT solver can keep its internal state from the previous call to the next call. *DPLL* with a guiding path can be easily modified to support such applications, treating the assumption as a special set of unit clauses which can be added into or removed from the solver. The guiding path allows us to backtrack to the new search space but never revisit the search space visited before.

To obtain a high-performance SAT solver, we may use the deletion techniques presented in section 3.3.6 as pre-processing techniques. They enable us to reduce the size of the formula before passing it to *DPLL*.

Inside the body of *DPLL*, it calls itself twice, once with  $S \cup \{(A)\}$  and once with  $S \cup \{(\bar{A})\}$ . A persistent data structure for  $S$  would quickly run out of memory as  $S$  is very large and the recursion tree is huge for practical SAT instances. A better approach is to use a destructive data structure for  $S$ : We remember all the operations performed onto  $S$  inside the first call, and then undo them if necessary when we enter the second call.

#### 4.1.3 *BCPw*: Implementation of Watch Literals

In Algorithm 3.4.16, we presented an efficient implementation of *BCP* as a decision procedure for Horn clauses. This version of *BCP*, called the procedure *BCPht*, uses the head/tail data structure for each non-unit clause. Naturally, *BCPht*

can be used to support *DPLL*. Better yet, with minor modification, *BCPht* can implement the idea of *watch literals*, which is based on head/tail literals and used in solver *zchaff* by Moskewicz and Zhang.

The function of watch literals is the same as head/tail literals: Use two literals to decide if a clause becomes unit or conflicting, as stated in Proposition 3.4.15. Head/tail literals are watch literals; when head/tail literals are allowed to appear anywhere in a clause, head/tail literals are identical to watch literals.

Since the idea of watch literals is based on that of head/tail literals, the implementation of watch literals is easy if we have an implementation of head/tail literals. In *BCPht*, when a head or tail literal of  $c$  becomes false, we look for its replacement at index  $x$ , where  $head(c) < x < tail(c)$ . There are two alternative ways to modify *BCPht*. One solution is to allow  $x \in \{0, 1, \dots, |c| - 1\} - \{head(c), tail(c)\}$ ; we will introduce an implementation of this solution shortly. The other solution is to fix  $head(c) = 0$  and  $tail(c) = 1$  or  $|c| - 1$ , and swap the elements of  $lits(c)$  if necessary.

The second solution saves the space for  $head(c)$  and  $tail(c)$ , but may cost more for visiting literals of  $lits(c)$ . For example, if a clause has  $k$  literals which become false in the order from left to right during *BCP*. The second solution may take  $O(k^2)$  time to visit the literals of the clause, while the first solution takes  $O(k)$  time. As a result, *BCP* cannot be a linear time decision procedure for Horn clauses if the second solution is used. Of course, we may introduce a pointer in the clause to remember the position of the last visited literal, so that the second solution still gives us a linear time algorithm.

**Algorithm 4.1.7.** *BCPw(U)* is an implementation of *BCP(C)* based on the head/tail data structure. *BCPw(U)* works the same way as *BCPht(U)*, where  $U$  is a stack of unit clauses and non-unit clauses are stored using the head/tail data structure. Like *BCPht*, *BCPw* returns a conflicting clause if an empty clause is generated; returns “SAT” otherwise. Unlike *BCPht*, values of  $head(c)$  and  $tail(c)$  can be any valid indices of  $lits(c)$ .

```

proc BCPw(U)
1   while  $U \neq \{\}$  do
2        $A := pop(U)$ ;
3       if  $val(A) = 0$  return  $reason(\bar{A})$ ; // an empty clause is found;
4       else if  $(val(A) \neq 1)$ 
5            $val(A) := 1; val(\bar{A}) := 0$ ;
6       for  $c \in cls(\bar{A})$  do //  $\bar{A}$  is either head or tail of  $c$ 
7           if  $(\bar{A} = lits(c)[head(c)])$ 

```

```

8            $e_1 := \text{head}(c); e_2 := \text{tail}(c); \text{step} := 1; //$  scan from left to right
9     else
10           $e_1 := \text{tail}(c); e_2 := \text{head}(c); \text{step} := -1; //$  scan from right to left
11     while true do
12           $x := (e_1 + \text{step}) \bmod |c|; //$   $x$  takes any valid index of  $\text{lits}(c)$ 
13          if  $x = e_1$  break;  $//$  exits the inner while loop
13'         if  $x = e_2$  continue;  $//$  go to line 12
14           $B := \text{lits}(c)[x];$ 
15          if  $(\text{val}(B) \neq 0)$ 
16              $\text{remove}(c, \text{cls}(\bar{A})); \text{insert}(c, \text{cls}(B));$ 
17             if  $(\text{step} = 1)$   $\text{head}(c) := x;$  else  $\text{tail}(c) := x;$ 
18             break;  $//$  exit the inner while loop
19     if  $(x = e_1)$   $//$  no new head or tail found
20          $A := \text{lits}(c)[e_2]; //$   $c$  is unit or conflicting.
21         if  $\text{val}(A) = 0$  return  $c;$   $//$   $c$  is conflicting.
22         else if  $(\text{val}(A) \neq 1)$   $//$   $c$  is unit.
23              $\text{push}(A, U); \text{reason}(A) := c; \text{val}(A) := 1; \text{val}(\bar{A}) := 0;$ 
24     return “SAT”;  $//$  no empty clauses are found;

```

The above algorithm comes from *BCPht* (Algorithm 3.4.16) by the following modifications, which allow the values of  $\text{head}(c)$  and  $\text{tail}(c)$  to be any valid indices of  $\text{lits}(c)$ .

- At line 12,  $x := x + \text{step}$  is replaced by  $x := (x + \text{step}) \bmod |c|$ .
- At line 13, the condition  $x = e_2$  is replaced by  $x = e_1$ .
- Line 13' is added to skip the literal watched by  $e_2$ ; this is the literal to be in the unit clause if the clause becomes unit.
- At line 19, the condition  $x = e_2$  is replaced by  $x = e_1$ ; when this is true, every literal of the clause is false, with the exception of the literal watched by  $e_2$ .
- At line 23, we record the reason for the assignment of a literal; this information is used at line 3 when a clause becomes conflicting.

In *DPLL*, the book-keeping required to detect when a clause becomes unit can involve a high computational overhead if implemented naively. Since it is sufficient to watch in each clause two literals that have not been assigned yet, assignments to the non-watched literals can be safely ignored. When a variable  $p$  is assigned 1, the SAT solver only needs to visit clauses watched by  $\bar{p}$ . Each time one of the watched

literals becomes false, the solver chooses one of the remaining unassigned literals to watch. If this is not possible, the clause is necessarily unit, or becomes a conflicting clause, or is already satisfied under the current partial assignment. Any sequence of assignments that makes a clause unit will include an assignment of one of the watched literals. The computational overhead of this strategy is relatively low: In a formula with  $m$  clauses and  $n$  variables,  $2m$  literals need to be watched, and  $m/n$  clauses are visited per assignment on average. This advantage is inherited from the idea of head/tail literals.

The key advantage of the watch literals over the head/tail literals is that the watched literals do not need to be updated upon backtracking in *DPLL*. That is, the key advantage of *BCP<sub>w</sub>* over *BCP<sub>ht</sub>* is its use inside *DPLL*: When you backtrack from a recursive call of *DPLL*, you do not need to undo the operations on *head*( $c$ ) and *tail*( $c$ ), because any two positions of  $c$  will be sufficient for checking if  $c$  becomes unit.

#### 4.1.4 Iterative Implementation of *DPLL*

Since the compiler will convert a recursion procedure into an iterative one using a stack, to avoid the expense of this conversion, in practice, the procedure *DPLL* is not implemented by means of recursion but in an iterative manner using a stack for the partial interpretation  $\sigma$ . If the head/tail data structure is used,  $\sigma$  can be obtained from *val*( $a$ ), but we still need a stack of literals to record the time when a literal is assigned. This stack is called *trail*, which can also serve as a partial interpretation. We keep track of the recursive case-splits and their implications using the trail. Every time a literal is assigned a value, either by a decision, or an implication inside *BCP*, the literal is pushed into the trail and its level in the recursion tree is remembered.

A trail may lead to a dead end, i.e., result in a conflicting clause, in which case we have to explore the alternative branch of one of the case splits previously made. This corresponds to backtracking which reverts one of the decisions. When we backtrack from a recursive call, we pop those literals off the trail, and make them as “unassigned” (denoted by  $\times$ ). When backtracking enough times, the search algorithm always yields a conflicting clause or a satisfying assignment and eventually exhausts all branches.

To implement the above idea, we use the trail  $\sigma$  which is a stack of literals and serves as a partial interpretation.  $\sigma$  will keep all the assigned literals. For each literal  $A$ , we use *lvl*( $A$ ) to record the level at which  $A$  is assigned a truth value.

This is actually done at line 5 of *BCPw*. The Line 5 of *BCPw* now is read as follows:

```
5 :   val(A) := 1; val( $\bar{A}$ ) := 0; lvl(A) := lvl( $\bar{A}$ ) := level; push(A,  $\sigma$ );
```

where *val*, *lvl*,  $\sigma$ , and *level* are global variables. From now on, we call the modified *BCPw* simply as *BCP*.

In practice, it is more convenient to set the level of a right child the same as the level of its parent in the recursion tree, to indicate that the right child has no alternatives to consider. In the following implementation of *DPLL*, we have used this idea.

**Algorithm 4.1.8.** Non-recursive *DPLL*(*C*) takes as input a set *C* of clauses and returns a model of *C* if *C* is satisfiable; returns 0 if *C* is unsatisfiable. Procedure *initialize*(*C*) will return the unit clauses of *C* and store the rest clauses in the head/tail data structure.

```
proc DPLL(C)
```

```
1   U := initialize(C); // initialize the head/tail data structure
2    $\sigma$  := {}; level := 0; // initialize  $\sigma$  and level
3   while (true) do
4     if (BCP(U) =  $\perp$ )
5       if (level = 0) return  $\perp$ ; // C is unsatisfiable
6       level := level - 1; // level of the second child
7       A := undo(); // A is the last literal picked at line 10
8       U := { $\bar{A}$ };
9     else
10      A := pickLiteral(); // pick an unassigned literal
11      if (A = nil) return  $\sigma$ ; // all literals assigned, a model is found
12      level := level + 1; // level of the first child
13      U := {A};
```

```
proc undo()
```

```
  // Backtrack to level, undo all assignments to literals of higher levels.
  // Return the last literal of level + 1 in the trail, or nil.
1  A := nil;
2  while ( $\sigma \neq \{\}$ )  $\wedge$  (lvl(top( $\sigma$ )) > level) do
3    A := pop( $\sigma$ );
4    val(A) := val( $\bar{A}$ ) :=  $\times$ ; //  $\times$  means “unassigned”
5  return A;
```

The procedure *undo()* will undo all the assignments (made at line 5 of *BCP*) to the literals of levels higher than *level* (including their complements), and return the last literal of *level* + 1 in  $\sigma$ , which must be the last decision literal *A* picked at line 10 of *DPLL*.

The procedure *pickLiteral(L)* (line 12) is the place for implementing various heuristics of best search strategies. It leaves the question of which literals to assign open. As we know from Theorem 4.1.5, this choice has no impact on the completeness of the search algorithm. It has, however, a significant impact on the performance of the solver, since this choice is instrumental in pruning the search space. We will address this issue later.

## 4.2 Conflict-Driven Clause Learning (CDCL)

This section introduces conflicting clauses as a source information to prevent the repeated exploration of futile search space. We start with a motivating example.

**Example 4.2.1.** Suppose *C* is a set of clauses over the variable set  $\{x_1, x_2, \dots, x_n\}$  and *C* has more than a million models. Let  $C' = C \cup S$ , where  $S = \{(y_1 \mid y_2), (y_1 \mid \bar{y}_2), (\bar{y}_1 \mid y_2), (\bar{y}_1 \mid \bar{y}_2)\}$ . Then *C'* has no models, because the added four clauses are unsatisfiable. Now we feed *C'* to *DPLL* and assume that *pickLiteral* will prefer  $x_i$  over  $y_1$  and  $y_2$ . Every time when *DPLL* picks  $y_1$  for case-split, a model of *C* was found at that point. If  $y_1$  is true, then one of  $(\bar{y}_1 \mid y_2)$  and  $(\bar{y}_1 \mid \bar{y}_2)$  will be conflicting. If  $y_1$  is false, then one of  $(y_1 \mid y_2)$  and  $(y_1 \mid \bar{y}_2)$  will be conflicting. Either value of  $y_1$  will lead to a conflicting clause in *S*, so *DPLL* will backtrack. Since *C* has more than a million models, it means *DPLL* will try to show that *S* is unsatisfiable more than a million times.

If we can learn from the conflicting clauses of *S*, we may avoid repeatedly doing the same work again and again. This is the motivation for the technique of conflict-driven clause learning (CDCL): Learn new information from conflicting clauses and use the new information to avoid futile search effort. Using CDCL, new information is obtained in the form of clauses which are then added back into the search procedure. For the above example, when  $y_1$  is true by case-split,  $(\bar{y}_1 \mid y_2)$  becomes unit, so  $y_2$  must be set to 1. Now  $(\bar{y}_1 \mid \bar{y}_2)$  becomes conflicting. A resolution between these two clauses on  $p_2$  will generate a new clause, i.e.,  $(\bar{y}_1)$ . The new clause has nothing to do with the assignments of  $x_i$ , so we can jump back to the level 0 of *DPLL* and the subsequent *BCP* will find the empty clause from  $S \wedge \bar{y}_1$  and *DPLL* returns  $\perp$ , thus avoiding finding models in *C* a million times.

### 4.2.1 Generating Clauses from Conflicting Clauses

Recall that in the last *DPLL* algorithm, when a non-unit clause  $c$  becomes unit, where  $c$  contains an unassigned literal  $A$ , we do  $reason(A) := c$ , i.e.,  $c$  is the reason for the unassigned literal  $A$  to be true. When a conflicting clause is found in *DPLL*, all of its literals are false, and we pick the literal  $A$ , which is assigned last in  $c$ , in the clause and do resolution on  $A$  between this clause and  $reason(\bar{A})$ , hoping the resolvent will be useful to cut the search space. Since resolution is sound, new clauses generated this way are logical consequences of the input and it is safe to keep them. For practical reason, we must selectively do resolution and keep at most one clause per conflicting clause.

**Algorithm 4.2.2.** *conflictAnalysis(c)* takes as input a conflicting clause  $c$  in the current partial interpretation at level  $level$ , generates a new clause by resolution, and returns it. Procedure *latestAssignedLiteral(c)* will return a literal  $A$  of  $c$  such that  $A$  is the last literal assigning a truth value in  $c$ . Procedure *countLevel(c, lvl)* will return the number of literals in  $c$  which are assigned at level  $lvl$ . Procedure *resolve(c<sub>1</sub>, c<sub>2</sub>)* returns the resolvent of  $c_1$  and  $c_2$ .

```

proc conflictAnalysis(c)
    //  $c$  is conflicting in the partial interpretation  $val()$ 
1    $\alpha := c$ ;
2   while ( $|\alpha| > 1$ ) do // as  $\alpha$  has more than one literal
3      $A := latestAssignedLiteral(\alpha)$ ; //  $A$  is assigned last in  $\alpha$ 
4      $\alpha := resolve(\alpha, reason(\bar{A}))$ ; // resolution on  $A$ 
5     if ( $countLevel(\alpha, level) \leq 1$ ) return  $\alpha$ ;
6   return  $\alpha$ ;

```

**Example 4.2.3.** Let the input clauses  $C = \{c_1 : (x_2 \mid x_3), c_2 : (\bar{x}_1 \mid \bar{x}_4), c_3 : (\bar{x}_2 \mid x_4), c_4 : (\bar{x}_1 \mid x_2 \mid \bar{x}_3)\}$ . If we pick  $x_1$  as the first decision literal in *DPLL*, it will result in the following assignments:

assignment	reason
$x_1@1$	$(x_1)$
$\bar{x}_4@1$	$c_2 : (\bar{x}_1 \mid \bar{x}_4)$
$\bar{x}_2@1$	$c_3 : (\bar{x}_2 \mid x_4)$
$x_3@1$	$c_1 : (x_2 \mid x_3)$

Clause  $c_4$  becomes conflicting, which contains three literals of level 1. Calling *conflictAnalysis(c<sub>4</sub>)*, we will have the following resolutions:

1.  $\alpha_0$  is  $c_4$ ;

2.  $\alpha_1 : (\overline{x_1} \mid x_2) = \text{resolve}(\alpha_0, c_1)$ .
3.  $\alpha_2 : (\overline{x_1} \mid x_4) = \text{resolve}(\alpha_1, c_3)$ .
4.  $\alpha_3 : (\overline{x_1}) = \text{resolve}(\alpha_2, c_2)$  and is returned.

**Proposition 4.2.4.** *Assuming  $c$  is conflicting in the current interpretation.*

- (a) *The procedure  $\text{conflictAnalysis}(c)$  will terminate and return a clause which is an entailment of the input clauses.*
- (b)  *$\text{conflictAnalysis}(c)$  will return a new conflicting clause in the current interpretation.*
- (c) *If we backtrack off the current level, the new clause returned by  $\text{conflictAnalysis}(c)$  will have at most one literal at level level  $- 1$ .*

*Proof.* (a): We assigned only a finite number of literals a truth value at the current level. The literals we chose to resolve off are in the backward order when they are assigned. Once resolved off, these literals can never come back into the new clause, thus there is no chance for a loop. Since the new clause is generated by resolution from the input clauses, it must be an entailment of the input clauses.

(b): In fact, any resolvent in  $\text{conflictAnalysis}(c)$  is a conflicting clause in the current interpretation. Since  $c$  is a conflicting clause, every literal is false in the current interpretation.  $c$  is resolved with a reason clause which has only one true literal in the interpretation and this true literal was resolved off during the resolution. The resulting resolvent will contain only false literals in the current interpretation.

(c): The procedure  $\text{conflictAnalysis}(c)$  will terminate when at most one literal, say  $B$ , of the current level remains in the new clause. If we backtrack off the current level,  $B$  will be unassigned and the rest literals are false by (b).  $\square$

#### 4.2.2 DPLL with CDCL

The new *DPLL* procedure with CDCL (conflict-driven clause learning) can be described as follows:

**Algorithm 4.2.5.** *DPLL(C)* takes as input a set  $C$  of clauses and returns a model of  $C$  if  $C$  is satisfiable; return  $\perp$  if  $C$  is unsatisfiable. Procedure  $\text{initialize}(C)$  will return the unit clauses of  $C$  and store the rest clauses in the head/tail data structure.

**proc** *DPLL(C)*

1     $U := \text{initialize}(C)$ ; // initialize the head/tail data structure

```

2    $\sigma := \{\}$ ;  $level := 0$ ; // initialize  $\sigma$  and  $level$ 
3   while (true) do
4      $res := BCP(U)$ ;
5     if ( $res \neq \text{"SAT"}$ ) //  $res$  contains a conflicting clause
6       if ( $level = 0$ ) return  $\perp$ ; //  $C$  is unsatisfiable
7        $U := insertNewClause(conflictAnalysis(res))$ ; // new level is set
8        $undo()$ ; // undo to the new level
9     else
10       $A := pickLiteral()$ ; // pick an unassigned literal
11      if ( $A = nil$ ) return  $\sigma$ ; // all literals assigned, a model is found
12       $level := level + 1$ ; // level of the first child
13       $U := \{A\}$ ;

proc  $insertNewClause(\alpha)$ 
1  if ( $|\alpha| = 1$ )
2     $level := 0$ ;
3     $A := literal(\alpha)$ ;  $reason(A) := \alpha$ ; // assume  $\alpha = (A)$ 
4    return  $\{A\}$ ;
5  else // insert  $\alpha$  into the clause set
6     $lits(\alpha) := makeArray(\alpha)$ ; // create an array of literals
7     $secondHigh := 0$ ; // look for the second highest level in  $\alpha$ 
8    for  $x := 0$  to  $|\alpha| - 1$  do
9      if ( $lvl(lits(\alpha)[x]) = level$ )  $head(\alpha) := x$ ;  $H := lits(\alpha)[x]$ ; //  $lvl(H) = level$ 
10     else if ( $lvl(lits(\alpha)[x]) > secondHigh$ )  $secondHigh := lvl(lits(\alpha)[x])$ ;  $t := x$ ;
11      $tail(\alpha) := t$ ;
12      $level := secondHigh$ ;
13      $reason(H) := \alpha$ ;
14     return  $\{H\}$ ; // the head literal of  $\alpha$ 

```

The procedure *insertNewClause* will return a singleton set of literals served as new unit clauses and reset the *level* variable of *DPLL*. If the new clause  $\alpha$  generated by *conflictAnalysis* is a unit clause, this clause will be returned and the level is set to 0. If the new clause is not unit, by Proposition 4.2.4, there is only one literal of the current level in  $\alpha$  and this literal will be the head literal of  $\alpha$  when the level is decreased. We then pick one literal of the highest level other than the head literal as the tail literal. The level will be reset to that of the tail literal. When we backtrack to this level, by Proposition 4.2.4,  $\alpha$  becomes unit and is the reason for the head literal, unassigned at this point, to be assigned true.

**Example 4.2.6.** Consider the clause set  $C = \{c_1 = (x_1 \mid x_2), c_2 = (x_1 \mid x_3 \mid \overline{x_7}), c_3 = (\overline{x_2} \mid \overline{x_3} \mid x_4), c_4 = (\overline{x_4} \mid x_5 \mid x_6), c_5 = (\overline{x_3} \mid \overline{x_4}), c_6 = (x_5 \mid \overline{x_6}), c_7 = (\overline{x_2} \mid x_3), c_8 = (\overline{x_1} \mid x_3), c_9 = (\overline{x_1} \mid x_2 \mid \overline{x_3})\}$ . If the first three literals chosen by *pickLiteral* in *DPLL* are  $\overline{x_7}, \overline{x_6}$ , and  $\overline{x_5}$ , in that order, then the only implied literal is  $\overline{x_4}@3$  by  $c_4$ . The next decision literal is  $\overline{x_1}$ , which implies the following assignments:  $x_2@4$  by  $c_1$  and  $\overline{x_3}@4$  by  $c_3$ .  $c_7$  becomes conflicting at level 4.

Calling *conflictAnalysis*( $c_7$ ), we obtain  $(\overline{x_2} \mid x_4)$  by resolution between  $c_7$  and  $c_3$ , which contains  $\overline{x_2}$  as the only literal at level 4. We add  $(\overline{x_2} \mid x_4)$  into *DPLL* as  $c_{10}$ .

The *level* of *DPLL* is set to 3 as  $lvl(x_4) = 3$ . So we backtrack to level 3, and  $c_{10} = (\overline{x_2} \mid x_4)$  becomes unit. The next call to *BGP* will have make the following assignments at level 3:  $\overline{x_2}@3$  by  $c_{10}$ ,  $x_1@3$  by  $c_1$ ,  $x_3@3$  by  $c_8$ , and  $c_9$  becomes conflicting.

Calling *conflictAnalysis*( $c_9$ ), we obtain  $c_{11} = (x_2)$  by two resolutions between  $c_8, c_9$ , and  $c_1$ . The new clause is unit and we set *level* = 0. Back jumping to level 0 by *undo*, the next call to *BGP* will imply the following assignments:  $x_2@0$  by  $c_{11}$ ,  $x_3@0$  by  $c_7$ ,  $x_4@0$  by  $c_{10}$ , and  $c_5$  becomes conflicting.

Calling *conflictAnalysis*( $c_5$ ), we obtain  $c_{12} = (\overline{x_2})$  by two resolutions between  $c_5, c_{10}$ , and  $c_7$ . Clauses  $c_{11}$  and  $c_{12}$  will generate the empty clause. Thus, the input clauses are unsatisfiable.

Clause learning with conflict analysis does not impair the completeness of the search algorithm: even if the learned clauses are dropped at a later point during the search, the trail guarantees that the solver never repeatedly enters a decision level with the same partial assignment. We have shown the correctness of clause learning by demonstrating that each conflicting clause is implied by the original formula.

The idea of CDCL was used first in solver *grasp* of Marques-Silva and Sakallah, and solver *relsat* of Bayardo and Schrag. In their solvers, they introduced independently a novel mechanism to analyze the conflicts encountered during the search for a satisfying assignment. There are many ways to generate new clauses from conflicting clauses and most of them are based on *implication graphs*. The method of learning new clauses through resolution is the easiest one to present.

CDCL brings a revolutionary change to DPLL so that DPLL is no longer a simple backtrack search procedure, because the level reset by *insertNewClause* may not be *level* - 1 and DPLL will backjumping to a higher level, thus avoiding unnecessary search space. If the new clause learned contains a single literal, the new level is set to 0. DPLL with CDCL can learn from failure and back jumping to the level where the source of failure is originated. Several new techniques are proposed based on CDCL, including generating a resolution proof when the input

clauses are unsatisfiable, or randomly restarting the search.

#### 4.2.3 Unsatisfiable Cores

The existence of a resolution proof for a set of unsatisfiable clauses is guaranteed by the completeness of resolution. How to find a resolution proof is a difficult job. It has been shown that CDCL as practiced in today's SAT solvers corresponds to a proof system exponentially more powerful than resolution. Specifically, each learning step is in fact a sequence of resolution steps, of which the learned clause is the final resolvent; conversely, a resolution proof can be obtained by regarding the learning procedure as a guided resolution procedure. In Example 4.2.6, we have seen how the empty clause is generated. The resolution proof leading to the empty clause can be displayed as follows:

$c_1$	$(x_1 \mid x_2)$	assumed
$c_3$	$(\overline{x_2} \mid \overline{x_3} \mid x_4)$	assumed
$c_5$	$(\overline{x_3} \mid \overline{x_4})$	assumed
$c_7$	$(\overline{x_2} \mid x_3)$	assumed
$c_8$	$(\overline{x_1} \mid x_3)$	assumed
$c_9$	$(\overline{x_1} \mid x_2 \mid \overline{x_3})$	assumed
$c_{10}$	$(\overline{x_2} \mid x_4)$	resolvent of $c_7, c_3$
$c_{11}$	$(\overline{x_1} \mid x_2)$	resolvent of $c_8, c_9$
$c_{12}$	$(x_2)$	resolvent of $c_{11}, c_1$
$c_{13}$	$(\overline{x_2} \mid \overline{x_3})$	resolvent of $c_5, c_{10}$
$c_{14}$	$(\overline{x_2})$	resolvent of $c_{13}, c_7$
$c_{15}$	$( )$	resolvent of $c_{12}, c_{14}$

This example shows the general idea of obtaining a resolution proof when *DPLL* finds a set of clauses unsatisfiable. During the search inside *DPLL*, conflicting clauses are generated due to decision literals. The procedure *conflictAnalysis* produces a new clause to show why the previous decisions were wrong and we go up in the decision tree according to this new clause. *DPLL* returns  $\perp$  when it finally finds a conflicting clause at level 0. In this case, a resolution proof is generated by *conflictAnalysis* and this proof uses either the input clauses or the clauses generated previously by *conflictAnalysis*.

Given an unsatisfiable set  $C$  of clauses, we can use all the resolution steps inside the procedure *conflictAnalysis* to construct a resolution proof. Obviously, such a proof provides an evidence for the unsatisfiability of  $C$ . The clauses used as the premises of the proof are a subset of the clauses of  $C$  and are called the unsatisfiable *core* of the proof. Note that a formula typically does not have a unique unsatisfiable core. Any unsatisfiable subset of  $C$  is an unsatisfiable core. Resolution

proofs and unsatisfiable cores have applications in hardware verification.

An unsatisfiable core is *minimal* if removing any clause from the core makes the remaining clauses in the core satisfiable. We may use this definition to check if every clause in the core is necessary, and this is obviously a difficult job when the core is huge. The core of the above example contains every input clauses except  $c_2$ ,  $c_4$ , and  $c_6$ ; it is a minimal core.

#### 4.2.4 Random Restart

Suppose we are looking for a model of  $C$  by  $DPLL(C)$  and the first decision literal is  $A$ . Unfortunately,  $C \wedge A$  is a hard unsatisfiable instance and your  $DPLL$  stuck there without success. Randomly restarting  $DPLL$  may be your only option. By “restart”, we mean that  $DPLL$  throws away all the previous decision literals (this can be easily done by  $level := 0$ ;  $undo()$  in  $DPLL$ ). By “random”, we mean that when you restart the search, the first decision literal will most likely be a different literal from  $A$ . Intuitively, randomly restarting means there is a chance of avoiding bad luck and getting luckier with guessing the right literal assignments that would lead to a quick solution.

Without CDCL, random restart makes  $DPLL$  incomplete. With CDCL, the input and generated clauses keep  $DPLL$  from choosing the same sequence of decision literals, and cannot generate the same clause again. Since the number of possible clauses is finite,  $DPLL$  cannot run forever.

The same intuition suggests random restart should be much more effective when the problem instance is in fact satisfiable than if it is not. Experimental results showed that random restart also helps when the input clauses are unsatisfiable. If restart is frequent, this assumes a deviation from standard practice, CDCL can be as powerful as general resolution, while  $DPLL$  without restart has been known to correspond to the weaker tree-like resolution. The theoretical justification for the speedup is that a restart allows the search to benefit from the knowledge gained about persistently troublesome conflicting variables sooner than backjumping would otherwise allow the partial assignment to be similarly reset. In effect, restarts may allow the discovery of shorter proofs of unsatisfiability.

To implement the restart strategy, when  $DPLL$  has found a certain number of failures (empty clauses) without success (a model), it is time for a restart. Today’s SAT solvers often use the following restart policy: Let  $r_i = r_0 \gamma^{i-1}$  be the number of failures for the  $i^{th}$  restart, where  $r_0$  is a given integer, and  $1 \leq \gamma < 2$ . If  $r_0 = 300$  and  $\gamma = 1.2$ , it means the first restart happens when  $DPLL$  has 300 failures, the second restart happens after another 360 failures, and so on. If  $\gamma = 1$ , it means  $DPLL$  restarts after a fix number of failures.

Given the common use of restarts in today's clause learning SAT solvers, the task of choosing a good restart policy appears appealing. While experimental results showed that no restart policy is better than others for a wide range of problems, a clause learning SAT solver could benefit substantially from a carefully designed restart policy. Specifically, experiments shows that nontrivial restart policies did significantly better than if restarts were disabled, and exhibited considerably different performance among themselves. This provides motivation for the design of better restart policies, particularly dynamic ones based on problem types and search statistics.

#### 4.2.5 Branching Heuristics for DPLL

Using CDCL or not, every implementation of DPLL needs a function for selecting literals for case-split. We assume that when a literal is chosen for case-split, the literal will be assigned true in the first branch of DPLL. Prior to the development of the CDCL techniques, branching heuristics were the primary method used to reduce the size of the search space. It seems likely, therefore, that the role of branching heuristics is likely to change significantly for algorithms that prune the search space using clause learning.

It is conventional wisdom that it is advantageous to assign the most tightly constrained variables, i.e., variables that occur in a large number of clauses. On representative of such a selection strategy is known as the MOMS rule, which branches on a literal which has the *Maximum Occurrences in clauses of Minimum Size*. If the clauses contain any binary clauses, then the MOMS rule will choose a literal in a binary clause. By assigning true to this literal, it is likely that a maximal number of binary clauses will become satisfied. On the other hand, if the primary criterion for the selection of a branch literal is to pick one that would enable a cascade of unit propagation (the result of such a cascade is a smaller subproblem), we would assign false to the literal chosen by the MOMS rule, in the presence of binary clauses, assigning this literal false will create likely a maximal number of unit clauses from all the binary clauses. MOMS provides a rough but easily computed approximation to the number of unit propagation that a particular variable assignment might cause.

Alternatively, one can call BCP multiple times on a set of promising variables in turn, and compute the exact number of unit clauses that would be caused by a branching choice. Each chosen variable is separately fixed to true and to false and BCP is executed for each choice. The precise number of unit propagation caused is then used to evaluate possible branching choices. Unlike the MOMS heuristic, this rule is obviously exact in its attempt to judge the number of unit clauses caused by a potential variable assignment. Unfortunately, it is also considerably more expensive

to compute because of the multiple executions of BCP and undoing them. It has shown that, using MOMS to choose a small number of promising candidates, each of which is then evaluated exactly using BCP, it outperforms other heuristics on randomly generated problems.

Another strategy is to branch on variables that are likely to be *backbone literals*. A *backbone literal* is one that must be true in all models of the input clauses. The likelihood that any particular literal is a backbone literal is approximated by counting the appearances of that literal in the satisfied clauses during the execution of DPLL. This heuristic outperforms those discussed in the previous paragraphs on many examples.

The development of CDCL techniques enabled solvers to attack more structured, realistic problems. There are no formal studies comparing the previously discussed heuristics on structured problems when CDCL techniques are used. CDCL techniques create many new clauses and make the occurrence counts of each literal difficult or less significant. Branching techniques and learning are deeply related, and the addition of learning to a DPLL implementation will have a significant effect on the effectiveness of any of these branching strategies. As new clauses are learned, the number of unit clauses an assignment will cause can be expected to vary; the reverse is also true in that the choice of branch literal can affect which clauses the algorithm learns.

Branching heuristics that are designed to function well in the context of clause learning generally try to branch on literals among the new clauses which have been learned recently. This tends to allow the execution of DPLL to keep “making progress” on a single section of the search space as opposed to moving from one area to another; an additional benefit is that existing learned clauses tend to remain relevant, avoiding the inefficiencies associated with losing the information present in learned clauses that become irrelevant and are discarded.

One of the popular branch heuristics for DPLL with clause learning is called “dynamic largest individual sum (DLIS) heuristic” and it behaves similarly as the MOMS rule. At each decision point, it chooses the assignment that satisfies the most unsatisfied clauses. Formally, let  $p_x$  be the number of unresolved clauses containing  $x$  and  $n_x$  be the number of unresolved clauses containing  $\bar{x}$ . Moreover, let  $x$  be the variable for which  $p_x$  is maximal, and let  $y$  be variable for which  $n_y$  is maximal. If  $p_x > n_y$ , choose  $x$  as the next branching literal; otherwise, choose  $\bar{y}$ . The disadvantage of this strategy is that the computational overhead is high: the algorithm needs to visit all clauses that contain a literal that has been set to true in order to update the values  $p_x$  and  $n_x$  for all variables contained in these clauses. Moreover, the process needs to be reversed upon backtracking.

A heuristic commonly used in contemporary SAT solvers favors literals in recently added conflict clauses. Each literal is associated with a count, which is initialized with the number of times the literal occurs in the clause set. When a learned clause is added, the count associated with each literal in the clause is incremented. Periodically, all counters divided by a constant greater than 1, resulting in a decay causing a bias towards branching on variables that appear in recently learned clauses. At each decision point, the solver then chooses the unassigned literal with the highest count (where ties are broken randomly by default). This approach, known as the variable state independent decaying sum (VSIDS) heuristics, was first implemented in `zchaff`. `Zchaff` maintains a list of unassigned literals sorted by count. This list is only updated when learned clauses are added, resulting in a very low overhead. Decisions can be made in constant time. The heuristic used in `berkmin` builds on this idea but responds more dynamically to recently learned clauses. The `berkmin` heuristic prefers to branch on literals that are unassigned in the most recently learned clause that is not yet satisfied.

The emphasis on variables that are involved in recent conflicting clauses leads to a locality based search, effectively focusing on a subspace. The subspaces induced by this decision strategy tend to coalesce, resulting in more opportunities for resolution of conflicting clauses, since most of the variables are common. Representing count using integer variables leads to a large number of ties. Solver `minisat` avoids this problem by using a floating point number to represent the weight. Another possible (but significantly more complex) strategy is to concentrate only on unresolved conflicting clauses by maintaining a stack of conflicting clauses.

When restart is supported, the branching heuristic should be a combination of random selection and other heuristics, so that DPLL will select a different branch literal after each restart. All told, there are many competing branching heuristics for satisfiability solvers, and there is still much to be done in evaluating their relative effectiveness with clause learning on realistic, structured problems.

### 4.3 Use of SAT Solvers

There exist excellent implementations of *DPLL* and you can find these SAT solvers at the International SAT competition web page ([satcompetition.org](http://satcompetition.org)). Most of them are available, free of charge, such as `minisat`, `glucose`, `lingeling`, `maple_LCM`, to name a few. These SAT solvers are also called general-purpose model generators, because they accept CNF as input and many problems can be specified in CNF as SAT is NP-complete.

If you have a hard search problem to solve, either solving a puzzle or finding

a solution under certain constraints, you may either write a special-purpose program or describe your problem in CNF and use one of these general-purpose model generators. While both special-purpose tools and general-purpose model generators rely on exhaustive search, there are fundamental differences. For the latter, every problem has a uniform internal representation (i.e., clauses for SAT solvers). This uniform representation may introduce redundancies and inefficiencies. However, since a single search engine is used for all the problems, any improvement to the search engine is significant. Moreover, we have an accumulated knowledge of three decades on how to make such a uniform search engine efficient. Using general-purpose model generators to solve hard combinatorial problems deserves the attention for at least two reasons.

- It is much easier to specify a problem in CNF than to write a special program to solve it. Similarly, fine-tuning a specification is much easier than fine-tuning a special program.
- General-purpose model generators can provide competitive and complementary results for certain combinatorial problems. Several examples will be provided in the following.

#### 4.3.1 Specify SAT Instances in DIMACS Format

Today SAT solvers use the format suggested by DIMACS (the Center for Discrete Mathematics and Theoretical Computer Science) many years ago for CNF formulas. The DIMACS format uses a text file: If the first character of a line is “c”, it means a comment. The CNF starts with a line “p cnf n m”, where  $n$  is the number of propositional variables and  $m$  is the number of clauses. Most SAT solvers only require  $m$  be an upper bound for the number of clauses. Each literal is represented by an integer: positive integers are positive literals and negative literals are negative literals. A clause is a sequence of integers ending with 0.

**Example 4.3.1.** To represent  $C = \{(x_1 \mid x_2 \mid \bar{x}_3), (\bar{x}_1 \mid x_2), (\bar{x}_2 \mid \bar{x}_3), (x_2 \mid x_3)\}$  in the DIMACS format, the text file will look as follows:

```
c a tiny example
c
p cnf 3 4
1 2 -3 0
-1 2 0
```

	2		5		1		9	
8			2		3			6
	3			6				7
		1				6		
5	4						1	9
		2				7		
	9			3				8
2			8		4			7
	1		9		7			6

Unsolved Sudoku

4	2	6	5	7	1	3	9	8
8	5	7	2	9	3	1	4	6
1	3	9	4	6	8	2	7	5
9	7	1	3	8	5	6	2	4
5	4	3	7	2	6	8	1	9
6	8	2	1	4	9	7	5	3
7	9	4	6	3	2	5	8	1
2	6	5	8	1	4	9	3	7
3	1	8	9	5	7	4	6	2

Solved Sudoku

Figure 4.3.1: A typical example of Sudoku puzzle and its solution.

```
-2 -3 0
2 3 0
```

To use today's SAT solvers, you need to prepare your SAT instances in the DIMACS format. This is done typically by writing a small program called *encoder*. Once a model is found by a SAT solver, you may need another program, called *decoder*, which translates the model into desired format. For example, the encoder for a Sudoku puzzle will generate a formula  $A$  in CNF from the puzzle and the decoder will take the model of  $A$  and generate the solution of the Sudoku.

#### 4.3.2 Sudoku Puzzle

The standard Sudoku puzzle is to fill a  $9 \times 9$  board with the digits 1 through 9, such that each row, each column, and each of nine  $3 \times 3$  blocks contain all the digits from 1 to 9. A typical example of Sudoku puzzle and its solution is shown in Figure 4.3.1.

To solve this problem by a SAT solver, we need to specify the puzzle in CNF. At first, we need to define the propositional variables. Let  $p_{i,j,k}$  be the propositional variable such that  $p_{i,j,k} = 1$  iff digit  $k$  is at row  $i$  and column  $j$  of the Sudoku board. Since  $1 \leq i, j, k \leq 9$ , there are  $9^3 = 729$  variables.

We may encode  $p_{i,j,k}$  as  $81(i-1) + 9(j-1) + k$ , then  $p_{1,1,1} = 1$  and  $p_{9,9,9} = 729$ . Another encoding is  $p_{i,j,k} = 100i + 10j + k$ , then  $p_{1,1,1} = 111$  and  $p_{9,9,9} = 999$ . The latter wasted some variables, but it allows us to see clearly the values of  $i, j$  and  $k$  from the integer. We will use the latter, since this is a very easy SAT problem, a tiny waste does not hurt. The first three lines of the CNF file may read as

```
c Sudoku puzzle
p cnf 999 1000000
```

where 1000000 is an estimate for the number of clauses from the constraints on the puzzle.

The general constraints of the puzzle can be stated as follows.

- Each cell contains exactly one digit of any value.
- Each row contains every digit exactly once, i.e., there are no duplicate copies of a digit in a row.
- Each column contains every digit exactly once.
- Each  $3 \times 3$  grid contains every digit exactly once.

A solution to a given board will say which digit should be placed in which cell.

The first constraint above can be divided into two parts: (1): Each cell contains at least one digit; and (2): Each cell contains at most one digit. If every cell in a row contributes one digit and each digit appears at most once, then it implies that every digit appears exactly once in this row. The same reasoning applies to columns and grids, so the next three constraints can be simplified by dropping the constraints like “every row contains every digit at least once”.

1. Each cell contains at least one digit.
2. Each cell contains at most one digit.
3. Each row contains every digit at most once.
4. Each column contains every digit at most once.
5. Each  $3 \times 3$  grid contains every digit at most once.

Now it is easy to convert the above constraints into clauses:

1. For  $1 \leq i, j \leq 9$ ,  $(p_{i,j,1} \mid p_{i,j,2} \mid \cdots \mid p_{i,j,9})$ .
2. For  $1 \leq i, j \leq 9, 1 \leq k < c \leq 9$ ,  $(\overline{p_{i,j,k}} \mid \overline{p_{i,j,c}})$ .
3. For  $1 \leq i, k \leq 9, 1 \leq j < b \leq 9$ ,  $(\overline{p_{i,j,k}} \mid \overline{p_{i,b,k}})$ .
4. For  $1 \leq j, k \leq 9, 1 \leq i < a \leq 9$ ,  $(\overline{p_{i,j,k}} \mid \overline{p_{a,j,k}})$ .
5. For  $0 \leq a, b \leq 2, 1 \leq i < i' \leq 3, 1 \leq j < j' \leq 3, 1 \leq k \leq 9$ ,  $(\overline{p_{3a+i,3b+j,k}} \mid \overline{p_{3a+i',3b+j',k}})$ .

Finally, we need to add the initial configuration of the board as clauses. If digit  $c$  is placed at row  $a$  and column  $b$ , we create a unit clause  $(p_{a,b,c})$  for it. Each Sudoku puzzle should have at least 17 digits placed initially on the board.

You may write a small program to convert the above clauses in the DIMACS format and use one of the SAT solvers available on the internet to solve the Sudoku puzzle. A SAT solver takes no time to find the solution. If you like to design a Sudoku puzzle for your friend, you may use a SAT solver to check if your puzzle has no solutions or has more than one solution.

### 4.3.3 Latin Square Problems

**Definition 4.3.2.** *Given a set  $S$  of  $n$  elements, a Latin square  $L$  over  $S$  is an  $n \times n$  matrix such that each row and each column of  $L$  is a permutation of  $S$ . The size of  $S$ ,  $|S|$ , is called the order of  $L$ .*

A Sudoku puzzle is a Latin square of order 9 over  $S = \{1, 2, \dots, 9\}$ . People usually use  $S = Z_n = \{0, 1, \dots, n-1\}$  for a Latin square of order  $n$ .

Every Latin square can be regarded as a definition of the function  $* : S \times S \rightarrow S$  such that  $(a * b) = c$  iff the cell at row  $a$  and column  $b$  contains  $c$ . The pair  $(S, *)$  is a special structure in abstract algebra, called *quasigroup*.

Latin squares of special properties have been extensively studied in design theory and SAT solvers have been used to solve many open problems there. For example, SAT solvers have been used successfully to show the non-existence of Latin squares of orders  $n = \{9, 10, 12, 13, 14, 15\}$  satisfying the constraint  $((y * x) * y) * y = x$  for  $x, y \in Z_n$ . It remains open if there exists such a Latin square of order 18.

To represent the constraint  $((y * x) * y) * y = x$  in CNF, we need to “flatten” the constraint by introducing new variables: the constraint  $((y * x) * y) * y = x$  is equivalent to

$$(y * x = z) \wedge (z * y = w) \rightarrow w * y = x.$$

If we replace  $y * x = z$  by  $p_{y,x,z}$ , etc., we have the following clause:

$$(\overline{p_{y,x,z}} \mid \overline{p_{z,y,w}} \mid p_{w,y,x})$$

for  $x, y, z, w \in Z_n$ . Since a Latin square can be specified in CNF as we did for a Sudoku puzzle, now searching for a Latin square of order  $n$  satisfying  $((y * x) * y) * y = x$  is the same as searching a model of the following clauses:

1. For  $0 \leq i, j < n$ ,  $(p_{i,j,0} \mid p_{i,j,2} \mid \dots \mid p_{i,j,n-1})$ .
2. For  $0 \leq i, j < n, 0 \leq k < c < n$ ,  $(\overline{p_{i,j,k}} \mid \overline{p_{i,j,c}})$ .

3. For  $0 \leq i, k < n, 0 \leq j < b < n, (\overline{p_{i,j,k}} \mid \overline{p_{i,b,k}})$ .
4. For  $0 \leq j, k < n, 0 \leq i < a < n, (\overline{p_{i,j,k}} \mid \overline{p_{a,j,k}})$ .
5. For  $0 \leq x, y, z, w < n, (\overline{p_{y,x,z}} \mid \overline{p_{z,y,w}} \mid p_{w,y,x})$ .

The first four sets of clauses are copied from those for the Sudoku puzzle, and the last set is the specification of  $((y * x) * y) * y = x$  in propositional logic.

#### 4.3.4 Graph Problems

Many application problems can be modeled as a graph problem, such as representing the exam scheduling problem as graph coloring problem. In the study of NP-completeness theory, we learned to reduce SAT to some graph problems. That is, we use the algorithms for the graph problems to solve SAT. Now, since we have a SAT solver, the reduction is reversed: graph problems are reduced to SAT, i.e., we use SAT solvers to solve graph problems. The challenge is to ensure that the reduction produces a CNF of minimal size, or at least, the reduction is polynomial.

**Definition 4.3.3. (graph coloring)** *Given an undirected simple graph  $G = (V, E)$  and a positive integer  $K$ , each number  $k, 1 \leq k \leq K$ , represents a color. The graph coloring problem is to assign a color  $k$ , to each  $v \in V$  such that the two ends of every edge of  $E$  receive different colors.*

Suppose  $V = \{v_1, v_2, \dots, v_n\}$ . We introduce  $kn$  propositional variables  $p_{i,k}$ ,  $1 \leq i \leq n, 1 \leq k \leq K$ , with the meaning that  $p_{i,k} = 1$  iff  $v_i$  receives color  $k$ . The constraints for the graph coloring are

1. Every vertex receives exactly one color.
2. The two ends of every edge receive different color.

These constraints translate into the following clauses:

1. For  $1 \leq i \leq |V|, (p_{i,1} \mid p_{i,2} \mid \dots \mid p_{i,K})$ .
2. For  $1 \leq i \leq |V|, 1 \leq j < k \leq K, (\overline{p_{i,j}} \mid \overline{p_{i,k}})$ .
3. For  $(v_i, v_j) \in E, 1 \leq k \leq K, (\overline{p_{i,k}} \mid \overline{p_{j,k}})$ .

The first two sets of clauses specify “at least one” and “at most one” for the constraint “exactly one”, as we did in the Sudoku puzzle. The size of the first set of clauses is  $K|V|$  (counting the number of literals); the size of the second set is

$(K - 1)K/2|V|$ ; the size of the third set is  $2K|E|$ . Strictly speaking, this is not a polynomial reduction if  $K$  is not a constant. If  $K$  is large, we may assume  $K = s * t$  for some integers  $s, t > 1$  and introduce  $s$  new variables  $q_{i,j}$ ,  $0 \leq j < s$ , such that  $q_{i,j}$  is true if node  $i$  receives color  $c$  for  $j * t < c \leq (j + 1) * t$ , and use these variables to reduce the size of the second set of clauses.

**Definition 4.3.4. (clique)** *Given an undirected simple graph  $G = (V, E)$  and a positive integer  $k$ , the clique problem is to find a subset  $S \subseteq V$  such that  $|S| = k$  and for any two vertices  $x, y$  of  $S$ ,  $(x, y) \in E$ .*

Suppose  $V = \{v_1, v_2, \dots, v_n\}$ . As a first try, let us define  $n$  propositional variables,  $x_i$ ,  $1 \leq i \leq n$ , with the understanding that  $x_i = 1$  iff  $v_i \in S$ . It is easy to specify if  $v_i, v_j \in S$ , then  $(v_i, v_j) \in E$ . That is, if  $(v_i, v_j) \notin E$ , then  $(\overline{x_i} \mid \overline{x_j})$ . However, it will be inefficient to specify that  $|S| = k$ . To ensure that  $|S| = k$ , we consider any subset  $T \subset V$  with  $n - k + 1$  vertices. One of them must be in  $S$ . To specify this constraint in logic, we have the following clauses:

For  $T \subset V$  and  $|T| = n - k + 1$ ,  $(\bigvee_{v_i \in T} x_i)$ .

This is a very large set of very long clauses, as evidently the reduction from the clique to CNF is not polynomial. For example, if  $n = 100$  and  $k = 5$ , then there are  $3,921,225 (= 100!/(4!96!))$  clauses of length 96. Obviously, any SAT solver will have difficulty to handle them.

An alternative solution is to consider that  $S$  is stored as a list  $L$  of  $k$  vertices. We introduce additionally  $k * n$  propositional variables,  $p_{i,j}$ ,  $1 \leq i \leq n, 1 \leq j \leq k$ , with the understanding that  $p_{i,j} = 1$  iff  $v_i$  is at position  $j$  of  $L$ . We then have the following constraints:

1. There is exactly one vertex at each position of  $L$ .
2. If  $v_i \in L$ , then  $v_i \in S$ .
3. If  $v_i \in S$ , then  $v_i \in L$  (optional).

The above constraints, plus the old constraint, can be specified in CNF as follows:

1. For  $1 \leq j \leq k$ ,  $(p_{1,j} \mid p_{2,j} \mid \dots \mid p_{n,j})$ .
2. For  $1 \leq j \leq k, 1 \leq i < i' \leq n$ ,  $(\overline{p_{i,j}} \mid \overline{p_{i',j}})$ .
3. For  $1 \leq j \leq k, 1 \leq i \leq n$ ,  $(\overline{p_{i,j}} \mid x_i)$ .
4. For  $1 \leq i \leq n$ ,  $(\overline{x_i} \mid p_{i,1} \mid p_{i,2} \mid \dots \mid p_{i,k})$ .

5. For  $1 \leq i, j \leq n$ , if  $(v_i, v_j) \notin E$ ,  $(\bar{x}_i \mid \bar{x}_j)$ .

Without loss of generality, we may assume that  $L$  is a sorted list by adding more clauses, to reduce some symmetric solutions.

If  $n = 100, k = 5$ , (1) produces 5 clauses of length 100; (2) produces 24,750 ( $=5 * 100 * 99/2$ ) clauses of length 2; (3) produces 500 clauses of length 2; and (4) produces 100 clauses of length 6, for a total of 25,355 clauses. This is a substantial reduction from 3,921,225 at the cost of introducing 500 new variables.

#### 4.4 Local Search Methods and MaxSAT

Given a set  $C$  of clauses and a full interpretation  $\sigma$ , we may count how many clauses are false under  $\sigma$ . Let

$$f(C, \sigma) = |\{c : c \in C, \sigma(c) = 0\}|,$$

then  $0 \leq f(C, \sigma) \leq |C|$ . When  $f(C, \sigma) = 0$ ,  $\sigma$  is a model of  $C$ . The maximum satisfiability problem *MAX-SAT* is thus the optimization problem where  $f$  serves as the objective function: find  $\sigma$  such that  $f(C, \sigma)$  is minimal (or  $|C| - f(C, \sigma)$  is maximal).

##### 4.4.1 Local Search Methods for SAT

Local search methods are widely used for solving optimization problems. In these methods, we first define a search space and for each point in the search space, we define the neighbors of the point. Starting with any point, we replace repeatedly the current point by one of its better neighbors under the objective function, until the point fits the search criteria or we run out of time.

For SAT or MAX-SAT, the search space of common local search methods is the set of full interpretations. If  $C$  contains  $n$  variables, then the search space size is  $O(2^n)$ , which is smaller than  $O(3^n)$ , the search space size for exhaustive search procedures like *DPLL*, where all partial interpretations are considered. For convenience, we assume that a full interpretation  $\sigma$  is represented by a set of  $n$  literals where each variable appears exactly once (or equivalently, a full minterm). The *neighbor* of  $\sigma$  is obtained by flipping the truth value of one variable in  $\sigma$ . Thus, every  $\sigma$  has  $n$  neighbors:

$$neighbors(\sigma) = \{\sigma - \{A\} \cup \{\neg A\} \mid A \in \sigma\}$$

Selman et al. proposed a greedy local search procedure called *GSAT* for SAT. The running time of *GSAT* is controlled by two parameters: *MaxTries* is the

maximal number of the starting points; and *MaxFlips* is the maximal number of flips allowed for any starting point.

**Algorithm 4.4.1.** *GSAT*( $C$ ) takes as input a set  $C$  of clauses and returns a model of  $C$  if it finds one; returns “unknown” if it terminates without finding a model.

```

proc GSAT( $C, \sigma$ )
1   for  $i := 1$  to MaxTries do
2        $\sigma :=$  a random full interpretation of  $C$ ; // a starting point
3       for  $j := 1$  to MaxFlips do
4           if  $\sigma(C) = 1$  return  $\sigma$ ; // a model is found
5           pick  $A \in \sigma$ ; // pick a literal of  $\sigma$  for flipping
6            $\sigma' = \sigma - \{A\} \cup \{\neg A\}$ ;
7            $\sigma = \sigma'$ ;
8   return “unknown”;

```

**Example 4.4.2.** Let  $C = \{c_1 : (x_1 \mid x_3 \mid x_4), c_2 : (\bar{x}_1 \mid \bar{x}_4), c_3 : (\bar{x}_2 \mid \bar{x}_4), c_4 : (\bar{x}_1 \mid x_2 \mid \bar{x}_3), c_5 : (x_2 \mid \bar{x}_4)\}$ . If the starting interpretation is  $\sigma_0 = \{x_1, x_2, x_3, x_4\}$ ,  $f(C, \sigma_0) = 2$  as  $\sigma_0(c_2) = \sigma_0(c_3) = 0$ . If  $x_1$  is chosen for flipping, we obtain  $\sigma_1 = \{\bar{x}_1, x_2, x_3, x_4\}$  and  $f(C, \sigma_1) = 1$  since  $\sigma_1(c_3) = 0$ . If  $x_2$  is chosen next for flipping, then  $\sigma_2 = \{\bar{x}_1, \bar{x}_2, x_3, x_4\}$  and  $f(C, \sigma_2) = 1$  since  $\sigma_2(c_5) = 0$ . If  $x_4$  is chosen next for flipping, then  $\sigma_3 = \{\bar{x}_1, \bar{x}_2, x_3, \bar{x}_4\}$  and  $f(C, \sigma_3) = 0$ , i.e.,  $\sigma_3$  is a model of  $C$ .

There exist many variations of *GSAT*, as there are many ways to pick a literal at line 5. In general, we need to pick a literal  $A \in neighbors(\sigma)$  such that the resulting interpretation  $\sigma'$ , here  $\sigma' = \sigma - \{A\} \cup \{\neg A\}$ , satisfies  $(f(C, \sigma') \leq f(C, \sigma))$ . We may randomly pick one of such literals, or pick the best one such that  $f(C, \sigma')$  is minimal (the original *GSAT* uses this strategy).

When  $f(C, \sigma') = f(C, \sigma)$ , *GSAT* still chooses to move from  $\sigma$  to  $\sigma'$ , and such moves are called “sideways” moves; they are moves that do not increase or decrease the total number of unsatisfied clauses. To avoid *GSAT* going back and forth between two neighboring interpretations of the same  $f$  value, the idea of *tabu search* is needed, so that recent visited interpretations are not allowed to be the next interpretation. The search of *GSAT* typically begins with a rapid greedy descent towards a better truth assignment, followed by long sequences of “sideways” moves. Experiments indicate that on many formulas, *GSAT* spends most of its time on sideways moves.

We can also allow moves that increase the value of  $f$  (called *upward moves*) occasionally, especially in the beginning of the execution. For example, using the idea of “simulated annealing”, upward moves can be allowed with a probability, and

this probability decreases as the number of tries increases. To implement this idea, we replace line 7 by the following two lines:

```
7      x := random(0, 1); // x is a random number in [0, 1)
7'     if (x < p ∨ (f(C, σ') ≤ f(C, σ)) σ := σ';
```

Selman et al. proposed a variation of *GSAT*, called **walksat**, which interleaves the greedy moves of *GSAT* with random walk moves of a standard Metropolis search. **Walksat**'s strategy for picking a literal to flip is as follows:

1. Randomly pick a conflicting clause, say  $\alpha$ ;
2. If there is a literal  $B \in \alpha$  such that the flipping of  $B$  does not turn any currently satisfied clauses to unsatisfied, return  $B$  as the picked literal;
3. With probability  $p$ , randomly pick  $B \in \alpha$  and return  $B$ ;
4. With probability  $1 - p$ , pick  $B \in \alpha$  such that the flipping of  $B$  turns the least number of currently satisfied clauses to unsatisfied, return  $B$ .

Since the early 1990's, there has been active research on designing, understanding, and improving local search methods for SAT. In the framework of *GSAT*, various ideas are proposed and they include the following:

- When picking a literal, take the literal that was flipped longest ago for breaking ties.
- A weight is given to each clause, incrementing the weight of unsatisfied clauses by one for each interpretation. One of the literal occurring in more clauses of maximum weight is picked.
- A weight is given to each variable, which is increased in variables that are often flipped. The variable with minimum weight is chosen.

Selman et al. showed through experiments that *GSAT* as well as **walksat** substantially outperformed the best DPLL-based SAT solvers on some classes of formulas, including randomly generated formulas and graph coloring problems.

*GSAT* returns either a model or “unknown”, as it can never declare that “the input clauses are unsatisfiable”. This is a typical feature of local search methods as they are incomplete in the sense that they do not provide the guarantee that it will

eventually either report a satisfying assignment or declare that the given formula is unsatisfiable.

There have also been attempts at hybrid approaches that explore the combination of ideas from DPLL methods and local search techniques. There has also been work on formally analyzing local search methods, yielding some of the best  $o(2^n)$  time algorithms for SAT. For instance, the expected running time, ignoring polynomial factors, of a simple local search method with restarts after every  $3n$  “flips” has been shown to  $(k + 1)/k)^n$  for  $k$ CNF instances of  $n$  variables, where each clause contains at most  $k$  literals. When  $k = 3$ , the result yields a complexity of  $(4/3)^n$  for 3SAT. This result has been derandomized to yield a deterministic algorithm with complexity  $1.481^n$  (up to polynomial factors) to solve 3SAT.

#### 4.4.2 2SAT versus Max2SAT

Let Max- $k$ SAT denote the decision version of MaxSAT which takes a positive integer  $m$  and a set  $C$  of clauses, where each clause has no more than  $k$  literals, as input and returns true iff there exists an interpretation  $\sigma$  such that at least  $m$  clauses are satisfied under  $\sigma$ . It is well-known that 2SAT can be solved in polynomial time, while Max2SAT is NP-complete.

**Theorem 4.4.3.** *There is a linear time decision procedure for 2SAT.*

*Proof.* This sketch of the proof comes from (Aspvall, Plass, Tarjan, 1979). Let  $C$  be in 2CNF without pure literals. We construct a directed graph  $G = (V, E)$  for  $C$  as follows:  $V$  contains the two literals for each variable appearing in  $C$ . For each clause  $(A \mid B) \in C$ , we add two edges,  $(\bar{A}, B)$  and  $(\bar{B}, A)$ , into  $E$ . Obviously, the edges represent the implication relation among the literals, because  $(A \mid B) \equiv (\bar{A} \rightarrow B) \wedge (\bar{B} \rightarrow A)$ . Since the implication relation is transitive and  $(A \rightarrow B) \equiv (\bar{B} \rightarrow \bar{A})$ , the following three statements are equivalent for any pair of literals  $x$  and  $y$  in  $G$ :

1.  $C \models (x \rightarrow y)$ ;
2. there is a path from  $x$  to  $y$  in  $G$ ; and
3. there is a path from  $\bar{y}$  to  $\bar{x}$  in  $G$ .

We then run the algorithm for strongly connected components (SCC), which can be done in linear time, on  $G$ . It is easy to see that  $C \models (x \leftrightarrow y)$  for any two literals  $x$  and  $y$  in the same component, because  $x$  and  $y$  share a cycle in  $G$ . If both a variable, say  $p$ , and its complement,  $\bar{p}$ , appear in the same component, then  $C$  is unsatisfiable, because  $C \models (p \leftrightarrow \bar{p})$ . If no such cases exist, then  $C$  must be

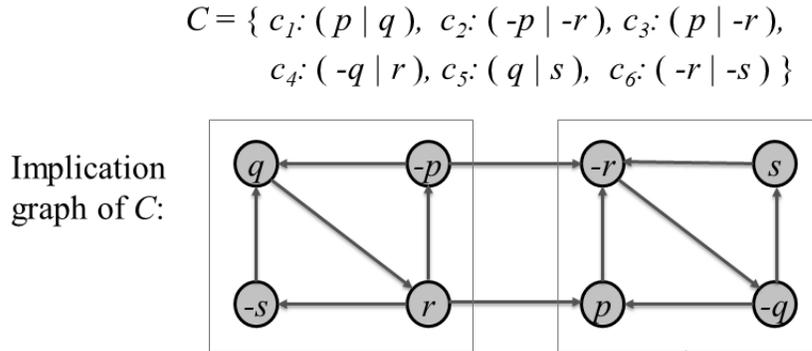


Figure 4.4.1: The implication graph for 2CNF  $C$  and the SCC.

satisfiable. A model can be obtained by assigning the same truth value to each node in the same component, following the topological order of the components induced by  $G$ : Always assign 0 to a component, unless one of the literals in the component had already received 1.  $\square$

**Example 4.4.4.** Given  $C = \{c_1 : (p \mid q), c_2 : (\bar{p} \mid \bar{r}), c_3 : (p \mid \bar{r}), c_4 : (\bar{q} \mid r), c_5 : (q \mid s), c_6 : (\bar{r} \mid \bar{s})\}$ , its implication graph is shown in Figure 4.4.1. There are two SCCs:  $S_1 = \{\bar{p}, q, r, \bar{s}\}$  and  $S_2 = \{p, \bar{q}, \bar{r}, s\}$ . We have to assign 0 to every node in  $S_1$  (and 1 to every node in  $S_2$ ) to obtain a model of  $C$ .

**Theorem 4.4.5.** *Max2SAT is NP-complete.*

*Proof.* This proof comes from (Garey, Johnson, and Stockmeyer, 1976). Max2SAT is in NP because we can guess non-deterministically an interpretation  $\sigma$  and check if there are at least  $K$  clauses are true under  $\sigma$  in polynomial time.

To show Max2SAT is NP-hard, we reduce 3SAT to Max2SAT. That is, if there exists an algorithm  $A$  to solve Max2SAT, we can construct algorithm  $B$  to solve 3SAT, and the computing times of  $A$  and  $B$  differ polynomially.

Given any instance  $C$  of 3SAT, we assume that every clause of  $C$  has exactly three literals. For every clause  $c_i = (l_1 \mid l_2 \mid l_3) \in C$ , where  $l_1, l_2$ , and  $l_3$  are arbitrary literals, we create 10 clauses in Max2SAT:

$$(l_1), (l_2), (l_3), (x_i), (\bar{l}_1 \mid \bar{l}_2), (\bar{l}_1 \mid \bar{l}_3), (\bar{l}_2 \mid \bar{l}_3), (l_1 \mid \bar{x}_i), (l_2 \mid \bar{x}_i), (l_3 \mid \bar{x}_i),$$

where  $x_i$  is a new variable. Let  $C'$  be the collection of these clauses created from the clauses in  $C$ , then  $|C'| = 10|C|$  and the following statements are true:

- If an interpretation satisfies  $c_i$ , then we can make seven of the ten clauses satisfied.

- If an interpretation falsifies  $c_i$ , then at most six of the ten clauses can be satisfied.

The above statements can be verified by a truth table on  $l_1, l_2, l_3$  and  $x_i$ :

$l_1$	$l_2$	$l_3$	$x_i$	$(\overline{l_1} \overline{l_2})$	$(\overline{l_1} \overline{l_3})$	$(\overline{l_2} \overline{l_3})$	$(l_1 \overline{x_i})$	$(l_2 \overline{x_i})$	$(l_3 \overline{x_i})$	sum
0	0	0	0	1	1	1	1	1	1	6
0	0	0	1	1	1	1	0	0	0	4
0	0	1	0	1	1	1	1	1	1	7
0	0	1	1	1	1	1	0	0	1	6
0	1	1	0	1	1	0	1	1	1	7
0	1	1	1	1	1	0	0	1	1	7
1	1	1	0	0	0	0	1	1	1	6
1	1	1	1	0	0	0	1	1	1	7

The first four columns serve both the truth values of  $\{l_1, l_2, l_3, x_i\}$  and the first four clauses of the ten clauses. The first two lines of the table shows that when  $(l_1 | l_2 | l_3)$  is false, the sum of the truth values of the ten clauses (the last column) is at most 6. The rest of the table shows the sums of the truth values of the ten clauses when one, or two, or three literals of  $\{l_1, l_2, l_3\}$  are true (some cases are ignored due to the symmetry of  $l_1, l_2$ , and  $l_3$ ). It is clear from the table that when  $(l_1 | l_2 | l_3)$  is true, we may choose the value of  $x_i$  to obtain 7 true clauses out of the ten clauses. Then  $C$  is satisfiable iff  $C'$  has an interpretation which satisfies  $K = 7|C|$  clauses of  $C'$ . If we have an algorithm to solve  $C'$  with  $K = 7|C|$ , then the output of the algorithm will tell if  $C$  is satisfiable or not.  $\square$

The above theorem shows that MaxSAT is a hard problem even for 2CNF.

## 4.5 Maximum Satisfiability

The algorithm *GSAT* can be used to solve the MaxSAT problem, which asks to find an interpretation for a set of clauses such that a maximal number of clauses are true. Since *GSAT* is incomplete, the answer found by *GSAT* is not guaranteed to be optimal. This is a typical feature of local search methods which always find local optimal solutions. *GSAT* can also be used to solve MaxSAT of general form: weighted MaxSAT and hybrid MaxSAT.

### 4.5.1 Weight MaxSAT and Hybrid MaxSAT

**Definition 4.5.1.** Given a set  $C$  of clauses and a function  $w : C \rightarrow \mathcal{N}$ ,  $w(c)$  is called the weight (or cost) of  $c \in C$ .  $(C, w)$  is called weighted CNF (WCNF).

Given  $(C, w)$ , the weighted MaxSAT problem is to find a full interpretation  $\sigma$  such that  $\sum_{c \in C, \sigma(c)=0} w(c)$  is minimal (or equivalently,  $\sum_{c \in C, \sigma(c)=1} w(c)$  is maximal). Any clause in  $C$  is called a weighted clause and denoted by  $(c; w(c))$ .

**Example 4.5.2.** Let  $C$  be a weighted MaxSAT instance, where  $C = \{c_1 : (x_1; 5), c_2 : (\overline{x_1} \mid x_2; 4), c_3 : (x_1 \mid \overline{x_2} \mid x_3; 3), c_4 : (\overline{x_1} \mid \overline{x_2}; 2), c_5 : (x_1 \mid x_2 \mid \overline{x_3}; 4), c_6 : (\overline{x_1} \mid x_3; 1), c_7 : (\overline{x_1} \mid \overline{x_2} \mid \overline{x_3}; 2)\}$ . Then  $\sigma = \{x_1, x_2, \overline{x_3}\}$  is the optimal interpretation for  $C$ , where  $c_4$  and  $c_6$  are false, with the total weights of falsified clauses being 3; the sum of weights of true clauses is 18.

**Definition 4.5.3.** Given a WCNF  $(C, w)$ , let  $(H, S)$  be a partition of the set  $C$ , where each clause  $c \in H$  has weight  $w(c) = \infty$ , a solution of  $(H, S)$  is a full interpretation  $\sigma$  such that every clause of  $H$  is true in  $\sigma$ . The hybrid MaxSAT problem is to find an optimal solution such that  $\sum_{c \in S, \sigma(c)=0} w(c)$  is minimal. The clauses in  $H$  are called hard clauses and the clauses in  $S$  are called soft clauses.

Hybrid MaxSAT is also called *weighted partial MaxSAT*. The definition of hybrid MaxSAT assumes that hybrid MaxSAT is a special case of weighted MaxSAT. Weighted MaxSAT can be also viewed as a special case of hybrid MaxSAT when  $H = \emptyset$ . In other words, the two definitions are equivalent in expressive power. In practice,  $\infty$  can be replaced by the sum of all weights in  $S$ . Separating  $H$  from  $S$  will allow us to handle  $H$  efficiently, using advanced techniques for SAT.

Since MaxSAT is a special case of weighted MaxSAT or hybrid MaxSAT, weighted MaxSAT or hybrid MaxSAT must be NP-complete, too. In fact, it is very easy to reduce NP-complete problems to weighted MaxSAT or hybrid MaxSAT.

The clique problem is a well-known NP-complete problem in graph theory. Given a graph  $G = (V, E)$ , let  $V = \{v_1, v_2, \dots, v_n\}$ . We define  $n$  propositional variables  $x_i$ ,  $1 \leq i \leq n$ , with the meaning that  $x_i = 1$  iff  $v_i$  is in the solution of a maximum clique. Then we may reduce the clique problem to the weighted MaxSAT problem by converting  $G = (V, E)$  into a formula of 2WCNF, which consists of the following two sets of weighted clauses.

- For  $1 \leq i \leq n$ , create  $(x_i; 1)$ .
- For  $v_i, v_j \in V$ ,  $i \neq j$  and  $(v_i, v_j) \notin E$ , create  $(\overline{x_i} \mid \overline{x_j}; n)$ .

The first set of clauses says that if a vertex is chosen in the solution, there is a reward of 1 (equivalently, if the vertex is not in the solution, there is a penalty of 1). The second set of clauses specifies the property of a clique: If two vertices are chosen in the solution and there is no edge between them, there is a penalty of  $|V|$ .

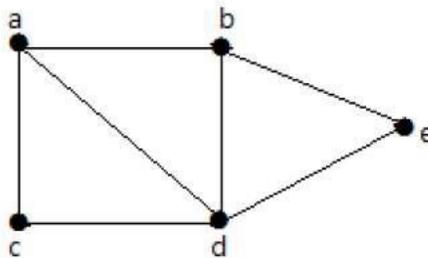


Figure 4.5.1: Finding a maximal clique in a graph.

**Example 4.5.4.** For the graph in Figure 4.5.1, the clauses are

$$C = \{(\overline{x_a} \mid \overline{x_e}; 5), (\overline{x_b} \mid \overline{x_c}; 5), (\overline{x_c} \mid \overline{x_e}; 5), (x_a; 1), (x_b; 1), (x_c; 1), (x_d; 1), (x_e; 1)\}.$$

One of the optimal solutions is  $\sigma = \{x_a \mapsto 1, x_b \mapsto 1, x_c \mapsto 0, x_d \mapsto 1, x_e \mapsto 0\}$ , under which, two clauses, i.e.,  $(x_c)$  and  $(x_e)$ , are false and the total penalty is 2.

Consider an interpretation  $\sigma$  in which one variable  $x_i$  is true; all the other variables are false. Given  $\sigma$ , all the clauses in the second set will be true and all the clauses in the first set will be false, except one. That is, the sum of weights of the falsified clauses is  $n - 1$  under  $\sigma$ . Thus, the total weights of all falsified clauses in any optimal solution cannot be greater than  $n - 1$ . In other words, all the clauses in the second set must be true in an optimal solution. This implies that there must be an edge between two vertices in the solution. The optimal solution picks a maximal number of  $x_i$ 's and thus gives us a maximal clique. It is natural to specify the second set of clauses as hard clauses in the hybrid MaxSAT format.

#### 4.5.2 The Branch-and-Bound Algorithm

To obtain an optimal solution of MaxSAT or hybrid MaxSAT, the traditional approach is to use the branch-and-bound algorithm. In these algorithms, we often cast MaxSAT as a minimization problem: find solution with minimum (optimal) cost. The branching part of the branch-and-bound algorithm is very much like *DPLL*, doing case-splits on selected literals. The bounding part of the algorithm checks if an estimation of the current solution is better or not than the solution found so far. If not better, the algorithm cuts this futile branch and backtracks. If we look for a minimal solution, the estimation produces a lower bound; for a maximal solution, the estimation produces an upper bound. For example, if we look for a minimal solution, and already have a solution of value 30, if the estimation says the lower bound is 32, then there is no point to continue, because the value of the solution from the current branch cannot be lower than 32.

Here we present the branch-and-bound algorithm for hybrid MaxSAT in the style of recursive *DPLL* (Algorithm 4.1.1).

**Algorithm 4.5.5.**  $HyMaxSATBB(H, F, \sigma, soln)$  takes as input a set  $H$  of hard clauses, a set  $F$  of weighted soft clauses, a partial interpretation  $\sigma$ , and the best solution  $soln$  found so far.  $H$  and  $\sigma$  do not share any variables. It looks for a full interpretation  $\sigma$  which is a model of  $H$  and the sum of the weights of the false clauses in  $F$  is minimal under  $\sigma$ ; this weight sum will be returned by the algorithm. Procedure  $pickLiteral(S)$  will pick a literal appearing in  $S$ .

```

proc  $HyMaxSATBB(H, F, \sigma, soln)$ 
1    $res := BCP(H)$ ;
2   if ( $res = \perp$ ) return  $\infty$ ; // hard causes violated
3   // Assume  $res = (U, S)$ 
4    $\sigma := \sigma \cup U$ ;
5    $F := simplifySoft(F, \sigma)$ ;
6   if  $soln \leq lowerBound(F)$  return  $soln$ ;
7    $A := pickLiteral(S \cup F)$ ;
8   if ( $A = nil$ ) return  $countFalseClauseWeight(F, \sigma)$ ;
9    $soln := min(soln, HyMaxSATBB(S \cup \{(A)\}, F, \sigma, soln))$ ;
10  return  $min(soln, HyMaxSATBB(S \cup \{(\bar{A})\}, F, \sigma, soln))$ ;

```

We have seen  $BCP(H)$  in Algorithm 4.1.1; the following four procedures are called inside  $HyMaxSATBB$ :

- $simplifySoft(F, \sigma)$  at line 5: Simplify soft clauses  $F$  using the partial interpretation  $\sigma$  and keep the total weights of false clauses unchanged under any interpretation. Details will follow.
- $lowerBound(F)$  at line 6: Estimate the minimal total weights of false clauses when  $\sigma$  is extended to a full interpretation in all possible ways. Details will be discussed later.
- $pickLiteral(S \cup F)$  at line 7: pick an unassigned literal for case-splits; returns  $nil$  if every literal has a truth value. Here is the place heuristics for branching are implemented.
- $countFalseClauseWeight(F, \sigma)$  at line 8: a simple procedure for computing the sum of all the weights of false clauses in  $F$ . Here,  $\sigma$  is a full interpretation and can be saved for late use.

For an efficient implementation of *HyMaxSATBB*, the first solution is always obtained by a local search algorithm such as *GSAT*. We can use all techniques used for *DPLL*, including iterative procedure, destructive data structure, and conflict-driven clause learning for hard clauses.

For soft clauses, the treatment is different. For example, unit resolution is a powerful simplification rule for hard clauses, but it is unsound for soft clauses. If  $C = \{(\bar{x}_1; 1), (x_1 \mid x_2; 1), (x_1 \mid \bar{x}_2; 1), (x_1 \mid x_3; 1), (x_1 \mid \bar{x}_3; 1)\}$ , the optimal solution has one false clause when  $x \mapsto 1$ . If we apply unit resolution to  $C$  to get  $C' = \{(\bar{x}_1; 1), (x_2; 1), (\bar{x}_2; 1), (x_3; 1), (\bar{x}_3; 1)\}$ , an optimal solution of  $C'$  will contain two false clauses. Of course, if a unit clause is hard, we may apply unit resolution to both hard and soft clauses.

#### 4.5.3 Simplification Rules and Lower Bound

**Definition 4.5.6.** *An inference rule or simplification rule applicable to a set  $C$  of weighted clauses is said to be sound, if the optimal solution of the resulting set  $C'$  of weighted clauses has the same value as the optimal solution of  $C$ .*

In other words, sound rules should preserve the value of optimal solutions. Most inference rules are unsound for weighted clauses because they generate extra clauses. One exception is when a hard clause can be generated from soft clauses. For example, given a soft clause  $(c; w)$ , if  $c$  becomes false,  $w$  will be added as a portion of the lower bound. If this addition of  $w$  causes the lower bound to exceed the value of the current solution, then we can generate the hard clause  $c$  to prevent  $c$  from being false. In the following, we focus on simplification rules, which replace some clauses by other clauses when we transform  $C$  to  $C'$ . Each simplification rule can be denoted by

$$C_1, C_2, \dots, C_m \implies C'_1, C'_2, \dots, C'_n$$

when we want to replace  $\{C_1, C_2, \dots, C_m\} \subseteq C$  by  $\{C'_1, C'_2, \dots, C'_n\}$ , where  $C_i, C'_j$  are weighted clauses.

In the following rules,  $\alpha$  and  $\beta$  denote arbitrary, possibly empty, disjunction of literals.

1. **zero elimination:** Weighted clauses with zero weights can be discarded.

$$(\alpha; 0) \implies 1$$

2. **tautology elimination:** Valid clauses are always true and can be discarded.

$$(p \mid \bar{p} \mid \alpha; w) \implies 1$$

3. **identical merge:** Identical clauses are merged into one.

$$(\alpha; w_1), (\alpha; w_2) \implies (\alpha; w_1 + w_2)$$

4. **unit clash:** Unit clauses with complement literals can be reduced to one.

$$(A; w_1), (\bar{A}; w_2) \implies (A; w_1 - w_2), (\perp; w_2)$$

where  $A$  is a literal and  $w_1 \geq w_2$ . For example, given  $(p; 2)$  and  $(\bar{p}; 3)$ , if  $p \mapsto 1$ , we get a cost of 3 from  $(\bar{p}; 3)$ ; if  $p \mapsto 0$ , we get a cost of 2 from  $(p; 2)$ . The same result can be obtained from  $(\bar{p}; 1)$  and  $(\perp; 2)$ ; the latter is added into the cost function.

5. **weighted resolution:**

$$(p \mid \alpha; w_1), (\bar{p} \mid \beta; w_2) \implies \begin{array}{l} (\alpha \mid \beta; u), (p \mid \alpha, w_1 - u), (\bar{p} \mid \beta; w_2 - u), \\ (p \mid \alpha \mid \bar{\beta}; u), (\bar{p} \mid \bar{\alpha} \mid \beta; u) \end{array}$$

where  $w_1 > 0$ ,  $w_2 > 0$ , and  $u = \min(w_1, w_2)$ . This rule is proposed by Larrosa et al. (2007). Since either  $w_1 - u = 0$  or  $w_2 - u = 0$ , the above rule cannot be used forever. For example, from  $(p \mid q; 2)$  and  $(\bar{p} \mid q; 3)$ ,  $\alpha = \beta = q$ ,  $u = \min(2, 3)$ , we obtain five clauses:  $(q; 2)$ ,  $(p \mid q; 0)$ ,  $(\bar{p} \mid q; 1)$ ,  $(p \mid q \mid \bar{q}; 2)$ , and  $(\bar{p} \mid \bar{q} \mid q; 2)$ . The second clause is removed by zero elimination, and the last two clauses by tautology elimination, we obtain two clauses:  $(q; 2)$ ,  $(\bar{p} \mid q; 1)$ . Note that  $\bar{\beta}$  and  $\bar{\alpha}$  in the last two clauses are not disjunction of literals in general. We need to convert  $(p \mid \alpha \mid \bar{\beta}; u)$  (or  $(\bar{p} \mid \bar{\alpha} \mid \beta; u)$ ) into a set of clauses, each having the weight  $u$ . In practice, this rule often applies to short clauses or when  $\alpha = \beta$ .

It is easy to check all of the above simplification rules are sound for weighted clauses, as they preserve optimal solutions.

To efficiently implement  $simplifySoft(F, \sigma)$ ,  $F$  will be simplified first by  $\sigma$ , as  $\sigma$  represents a set of hard unit clauses. Every true clause under  $\sigma$  will be removed from  $F$  and every false literal under  $\sigma$  will be removed from the clauses in  $F$ . After this step,  $F$  and  $\sigma$  do not share any variable. The second step is to apply the above simplification rules. Some rules like zero elimination and tautology elimination, are applied eagerly. Some rules like weighted resolution are applied selectively.

To implement  $lowerBound(F, soln)$ , at first, we assumed that  $F$  is simplified and we count the total weights of all false clauses. If the sum exceeds the current solution, we exit.

For non-false clauses of  $F$ , there exist various techniques to estimate the lower bound. One technique is to formulate  $F$  as an instance of Integer Linear Programming (ILP). ILP Solvers are common optimization tools for operation research. ILP solvers solve problems with linear constraints and linear objective function where some variables are integers. State-of-the-art ILP solvers are very powerful and effective and can attack MaxSAT directly. In practice, ILP solvers are effective on many standard optimization problems, e.g., vertex cover, but for problems where there are many boolean constraints ILP is not as effective.

To obtain a lower bound of ILP, the constraints on integer solutions are relaxed so that ILP becomes an instance of the Linear Programming (LP), which is known to have polynomial-time solutions.

If we view  $\top$  as 1 and  $\perp$  as 0, it is well known that a clause can be converted into a linear inequation, such as  $(\bar{x}_1 \mid x_2 \mid \bar{x}_3)$  becomes  $(1 - x_1) + x_2 + (1 - x_3) \geq 1$ , so that the clauses are satisfiable iff the set of inequations has a 0/1 solution for all variables.

For any clause  $c$ , let  $P(c)$  and  $N(c)$  be the sets of positive literals and negative literals in  $c$ , respectively. Let  $F = \{(c_i; w_i) : 1 \leq i \leq m\}$  over  $n$  variables, say  $\{x_1, x_2, \dots, x_n\}$ .  $F$  can be formulated as the following instance of ILP:

$$\begin{aligned} & \text{minimize } \sum_{1 \leq i \leq m} w_i(1 - y_i) && // \text{ minimize the weights of falsified clauses} \\ \text{subject to } & \sum_{x \in P(c_i)} x + \sum_{\bar{x} \in N(c_i)} (1 - x) \geq y_i, 1 \leq i \leq m, && // c_i \text{ is true iff } y_i = 1 \\ & y_i \in \{0, 1\}, 1 \leq i \leq m, && // \text{ every clause is falsified or satisfied} \\ & x_j \in \{0, 1\}, 1 \leq j \leq n. && // \text{ every variable is false or true} \\ & \text{where } y_i \text{ is a new variable for each clause } c_i \text{ in } F. \end{aligned}$$

Intuitively, the more  $y_i$ 's with  $y_i = 1$ , the less the value of  $\sum_{1 \leq i \leq m} w_i(1 - y_i)$ . And  $y_i = 1$  implies that  $\sum_{x \in P(c_i)} x + \sum_{\bar{x} \in N(c_i)} (1 - x) \geq 1$ , i.e.,  $c_i$  is true.

The above ILP instance can be relaxed to an instance  $L$  of LP:

$$\begin{aligned} & \text{minimize } \sum_{1 \leq i \leq m} w_i(1 - y_i) \\ \text{subject to } & \sum_{x \in P(c_i)} x + \sum_{\bar{x} \in N(c_i)} (1 - x) \geq y_i, 1 \leq i \leq m, \\ & 0 \leq y_i \leq 1, 1 \leq i \leq m, \text{ and} \\ & 0 \leq x_j \leq 1, 1 \leq j \leq n. \end{aligned}$$

If  $L$  has a solution, then  $\sum_{1 \leq i \leq m} w_i(1 - y_i)$  is a lower bound for the solutions of  $F$ . The ILP can be used to obtain an approximate solution for  $F$  when we round up/down the values of  $x_i$  from reals to integers.

Other techniques for lower bounds look for inconsistencies that force some soft clause to be falsified. Some suggested learning falsified soft clauses, some people

use the concepts of *clone*, *minibuckets*, or width-restricted BDD in relaxation. For instance, we may treat the soft clauses as if they were hard and then run BCP to locate falsified clauses. This can help us to find unsatisfiable subsets of clauses quickly, and use them to estimate the total weights of falsified clauses. These techniques can be effective on small combinatorial problems. Once the number of variables in  $F$  gets to 1,000 or more, these lower bound techniques become weak or too expensive. Strategies for turning on/off these techniques in *HyMaxSATBB* should be considered.

Besides the algorithm based on branch-and-bound, some MaxSAT solvers convert a MaxSAT instance into a sequence of SAT instances where each instance encodes a decision problem of the form, for different values of  $k$ , “is there an interpretation that falsifies soft clauses of total weights at most  $k$ ?”

For example, if we start with a very small value of  $k$ , the SAT instance will be unsatisfiable. By increasing  $k$  gradually, the first  $k$  when the SAT instance becomes satisfiable will be the value of an optimal solution. Of course, this linear strategy on  $k$  does not provide necessarily an effective MaxSAT solver.

There are many ongoing research in this direction. The focus of this approach is mainly doing two things:

1. Develop economic ways to encode the decision problem at each stage of the sequence.
2. Exploit information obtained from the SAT solver at each stage in the next stage.

Some MaxSAT solvers use innovation techniques to obtain more efficient ways to encode and solve the individual SAT decision problems. Some solvers use unsatisfiable cores to improve SAT solving efficiency. These solvers appear promising as they are effective on some large MaxSAT problems, especially those with many hard clauses.

#### 4.5.4 Use of Hybrid MaxSAT Solvers

Most real-world problems involve an optimization component and there are high demand for automated approaches to finding good solutions to computationally hard optimization problems. We have seen how to encode the clique problem easily and compactly into a MaxSAT instance. MaxSAT allows for compactly encoding various types of high-level finite-domain soft constraints due to Cook-Levin Theorem: Any NP constraint can be polynomially represented as a weighted SAT instance.

There exists an extension of DIMACS format for WCNF. To specify an instance of hybrid MaxSAT, the line starting with `p` will contain the keyword `wcnf`, followed by the number of variables, the number of clauses, and the label for hard clauses. That is,

```
p wcnf <# of variables> <# of clauses> <hard clause label>
```

The hard clause label is an integer large than any weight of soft clauses. Each clause is started with the weight value, followed by literals and ending with 0. For the small clique example of Example 4.5.4, the input file to a MaxSAT solver will look like:

```
c small clique example
c xa = 1, xb = 2, xc = 3, xd = 4, xe = 5
p wcnf 5 8 9
9 -1 -5 0
9 -2 -3 0
9 -3 -5 0
1 1 0
1 2 0
1 3 0
1 4 0
1 5 0
```

A clause is hard if the weight of the clause is the label for hard clauses.

The standard format for MaxSAT instances makes the annual evaluation of MaxSAT solvers possible. The latest website for MaxSAT evaluation can be found on the internet and the site for the year 2020 is given below:

```
https://maxsat-evaluations.github.io/2020/
```

Researchers on MaxSAT can assess the state of the art in the field of MaxSAT solvers, create a collection of publicly available MaxSAT benchmark instances, and use tens of solvers from various research groups.

There MaxSAT solvers are built on the successful techniques for SAT, evolve constantly in practical solver technology, and offer an alternative to traditional approaches, e.g., integer programming. MaxSAT solvers use propositional logic as the underlying declarative language and are especially suited for inherently “very Boolean” optimization problems. The performance of MaxSAT solvers have surpassed that of specialized algorithms on several problem domains:

- correlation clustering
- probabilistic inference
- maximum quartet consistency
- software package management
- fault localization
- reasoning over bionetwork
- optimal covering arrays
- treewidth computation
- Bayesian network structure learning
- causal discovery
- cutting planes for IP (integer programming)
- argumentation dynamics

The advances of contemporary SAT and MaxSAT solvers have transformed the way we think about NP complete problems. They have shown that, while these problems are still unmanageable in the worst case, many instances can be successfully tackled.

#### 4.6 Exercise Problems

1. Draw the four complete decision trees (trees of recursive calls) of *DPLL* on the following set of clauses,

$$\{(p \mid q), (p \mid \bar{r}), (\bar{p} \mid r), (\bar{p} \mid q \mid \bar{r}), (\bar{q} \mid \bar{r})\}$$

using the orders of  $p, q, r$  and  $r, q, p$ , respectively, for splitting (either value can be chosen first) and with or without BCP.

2. Draw the complete decision tree of *DPLL* on the following set of clauses,

$$\{(p \mid q \mid s), (p \mid \bar{r} \mid \bar{s}), (q \mid \bar{r} \mid s), (\bar{p} \mid q \mid \bar{r}), (\bar{q} \mid \bar{r} \mid \bar{s})\}$$

using the orders of  $p, q, r, s$  for splitting (the value 1 is used first for each splitting variable) and report the truth values of variables at various levels. How many models are found by *DPLL*? What guiding paths are found when each model is found?

3. If we feed the clause set in the previous problem to a SAT solver which works as a black box and can produce only one model at a time, what clause should be added into the clause set so that the SAT solver will produce a new model?
4. Estimate the upper bound of the size of the inference graph for resolution, if the input clauses contain  $n$  variables.
5. Provide the pseudo-code for implementing *BCPw* based on *BCPht* such that the space for  $head(c)$  and  $tail(c)$  are freed by assuming  $head(c) = 0$  and  $tail(c) = |c| - 1$  for any non-unit clause  $c$ , and no additional space can be used for  $c$ . In your code, you may call  $swap(lits(c), i, j)$ , which swaps the two elements at indices  $i$  and  $j$  in the array  $lits(c)$ .

6. Given a set  $H$  of  $m + 1$  Horn clauses over  $m$  variables:

$$H = \{c_1 = (\bar{x}_1 \mid \bar{x}_2 \mid \cdots \mid \bar{x}_{m-1} \mid \bar{x}_m), c_2 = (\bar{x}_1 \mid x_2), c_3 = (\bar{x}_2 \mid x_3), \dots, c_m = (\bar{x}_{m-1} \mid x_m), c_{m+1} = (x_1)\}$$

Show that the size of  $H$  is  $O(m)$ , and the algorithm *BCPw* in the previous problem will take  $\Omega(m^2)$  time on  $H$ , if all the clauses in  $cls(A)$  must be processed before considering  $cls(B)$ , where  $\neg A$  is assigned true before  $\neg B$ .

7. We apply *DPLL* (without CDCL) to the following set  $S$  of clauses:  $S = \{1. (a \mid c \mid \bar{d} \mid f), 2. (a \mid d \mid \bar{e}), 3. (\bar{a} \mid c), 4. (b \mid \bar{d} \mid e), 5. (b \mid d \mid \bar{f}), 6. (\bar{b} \mid \bar{f}), 7. (c \mid \bar{e}), 8. (c \mid e), 9. (\bar{c} \mid \bar{d})\}$ , assuming that *pickLiteral* picks a literal in the following order  $\bar{a}, \bar{b}, \bar{c}, \bar{d}, \bar{e}, \bar{f}$  (false is tried first for each variable). Please (a) draw the decision tree of DPLL; (b) provide the unit clauses derived by BCP at each node of the decision tree; and (c) identify the head/tail literals of each clause at the end of *DPLL* (assuming initially the first and last literals are the head/tail literals).
8. We apply *GSAT* to the set  $S$  of clauses in the previous problem. Assume that the initial node (i.e., interpretation) is  $\sigma_0 = \{\bar{a}, b, \bar{c}, d, \bar{e}, f\}$ , and the strategy for picking a neighbor to replace the current node is (a) the best among all the neighbors; (b) the first neighbor which is better than the current node (flipping variables in the order of  $a, b, \dots, f$ . Sideways moves are allowed if no neighbors are better than the current node. If no neighbors are better than the current node, a local optimum is found and you need to start with the node  $\sigma_0 = \{a, \bar{b}, c, \bar{d}, e, \bar{f}\}$ .
9. We apply DPLL with CDCL to the clause set  $C : \{c_1 : (x_1 \mid x_2 \mid x_5), c_2 : (x_1 \mid x_3 \mid \bar{x}_4 \mid x_5), c_3 : (\bar{x}_1 \mid x_2 \mid x_3 \mid x_4), c_4 : (\bar{x}_3 \mid x_4 \mid x_5), c_5 : (\bar{x}_2 \mid x_3 \mid x_5), c_6 :$

$(\overline{x_3} \mid \overline{x_4} \mid \overline{x_5}), c_7 : (\overline{x_3} \mid x_4 \mid \overline{x_5})\}$ . If *pickLiteral* picks a literal for case-split in the following order  $\overline{x_5}, \overline{x_4}, \dots, \overline{x_1}$  ( $x_i = 0$  is tried first), please answer the following questions:

- (a) For each conflicting clause found in  $DPLL(C)$ , what new clause will be generated by *conflictAnalysis*?
  - (b) What will be the head and tail literals if the new clause generated in (a) is not unit? And what will be the value of *level* at the end of the call to *insertNewClause*?
  - (c) Draw the DPLL decision tree until either a model is found or  $DPLL(C)$  returns *false*.
10. The pigeon hole problem is to place  $n + 1$  pigeons into  $n$  holes such that every pigeon must be in a hole and no holes can hold more than one pigeon. Please write an encoder in any programming language that reads in an integer  $n$  and generates the CNF formula for the pigeon hole problem with  $n + 1$  pigeons and  $n$  holes in the DIMACS Format. Store the output of your encoder in a file named `pigeonN.cnf` and feed `pigeonN.cnf` to a SAT solver which accepts CNF in the DIMACS format. Turn in both the encoder, the result of the SAT solver for  $n = 4, 6, 8, 10$ , and the CNF file for  $n = 4$  (i.e., `pigeon4.cnf`), plus a summary of the outputs of SAT solvers on the sizes of the input, the computing times, numbers of conflicts. What relation between the size of the input and the computing time?
  11. The  $n$ -queen problem is to place  $n$  queens on an  $n \times n$  board such that no two queens are in the same row, the same column, or the same diagonal. Please write a smaller program called *encoder* in any programming language that reads in an integer  $n$  and generates the CNF formula for the  $n$ -queen problem in the DIMACS Format. Store the output of your encoder in in a file named `queenN.cnf` and feed `queenN.cnf` to a SAT solver which accepts CNF in the DIMACS format. Write another smaller program called *decoder* which reads the output of the SAT solver and displays a solution of  $n$ -queen on your computer screen. Turn in both the encoder and decoder, the result of the SAT solver for  $n = 5, 10, 15$ , and the CNF file for  $n = 5$  (i.e., `queen5.cnf`), and the output of your decoder for  $n = 10$ .
  12. You are asked to write two small programs, one called *encoder* and the other called *decoder*. The encoder will read a sudoku puzzle from a file (some examples are provided online) and generates the CNF formula for this puzzle

in the DIMACS Format. Store the output of your encoder in a file named `sudokuN.cnf` and feed `sudokuN.cnf` to a SAT solver which accepts CNF in the DIMACS format. The decoder will read the output of the SAT solver and displays a solution of the Sudoku puzzle on your computer screen. Turn in both the encoder and decoder, the result of the SAT solver, and the result of your decoder for the provided Sudoku puzzles.

13. A Latin square over  $N = \{0, 1, \dots, n-1\}$  is an  $n \times n$  square where each row and each column of the square is a permutation of the numbers from  $N$ . Every Latin square defines a binary function  $* : N, N \rightarrow N$  such that  $(x*y) = z$  iff the entry at the  $x^{\text{th}}$  row and the  $y^{\text{th}}$  column is  $z$ . You are asked to write a small program called *encoder*. The encoder will read an integer  $n$  and generates the CNF formula for the Latin square of size  $n$  satisfying the constraints  $x*x = x$  and  $(x*y)*(y*x) = y$  (called *Stein's third law*) in the DIMACS Format. Store the output of your encoder in a file named `latinN.cnf` and feed `latinN.cnf` to a SAT solver which accepts CNF in the DIMACS format. Turn in the encoder, and a summary of the results of the SAT solver for  $n = 4, 5, \dots, 9$ .
14. Given an undirected simple graph  $G = (V, E)$ , an *independent set* of  $G$  is a subset  $X \subseteq V$  such that for any two vertices of  $X$ ,  $(x, y) \notin E$ . The *independent set problem* is to find a maximum independent set of  $G$ . Specify the independent set problem in hybrid MaxSAT; and (b) specify the decision version of the independent set problem in SAT.
15. Given an undirected simple graph  $G = (V, E)$ , a *vertex cover* of  $G$  is a subset  $X \subseteq V$  such that for every edge  $(x, y) \in E$ ,  $\{x, y\} \cap X \neq \emptyset$ . The *vertex cover problem* is to find a minimal vertex cover of  $G$ . (a) Specify the vertex cover problem in hybrid MaxSAT; and (b) specify the decision version of the vertex cover problem in SAT.
16. Given an undirected simple graph  $G = (V, E)$ , the *longest path problem* is to find the longest simple path (no repeated vertices) in  $G$ . (a) Specify the longest path problem in hybrid MaxSAT; and (b) specify the decision version of the longest path problem in SAT.
17. Given a set  $S$  of binary clauses,

$$S = \{(\bar{a} \mid \bar{b}), (\bar{a} \mid \bar{c}), (a \mid \bar{d}), (b \mid \bar{c}), (b \mid d), (c \mid \bar{d})\}$$

- (a) draw the implication graph  $G = (V, E)$  of  $S$ ; (b) find the strongly connected components of  $G$ ; (c) From  $G$ , decide if  $S$  is satisfiable or not. If satisfiable, construct a model of  $S$  from  $G$ .

## CHAPTER 5

### FIRST ORDER LOGIC

Propositional logic provides a good start at describing the general principles of logical reasoning, but it is not possible to express general properties of many important sets, especially when a set is infinite. For example, how to specify that every natural number is either even or odd, a natural number is even iff its successor is odd, or the relation  $<$  is transitive over the natural numbers?

The weak expressive power of propositional logic accounts for its relative mathematical simplicity, but it is a very severe limitation, and it is desirable to have more expressive logics. Mathematics and some other disciplines such as computer science often consider sets of elements in which certain relations and operations are singled out. When using the language of propositional logic, our ability to talk about the properties of such relations and operations is very limited. Thus, it is necessary to refine the logic language, in order to increase the expressive power of logic. These are exactly what is provided by a more expressive logical framework known as first-order logic, which will be the topic of the next three chapters. First-order logic is a considerably richer logic than propositional logic, but yet enjoys many nice mathematical properties.

In first-order logic, many interesting and important properties about various sets can be expressed. Technically, this is achieved by allowing the propositional symbols to have arguments ranging over elements of sets, to express the relations and operations in question. These propositional symbols are called *predicate* symbols. This is the reason that first-order logic is often called *predicate calculus*.

#### 5.1 Syntax of First Order Languages

A logic is a collection of closely related formal languages. There are certain languages called first-order languages, and together they form first-order logic. Our study of first-order logic will parallel the study of propositional logic: We will at first define the syntax and semantics of first-order languages as an extension of propositional languages, and then present various reasoning tools for these languages.

##### 5.1.1 Terms and Formulas

Consider some general statements about the natural numbers:

- Every natural number is even or odd, but not both.
- A natural number  $x$  is even iff  $x + 1$  is odd.
- For any natural number  $x$ ,  $x < x + 1$  and  $\neg(x < x)$ .
- For any three natural numbers  $x$ ,  $y$ , and  $z$ , if  $x < y$  and  $y < z$ , then  $x < z$ .
- For every natural number  $x$ , there exists a natural number  $y$ , such that  $x < y$ .

Before any attempt to prove these statements, let us consider how to write them in a formal language. At first, we need a way to describe the set of natural numbers. Like propositional logic, first-order logic does not assume the existence of any special sets, except some symbols, and everything needs to be constructed from scratch. Fortunately, from Example 1.3.11, we know that the set of natural numbers can be constructed recursively from two symbols: the constant 0 and the unary function symbol  $s$ : Natural number  $i$  is uniquely represented by the term  $s^i(0)$ , which means there are  $i$  copies of  $s$  in  $s(s(\dots s(0)\dots))$ . Hence, our first-order language will have 0 and  $s$ .

Next, to express “even” or “odd” properties, we may introduce the predicate symbols *even* and *odd*, which take a natural number and return a Boolean value. The relation  $<$  is used in the statement and we take it in the first-order language. Since some variables like  $x$ ,  $y$ , and  $z$  are already used in the statements, we borrow them, too.

Finally, to express “for every” and “there exists” in natural languages, we add two important symbols,  $\forall$  and  $\exists$ , respectively, in a first-order language. The symbol  $\forall$  is called *universal quantifier* and  $\exists$  is called *existential quantifier*. These two symbols will be discussed extensively in the next subsection.

**Definition 5.1.1.** *A first-order language  $L$  is built up by an alphabet that consists of four sets of symbols:*

- A set  $P$  of predicate symbols, including propositional variables, usually denoted by  $p$ ,  $q$ ,  $r$ , with or without subscripts, or some popular binary relation symbols, such as  $<$ ,  $\leq$ ,  $\geq$ ,  $>$ , etc.
- A set  $F$  of function symbols, including constants, usually denoted by  $a$ ,  $b$ ,  $c$  (constants),  $f$ ,  $g$ ,  $h$ , with or without subscripts.
- A set  $X$  of variables, usually denoted by  $x$ ,  $y$ , and  $z$ , with or without subscripts.
- A set  $Op$  of logical operators such as  $\top$ ,  $\perp$ ,  $\neg$ ,  $\wedge$ ,  $\vee$ ,  $\rightarrow$ , and  $\leftrightarrow$ , plus  $\forall$  and  $\exists$ .

For each symbol in  $P \cup F$ , the number of arguments for that symbol is fixed and is dictated by the arity function  $\text{arity} : P \cup F \rightarrow \mathcal{N}$ . If  $\text{arity}(p) = 0$  for  $p \in P$ ,  $p$  is the same as a propositional variable; if  $\text{arity}(f) = 0$  for  $f \in F$ ,  $f$  is a constant.

The triple  $(P, F, \text{arity})$  is called the signature of the language. We denote the language  $L$  by  $L = (P, F, X, Op)$ .

In general,  $P$  or  $F$  could be empty;  $X$  and  $Op$  can never be empty. In theory,  $P$  or  $F$  or  $X$  could be infinite, but in practice, they are always finite. We also assume that no symbols appear more than once in these four sets. Rather than fixing a single language once and for all, different signatures generate different first-order languages and allow us to specify the symbols we wish to use for any given domain of interest.

**Example 5.1.2.** To specify the properties of the natural numbers listed at the beginning of this section, we use the signature  $(P, F, \text{arity})$ , where  $P = \{\text{even}/1, \text{odd}/1, </2\}$ ,  $F = \{0/0, s/1\}$ , the arity is given after the symbol / following each predicate or function symbol.

Once a signature is given, we add some variables  $X$  and logical operators  $Op$ , to obtain the alphabet for a first-order language. First-order logic allows us to build complex expressions out of the alphabet. Starting with the variables and constants, we can use the function symbols to build up compound expressions which are called *terms*.

**Definition 5.1.3.** Given a first-order language  $L = (P, F, X, Op)$ , the set of terms of  $L$  are the set of strings built up by  $F$  and  $X$ , which can be formally defined by the following BNF grammar:

$$\begin{aligned} \langle \text{Constants} \rangle & ::= a \text{ if } a \in F, \text{arity}(a) = 0 \\ \langle \text{Variables} \rangle & ::= x \text{ if } x \in X \\ \langle \text{Terms} \rangle & ::= \langle \text{Constants} \rangle \mid \langle \text{Variables} \rangle \mid \\ & \quad f(\langle \text{Terms} \rangle_1, \langle \text{Terms} \rangle_2, \dots, \langle \text{Terms} \rangle_k) \text{ if } f \in F, \text{arity}(a) = k \end{aligned}$$

The set of terms in  $\langle \text{Terms} \rangle$  is often denoted by  $T(F, X)$ . A term is said to be ground if it does not contain any variable. The set of all ground terms is denoted by  $T(F)$ . Note that  $T(F) = \emptyset$  if  $F$  contains no constants.

For  $F = \{0, s\}$  in the previous example,  $T(F) = \{0, s(0), s(s(0)), \dots\}$  and  $T(F, X) = \{0, s(0), s(x), s(y), s(s(0)), s(s(x)), s(s(y)), \dots\}$  for  $X = \{x, y\}$ .

Intuitively, the terms name objects in the intended domain of interest. In a term, besides the symbols from  $F$  and  $X$ , the parentheses “(” and “)” and the comma “,” are also used to identify its structure. Like every propositional formula

has a formula tree (Definition 2.1.2), for each term we can have a term tree in which parentheses and commas are omitted. We may assign a *position* to each node of the term tree and use this position to identify the subterm denoted by that node.

**Definition 5.1.4.** *Given a term  $t$ , a position  $p$  of  $t$ , with the subterm at  $p$ , denoted by  $t/p$ , is recursively defined as follows:*

- *If  $t$  is a constant or a variable, then  $p = \epsilon$ , the empty sequence, and  $t/p = t$ .*
- *If  $t = f(t_1, \dots, t_k)$ , then either  $p = \epsilon$  and  $t/p = t$ , or  $p = i.q$ , where  $1 \leq i \leq k$ , and  $t/p = t_i/q$ .*

For example, if  $t = f(x, g(g(y)))$ , then the legal positions of  $t$  are  $\epsilon, 1, 2, 2.1$ , and  $2.1.1$ , and the corresponding subterms are  $t, x, g(g(y)), g(y)$ , and  $y$ , respectively.

Now adding the usual logical operators from propositional logic, we have pretty much all the symbols needed for a first-order language.

**Definition 5.1.5.** *Given a first-order language  $L = (P, F, X, Op)$ , the formulas of  $L$  are the set of strings built up by terms, predicate symbols and logical operators. Let  $op$  be the binary operators used in the current application, then the formulas for this application can be defined by the following BNF grammar:*

$$\begin{aligned}
 \langle op \rangle &::= \wedge \mid \vee \mid \rightarrow \mid \oplus \mid \leftrightarrow \\
 \langle Atoms \rangle &::= p(\langle Terms \rangle_1, \langle Terms \rangle_2, \dots, \langle Terms \rangle_k) \text{ if } p \in P, \text{arity}(a) = k \\
 \langle Literals \rangle &::= \langle Atoms \rangle \mid \neg \langle Atoms \rangle \\
 \langle Formulas \rangle &::= \top \mid \perp \mid \langle Atoms \rangle \mid \neg \langle Formulas \rangle \mid \\
 &\quad (\langle Formulas \rangle \langle op \rangle \langle Formulas \rangle) \mid \\
 &\quad (\forall \langle Variables \rangle \langle Formulas \rangle) \mid (\exists \langle Variables \rangle \langle Formulas \rangle)
 \end{aligned}$$

A formula (clause, literal, atom) is said to be ground if it does not contain any variable.

Atoms are a shorthand for *atomic formulas*. We can have a formula tree for each first-order formula.

Some parentheses are unnecessary if we use a precedence relation among the logical operators. When displaying a formula, it is common to drop the outmost parentheses or those following a predicate symbol of zero arity, e.g. we write  $p$  instead of  $p()$ . We may also define subformulas of a formula as we did for propositional formulas and terms.

**Example 5.1.6.** Given a signature  $(P, F, \text{arity})$ , where  $P = \{\text{child}/2, \text{love}/2\}$  and  $F = \{a_i/0 : 1 \leq i \leq 100\}$ , then  $T(F) = F$  and  $T(F, X) = F \cup X$ , because  $F$  does not have any non-constant function symbols. If  $X = \{x, y, z\}$ , the following are some examples of formulas:

- **ground atoms:**  $child(a_2, a_1)$ ,  $child(a_3, a_1)$ ,  $love(a_1, a_2)$ , ...
- **non-ground atoms:**  $child(x, y)$ ,  $loves(y, x)$ ,  $love(z, x)$ , ...
- **ground non-atom formulas:**  $\neg child(a_3, a_4)$ ,  $child(a_2, a_1) \vee child(a_3, a_1)$ ,  $child(a_4, a_1) \rightarrow love(a_1, a_4)$ , ...
- **general formulas:**  $\neg child(a_1, y)$ ,  $\neg child(x, y) \vee love(y, x)$ ,  $love(x, y) \vee \neg love(y, x)$ , ...

In these examples, if we replace each atom by a propositional variable, they become propositional formulas. This is true for every ground formulas. Examples of formulas with quantifiers are given in the next subsection.

### 5.1.2 The Quantifiers

What makes first-order logic powerful is that it allows us to make general assertions using quantifiers: Universal quantifier  $\forall$  and existential quantifier  $\exists$ . In the definition of  $\langle Formulas \rangle$ , every occurrence of quantifiers is followed by a variable, say  $x$ , and then followed by a formula, say  $A$ . That is, the quantifiers appear in a formula like either  $\forall x A$  or  $\exists x A$ ; the former says  $A$  is hold for every value of  $x$  and the latter says  $A$  is hold for some value of  $x$ .

In middle school algebra, when we say  $+$  is commutative, we use formula  $x + y = y + x$ , here  $x$  and  $y$  are meant for every values. When solving an equation like  $x^2 - 3x + 2 = 0$ , we try to find some value of  $x$  so that the equation holds. In first-order logic, the two cases are expressed as  $\forall x(\forall y(x + y = y + x))$  and  $\exists x(x^2 - 3x + 2 = 0)$ , respectively.

**Example 5.1.7.** Continue from Example 5.1.6, if  $child(x, y)$  means  $x$  is a child of  $y$ ,  $parent(x, y)$  means  $x$  is a parent of  $y$ ,  $descen(x, y)$  means  $x$  is a descendant of  $y$ , and  $love(x, y)$  means  $x$  loves  $y$ , then we can write many assertions in first-order logic.

- $x$  is a child of  $y$  iff  $y$  is a parent of  $x$ :  $\forall x \forall y child(x, y) \leftrightarrow parent(y, x)$ .
- If  $x$  is a child of  $y$ , then  $x$  is a descendant of  $y$ :  $\forall x \forall y child(x, y) \rightarrow descen(x, y)$ .
- Every parent loves his child:  $\forall x \forall y (parent(x, y) \rightarrow love(x, y))$ .
- Everybody is loved by somebody:  $\forall x \exists y love(y, x)$ .
- The descendant relation is transitive:  $\forall x \forall y \forall z (descen(x, y) \wedge descen(y, z) \rightarrow descen(x, z))$ .

When using  $\forall$  to describe a statement under various conditions, these conditions usually should be the first argument of  $\rightarrow$ . When using  $\exists$  to describe a statement under various conditions, we cannot use *imply*. For example, let  $own(x, y)$  means “ $x$  owns  $y$ ”,  $red(x)$  means “ $x$  is red”,  $toy(x)$  means “ $x$  is a toy”,  $car(x)$  means “ $x$  is a car”. Then the statement “Tom loves every red toy he owns” can be expressed as

$$\forall x (own(Tom, x) \wedge red(x) \wedge toy(x) \rightarrow love(Tom, x)).$$

It is natural to read it as a restricted quantifier: “for everything  $x$ , if Tom owns  $x$ ,  $x$  is red and  $x$  is a toy, then Tom loves  $x$ .” Of course, if we get rid of  $\rightarrow$  in the above formula, we have

$$\forall x \neg own(Tom, x) \vee \neg red(x) \vee \neg toy(x) \vee love(Tom, x),$$

which can be read as ““for everything  $x$ , either Tom doesn’t own  $x$ ,  $x$  is not red,  $x$  is not a toy, or Tom loves  $x$ .” Note that if we replace  $\rightarrow$  by  $\wedge$  in the above formula, then it claims that “Tom owns everything”, “everything is red”, “everything is a toy”, and “Tom loves everything”.

On the other hand, to express “Tom loves one of his red toy cars”, the formula is

$$\exists x (own(Tom, x) \wedge toy(x) \wedge red(x) \wedge car(x) \wedge love(Tom, x)).$$

If we replace the last  $\wedge$  by  $\rightarrow$  in the above formula, then it claims that there exists something, denoted by “it”, and one (or more) of the following is true: “Tom doesn’t own it”, “it is not a toy”, “it is not red”, “it is not a car”, or “Tom loves it”. Since there are many things that “Tom doesn’t own it”, “it is not a toy”, “it is not red”, or “it is not a car”, the above formula does not catch the meaning that “Tom loves one of his red toy cars.”

**Definition 5.1.8.** In  $\forall x A$  or  $\exists x A$ , the scope of the variable  $x$  is  $A$  and every occurrence of  $x$  in  $A$  is said to be bounded. A variable in  $A$  is said to be free if it is not bounded. A formula  $A$  is said to be closed if it has no free variables. A sentence is a closed formula.

In  $love(x, y)$ , both  $x$  and  $y$  are free; in  $\exists y p(x, y)$ ,  $x$  is free,  $y$  is not free; in  $\forall x \exists y p(x, y)$ , neither  $x$  nor  $y$  are free; in  $(\forall x q(x)) \vee r(x)$ , the first occurrence of  $x$  is bounded and the second occurrence of  $x$  is free. For the last example, we cannot say “ $x$  in  $q(x)$  is bounded” as “ $x$  is free in  $q(x)$ ” but we can say “ $x$  is bounded in  $\forall x q(x)$ ”.

When  $A$  represents any formula, we write  $A(x_1, x_2, \dots, x_n)$  to indicate that the free variables of  $A$  are among  $x_1, x_2, \dots, x_n$ .

Like propositional formulas, we may use the precedence relation on the logical operators to avoid writing out all the parentheses in a formula. The precedence relation used in Chapter 2 is

$$\neg, \wedge, \vee, \rightarrow, \{\oplus, \leftrightarrow\}.$$

Where to place  $\forall$  and  $\exists$ ? According to *Wikipedia* on first-order logic, the quantifiers can be placed right before  $\rightarrow$ . However, as a common practice, we will place them at the end, after  $\leftrightarrow$ . That is,

$$\neg, \wedge, \vee, \rightarrow, \{\oplus, \leftrightarrow\}, \{\forall, \exists\}.$$

For example,  $\forall x A \vee B$  is interpreted as  $\forall x (A \vee B)$ , and we would write  $(\forall x A) \vee$  to limit the scope of  $x$  on  $A$  only.

Now, using the predicates and relation symbols, we can create formulas about the statements mentioned in the beginning of this section:

- $\forall x \text{odd}(x) \leftrightarrow \neg \text{even}(x)$ : Every natural number is even or odd, but not both.
- $\forall x \text{even}(s(x)) \leftrightarrow \text{odd}(x)$ : A natural number  $x$  is odd iff  $x + 1$  is even.
- $\forall x x < s(x) \wedge \neg(x < x)$ : For any natural number  $x$ ,  $x < x + 1$  and  $\neg(x < x)$ .
- $\forall x \exists y x < y$ : For every natural number  $x$ , there exists a natural number  $y$ , such that  $x < y$ .
- $\forall x \forall y \forall z (x < y) \wedge (y < z) \rightarrow (x < z)$ : For any three natural numbers  $x$ ,  $y$ , and  $z$ , if  $x < y$  and  $y < z$ , then  $x < z$ , that is,  $<$  is transitive.

### 5.1.3 Unsorted and Many-Sorted Logics

First-order languages are generally regarded as unsorted languages. If we look at its syntax carefully, two sets of objects are defined:  $\langle \text{Terms} \rangle$  and  $\langle \text{Formulas} \rangle$ . Obviously,  $\langle \text{Formulas} \rangle$  are of sort boolean, denoted by  $\text{Bool} = \{0, 1\}$ .  $\langle \text{Terms} \rangle$  represent objects of interests, and the set of objects is commonly called the *universe* (also called the *domain* of discourse) of the language. For example, the universe could use  $T(F)$ , the set of ground terms, to represent the universe.

If we use  $U$  to denote the sort of  $\langle \text{Terms} \rangle$ , for any signature  $(P, F, \text{arity})$ , predicate symbol  $p \in P$  represents a function  $p : U^{\text{arity}(p)} \rightarrow \text{Bool}$  and function symbol  $f \in F$  represents a function  $f : U^{\text{arity}(f)} \rightarrow U$ . The universal and existential quantifiers range over that universe  $U$ . For example, the first-order language

in Example 5.1.6 talks about 100 people living in a certain town, with a relation  $love(x, y)$  to express that  $x$  loves  $y$ . In such a language, we might express the statement that “everyone loves someone” by writing  $\forall x \exists y love(x, y)$ .

For many applications,  $Bool$  and  $U$  are disjoint, like in Examples 5.1.2 and 5.1.6. For some applications,  $Bool$  can be regarded as a subsort of  $U$ . For example, to model the if-then-else command in a programming language, we may use the symbol  $ite : Bool \times U \times U \rightarrow U$ . Another example is the function symbol  $says : U \times Bool \rightarrow Bool$  in the Knight and Knave puzzles (Example 2.5.4).

It is not difficult to introduce sorts in first-order logic and this is done by using unary predicate symbols ( $arity = 1$ ). In Example 5.1.6, we use  $a_i$  to stand for person  $i$ . If we need function  $age(x)$  to tell the age of person  $x$ , we need the function symbols for natural numbers. If we use those in Example 5.1.2, then  $F = \{0, s, a_i \mid 1 \leq i \leq 100\}$ . In this application, we do need to tell which are persons and which are numbers by using the following predicate symbols:

- $person : U \rightarrow Bool$ :  $person(x)$  is true iff  $x$  is a person.
- $nat : U \rightarrow Bool$ :  $nat(x)$  is true iff  $x$  is a natural number.

$person(x)$  can be easily defined by asserting  $person(a_i)$  for  $1 \leq i \leq 100$ .  $nat(x)$  can also be defined by asserting  $nat(0)$  and  $nat(s(x)) \equiv nat(x)$ .

Now, to state general properties as in Examples 5.1.2 and 5.1.6, we have to put restrictions on the variables. For example, to express “everybody is loved by somebody”, the formula now is

$$\forall x \exists y person(x) \wedge person(y) \rightarrow love(y, x).$$

To say “ $<$ ” is transitive on natural numbers, we use

$$\forall x, y, z nat(x) \wedge nat(y) \wedge nat(z) \wedge (x < y) \wedge (y < z) \rightarrow (x < z).$$

Dealing with such unary predicates is tedious. To get rid of such clumsy, people introduce “*sorts*”. Let

$$\begin{aligned} Person &= \{x \mid person(x) \wedge x \in U\} \\ Nat &= \{x \mid nat(x) \wedge x \in U\} \end{aligned}$$

Then define  $love : Person^2 \rightarrow Bool$  and  $< : Nat^2 \rightarrow Bool$  and the formulas become neat again: We do not need  $person(x)$ ,  $nat(x)$ , etc., in the formula. For the  $age$  function, we have  $age : Person \rightarrow Nat$ . That is the birth of many-sorted logic, which is a mild extension of first-order logic.

The idea of many-sorts is based on set theory, where each sort is a set of similar objects. Sets are usually defined by unary predicates. From a logician’s view point, a set  $S = \{x \mid p(x)\}$  and a predicate  $p(x)$  have the same meaning. In many-sorted logic, one can have different sorts of objects — such as *persons* and *natural numbers* — and a separate set of variables ranging over each. Moreover, the specification of function symbols and predicate symbols indicates what sorts of arguments they expect (the *domain*), and, in the case of function symbols, what sort of value they return (the *range*). The concept of *signature* is extended to contain the domain and range information of each function and predicate symbol.

Many-sorted logic can reflect formally our intention not to handle the universe as a collection of objects, but to partition it in a way that is similar to types in programming languages. This partition can be carried out on the syntax level: substitution of variables can be done only accordingly, respecting the “*sorts*”.

## 5.2 Semantics

We should be aware that, at this stage, all the function and predicate symbols are just symbols. In propositional logic, a symbol  $p$  can represent “Tom runs fast”, or “Mary is tall”. The same holds true in first-order logic as we may give different meaning to a symbol.

In Example 5.1.7, we have seen the following two formulas:

- $\forall x \forall y (child(x, y) \rightarrow descen(x, y))$ .
- $\forall x \forall y (parent(x, y) \rightarrow love(x, y))$ .

The first formula says “if  $x$  is a child of  $y$ , then  $x$  is a descendant of  $y$ ”; the second says “every parent loves his child.” What is the difference between them? They differ only on predicate symbols. If we give different meanings to the predicate symbols, one can replace the other. That is, we have designed the language with a certain interpretation in mind, but one could also interpret the same formula differently.

In this section, we will spell out the concepts of interpretations, models, satisfiable or valid formulas, as an extension of these concepts of propositional logic.

### 5.2.1 Interpretation

For propositional logic, an interpretation is an assignment of propositional variables to truth values. For a first-order language  $L = (P, F, X, Op)$ , we need an interpretation for every symbol of  $L$ .

**Definition 5.2.1.** Given a closed formula  $A$  of  $L = (P, F, X, Op)$ , an interpretation of  $A$  is a triple  $I = (D, R, G)$ , where

- $D$  is a non-empty set called the universe of  $I$ , often denoted by  $D^I$ .
- For each predicate symbol  $p \in P$  of arity  $k$ , there is a  $k$ -ary relation  $r_p \in R$  over  $D$ ,  $r_p \subseteq D^k$ , denoted by  $r_p = p^I$ . If  $k = 0$ , then  $p$  is a propositional variable and  $p^I$  is either 1 or 0, a truth assignment to  $p$ .
- For each function symbol  $f \in F$  of arity  $k$ , there is a function  $g : D^k \rightarrow D \in G$ , denoted by  $g = f^I$ . If  $k = 0$ , then  $f^I \in D$ .

**Example 5.2.2.** Consider  $L = (P, F, X, Op)$ , where  $P = \{p/2\}$ ,  $F = \{a/0\}$ , and  $X = \{x\}$ . Then for the formula  $A = \forall x p(a, x)$ , we may have the following interpretations:

1.  $I_1 = (\mathcal{N}, \{\leq\}, \{0\})$  and the meaning of  $A$  is “for every natural number  $x$ ,  $0 \leq x$ .”
2.  $I_2 = (\mathcal{N}, \{|\}, \{1\})$ ,  $x|y$  returns true iff  $x$  divides  $y$ , and the meaning of  $A$  is “for every natural number  $x$ ,  $1 | x$ .”
3.  $I_3 = (\{0, 1\}^3, \{prefix\}, \{\epsilon\})$ ,  $prefix(x, y)$  returns true iff  $x$  is a prefix of  $y$ , the meaning of  $A$  is “for every binary string  $x$  of length 3, the empty string  $\epsilon$  is a prefix of  $x$ .”
4.  $I_4 = (\mathcal{P}(\mathcal{N}), \{\subseteq\}, \{\emptyset\})$ ,  $\mathcal{P}(\mathcal{N})$  is the power set of  $\mathcal{N}$ , and the meaning of  $A$  is “for every subset  $x$  of natural numbers,  $\emptyset \subseteq x$ .”
5.  $I_5 = (V, \{E\}, \{a\})$ , where  $V = \{a, b, c\}$ ,  $E = \{(a, a), (a, b), (a, c), (c, c)\}$ ,  $G = (V, E)$  is a graph, and the meaning of  $A$  is “for every vertex  $x \in V$ , there is an edge  $(a, x) \in E$ .”

In propositional logic, once we have an interpretation, we may use a procedure like *eval* (Algorithm 2.2.4) to obtain the truth value of a formula. We can do the same in first-order logic, provided that we have a way of handling free variables and quantifiers in the formula.

Let  $free(A)$  be the free variables of formula  $A$ . Given a universe  $D$ , an *assignment* is a mapping  $\theta : free(A) \rightarrow D$ , which maps each free variable of  $A$  to an element of  $D$ . The purpose of  $\theta$  is that, during the evaluation of a formula under an interpretation, if we see a free variable  $x$ , we will use the value  $\theta(x)$  for  $x$ .

**Procedure 5.2.3.** The procedure *eval* takes a formula or a term  $A$ , with or without free variables, in  $L = (P, F, X, Op)$ , an interpretation  $I$  for  $L$ , and an assignment  $\theta : free(A) \rightarrow D^I$ , and returns a Boolean value for a formula and an element of  $D^I$  for a term.

```

proc eval( $A, I, \theta$ )
  if  $A = \top$  return 1; // 1 means true
  if  $A = \perp$  return 0; // 0 means false
  if  $A \in X$  return  $\theta(A)$ ; //  $A$  is a free variable
  if  $A = f(t_1, t_2, \dots, t_k)$  return  $f^I(t'_1, t'_2, \dots, t'_k)$ , where  $t'_i = eval(t_i, I, \theta)$ ;
  if  $A = \neg B$  return  $\neg eval(B, I, \theta)$ ;
  if  $A = (B \text{ op } C)$  return  $eval(B, I, \theta) \text{ op } eval(C, I, \theta)$ , where  $op \in Op$ ;
  if  $A = (\forall x B)$  return allInD( $B, I, \theta, x, D^I$ );
  if  $A = (\exists x B)$  return someInD( $B, I, \theta, x, D^I$ );
  else return “unknown”;

```

Procedure *allInD*( $B, I, \theta, x, S$ ) takes formula  $B$ , interpretation  $I$ , assignment  $\theta$ , free variable  $x$  in  $B$ , and  $S \subseteq D$ , evaluates the value of  $B$  under  $I$  and  $\theta \cup \{x \leftarrow d\}$  for every value  $d \in S$ . If *eval*( $B, I, \theta \cup \{x \leftarrow d\}$ ) returns 0 for one  $d \in D^I$ , return 0; otherwise, return 1.

```

proc allInD( $B, I, \theta, x, S$ )
  if  $S = \emptyset$  return 1;
  pick  $d \in S$ ;
  if ( $eval(B, I, \theta \cup \{x \leftarrow d\}) = 0$ ) return 0;
  return allInD( $B, I, \theta, x, S - \{d\}$ );

```

Procedure *someInD*( $B, I, \theta, x, S$ ) works in the same way as *allInD*, except that if one of the evaluations, *eval*( $B, I, \theta \cup \{x \leftarrow d\}$ ), return 1, return 1; otherwise, return 0.

```

proc someInD( $B, I, \theta, x, S$ )
  if  $S = \emptyset$  return 0;
  pick  $d \in S$ ;
  if ( $eval(B, I, \theta \cup \{x \leftarrow d\}) = 1$ ) return 1;
  return someInD( $B, I, \theta, x, S - \{d\}$ );

```

**Proposition 5.2.4.** Given a formula  $A$  of  $L = (P, F, X, Op)$ , an interpretation  $I$ , and an assignment  $\theta : free(A) \rightarrow D^I$ , *eval*( $A, I, \theta$ ) will return a boolean value; if  $t$  is a term appearing in  $A$ , then *eval*( $t, I, \theta$ ) will return a value of  $D^I$ .

*Proof.* (sketch) Let  $D$  be  $D_I$ . We do induction on the tree structure of  $A$ . The base cases when  $A$  or  $t$  are symbols of arity 0 are trivial. If  $t = f(t_1, t_2, \dots, t_k)$  is a term in  $A$ , the returned value is  $f^I(t'_1, t'_2, \dots, t'_k) \in D^I$  because  $f^I : D^k \rightarrow D$  is a function over  $D$ . If  $A$  is an atom  $p(t_1, t_2, \dots, t_k)$ , where  $p \in P$ , then  $p^I$  is a relation and  $p^I(t'_1, t'_2, \dots, t'_k) = 1$  iff  $(t'_1, t'_2, \dots, t'_k) \in p^I$  by our assumption. That is,  $p^I(t'_1, t'_2, \dots, t'_k)$  is boolean. For non-atom formulas, the returned values are always boolean as shown by the next proposition.  $\square$

From now on, we will denote  $eval(A, I, \theta)$  by  $I(A, \theta)$  and  $eval(A, I, \emptyset)$  by  $I(A)$  for brevity.

**Proposition 5.2.5.** *For formulas  $A$  and  $B$ , interpretation  $I$ , and assignment  $\theta$ ,*

- $I(\neg A, \theta) = \neg I(A, \theta)$ ;
- $I(A \vee B, \theta) = I(A, \theta) \vee I(B, \theta)$ ;
- $I(A \wedge B, \theta) = I(A, \theta) \wedge I(B, \theta)$ ;
- $I(A \oplus B, \theta) = I(A, \theta) \oplus I(B, \theta)$ ;
- $I(A \leftrightarrow B, \theta) = I(A, \theta) \leftrightarrow I(B, \theta)$ ;
- $I(\forall x B(x), \theta) = allInD(B(x), I, \theta, x, D^I)$ ;
- $I(\exists x B(x), \theta) = someInD(B(x), I, \theta, x, D^I)$ .

From the definition of *eval*, it is easy to check that the above proposition holds. Obviously, *eval* cannot be an algorithm because *allInD* will not terminate when  $D^I$  is infinite. We present *eval* as a procedure because we assume that the reader is familiar with algorithms. We need other tools such as induction to evaluate a formula when the universe of the interpretation is infinite. When  $D^I$  is finite, the termination of *eval* is easy to establish.

**Example 5.2.6.** For  $I_5$  in Example 5.2.2,  $I_5(\forall x p(a, x))$  will call *allInD* with parameters  $(p(a, x), I_5, \emptyset, x, \{a, b, c\})$ , which will call *eval* three times, with  $\theta = \{x \leftarrow a\}$ ,  $\{x \leftarrow b\}$ , and  $\{x \leftarrow c\}$ , respectively. Since  $E$  contains  $(a, x)$  for  $x \in \{a, b, c\}$ , all the evaluations return 1. The returned result can be expressed as

$$I_5(\forall x p(a, x)) = \bigwedge_{d \in \{a, b, c\}} I_5(p(a, x), \{x \leftarrow d\})$$

We can do the same for all the interpretations in Example 5.2.2. In general,

$$I(\forall x B(x), \theta) = \bigwedge_{d \in D^I} I(B(x), \theta \cup \{x \leftarrow d\})$$

which catches the meaning that “ $I(\forall x B(x), \theta)$  returns 1 iff  $B(x)$  is evaluated to be 1 under  $I$  and  $\theta \cup \{x \leftarrow d\}$  for every value  $d \in D^I$ ”.

Similarly, we have the following notation for  $I(\exists x B(x), \theta)$ , which is equal to *someInD*( $B(x), I, \theta, x, D$ ):

$$I(\exists x B(x), \theta) = \bigvee_{d \in D^I} I(B(x), \theta \cup \{x \leftarrow d\})$$

which catches the meaning that “ $I(\exists x B(x), \theta)$  returns 1 iff  $B(x)$  is evaluated to be true under  $I$  and  $\theta \cup \{x \leftarrow d\}$  for some value  $d \in D^I$ ”. These notations will help us to understand quantified formulas and establish their properties. Note that if the domain  $D^I$  is empty, then  $I(\forall x B(x), \theta) = 1$  and  $I(\exists x B(x), \theta) = 0$ .

### 5.2.2 Models, Satisfiability, and Validity

Here are just straightforward extensions from propositional logic to first-order logic. Recall that a sentence is a closed formula.

**Definition 5.2.7.** *Given a sentence  $A$  and an interpretation  $I$ ,  $I$  is said to be a model of  $A$  if  $I(A) = 1$ . If  $A$  has a model,  $A$  is said to be satisfiable.*

For a propositional formula, the number of models is finite; for a first-order formula, the number of models is infinite in general, because we have an infinite number of choices for domains, relations, and functions in an interpretation. However, we can still borrow the notation  $\mathcal{M}(A)$  from propositional logic.

**Definition 5.2.8.** *Given a closed formula  $A$ , let  $\mathcal{M}(A)$  be the set of all models of  $A$ . If  $\mathcal{M}(A) = \emptyset$ ,  $A$  is said to be unsatisfiable; if  $\mathcal{M}(A)$  contains every interpretation, i.e., every interpretation is a model of  $A$ , then  $A$  is said to be valid.*

With the identical definitions from propositional logic to first-order logic, the following result is expected.

**Proposition 5.2.9.** *Every valid propositional formula is a valid formula in first-order logic, and every unsatisfiable propositional formula is unsatisfiable in first-order logic.*

Of course, there are more valid formulas in first-order logic and our attention is on formulas with quantifiers.

**Example 5.2.10.** Let us check some examples:

1.  $A_1 = \forall x (p(a, x) \rightarrow p(a, a))$ .

Consider  $I = (\{a, b\}, \{p\}, \{a\})$ , where  $p^I = p = \{\langle a, a \rangle\}$ ,  $a^I = a$ ,  $I(A_1) = 1$  because  $p^I(a, a) = 1$  and

$$\begin{aligned} I(A_1) &= \bigwedge_{d \in \{a, b\}} I(p(a, x) \rightarrow p(a, a), \{x \leftarrow d\}) \\ &= I(p(a, a) \rightarrow p(a, a)) \wedge I(p(a, b) \rightarrow p(a, a)) \\ &= 1 \wedge (p^I(a, b) \rightarrow p^I(a, a)) \\ &= 1 \end{aligned}$$

Let  $I' = (\{a, b\}, \{\{\langle a, b \rangle\}\}, \{a\})$ , then  $I'(A_1) = 0$  because  $p^{I'}(a, a) = 0$  and  $p^{I'}(a, b) = 1$  imply that  $I'(p(a, b) \rightarrow p(a, a)) = 0$ .  $I'(A_1) = I'(p(a, a) \rightarrow p(a, a)) \wedge I'(p(a, b) \rightarrow p(a, a)) = 1 \wedge 0 = 0$ . Thus  $A_1$  is satisfiable but not valid.

2.  $A_2 = (\forall x p(a, x)) \rightarrow p(a, a)$ .

$A_2$  is different from  $A_1$  because the scopes of  $x$  are different in  $A_1$  and  $A_2$ . For any interpretation  $I$ , we consider the truth value of  $p^I(a^I, a^I)$ . If  $p^I(a^I, a^I) = 1$ , then  $I(A_2, \emptyset) = 1$ . If  $p^I(a^I, a^I) = 0$ , then  $I(\forall x p(a, x)) = 0$  because  $x$  will take on value  $a^I$ . So  $I(A_2) = 1$  in any case. Since  $I$  is arbitrary,  $A_2$  must be valid.

3.  $A_3 = \exists x \exists y (p(x) \wedge \neg p(y))$ .

Consider  $I = (\{a, b\}, \{p\}, \emptyset)$ , where  $p(a) = 1$  and  $p(b) = 0$ , then  $p(a) \wedge \neg p(b) = 1$ . From  $I(p(x) \wedge \neg p(y), \{x \leftarrow a, y \leftarrow b\}) = 1$ , we conclude that  $I$  is a model of  $A_3$ .  $A_3$  is false in any interpretation whose domain is empty or contains a single element. Thus,  $A_3$  is not valid.

**Definition 5.2.11.** Let  $A(x)$  be a formula where  $x$  is a free variable of  $A$ . Then  $A(t)$  denotes the formula  $A(x)[x \leftarrow t]$ , where every occurrence of  $x$  is replaced by term  $t$ . The formula  $A(t)$  is called an instance of  $A(x)$  by substituting  $x$  for  $t$  (ref. Definition 2.1.4).

We will use substitutions extensively in the next chapter. The substitution  $\theta : X \rightarrow D$  used in the evaluation procedure *eval* coincides with the above definition if  $D = T(F)$  (the set of ground terms built on  $F$ ). That is, for any interpretation  $I = (D, R, G)$  of  $L = (P, F, X, Op)$ , if  $D = T(F)$ , then  $I(A(x), \{x \leftarrow d\}) = I(A(d))$ , i.e.,  $eval(A(x), I, \{x \leftarrow d\}) = eval(A(d), I, \emptyset)$ , where  $x \in X, d \in T(F)$ , and  $A(d) = A(x)[x \leftarrow d]$ .

**Proposition 5.2.12.** For any formula  $A(x)$ , where  $x$  is a free variable of  $A$  and variable  $y$  does not appear in  $A(x)$ ,  $\forall x A(x) \equiv \forall y A(y)$  and  $\exists x A(x) \equiv \exists y A(y)$ .

*Proof.* For any interpretation  $I$ ,

$$\begin{aligned} I(\forall x A(x)) &= \bigwedge_{d \in D} I(A(x), \{x \leftarrow d\}) \\ I(\forall y A(y)) &= \bigwedge_{d \in D} I(A(y), \{y \leftarrow d\}) \end{aligned}$$

The variable  $x$  serves as a placeholder in  $A(x)$ , just like  $y$  serves as the same placeholder in  $A(y)$ , to tell *eval* when to replace  $x$  or  $y$  by  $d$ . Any distinct name can be used for this placeholder, thus  $I(\forall x A(x), \theta) = I(\forall y A(y), \theta)$ . Since  $I$  is arbitrary, it must be the case that  $\forall x A(x) \equiv \forall y A(y)$ . The proof for the second part is identical.

□

The above proposition allows us to change safely the names of bounded variables, so that different variables are used for each occurrence of quantifiers. For example, it is better to replace  $\forall x p(x) \vee \exists x q(x, y)$  by  $\forall x p(x) \vee \exists z q(z, y)$ , to improve the readability.

We will keep the convention that a set  $S$  of formulas denotes the conjunction of the formulas appearing in the set. If every formula is true in the same interpretation, then the set  $S$  is said to be satisfiable and the interpretation is its model.

**Example 5.2.13.** Let  $S = \{\forall x \exists y p(x, y), \forall x \neg p(x, x), \forall x \forall y \forall z p(x, y) \wedge p(y, z) \rightarrow p(x, z)\}$ . The second formula states that “ $p$  is irreflexive”; the third formula says “ $p$  is transitive”.

Consider the interpretation  $I = (\{\mathcal{N}\}, \{<\}, \emptyset)$ , it is easy to check that  $I$  is a model of  $S$ , because  $<$  is transitive and irreflexive, i.e.,  $\neg(x < x)$ , and for every  $x \in \mathcal{N}$ , there exists  $x + 1 \in \mathcal{N}$  such that  $x < x + 1$ .

Does  $S$  has a finite model, i.e., a model in which the universe is finite? If  $S$  has a finite model, say  $I = (D, \{r_p\}, \emptyset)$ , where  $|D| = n$ , consider the directed graph  $G = (D, r_p)$ , then  $G$  cannot have any cycle, because, if there exists a cycle, then there exists a path from a vertex in the cycle to itself. If there exists a path from  $a$  to  $b$  in  $G$ , there must exist an edge from  $a$  to  $b$ , because  $r_p$  is transitive. We cannot have an edge from a vertex to itself because  $r_p$  is irreflexive. For a vertex  $x$  at the end of a path, we cannot find another vertex  $y$  satisfying  $\forall x \exists y p(x, y)$ . Hence,  $S$  cannot have any finite model.

### 5.2.3 Entailment and Equivalence

The following two definitions are copied from Definitions 2.2.13 and 2.2.17.

**Definition 5.2.14.** Given two formulas  $A$  and  $B$ , we say  $A$  entails  $B$ , or  $B$  is a logical consequence of  $A$ , denoted by  $A \models B$ , if  $\mathcal{M}(A) \subseteq \mathcal{M}(B)$ .

To denote that  $A$  is valid, we may write simply  $\models A$ .

**Definition 5.2.15.** *Given two formulas  $A$  and  $B$ ,  $A$  and  $B$  are said to be logically equivalent if  $\mathcal{M}(A) = \mathcal{M}(B)$ , denoted by  $A \equiv B$ .*

With the same definitions, we have the same results for first-order logic.

**Theorem 5.2.16.** *For any two formulas  $A$  and  $B$ , (a)  $A \models B$  iff  $A \rightarrow B$  is valid; and (b)  $A \equiv B$  iff  $A \leftrightarrow B$  is valid.*

The substitution theorem from propositional logic still holds in first-order logic.

**Theorem 5.2.17.** *For any formulas  $A$ ,  $B$ , and  $C$ , where  $B$  is a subformula of  $A$ , and  $B \equiv C$ , then  $A \equiv A[B \leftarrow C]$ .*

**Proposition 5.2.18.** *Let  $A(x)$ ,  $B(x)$ , and  $C(x, y)$  be the formulas with free variables  $x$  and  $y$ .*

1.  $\neg \forall x A(x) \equiv \exists x \neg A(x)$
2.  $\neg \exists x A(x) \equiv \forall x \neg A(x)$
3.  $\forall x A(x) \wedge B(x) \equiv (\forall x A(x)) \wedge (\forall x B(x))$
4.  $\exists x A(x) \vee B(x) \equiv (\exists x A(x)) \vee (\exists x B(x))$
5.  $(\forall x A(x)) \vee B \equiv \forall x A(x) \vee B$  if  $x$  is not free in  $B$
6.  $(\exists x A(x)) \wedge B \equiv \exists x A(x) \wedge B$  if  $x$  is not free in  $B$
7.  $\forall x \forall y C(x, y) \equiv \forall y \forall x C(x, y)$
8.  $\exists x \exists y C(x, y) \equiv \exists y \exists x C(x, y)$

*Proof.* We give the proofs of 1, 3, 5, and 7, and leave the rest as exercises.

1. Consider any interpretation  $I$ , using de Morgan's law,  $\neg(A \wedge B) \equiv \neg A \vee \neg B$ ,

$$\begin{aligned}
 & I(\neg \forall x A(x)) \\
 = & \neg I(\forall x A(x)) \\
 = & \neg \bigwedge_{d \in D^I} I(A(x), \{x \leftarrow d\}) \\
 = & \bigvee_{d \in D^I} \neg I(A(x), \{x \leftarrow d\}) \\
 = & \bigvee_{d \in D^I} I(\neg A(x), \{x \leftarrow d\}) \\
 = & I(\exists x \neg A(x)).
 \end{aligned}$$

Because  $I$  is arbitrary, the equivalence holds.

3. Consider any interpretation  $I = (D, R, G)$ , using the commutativity and associativity of  $\wedge$ ,

$$\begin{aligned}
& I(\forall x A(x) \wedge B(x)) \\
&= \bigwedge_{d \in D} I(A(x) \wedge B(x), \{x \leftarrow d\}) \\
&= \bigwedge_{d \in D} I(A(x), \{x \leftarrow d\}) \wedge I(B(x), \{x \leftarrow d\}) \\
&= (\bigwedge_{d \in D} I(A(x), \{x \leftarrow d\})) \wedge (\bigwedge_{d \in D} I(B(x), \{x \leftarrow d\})) \\
&= I(\forall x A(x)) \wedge I(\forall x B(x)) \\
&= I((\forall x A(x)) \wedge (\forall x B(x)))
\end{aligned}$$

Because  $I$  is arbitrary, the equivalence holds.

5. In the following proof, we assume that  $x$  does not appear in  $B$ . Consider any interpretation  $I = (D, R, G)$ , using the distribution law of  $\vee$  over  $\wedge$ ,

$$\begin{aligned}
& I((\forall x A(x)) \vee B) \\
&= I(\forall x A(x)) \vee I(B) \\
&= (\bigwedge_{d \in D} I(A(x), \{x \leftarrow d\})) \vee I(B) \\
&= \bigwedge_{d \in D} (I(A(x), \{x \leftarrow d\}) \vee I(B)) \\
&= \bigwedge_{d \in D} (I(A(x), \{x \leftarrow d\}) \vee I(B, \{x \leftarrow d\})) \\
&= \bigwedge_{d \in D} I(A(x) \vee B, \{x \leftarrow d\}) \\
&= I(\forall x A(x) \vee B)
\end{aligned}$$

Because  $I$  is arbitrary, the equivalence holds.

7. Consider any interpretation  $I = (D, R, G)$ , using the commutativity and associativity of  $\wedge$ ,

$$\begin{aligned}
& I(\forall x \forall y C(x, y)) \\
&= \bigwedge_{d \in D} I(\forall y C(x, y), \{x \leftarrow d\}) \\
&= \bigwedge_{d \in D} (\bigwedge_{e \in D} I(C(x, y), \{x \leftarrow d, y \leftarrow e\})) \\
&= \bigwedge_{e \in D} (\bigwedge_{d \in D} I(C(x, y), \{x \leftarrow d, y \leftarrow e\})) \\
&= \bigwedge_{e \in D} I(\forall x C(x, y), \{y \leftarrow e\}) \\
&= I(\forall y \forall x C(x, y))
\end{aligned}$$

Because  $I$  is arbitrary, the equivalence holds. □

**Example 5.2.19.** Consider  $\exists x \forall y p(x, y)$  and  $\forall y \exists x p(x, y)$ . If  $p(x, y)$  stands for “ $x$  loves  $y$ ”, then the former means “there exists somebody who loves everyone”, but the latter means “everyone is loved by somebody”. Obviously, the two are not equivalent. To show

$$\exists x \forall y p(x, y) \models \forall y \exists x p(x, y),$$

consider any model  $I = (D, R, G)$  of  $\exists x \forall y p(x, y)$ . There must exist  $d \in D$  such that  $I(\forall y p(x, y), \{x \leftarrow d\}) = 1$ . Then  $\forall y \exists x p(x, y)$  must be true in  $I$ :

$$\begin{aligned}
 & I(\forall y \exists x p(x, y), \emptyset) \\
 &= \bigwedge_{e \in D} I(\exists x p(x, y), \{y \leftarrow e\}) \\
 &= \bigwedge_{e \in D} I(p(x, y), \{y \leftarrow e, x \leftarrow d\}) \\
 &= I(\forall y p(x, y), \{x \leftarrow d\}) \\
 &= 1
 \end{aligned}$$

To formally show that the two formulas are not equivalent, we need to provide a model in which one is true and the other is false.

**Example 5.2.20.** Let us prove that  $(\exists x p(x)) \rightarrow \forall y q(y)$  and  $\forall x p(x) \rightarrow q(x)$  are not equivalent. Let  $I = (\{a, b\}, \{r_p, r_q\}, \emptyset)$  where  $r_p(a) = r_q(a) = 1$  and  $r_p(b) = r_q(b) = 0$ . Then  $I(\exists x p(x)) = 1$  and  $I(\forall y q(y)) = 0$ . So  $I((\exists x p(x)) \rightarrow \forall y q(y)) = 0$ . On the other hand,  $I(\forall x p(x) \rightarrow q(x)) = 1$  because  $r_p = r_q$ , thus  $I(p(x) \rightarrow q(x), \{x \leftarrow d\}) = 1$ . Hence,  $(\exists x p(x)) \rightarrow \forall y q(y)$  and  $\forall x (p(x) \rightarrow q(x))$  are not equivalent.

**Example 5.2.21.** We may use known equivalence relations to show  $\exists x A(x) \rightarrow B(x) \equiv (\forall x A(x)) \rightarrow \exists x B(x)$ .

$$\begin{aligned}
 & \exists x A(x) \rightarrow B(x) \\
 \equiv & \exists x \neg A(x) \vee B(x) && // A \rightarrow B \equiv \neg A \vee B \\
 \equiv & (\exists x \neg A(x)) \vee \exists x B(x) && // \text{Proposition 5.2.18 4.} \\
 \equiv & \neg(\forall x A(x)) \vee \exists x B(x) && // \text{Proposition 5.2.18 1.} \\
 \equiv & (\forall x A(x)) \rightarrow \exists x B(x) && // A \rightarrow B \equiv \neg A \vee B
 \end{aligned}$$

### 5.3 Proof Methods

In Chapter 3, we introduced various proof systems for propositional logic: semantic tableau, natural deduction, resolution, etc. All these systems can be extended to first-order logic. In this section, we briefly introduce the extensions of semantics tableau and natural deduction to first-order logic and leave resolution to the next chapter.

#### 5.3.1 Semantic Tableau

Semantic tableau as given in Chapter 3 transforms a propositional formula into DNF by a set of  $\alpha$ -rules (which handle conjunctions and generate one child in the tableau) and  $\beta$ -rules (which handle disjunctions and generate two children in the tableau). These rules will be inherited for first-order logic. Keep in mind that the rules apply to the top logical operators in a formula in a leaf node of the tableaux.

$\alpha$	$\alpha_1, \alpha_2$
$A \wedge B$	$A, B$
$\neg(A \vee B)$	$\neg A, \neg B$
$\neg(A \rightarrow B)$	$A, \neg B$
$\neg(A \oplus B)$	$(A \vee \neg B), (\neg A \vee B)$
$A \leftrightarrow B$	$(A \vee \neg B), (\neg A \vee B)$
$\neg(A \uparrow B)$	$A, B$
$A \downarrow B$	$\neg A, \neg B$
$\neg\neg A$	$A$

$\beta$	$\beta_1$	$\beta_2$
$\neg(A \wedge B)$	$\neg A$	$\neg B$
$A \vee B$	$A$	$B$
$A \rightarrow B$	$\neg A$	$B$
$A \oplus B$	$A, \neg B$	$\neg A, B$
$\neg(A \leftrightarrow B)$	$A, \neg B$	$\neg A, B$
$A \uparrow B$	$\neg A$	$\neg B$
$\neg(A \downarrow B)$	$A$	$B$

To handle the quantifiers, we introduce two rules, one for the universal quantifier, called the  $\forall$ -rule, and one for the existential quantifier, called the  $\exists$ -rule.

**Definition 5.3.1.** *The  $\forall$ -rule is an inference rule:*

$$(\forall) \frac{\forall x A(x)}{A(t)} \quad t \text{ is a ground term}$$

The  $\forall$ -rule works as follows: if  $\forall x A(x)$  appears in a leaf node  $n$ , we create a child node  $n'$  of  $n$ , copy every formula of  $n$  into  $n'$ , and add  $A(t)$  into  $n'$ , where  $t$  is a ground term (without variables) appearing in  $n$ . If  $n'$  does not have any ground term, we may introduce a new constant as  $t$ .

We emphasize that the  $\forall$ -rule is an inference rule, not a simplification rule, because it creates and adds something, not replacing  $\forall x A(x)$  by something. Because of the  $\forall$ -rule, we may not be able to show the termination of the semantic tableau method.

**Definition 5.3.2.** *The  $\exists$ -rule is a simplification rule:*

$$(\exists) \frac{\exists x A(x)}{A(c)} \quad c \text{ is a new constant}$$

The  $\exists$ -rule works as follows: if  $\exists x A(x)$  appears in a leaf node  $n$ , we create a child node  $n'$  of  $n$ , copy every formula of  $n$  into  $n'$ , and replace  $\exists x A(x)$  in  $n'$  by  $A(c)$ . Later after you finish reading this chapter, you will know that the new constant is called *Skolem constant*, to remember Skolem's contribution to logic.

If the formulas are not in NNF (negation normal form), we also need the  $\neg\exists$ -rule, which is also an inference rule:

$$(\neg\exists) \frac{\neg\exists x A(x)}{\neg A(t)} \quad t \text{ is a ground term}$$

and the  $\neg\forall$ -rule, which is a simplification rule:

$$(\neg\forall) \frac{\neg\forall x A(x)}{\neg A(c)} \quad c \text{ is a new constant}$$

**Example 5.3.3.** To show  $(\exists x\forall y p(x, y)) \rightarrow \forall y\exists x p(x, y)$  is valid by semantic tableaux, we show that the negation of the formula has a closed tableau.

$\epsilon :$	$\neg((\exists x\forall y p(x, y)) \rightarrow \forall y\exists x p(x, y))$	$\alpha \neg \rightarrow$
$1 :$	$\exists x\forall y p(x, y), \neg\forall y\exists x p(x, y)$	$\exists$
$11 :$	$\forall y p(c_1, y), \neg\forall y\exists x p(x, y)$	$\neg\forall$
$111 :$	$\forall y p(c_1, y), \neg\exists x p(x, c_2)$	$\forall$
$1111 :$	$\forall y p(c_1, y), \neg\exists x p(x, c_2), p(c_1, c_2)$	$\neg\exists$
$11111 :$	$\forall y p(c_1, y), \neg\exists x p(x, c_2), p(c_1, c_2), \neg p(c_1, c_2)$	<i>closed</i>

**Example 5.3.4.** Let  $A$  be  $B \wedge C \wedge \wedge(\forall x\forall y \neg q(x) \vee \neg q(y))$  where  $B = (\forall x p(x, s(x)))$  and  $C = (\forall x\forall y q(x) \vee q(y))$ . If we choose to work on  $q$  first, we will find a closed tableau for  $A$ , thus  $A$  is unsatisfiable.

$\epsilon :$	$B \wedge C \wedge (\neg\exists x q(x))$	$\alpha \wedge$
$1 :$	$B, C \wedge (\neg\exists x q(x))$	$\alpha \wedge$
$11 :$	$B, C, (\neg\exists x q(x))$	$\forall$
$111 :$	$B, C, (\neg\exists x q(x)), (\forall y q(c) \vee q(y))$ ,	$\forall$
$1111 :$	$B, C, (\neg\exists x q(x)), (\forall y q(c) \vee q(y)), q(c) \vee q(c)$ ,	$\neg\exists$
$11111 :$	$B, C, (\neg\exists x q(x)), (\forall y q(c) \vee q(y)), q(c) \vee q(c), \neg q(c)$ ,	$\beta \vee$
$111111 :$	$B, C, (\neg\exists x q(x)), (\forall y q(c) \vee q(y)), q(c), \neg q(c)$ ,	<i>closed</i>
$111112 :$	$B, C, (\neg\exists x q(x)), (\forall y q(c) \vee q(y)), q(c), \neg q(c)$ ,	<i>closed</i>

If we choose to work on  $p$  first, then the tableau will never terminate:

$\epsilon :$	$(\forall x p(x, s(x))), (\forall x\forall y q(x) \vee q(y)), (\neg\exists x q(x))$	$\forall$
$1 :$	$(\forall x p(x, s(x))), \dots, p(c, s(c))$ ,	$\forall$
$11 :$	$(\forall x p(x, s(x))), \dots, p(c, s(c)), p(s(c), s(s(c)))$	$\forall$
$111 :$	$(\forall x p(x, s(x))), \dots, p(c, s(c)), p(s(c), s(s(c))), p(s(s(c)), s(s(s(c))))$	$\forall$
$1\dots 1 :$	$\dots$	

This example shows clearly that the nondeterminism in the application of the rules (i.e., choosing a leaf node, choosing a rule, or choosing a term in the  $\forall$  rule) does matter for obtaining a closed tableau. If we use a fair strategy for the application of the rules, then a closed tableau is guaranteed to be found when the input formula is unsatisfiable, as indicated by the following result (without a proof).

**Theorem 5.3.5.** *A formula  $A$  is unsatisfiable iff  $A$  has a closed tableau, which can be found by a fair strategy.*

### 5.3.2 Natural Deduction

In natural deduction for propositional logic, we have introduction rule and elimination rule for each logical operator. These rules will be inherited, plus the new rules for the quantifiers over the sequents:

<i>op</i>	<b>introduction</b>	<b>elimination</b>
$\forall$	$(A(x) \mid \alpha) \vdash (\forall y A(x) \mid \alpha)$	$(\forall x A(x) \mid \alpha) \vdash (A(t) \mid \alpha)$
$\exists$	$(A(t) \mid \alpha) \vdash (\exists x A(x) \mid \alpha)$	$(\exists x A(x) \mid \alpha) \vdash (A(c) \mid \alpha)$

where  $x$  and  $y$  are variables,  $t$  is a term and  $c$  is a new constant. From these rules, we can see that free variables have the same function as universally quantified variables. We use two examples to illustrate these rules.

**Example 5.3.6.** Let  $A = \{\forall x p(x), \forall x q(x)\}$ . A proof of  $A \models \forall y p(y) \wedge q(y)$  in natural deduction is given below.

1.  $(\forall x p(x))$             assumed
2.  $(\forall x q(x))$             assumed
3.  $(p(y))$                  $\forall_E$  from 1
4.  $(q(y))$                  $\forall_E$  from 2
5.  $(p(y) \wedge q(y))$         $\wedge_I$  from 3, 4
6.  $(\forall y p(y) \wedge q(y))$   $\forall_I$  from 5

By Theorem 5.2.16,  $(\forall x p(x) \wedge (\forall x q(x)) \rightarrow \forall y p(y) \wedge q(y))$  is valid.

**Example 5.3.7.** Show that  $(\exists x p(x) \vee q(x)) \models (\exists x p(x)) \vee (\exists x q(x))$  in natural deduction.

1.  $(\exists x p(x) \vee q(x))$        assumed
2.  $(p(c) \vee q(c))$          $\exists_E$  from 1
3.  $(p(c) \mid q(c))$          $\vee_E$  from 2
4.  $(\exists x p(x) \mid q(c))$      $\exists_I$  from 3
5.  $(\exists x p(x) \mid (\exists x q(x)))$   $\exists_I$  from 4
6.  $(\exists x p(x) \vee (\exists x q(x)))$   $\vee_I$  from 5

By Theorem 5.2.16,  $(\exists x p(x) \vee q(x)) \rightarrow ((\exists x p(x)) \vee (\exists x q(x)))$  is valid.

## 5.4 Conjunctive Normal Form

We just introduced an extension of semantic tableau and natural deduction to first-order logic and illustrated the use of these two methods for proving theorems in

first-order logic. In practice, these two methods are not as effective as the resolution method, which is the focus of the next chapter.

Recall that in propositional logic, the resolution is a refutational method and the input formula must be in clausal form. This is still true for first-order logic. In this section, we discuss how to convert first-order formulas into clausal form.

**Definition 5.4.1.** *In first-order logic, a literal is either an atom (positive literal) or the negation of an atom (negative literal); a clause is a disjunction of literals; a formula is in CNF if it is a conjunction of clauses.*

From the definition, if we replace each atom in a clause by a propositional variable, we obtain a clause in propositional logic. Since clauses use only three logical operators, i.e.,  $\neg$ ,  $\vee$  and  $\wedge$ , we need to get rid of  $\rightarrow$ ,  $\leftrightarrow$ ,  $\oplus$ , etc., as we did in propositional logic. We also need to get rid of quantifiers and we will talk about it now.

#### 5.4.1 Prenex Normal Form

**Definition 5.4.2.** *A formula  $A$  is in prenex normal form (PNF) if*

$$A = Q_1x_1Q_2x_2 \cdots Q_kx_kB(x_1, x_2, \dots, x_k),$$

where  $Q_i \in \{\forall, \exists\}$ , and  $B$  does not contain any quantifiers.

The idea of prenex normal form is to move all the quantifiers to the top of the formula. For example,  $\forall x \forall y p(x, y) \vee q(x)$  is a PNF but  $\forall x (\forall y p(x, y)) \vee q(x)$  is not.

To make the discussion easier, we assume that the formulas contain only  $\neg$ ,  $\wedge$ ,  $\vee$ , plus the two quantifiers, as operators. A quantifier which is not at the top must have another operator immediately above it: in this case the operator will be one of the three operators:  $\neg$ ,  $\wedge$  or  $\vee$ . Since the quantifier can be in either parameter position of  $\wedge$  or  $\vee$ , and there are two quantifiers, there are 10 cases to consider. If we consider  $\wedge$  and  $\vee$  are commutative, we need six rules, which are provided in Proposition 5.2.18:

Equivalence Rules for PNF

1.  $\neg \forall x A(x) \equiv \exists x \neg A(x)$
2.  $\neg \exists x A(x) \equiv \forall x \neg A(x)$
3.  $(\forall x A(x)) \wedge B \equiv (\forall x A(x) \wedge B)$  if  $x$  is not free in  $B$

4.  $(\forall x A(x)) \vee B \equiv (\forall x A(x) \vee B)$  if  $x$  is not free in  $B$
5.  $(\exists x A(x)) \wedge B \equiv (\exists x A(x) \wedge B)$  if  $x$  is not free in  $B$
6.  $(\exists x A(x)) \vee B \equiv (\exists x A(x) \vee B)$  if  $x$  is not free in  $B$

where the arguments of  $\vee$  and  $\wedge$  can be switched.

To meet the condition that “ $x$  is not free in  $B$ ”, it is necessary first to rename all the quantified variables so that no variable appears more than once following a quantifier.

In Chapter 2, we introduced the concept of “negation normal form”, where the negation appears only in literals. Using the first two of the above rules plus de Morgan’s laws and the double negation law, we may push the negation down, until all negations appear in literals. The resulting formulas are also said to be in *negation normal form* (NNF), a straightforward extension from propositional logic.

**Proposition 5.4.3.** *Every first-order formula can be transformed into an equivalent formula in NNF.*

**Example 5.4.4.** Given the following formula

$$(\forall x \exists y p(x, y)) \wedge \neg(\forall x \forall y q(x, y) \wedge q(y, x)),$$

we may choose to rename the variable names first so that every quantified variable has a different name:

$$(\forall x \exists y p(x, y)) \wedge \neg(\forall z \forall w q(z, w) \wedge q(w, z))$$

Applying equivalence rule 1 twice,

$$(\forall x \exists y p(x, y)) \wedge \exists z \exists w \neg(q(z, w) \wedge q(w, z))$$

Now we may apply the deMorgan’s law to get

$$(\forall x \exists y p(x, y)) \wedge \exists z \exists w \neg q(z, w) \vee \neg q(w, z))$$

which is in NNF; this step is not required for obtaining PNF. We have two options now: either choose equivalence rule 3 or rule 5. To use rule 3, we have

$$\forall x (\exists y p(x, y)) \wedge \exists z \exists w \neg q(z, w) \vee \neg q(w, z))$$

To apply rule 5 three times, we have

$$\forall x \exists y \exists z \exists w p(x, y) \wedge (\neg q(z, w) \vee \neg q(w, z)),$$

which is a PNF.

Note that the order of quantified variables in this formula is  $xyzw$ . There are five other possible outcomes where the order is  $xzyw$ ,  $xzwy$ ,  $zxyw$ ,  $zxwy$ , or  $zwx y$ . If we are allowed to switch the quantifiers of the same type (the last two in Proposition 5.2.18), we obtain more equivalent formulas where the order is  $xywz$ ,  $xwyz$ ,  $xwzy$ ,  $wxyz$ ,  $wxzy$ , or  $wzxy$ . These formulas are all equivalent as we use the equivalence relations in obtaining PNF.

**Proposition 5.4.5.** *Every first-order formula can be transformed into an equivalent PNF.*

*Proof.* By Proposition 5.2.12, renaming quantifiers variable names preserve the equivalence. The equivalence rules for PNF are terminating (an exercise problem) and preserve the equivalence. Apply these rules repeatedly, we will obtain the desired result.  $\square$

#### 5.4.2 Skolemization

Skolemization is the process in which we remove all the quantifiers, both universal and existential, leaving a formula with only free variables. This process is made easier by having the original formula in prenex normal form, and we assume this in this section.

Removing the universal quantifiers while preserving meaning is easy since we assume that the meaning of a free variable is the same as a universally quantified variable: The formula will be true for any value of a free variable. We need only to discard the universal quantifiers from a PNF.

Doing the same for an existentially quantified variable is not possible because a free variable can have only one meaning. The key idea is to introduce a new function to replace that variable. These are called *Skolem functions*. When the function are nullary they are called *Skolem constants*.

**Example 5.4.6.** Consider the formula  $\forall x \exists y y^2 \leq x \wedge \neg(y+1)^2 \leq x$ . This asserts the existence of a maximal number  $y$  whose square is no more than  $x$  for any value of  $x$ . We may drop the quantifier  $\forall$  so that  $x$  becomes free and keep the meaning of the formula. The same effect could be achieved by asserting the existence of a function  $f$ , depending on  $x$ , satisfying

$$f(x)^2 \leq x \wedge \neg(f(x)+1)^2 \leq x.$$

The function  $f$  is indeed a Skolem function. We might then use our knowledge of algebra to prove that  $f(x) = \lfloor \sqrt{x} \rfloor$ , but this is beyond the scope of our present concerns.

In general, the parameters to a Skolem function will be free variables, which have the same meaning as universally quantified variables, thus allowing a different value for each different values of free variables.

**Algorithm 5.4.7.**  $Sko(A)$  takes as input a formula  $A$  and returns a formula without a formula without quantifiers. Procedure  $PrenexNF(A)$  will return a formula in PNF and equivalent to  $A$ .  $free(B)$  returns the list of free variables in  $B$ .

```

proc  $Sko(A)$ 
1    $A := PrenexNF(A)$ ;
2   while  $A$  has a quantifier do
3     if ( $A = \exists x B$ )  $A := B[x \leftarrow f(free(B))]$ ; //  $f$  is a new Skolem function
4     else if ( $A = \forall x B$ )  $A := B$ ;
5   return  $A$ ;

```

In the algorithm  $Sko(A)$ , we see that a formula is first converted into PNF, then each quantifier is peeled off from left to right, until they are all eliminated.

**Example 5.4.8.** In Example 5.4.4, let

$$A = \forall x \exists y \exists z \exists w p(x, y) \wedge (\neg q(z, w) \vee \neg q(w, z)),$$

then  $Sko(A) = p(x, f_1(x)) \wedge (\neg q(f_2(x), f_3(x)) \vee \neg q(f_3(x), f_2(x)))$ . For

$$B = \exists z \exists w \forall x \exists y p(x, y) \wedge (\neg q(z, w) \vee \neg q(w, z)),$$

$Sko(B) = p(x, f_3(x)) \wedge (\neg q(c_1, c_2) \vee \neg q(c_2, c_1))$ . We know that  $A \equiv B$ .  $Sko(A) \not\equiv Sko(B)$ , however, because they contain different function symbols.

The above example illustrates that Skolemization does not preserve the equivalence. The Skolemization process, applied to a formula in prenex normal form, terminates with a weakly equivalent formula. To see that the process terminates is easy, since the process of converting a formula into PNF is terminating and the number of quantifiers is reduced in the body of the while loop of  $Sko(A)$ . The loop continues while a quantifier remains at the top of the formula, so the final formula will have no quantifiers. The input and output formulas are not necessarily equivalent, but their meanings are closely related, and we will capture this relationship with the notion “equally satisfiable”.

Recall that in propositional logic (Definition 3.3.6), we say two formulas,  $A$  and  $B$ , are equally satisfiable, denoted by  $A \approx B$ , if  $\mathcal{M}(A) = \emptyset$  whenever  $\mathcal{M}(B) = \emptyset$  and vice versa. This relation is weaker than the logical equivalence, which requires  $\mathcal{M}(A) = \mathcal{M}(B)$ . We will reuse  $\approx$  for first-order formulas.

It is easy to see that elimination of universal quantifiers preserves equivalence. The difficulties arise with the introduction of Skolem functions/constants.

**Lemma 5.4.9.** *Let  $A$  be  $\forall x \exists y B(x, y)$ . Then there exists a function  $f(x)$  such that (a)  $\forall x B(x, f(x)) \models A$ ; and (b)  $\forall x B(x, f(x)) \approx A$ .*

*Proof.* (a): For any model  $I$  of  $\forall x B(x, f(x))$ , and for every  $d \in D^I$ , it must be the case that  $I(B(x, f(x)), \{x \leftarrow d\}) = 1$ . Let  $f^I(d) = e \in D^I$ , where  $f^I : D^I \rightarrow D^I$  is the interpretation of  $f$  in  $I$ , then  $I(B(x, y), \{x \leftarrow d, y \leftarrow e\}) = 1$ . So we have

$$I(\exists y B(x, y), \{x \leftarrow d\}) = 1$$

Since  $d$  is arbitrary in  $D^I$ ,  $I(\forall x \exists y B(x, y)) = 1$ . In other words,  $I$  is a model of  $A$ .

(b): Since (a) says “if  $\forall x B(x, f(x))$  is satisfiable, so is  $A$ ”, we need only to show that if  $A$  is satisfiable, so is  $\forall x B(x, f(x))$ . Suppose  $I = (D, R, G)$  is a model of  $A$ . Then for any  $d \in D$ , there exists  $e \in D$  such that

$$I(B(x, y), \{x \leftarrow d, y \leftarrow e\}) = 1$$

Let  $g : D \rightarrow D$  such that  $f(d) = e$ , and  $I' = (D, R, G \cup \{g\})$  with  $g = f^I$ . Then  $I'$  is a model of  $\forall x B(x, f(x))$ .  $\square$

**Corollary 5.4.10.** *Let  $A = \forall x_1 \forall x_2 \cdots \forall x_k \exists y B(x_1, x_2, \dots, x_k, y)$ . Then there exists a function  $f(x_1, x_2, \dots, x_k)$  such that*

- (a)  $\forall x_1 \forall x_2 \cdots \forall x_k B(x_1, x_2, \dots, x_k, f(x_1, x_2, \dots, x_k)) \models A$ ; and  
 (b)  $\forall x_1 \forall x_2 \cdots \forall x_k B(x_1, x_2, \dots, x_k, f(x_1, x_2, \dots, x_k)) \approx A$ .

*Proof.* Replace  $x$  by  $x_1, x_2, \dots, x_k$  in the proof of Lemma 5.4.9.  $\square$

**Theorem 5.4.11.** (a)  $Sko(A) \models A$ ; and (b)  $Sko(A) \approx A$ .

*Proof.* Since a free variable and a universally quantified variable have the same interpretation, drop a universal quantifier does not change the semantics of a formula. When we remove an existential quantifier, we introduce a Skolem function. By Corollary 5.4.10, the resulting formula is equally satisfiable to the original formula. Applying Corollary 5.4.10 to every existential quantifier, we have the desired result.  $\square$

By Corollary 5.4.10 (a),  $A$  is a logical consequence of  $Sko(A)$ , i.e.,  $Sko(A) \models A$ , but not  $A \models Sko(A)$ . However, we have  $A \approx Sko(A)$ , this is what we mean by “weakly equivalent”, and we will see that it is sufficient for our purposes.

### 5.4.3 Skolemizing Non-Prenex Formulas

Although Skolemization is easy to explain and justify on formulas in PNF, it can actually be applied directly to any formula. The process involved is not easily explained as in the case of PNF. Instead we will provide a recursive algorithm for Skolemization which takes two parameters: The first parameter is a formula to be Skolemized, which must be in NNF (negation normal form), and the second parameter is the list of free variables collected so far.

**Algorithm 5.4.12.**  $Skol(A, V)$  takes as input a formula  $A$  which is in NNF and a list  $V$  of variables which are free in  $A$ , and returns a formula without quantifiers.

```

proc  $Skol(A, V)$ 
1   if  $A$  is a literal or  $A = \top$  or  $A = \perp$  return  $A$ ; // base case
2   if  $(A = B \text{ op } C)$  return  $Skol(B, V) \text{ op } Skol(C, V)$ ; //  $op$  is either  $\wedge$  or  $\vee$ 
3   if  $(A = \exists x B)$  return  $Skol(B[x \leftarrow f(V)], V)$ ; //  $f$  is a new Skolem function
4   if  $(A = \forall x B)$  return  $Skol(B, V \cup \{x\})$ ;

```

When  $A$  is a literal, i.e., either an atom, or the negation of an atom,  $A$  is returned. When the topmost operator is  $\vee$  or  $\wedge$ , the situation is simple:  $Skol$  is merely called recursively on each of the arguments of the operator. If the topmost symbol is  $\exists$ ,  $Skol$  introduces a new Skolem function  $f$  (a Skolem constant if  $V$  is empty) and replace every occurrence of  $x$  by  $f(V)$ . If the topmost symbol is  $\forall$ , the variable  $x$  is collected in the second parameter for the recursive call on  $B$ . Assuming computing  $B[x \leftarrow f(V)]$  takes linear time,  $Skol$  will take quadratic time in terms of the size of the input formula.

**Example 5.4.13.** Let  $A$  be the first NNF in Example 5.4.4

$$A = (\forall x \exists y p(x, y)) \wedge (\exists z \exists w \neg q(z, w) \vee \neg q(w, z))$$

Then  $Skol(A, ( ))$  calls  $Skol$  on the two arguments of  $\wedge$ :

- $Skol(\forall x \exists y p(x, y), \emptyset)$  calls  $Skol(\exists y p(x, y), (x))$  and returns  $p(x, f(x))$ ;
- $Skol(\exists z \exists w \neg q(z, w) \vee \neg q(w, z), \emptyset)$  returns  $\neg q(c_1, c_2) \vee \neg q(c_2, c_1)$ .

Hence  $Skol(A)$  returns  $p(x, f(x)) \wedge (\neg q(c_1, c_2) \vee \neg q(c_2, c_1))$ .

**Theorem 5.4.14.** (a)  $Skol(A) \models A$ , and (b)  $Skol(A) \approx A$ .

The proof of the above theorem is similar to that of Theorem 5.4.11 and is omitted here. As pointed out in Example 5.4.4, we have many different PNFs for the same formula and each PNF may give us a different result after Skolemization. All PNFs of the same formula are logically equivalent and all Skolemized formulas of the same formula are equally satisfiable.

In practice, we wish to produce Skolemized formulas where Skolem functions have minimal number of arguments, so that the ground terms of these formulas will be simpler to deal with. Note that the equivalence rules for obtaining PNFs try to extend the scope of each quantified variable. To obtain Skolem functions with minimal number of arguments, we may need to use these equivalence rules backward.

#### 5.4.4 Clausal Form

By CNF, we meant that the formula is a conjunction of clauses, and a clause is a disjunction of literals. To obtain CNF, a formula  $A$  goes through the following stages:

- stage 1: Convert  $A$  into NNF  $A'$ ;
- stage 2: Skolemize  $A'$  into quantifier-free  $B$ ;
- stage 3: Convert  $B$  into CNF  $C$ .

Stage 3 converts Skolemized NNF  $B$  into clauses, as we did in propositional logic to convert NNF into clauses. That is, use the distribution law  $A \vee (B \wedge C) \equiv (A \vee B) \wedge (A \vee C)$  repeatedly, until it no longer applies.

Stages 1 and 3 preserve the logical equivalence,  $\equiv$ , and stage 2 preserves only a weak equivalence,  $\approx$ . Thus, we have  $A \equiv A'$ ,  $A' \approx B$ , and  $B \equiv C$ . As a result,  $A \approx C$ . So  $A$  is satisfiable iff  $C$  is satisfiable.

**Example 5.4.15.** Consider the following formulas and their intended meanings:

- $A = \forall x \text{ human}(x) \rightarrow \text{mortal}(x)$ : every human is mortal.
- $B = \forall x \text{ iowan}(x) \rightarrow \text{human}(x)$ : every Iowan is human.
- $C = \forall x \text{ iowan}(x) \rightarrow \text{mortal}(x)$ : every Iowan is mortal.

To show  $A \wedge B \models C$ , we convert  $A \wedge B \wedge \neg C$  into a set  $S$  of four clauses, where  $c_1$  comes from  $A$ ,  $c_2$  from  $B$ , and  $c_3$  and  $c_4$  from  $\neg C$ :

$$\begin{aligned} c_1 &: (\neg human(x) \mid mortal(x)), \\ c_2 &: (\neg iowan(y) \mid human(y)), \\ c_3 &: (iowan(c)), \\ c_4 &: (\neg mortal(c)) \end{aligned}$$

An instance of  $c_1$  is  $c'_1 = (\neg human(c) \mid mortal(c))$ , and an instance of  $c_2$  is  $c'_2 = (\neg iowan(c) \mid human(c))$ . Resolution between  $c_3$  and  $c'_2$ , we obtain  $c_5 = (human(c))$ . Resolution between  $c_5$  and  $c'_1$ , we obtain the empty clause  $(\ )$ .

Our formulas are now in clausal form, and this is the normal form required by the resolution rule. Keep in mind the variables in a clause are assumed to be free and have the same semantics as the universally quantified variables, so that they can be renamed. These variables are different from clauses to clauses, even though they use the same name. A practice is to give different names for variables in different clauses. For example, in clause  $C_i$ , we use  $x_i, y_i, z_i$ , etc. This way, we ensure that no variable appears in more than one clause.

Stage 3 may blow up the sizes of the output formulas, while stage 2 outputs a formula of the linear size of the input formula. In propositional logic, to obtain a set of clauses of linear size of the original formula, we may introduce new variables for each occurrence of binary operators and the resulting clauses are 3CNF, which are not equivalent, but are equally satisfiable, to the original formula. We may do the same in stage 3 so that the final clauses are of linear size of the original formula.

Clausal form might seem a highly constrained and unnatural representation for mathematical knowledge, and it has been criticized on those grounds. We will see that it is in fact surprisingly natural in many cases. In fact, a problem can be expressed as a set of constraints over some variables, such as constraint satisfaction problems (CSP). A constraint can be expressed in CNF and a set of constraints give us a set of clauses.

We have seen that any predicate calculus formula can be converted into an equally satisfiable formula in CNF, which is useful for proving theoretical results by resolution.

#### 5.4.5 Herbrand Models for CNF

When we consider interpretations for a set of formulas, there are simply too many ways to choose a set as the universe, too many ways to choose a relation for a predicate symbol, too many ways to choose a function for a function symbol, and

use the elements of each universe in all possible ways for variables. We can improve on this considerably because there is a class of interpretations, called *Herbrand interpretations*, which can stand for all the others. The Herbrand interpretations share the same universe, called the *Herbrand universe*. What must this universe contain? For each  $n$ -ary function,  $f$ , in the formula (including constants) and  $n$  elements,  $t_1, t_2, \dots, t_n$ , in the universe, it must contain the result of applying  $f$  to these  $n$  elements. This can be achieved by a simple syntactic device; we put all the ground terms in the universe and let each term denote the result of applying its function to its parameters. We will call this syntactic device “self denotation”. This way, function symbols are bounded to a unique interpretation over the universe.

**Definition 5.4.16.** *Given a CNF formula  $A$  in a first-order language  $L = (P, F, X, Op)$ , the Herbrand universe of  $A$  consists of all the ground terms, i.e., the set  $T(F)$ ; the Herbrand base is the set of all ground atoms, denoted by  $B(P, F)$ . A Herbrand interpretation is a mapping  $I : B(P, F) \rightarrow \{1, 0\}$ .*

Note that  $T(F)$  is the set of variable-free terms that we can make from the constants and functions in the formula. If there are no constants in  $F$ ,  $T(F)$  would be empty. To avoid this case, we assume that there must exist at least one constant in  $F$ , say  $a$ . From now on it will be useful to distinguish the constants from non-constant functions, so by “function” we mean non-nullary function. If  $F$  does not contain function symbols, then  $T(F)$  will be a finite set of constants.

The Herbrand base  $B(P, F)$  consists of all the atoms like  $p(t_1, t_2, \dots, t_n)$ , where  $p$  is an  $n$ -ary predicate symbol in  $P$  and  $t_i \in T(F)$ . If  $n = 0$  for all predicates, then  $B(P, F) = P$ , reduced to propositional logic. Otherwise,  $B(P, F)$  is finite iff  $T(F)$  is finite. For convenience, we may represent  $I$  by a set  $H$  of ground literals such that for any  $g \in B(P, F)$ ,  $g \in H$  iff  $I(g) = 1$  and  $\bar{g} \in H$  iff  $I(g) = 0$ .

**Example 5.4.17.** In Example 5.2.13, we showed that the set  $S$  has no finite models:

$$S = \{\forall x \exists y p(x, y), \forall x \neg p(x, x), \forall x \forall y \forall z p(x, y) \wedge p(y, z) \rightarrow p(x, z)\}.$$

Note that  $S$  has no function symbols. Converting  $S$  to CNF, we obtain the following set  $S'$  of clauses:

$$S' = \{c_1 : (p(x, f(x))), c_2 : (\overline{p(x, x)}), c_3 : (\overline{p(x, y)} \mid \overline{p(y, z)} \mid p(x, z))\}.$$

Adding a constant  $a$ ,  $F = \{f, a\}$  and  $P = \{p\}$ , then the Herbrand universe of  $S'$  is

$$T(F) = \{a, f(a), f(f(a)), \dots, f^i(a), \dots\},$$

where  $f^0(a) = a$  and  $f^{i+1}(a) = f(f^i(a))$ . The Herbrand base of  $S'$  is

$$B(P, F) = \{p(a, a), p(a, f(a)), p(f(a), a), p(a, f(f(a))), \dots, p(f^i(a), f^j(a)), \dots\}.$$

A Herbrand interpretation of  $S'$  is

$$H = \{p(f^i(a), f^j(a)) \mid 0 \leq i < j\} \cup \{\overline{p(f^i(a), f^j(a))} \mid 0 \leq j \leq i\}.$$

**Proposition 5.4.18.** *A Herbrand interpretation defines a unique interpretation for a formula in CNF.*

*Proof.* Let  $A$  be a formula in CNF and  $H$  is a Herbrand interpretation. Define interpretation  $I = (D, R, G)$ , where  $D = T(F)$  (the Herbrand universe),

$$R = \{r_p \mid p \in P, r_p = \{\langle t_1, \dots, t_n \rangle \mid H(p(t_1, \dots, t_n)) = 1\}\}$$

and  $G = \{g_f \mid f \in F, g_f(t_1, \dots, t_n) = f(t_1, \dots, t_n)\}$ , then  $I$  is an interpretation as defined in Definition 5.2.1 and is uniquely defined by  $H$ .  $\square$

Since  $D$  and  $G$  are fixed in  $I$ ,  $H$  is the only part which affects  $I$ . That is, the only way in which Herbrand interpretations differ from one to another is over the meanings of predicate symbols. The above proposition gives us the justification for calling  $H$  “an interpretation”.

Once we have an interpretation, we can compute the truth value of a formula  $A$  in this interpretation using a procedure like *eval* (Procedure 5.2.3). If  $A$  is true under a Herbrand interpretation  $H$ , we say  $H$  is a *Herbrand model*. Since a CNF  $A$  can be represented by a set  $S$  of clauses, and the variables in each clause have the same meaning as universally quantified variables, for any interpretation  $I$ ,  $I(A) = I(S) = \bigcap_{C \in S} I(C)$ , and for each clause  $C$ ,  $I(C) = I(\forall x_1 \dots \forall x_k C)$ , where  $x_1, \dots, x_k$  are the free variables of  $C$ .

The definition of Herbrand interpretation gives us a convenient way of expressing these meanings by assigning the truth value to each ground atom in the Herbrand base. Of course, the truth values of some atoms are irrelevant if these atoms cannot be instances of any predicate in the formula, as the truth of the formula does not depend on these atoms. From now on, when we say “a set of clauses has a model” we meant “it has a Herbrand model”.

## 5.5 Exercise Problems

- Using a first-order language with variables ranging over people, and predicates  $trusts(x, y)$ ,  $politician(x)$ ,  $crazy(x)$ ,  $know(x, y)$ , and  $related(x, y)$ , and  $rich(x)$ , write down first-order sentences asserting the following:

- (a) Nobody trusts a politician.
- (b) Anyone who trusts a politician is crazy.
- (c) Everyone knows someone who is related to a politician.
- (d) Everyone who is rich is either a politician or knows a politician.
2. Let the meaning of  $taken(x, y)$  be “student  $x$  has taken CS course  $y$ ”. Please write the meaning of the following formulas:
- (a)  $\exists x \exists y taken(x, y)$
- (b)  $\exists x \forall y taken(x, y)$
- (c)  $\forall x \exists y taken(x, y)$
- (d)  $\exists y \forall x taken(x, y)$
- (e)  $\forall y \exists x taken(x, y)$
3. Let  $c(x)$  mean “ $x$  is a criminal” and let  $s(x)$  mean “ $x$  is sane”. Then formalize each of the following sentences.
- (a) All criminals are sane.
- (b) Every criminal is sane.
- (c) Only criminals are sane.
- (d) Some criminals are sane.
- (e) There is a criminal that is sane.
- (f) No criminal is sane.
- (g) Not all criminals are sane.
4. Given the predicates  $male(x)$ ,  $female(x)$ ,  $parent(x, y)$  ( $x$  is a parent of  $y$ ), define the following predicates in terms of  $male$ ,  $female$ ,  $parent$ , or other predicates defined from these three predicates, in first-order logic:
- (a)  $father(x, y)$ :  $x$  is the father of  $y$ .
- (b)  $mother(x, y)$ :  $x$  is the mother of  $y$ .
- (c)  $ancesotor(x, y)$ :  $x$  is an ancestor of  $y$ .
- (d)  $descendant(x, y)$ :  $x$  is a descendant of  $y$ .
- (e)  $son(x, y)$ :  $x$  is a son of  $y$ .
- (f)  $daughter(x, y)$ :  $x$  is a daughter of  $y$ .

(g)  $sibling(x, y)$ :  $x$  is a sibling of  $y$ .

5. Let  $Set$  be defined by the following BNF:

$$\langle Set \rangle ::= \emptyset \mid adjoin(\langle Any \rangle, \langle Set \rangle)$$

where  $adjoin(e, s)$  adjoins an element  $e$  into the set  $s$ . Please define the following predicates in terms of  $\emptyset$  and  $adjoin$  in first-order logic:

- (a)  $isSet(x)$ :  $x$  is an object of  $Set$ .
- (b)  $member(x, s)$ :  $x$  is a member of set  $s$ .
- (c)  $subset(s, t)$ :  $s$  is a subset of set  $t$ .
- (d)  $equal(s, t)$ :  $s$  and  $t$  contain the same elements.

6. Decide if the following formulas are valid, not valid but satisfiable, or unsatisfiable, using the definition of interpretations:

- (a)  $p(a, a) \vee \neg \exists x p(x, a)$
- (b)  $\neg p(a, a) \vee \exists x p(x, a)$
- (c)  $\neg q(a) \vee \forall x q(x)$
- (d)  $\neg(\exists x q(x)) \vee \forall x q(x)$
- (e)  $\exists x q(x) \vee \neg \forall x q(x)$

7. Prove the following equivalence relations:

- (a)  $\neg \exists x p(x) \equiv \forall x \neg p(x)$
- (b)  $\exists x p(x) \vee q(x) \equiv (\exists x p(x)) \vee (\exists x q(x))$
- (c)  $(\exists x p(x)) \wedge q \equiv \exists x (p(x) \wedge q)$
- (d)  $\exists x \exists y r(x, y) \equiv \exists y \exists x r(x, y)$
- (e)  $\exists x p(x) \rightarrow q(x) \equiv (\forall x p(x)) \rightarrow (\exists x q(x))$

8. Prove that the following entailment relations are true, but the converse of the relation is not true.

- (a)  $(\forall x p(x)) \vee (\forall x q(x)) \models \forall x p(x) \vee q(x)$
- (b)  $\exists x p(x) \wedge q(x) \models (\exists x p(x)) \wedge (\exists x q(x))$
- (c)  $(\forall x p(x)) \rightarrow (\forall x q(x)) \models \forall x p(x) \rightarrow q(x)$
- (d)  $(\exists x p(x)) \rightarrow (\exists x q(x)) \models \exists x p(x) \rightarrow q(x)$

- (e)  $(\exists x p(x)) \rightarrow (\forall x q(x)) \models \forall x p(x) \rightarrow q(x)$
- (f)  $\forall x p(x) \rightarrow q(x) \models (\exists x p(x)) \rightarrow (\exists x q(x))$
- (g)  $\forall x p(x) \rightarrow q(x) \models (\forall x p(x)) \rightarrow (\exists x q(x))$
- (h)  $\forall x p(x) \leftrightarrow q(x) \models (\forall x p(x)) \leftrightarrow (\forall x q(x))$
- (i)  $\forall x p(x) \leftrightarrow q(x) \models (\exists x p(x)) \leftrightarrow (\exists x q(x))$

9. Use the semantic tableau method to show the entailment relations are true in the previous problem.
10. Provide the pseudo-code for the algorithm  $PrenexNF(A)$  and show that it runs in  $O(n)$  time when  $A$  is in NNF, where  $n$  is the size of  $A$ .
11. Show that the algorithm  $Sko(A)$  runs in  $O(n^2)$  time, where  $A$  is in NNF and  $n$  is the size of  $A$ .
12. Show that the algorithm  $Skol(A)$  runs in  $O(n^2)$  time, where  $A$  is in NNF and  $n$  is the size of  $A$ .
13. Convert each one of the following formulas into a set of clauses:

- (a)  $\neg((\forall x p(x)) \vee (\forall x q(x)) \rightarrow \forall x p(x) \vee q(x))$
- (b)  $\neg((\exists x p(x) \wedge q(x)) \rightarrow (\exists x p(x)) \wedge (\exists x q(x)))$
- (c)  $\neg(((\forall x p(x)) \rightarrow (\forall x q(x))) \rightarrow (\forall x p(x) \rightarrow q(x)))$
- (d)  $\neg(((\exists x p(x) \rightarrow q(x)) \rightarrow (\exists x p(x))) \rightarrow (\exists x q(x)))$
- (e)  $\neg((\exists x p(x)) \rightarrow (\forall x q(x)) \rightarrow (\forall x p(x) \rightarrow q(x)))$
- (f)  $\neg(((\forall x p(x) \rightarrow q(x)) \rightarrow (\exists x p(x))) \rightarrow (\exists x q(x)))$
- (g)  $\neg((\forall x p(x) \rightarrow q(x)) \rightarrow ((\forall x p(x)) \rightarrow (\exists x q(x))))$
- (h)  $\neg(((\forall x p(x) \leftrightarrow q(x)) \rightarrow (\forall x p(x))) \leftrightarrow (\forall x q(x)))$
- (i)  $\neg(((\forall x p(x) \leftrightarrow q(x)) \rightarrow (\exists x p(x))) \leftrightarrow (\exists x q(x)))$

14. Given a set  $S$  of clauses,  $S = \{(p(f(x), c)), (\neg p(f(w), z) \mid p(f(f(w)), f(z))), (\neg p(y, f(y)))\}$ , please answer the following questions: (a) What is the Herbrand universe of  $S$ ? (b) What is the Herbrand base of  $S$ ? (c) What is the minimum Herbrand model of  $S$  (no subset of this model is a model)? (d) What is the maximum Herbrand model of  $S$  (no superset of this model is a model)?

## CHAPTER 6

### UNIFICATION AND RESOLUTION

Resolution is known to be a refutational proof method from Chapter 3. By refutational method, we meant that the method will show that the input formula is unsatisfiable. If we like to show formula  $A$  is valid, we need to convert  $\neg A$  into CNF and show that the clauses are unsatisfiable. To show that  $A \models B$ , we need to convert  $A \wedge \neg B$  into CNF and show that the clauses are unsatisfiable. Note that  $A$  and  $B$  must be closed formulas. For example, if we want to express that  $p(x, y)$  is a commutative relation, we may use the formula  $p(x, y) \rightarrow p(y, x)$ , where  $x, y$  are assumed to be arbitrary values. The negation of this formula is  $\neg(\forall x \forall y p(x, y) \rightarrow p(y, x))$ , not  $\neg(p(x, y) \rightarrow p(y, x))$ . The CNF of the former is  $p(a, b) \wedge \neg p(b, a)$ , where  $a$  and  $b$  are Skolem constants; the CNF of the latter is  $p(x, y) \wedge \neg p(y, x)$ . We can substitute  $y$  by  $x$  in this CNF and obtain two unit clauses  $(p(x, x))$  and  $(\overline{p(x, x)})$ . These two clauses will generate the empty clause by resolution. The process of finding a substitution to make two atoms identical is called “unification”, which is needed in the definition of the resolution rule for first-order formulas.

#### 6.1 Unification

In Example 5.4.15, when we apply the resolution rule on two clauses, we need to find their instances such that one instance contains  $A$  and the other instance contains  $\neg A$ , where  $A$  is an atom. A clause has infinite many instances and the unification process helps us to find the right instances for resolution.

##### 6.1.1 Substitutions and Unifiers

Given a first-order language  $L = (P, F, X)$ , a mapping  $\sigma : X \rightarrow T(F, X)$  is said to be a *substitution*. If  $\sigma(x) \neq x$ , we say  $x$  is *affected* by  $\sigma$ . If  $\sigma(x)$  is a distinct variable for every affected  $x$ ,  $\sigma$  is said to be a *renaming*.

When displaying  $\sigma$ , we only display the pair  $x \leftarrow \sigma(x)$  for affected variable  $x$ , i.e., when  $x \neq \sigma(x)$ . For example, if  $X = \{x, y\}$ , then  $x$  is affected by  $\sigma_1 = [x \leftarrow f(y)]$  and  $\sigma_2 = [x \leftarrow f(x)]$ ; both  $x$  and  $y$  are affected by  $\sigma_3 = [x \leftarrow f(y), y \leftarrow a]$ .

A substitution  $\sigma : X \rightarrow T(F, X)$  can be extended to be a mapping  $\sigma : T(F, X) \rightarrow T(F, X)$  such that if  $t = x$ , then  $\sigma(t) = \sigma(x)$ ; if  $t = f(t_1, t_2, \dots, t_k)$ ,

then  $\sigma(t) = f(\sigma(t_1), \sigma(t_2), \dots, \sigma(t_k))$ . In a similar way,  $\sigma$  can be extended to be a mapping from formulas to formulas.

We used  $p(t)$  to denote an instance of  $p(x)$ , i.e.,  $p(t) = p(x)[x \leftarrow t]$ , where every occurrence of  $x$  is replaced by term  $t$ . Here,  $\sigma = [x \leftarrow t]$  is a substitution and  $p(t) = \sigma(p(x)) = p(x)\sigma$ , as we are free to write either  $p(x)\sigma$  or  $\sigma(p(x))$ .

**Definition 6.1.1.** (idempotent substitution)  $\sigma$  is said to be idempotent if for any term  $t$ ,  $t\sigma = (t\sigma)\sigma$ .

For example,  $\sigma_1 = [x \leftarrow f(y)]$  is idempotent;  $\sigma_2 = [x \leftarrow f(x)]$  is not idempotent, because  $x\sigma_2 = f(x)$  and  $(x\sigma_2)\sigma_2 = f(f(x))$ . Neither  $\sigma_3 = [x \leftarrow f(y), y \leftarrow a]$ , because  $x\sigma_3 = f(y)$  and  $(x\sigma_3)\sigma_3 = f(a)$ . Note that  $\sigma_3$  can be modified to make it idempotent:  $\sigma'_3 = [x \leftarrow f(a), y \leftarrow a]$ , while  $\sigma_2$  cannot.

**Proposition 6.1.2.** A substitution  $\sigma$  is idempotent iff for every  $x \in X$ , either  $\sigma(x) = x$  or  $x$  does not appear in  $\sigma(y)$  for any  $y \in X$ .

*Proof.* Given any  $\sigma$ , if an affected  $x$  appears in  $\sigma(y) = t(x)$ , suppose  $\sigma(x) = c$ , then  $y\sigma = t(x)$  and  $(y\sigma)\sigma = t(c)$ , so  $\sigma$  is not idempotent. Otherwise, for any term  $s$ ,  $s\sigma$  does not contain any variable affected by  $\sigma$ , so  $(s\sigma)\sigma = s\sigma$  and  $\sigma$  is idempotent.  $\square$

**Definition 6.1.3.** (unifier) Let  $s, t$  be either two terms of  $T(F, X)$  or two atoms. We say  $s$  and  $t$  are unifiable, if there exists an idempotent substitution  $\sigma : X \rightarrow T(F, X)$  such that  $s\sigma = t\sigma$ ;  $\sigma$  is said to be a unifier of  $s$  and  $t$ .

A set  $S$  of pairs of terms or atoms is unifiable if there exists a substitution  $\sigma$  such that for every pair  $(s_i, t_i) \in S$ ,  $s_i\sigma = t_i\sigma$  and  $\sigma$  is said to be a unifier of  $S$ .

**Example 6.1.4.** Let  $s = f(x, g(y))$  and  $t = f(g(z), x)$ , then  $s$  and  $t$  are unifiable with unifier  $\sigma = [x \leftarrow g(z), y \leftarrow z]$ , because

$$s\sigma = f(x, g(y))[x \leftarrow g(z), y \leftarrow z] = f(g(z), g(z)) = f(g(z), x)[x \leftarrow g(z), y \leftarrow z] = t\sigma.$$

In fact,  $s$  and  $t$  have infinite many unifiers like

$$[x \leftarrow g(k), y \leftarrow k, z \leftarrow k]$$

for any  $k \in T(F, X)$ .

In the above example, we see that a pair of terms may have multiple unifiers. How to compare these unifiers? We need some relation over unifiers and this relation is based on the composition of substitutions.

### 6.1.2 Combining Substitutions

**Definition 6.1.5.** (composition of substitution) *Given two substitutions  $\sigma$  and  $\theta$ , the composition of  $\sigma$  and  $\theta$  is the substitution  $\gamma$ , written  $\gamma = \sigma \cdot \theta$ , or simply  $\sigma\theta$ , where*

$$\gamma = [x \leftarrow t\theta \mid (x \leftarrow t) \in \sigma] \cup [y \leftarrow s \mid y = y\sigma, (y \leftarrow s) \in \theta]$$

We will assume that  $t\sigma\theta = (t\sigma)\theta = \theta(\sigma(t))$ . To denote the composition  $\sigma\theta$  on  $t$ , we write  $t(\sigma\theta)$ .

**Example 6.1.6.** Given  $\sigma = [x \leftarrow f(y), z \leftarrow f(u), v \leftarrow u]$  and  $\theta = [y \leftarrow g(a), u \leftarrow z, v \leftarrow f(f(a))]$ , then

$$\begin{aligned} \sigma\theta &= [x \leftarrow f(g(a)), y \leftarrow g(a), z \leftarrow f(z), u \leftarrow z, v \leftarrow z] \\ p(x, y, z, u, v)\sigma\theta &= p(f(y), y, f(u), u, u)\theta \\ &= p(f(g(a)), g(a), f(z), z, z) \\ &= p(x, y, z, u, v)(\sigma\theta) \end{aligned}$$

Note that the composition does not necessarily preserve the idempotent property as both  $\sigma$  and  $\theta$  are idempotent but  $z \leftarrow f(z)$  appears in  $\sigma\theta$  in the above example.

**Proposition 6.1.7.** *For any substitutions  $\sigma$  and  $\theta$ , and any term  $t$ ,  $t\sigma\theta = t(\sigma\theta)$ .*

*Proof.* For any  $x \in X$ , if  $x$  affected by  $\sigma$ , say  $\sigma(x) = s$ , then  $x\sigma\theta = s\theta = x(\sigma\theta)$ . If  $x$  is unaffected by  $\sigma$ , then  $x\sigma\theta = x\theta = x(\sigma\theta)$ .  $\square$

The above proposition states the property that applying two substitutions to a term in succession produces the same result as applying the composition to the term.

The composition operation is not commutative in general. For example, if  $\sigma = [x \leftarrow a]$  and  $\theta = [x \leftarrow b]$ , then  $\sigma\theta = \sigma$  and  $\theta\sigma = \theta$ . Thus,  $\sigma\theta \neq \theta\sigma$ . By definition, if  $x$  is affected by both  $\sigma$  and  $\theta$ , then the pair  $x \leftarrow \theta(x)$  is ignored in  $\sigma\theta$  and the pair  $x \leftarrow \sigma(x)$  is ignored in  $\theta\sigma$ .

**Proposition 6.1.8.** *The composition is associative, i.e., for any substitutions  $\sigma$ ,  $\theta$ , and  $\gamma$ ,  $(\sigma\theta)\gamma = \sigma(\theta\gamma)$ .*

The proof is left as an exercise.

### 6.1.3 Rule-Based Unification

**Definition 6.1.9. (most general unifier)** Let  $\sigma$  and  $\theta$  be unifiers of  $s$  and  $t$ . We say  $\sigma$  is more general than  $\theta$ , if there exists a substitution  $\gamma$  such that  $\theta = \sigma\gamma$ .  $\sigma$  is a most general unifier (mgu) of  $s$  and  $t$ , if  $\sigma$  is more general than any unifier of  $s$  and  $t$ .

**Example 6.1.10.** In Example 6.1.4,  $\sigma = [x \leftarrow g(z), y \leftarrow z]$  is a mgu of  $s = f(x, g(y))$  and  $t = f(g(z), x)$ . For any unifier  $\theta = [x \leftarrow g(t'), y \leftarrow t', z \leftarrow t']$ ,  $\theta = \sigma\gamma$ , where  $\gamma = [z \leftarrow t']$ . Another mgu of  $s$  and  $t$  is  $\sigma' = [x \leftarrow g(y), z \leftarrow y]$ ,  $\sigma = \sigma'[y \leftarrow z]$  and  $\sigma' = \sigma[z \leftarrow y]$ .

As illustrate by the above example, there are more than one mgu for a pair of terms; these mgus differ by the renaming of variables. We can say *the* most general unifier if the renaming of variables are considered.

Given a pair  $s$  and  $t$  of terms, how to decide if  $s$  and  $t$  are unifiable? If they are, how to compute their mgu? These questions will be answered by unification algorithms.

Let  $s \doteq t$  denote that “ $s$  and  $t$  are unifiable”. Suppose  $s = f(s_1, s_2, \dots, s_m)$  and  $t = g(t_1, t_2, \dots, t_n)$ . If  $f \neq g$  or  $m \neq n$ , then  $s \doteq t$  is false, i.e.,  $s$  and  $t$  are not unifiable. This is called *clash* failure.

If  $s \doteq t$  is true, so are  $s_i \doteq t_i$  for  $1 \leq i \leq m = n$ . In other words, if  $s_i \doteq t_i$  is false for some  $i$ , then  $s \doteq t$  is false. If for some  $i$ ,  $s_i = x$  and  $t_i = f(x)$  ( $x$  occurs in  $t_i$ ) then there exists no substitution  $\sigma$  such that  $x\sigma = f(x)\sigma$ . Hence  $x \doteq f(x)$  must be false. This is called *occur-check* failure, as  $x$  occurs in  $f(x)$ . If we accept  $x \leftarrow f(x)$ , then the unifier cannot be idempotent.

Rule-based unification algorithm uses a set of rules which transform  $s \doteq t$  into a substitution  $\sigma$  through a set of pairs  $\{s_i \doteq t_i\}$  and each rule preserves the unifiability of  $s \doteq t$ . These rules are give below.

**Definition 6.1.11.** *The unification rules for the rule-based unification algorithm are given below:*

- **decompose:**  $\langle S \cup \{f(s_1, s_2, \dots, s_m) \doteq f(t_1, t_2, \dots, t_m)\}, \sigma \rangle \longrightarrow \langle S \cup \{s_i \doteq t_i \mid 1 \leq i \leq m\}, \sigma \rangle$ .
- **clash:**  $\langle S \cup \{f(s_1, s_2, \dots, s_m) \doteq g(t_1, t_2, \dots, t_n)\}, \sigma \rangle \longrightarrow \langle S, fail \rangle$ .
- **occur-check:**  $\langle S \cup \{x \doteq t \mid x \neq t, x \text{ occurs in } t\}, \sigma \rangle \longrightarrow \langle S, fail \rangle$ .
- **redundant:**  $\langle S \cup \{t \doteq t\}, \sigma \rangle \longrightarrow \langle S, \sigma \rangle$ .
- **orient:**  $\langle S \cup \{t \doteq x \mid t \text{ is not variable}\}, \sigma \rangle \longrightarrow \langle S \cup \{x \doteq t\}, \sigma \rangle$ .

- **substitute**:  $\langle S \cup \{x \doteq t \mid x \text{ does not occur in } t\}, \sigma \rangle \longrightarrow \langle S[x \leftarrow t], \sigma \cdot [x \leftarrow t] \rangle$

**Algorithm 6.1.12.** The algorithm  $unify(s, t)$  takes two terms  $s$  and  $t$ , and returns either  $fail$  or  $\sigma$ , a mgu of  $s$  and  $t$ .

```

proc  $unify(s, t)$ 
1    $S := \{s \doteq t\}; \sigma = \emptyset;$ 
2   while (true) do
3        $\langle S, \sigma \rangle :=$  apply one unification rule on  $\langle S, \sigma \rangle;$ 
4       if ( $\sigma = fail$ ) return  $fail;$ 
5       if ( $S = \emptyset$ ) return  $\sigma;$ 

```

**Example 6.1.13.** Given  $s = p(f(x, h(x), y), g(y))$  and  $t = p(f(g(z), w, z), x)$ , then  $unify(s, t)$  will transform  $\langle \{s \doteq t\}, \emptyset \rangle$  into  $\langle \emptyset, \sigma \rangle$  as follows:

		$\{s \doteq t\}, \emptyset$
<b>decompose</b>	$\longrightarrow$	$\{f(x, h(x), y) \doteq f(g(z), w, z), g(y) \doteq x\}, \emptyset$
<b>orient</b>	$\longrightarrow$	$\{f(x, h(x), y) \doteq f(g(z), w, z), x \doteq g(y)\}, \emptyset$
<b>substitute</b>	$\longrightarrow$	$\{f(g(y), h(g(y)), y) \doteq f(g(z), w, z)\}, [x \leftarrow g(y)]$
<b>decompose</b>	$\longrightarrow$	$\{g(y) \doteq g(z), h(g(y)) \doteq w, y \doteq z\}, [x \leftarrow g(y)]$
<b>substitute</b>	$\longrightarrow$	$\{g(z) \doteq g(z), h(g(z)) \doteq w\}, [x \leftarrow g(z), y \leftarrow z]$
<b>redundant</b>	$\longrightarrow$	$\{h(g(z)) \doteq w\}, [x \leftarrow g(z), y \leftarrow z]$
<b>orient</b>	$\longrightarrow$	$\{w \doteq h(g(z))\}, [x \leftarrow g(z), y \leftarrow z]$
<b>substitute</b>	$\longrightarrow$	$\emptyset, [x \leftarrow g(z), y \leftarrow z, w \leftarrow h(g(z))]$

The algorithm  $unify(s, t)$  does not specify the order of rules, so we may use these rules in any order when they are applicable. However, none of these rules can be applied forever. Hence  $unify(s, t)$  will terminate eventually.

**Lemma 6.1.14.** *Given any pair  $s$  and  $t$ ,  $unify(s, t)$  will terminate with either  $fail$  or  $\sigma$ .*

*Proof.* At first, we like to point out that the **substitute** rule may increase the sizes of terms in  $S$  but it moves variable  $x$  from  $S$  into  $\sigma$ , that is,  $x$  disappears from  $S$ , so this rule can be applied no more than the number of variables.

For the rules used between two applications of the **substitute** rule, if the **clash** rule or the **occur-check** rule apply, the process ends with  $fail$  immediately. Otherwise, the **decompose** rule removes two occurrences of  $f$ . The **orient** rule switches  $t \doteq x$  to  $x \doteq t$ , where  $t$  is not a variable. The **redundant** rule removes  $t \doteq t$ . None of these rule can be used forever. In the end,  $S$  becomes empty and  $\sigma$  will be returned.  $\square$

**Lemma 6.1.15.** *If  $\langle S', \sigma' \rangle$  is the result of applying the **substitute** rule on  $\langle S, \sigma \rangle$  and  $\theta'$  is a mgu of  $S'$ , then  $\theta = [x \leftarrow t] \cdot \theta'$  is a mgu of  $S$ .*

*Proof.* By the definition of the **substitute** rule,  $S' = (S - \{x \doteq t\})[x \leftarrow t]$  and  $\sigma' = \sigma \cdot [x \leftarrow t]$ , where  $x$  is not affected by  $\sigma$ . By the assumption,  $S'$  is unifiable with a mgu  $\theta'$ . Since  $x$  does not appear in  $S'$ ,  $\theta = [x \leftarrow t] \cdot \theta'$  is a mgu of  $S' \cup \{x \doteq t\}$ . It is easy to check that  $\theta$  is also a mgu of  $S$ .  $\square$

**Theorem 6.1.16.** *The algorithm  $\text{unify}(s, t)$  will return a mgu of  $s$  and  $t$  iff  $s$  and  $t$  are unifiable.*

*Proof.* Lemma 6.1.14 ensures that  $\text{unify}(s, t)$  will terminate. Inside  $\text{unify}(s, t)$ , when  $\langle S, \sigma \rangle$  is transformed to  $\langle S', \sigma' \rangle$ , we can check for each rule that  $S$  is unifiable ( $\sigma$  does not affect the unifiability of  $S$ ) iff  $\sigma' \neq \text{fail}$  and  $S'$  is unifiable. That is, when  $\sigma'$  is *fail*,  $S$  cannot be unifiable. When  $S' = \emptyset$ ,  $S'$  is trivially unifiable, so  $S$  must be unifiable.

Before the **substitute** rule moves  $x \doteq t$  from  $S$ ,  $x$  is not affected by  $\sigma$ ;  $x$  is affected by  $\sigma' = \sigma \cdot [x \leftarrow t] = \sigma[x \leftarrow t] \cup [x \leftarrow t]$ , where  $t$  does not contain  $x$ . Since  $S[x \leftarrow t]$  and  $\sigma[x \leftarrow t]$  remove all the occurrences of  $x$  in  $S$  and  $\sigma$ ,  $x$  has a unique occurrence in  $\sigma'$  and that is also true for all affected variables in  $\sigma$ . Thus,  $\sigma$  is idempotent by Proposition 6.1.2.

Assume  $s$  and  $t$  are unifiable. Let  $S_0 = \{s \doteq t\}$ ,  $\sigma_0 = \emptyset$ , and  $\sigma_i = \sigma_{i-1} \cdot [x_i \leftarrow t_i]$  records every application of the **substitute** rule for  $1 \leq i \leq m$ . Let the corresponding  $S$  be  $S_1, S_2, \dots, S_m$ , respectively, where  $S_m = \emptyset$ . Since the mgu of  $S_m$  is  $\theta_m = \emptyset$ . So by Lemma 6.1.15,

$$\theta_{m-1} = [x_m \leftarrow t_m]\theta_m = [x_m \leftarrow t_m]$$

is a mgu of  $S_{m-1}$ . Again by Lemma 6.1.15,

$$\theta_{m-2} = [x_{m-1} \leftarrow t_{m-1}]\theta_{m-1} = [x_{m-1} \leftarrow t_{m-1}, x_m \leftarrow t_m],$$

is a mgu of  $S_{m-2}$ , and so on. Finally,

$$\theta_0 = [x_1 \leftarrow t_1]\theta_1 = [x_1 \leftarrow t_1, x_{m-1} \leftarrow t_{m-1}, x_m \leftarrow t_m] = \sigma_m$$

is a mgu of  $S_0 = \{s \doteq t\}$ .  $\square$

**Example 6.1.17.** Consider the following pair of terms:

$$\begin{aligned} s &= f(f(f(f(a, x_1), x_2), x_3), x_4), \\ t &= f(x_4, f(x_3, f(x_2, f(x_1, a))))). \end{aligned}$$

		$\{s \doteq t\}, \emptyset$
<b>decompose</b>	$\longrightarrow$	$\{f(f(f(a, x_1), x_2), x_3) \doteq x_4, x_4 \doteq f(x_3, f(x_2, f(x_1, a))))\}, \emptyset$
<b>substitute</b>	$\longrightarrow$	$\{f(f(f(a, x_1), x_2), x_3) \doteq f(x_3, f(x_2, f(x_1, a)))\},$ $[x_4 \leftarrow f(x_3, f(x_2, f(x_1, a)))]$
<b>decompose</b>	$\longrightarrow$	$\{f(f(a, x_1), x_2) \doteq x_3, x_3 \doteq f(x_2, f(x_1, a))\},$ $[x_4 \leftarrow f(x_3, f(x_2, f(x_1, a)))]$
<b>substitute</b>	$\longrightarrow$	$\{f(f(a, x_1), x_2) \doteq f(x_2, f(x_1, a))\},$ $[x_4 \leftarrow f(f(x_2, f(x_1, a)), f(x_2, f(x_1, a))), x_3 \leftarrow f(x_2, f(x_1, a))]$
<b>decompose</b>	$\longrightarrow$	$\{f(a, x_1) \doteq x_2, x_2 \doteq f(x_1, a)\}, [x_4 \leftarrow \dots, x_3 \leftarrow f(x_2, f(x_1, a))]$
<b>substitute</b>	$\longrightarrow$	$\{f(a, x_1) \doteq f(x_1, a)\}, [x_4 \leftarrow \dots, x_3 \leftarrow \dots, x_2 \leftarrow f(x_1, a)]$
<b>decompose</b>	$\longrightarrow$	$\{a \doteq x_1, x_1 \doteq a\}, [x_4 \leftarrow \dots, x_3 \leftarrow \dots, x_2 \leftarrow f(x_1, a)]$
<b>substitute</b>	$\longrightarrow$	$\{a \doteq a\}, [x_4 \leftarrow \dots, x_3 \leftarrow \dots, x_2 \leftarrow f(a, a), x_1 \leftarrow a]$
<b>redundant</b>	$\longrightarrow$	$\emptyset, [x_4 \leftarrow \dots, x_3 \leftarrow \dots, x_2 \leftarrow f(a, a), x_1 \leftarrow a]$

Let  $t_1 = a$ ,  $t_2 = f(t_1, t_1)$ ,  $t_3 = f(t_2, t_2)$ ,  $t_4 = f(t_3, t_3)$ , then the mgu of  $s \doteq t$  is  $[x_i \leftarrow t_i \mid 1 \leq i \leq 4]$ .

In general, if  $s$  and  $t$  are the following pair of terms,

$$s = f(f(f(\dots f(f(a, x_1), x_2), \dots), x_{n-1}), x_n)$$

$$t = f(x_n, f(x_{n-1}, f(\dots, f(x_2, f(x_1, a))\dots)))$$

then the mgu of  $s$  and  $t$  is  $\sigma = [x_i \leftarrow t_i \mid 1 \leq i \leq n]$ , where  $t_1 = a$  and  $t_{i+1} = f(t_i, t_i)$ . The sizes of  $s$  and  $t$  are the same, i.e.,  $2n + 1$ . However, the size of  $t_n$  is  $2^n - 1$ .

The above example illustrates that  $unify(s, t)$  will take  $O(2^n)$  time and space for the input terms of size  $O(n)$ . This high cost cannot be avoided if we use trees to represent terms. Using the data structure of directed graphs for terms, identical subterms can be shared, so we do not need a space of exponential size.

#### 6.1.4 Practically Linear Time Unification

Corbin and Bidoit (1983) improved Robinson's unification algorithm to a quadratic one by using directed graphs. Ruzicka and Privara (1989) further improved it to an almost linear one, by using the union/find data structures and postponing the occur-check at the end of unification. The union/find data structure stores a collection of disjoint sets, each of the sets represents an equivalence class. Thus, the disjoint sets represent an equivalence relation and are a partition of the involved elements. The *Union* operation will merge two sets into one, and the *Find* operation will return a representative member of a set, so that it allows to find out efficiently if any two elements are in the same or different sets.

In the following, we introduce an almost-linear time unification algorithm different from that of Ruzicka and Privara.

Given two terms  $s$  and  $t$ , the *term graph* for  $s$  and  $t$  is the directed acyclic graph  $G = (V, E)$  obtained by combining the formula trees of  $s$  and  $t$  with the following operation: All the nodes labeled with the same variable in  $s$  and  $t$  is merged into a single node of  $G$ .

The graph  $G$  has the following features:

- Each node  $v$  of  $G$  is labeled with a symbol in  $s$  or  $t$ , denoted by  $label(v)$ .
- The roots of the formula trees are the two nodes in  $G$  with zero in-degree; the other non-variable nodes' in-degrees are one. The in-degree of variable nodes is the sum of the occurrences of that variable in  $s$  and  $t$ .
- The leaf nodes of the formula trees are the nodes with zero out-degree in  $G$ , which are labeled with variables or constants.
- If a node is labeled with a function symbol  $f$  whose arity is  $k$ , then this node has  $k$  ordered successors. We will use  $child(v, i)$  to refer to the  $i^{th}$  successor of  $v$ .

The term associated with each node  $n$  of  $G$  can be recursively defined as follows:

1. If  $label(n)$  is a variable or constant,  $term(n) = label(n)$ ;
2. If  $label(n)$  is a function  $f$  other than a constant,  $term(n) = f(s_1, s_2, \dots, s_k)$ , where  $s_i = term(child(n, i))$  for  $1 \leq i \leq arity(f)$ .

We say  $term(n)$  is the term represented by  $n$  in  $G$ .

The unification algorithm will work on  $G = (V, E)$  and is divided into two stages: In stage 1, *unify1* checks if  $s$  and  $t$  have no clash failure; in stage 2, *postOccurCheck* returns true iff they have no occur-check failure.

In stage 1, the algorithm *unify1* takes the two nodes representing  $s$  and  $t$ , respectively, and visits  $G$  in tandem by depth-first search. During the search, we will assign to each node  $v$  of  $G$  another node  $u$  as the *representative* of  $v$ , denoted by  $up(v) = u$ . Initially,  $up(v) = v$  for every node of  $G$ . If  $up(v) = u$  and  $v \neq u$ , it means the terms represented by  $v$  and  $u$ , respectively, are likely unifiable (no clash failure) and we use  $u$  as the representative of  $v$  in the rest of the search. At the end of the search, *unify1* returns true if no clash failure was found. In this case, the relation  $R$  defined by  $up$ , i.e.,  $R = \{(u, up(u)) \mid u \in V\}$ , can generate an equivalence relation  $E$  (i.e.,  $E$  is the reflexive, symmetric and transitive closure of  $R$ ). The relation  $E$  partitions  $V$  into equivalence classes and the nodes in the same equivalent class represent the terms which must be unifiable if  $s$  and  $t$  are unifiable.

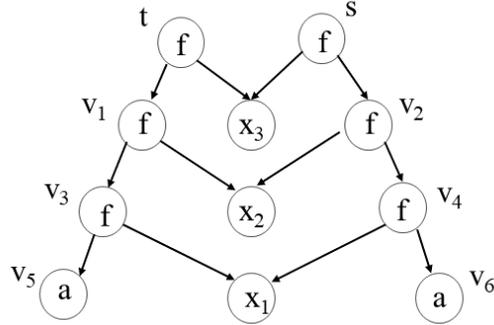


Figure 6.1.1: The term graph for  $t = f(f(f(a, x_1), x_2), x_3)$  and  $s = f(x_3, f(x_2, f(x_1, a)))$ .

**Example 6.1.18.** Given  $t = f(f(f(a, x_1), x_2), x_3)$  and  $s = f(x_3, f(x_2, f(x_1, a)))$ , their term graph is shown in Figure 6.1.1. The *up* relation is initialized as  $up(v) = v$ .

The algorithm  $unify1(t, s)$  starts the depth-first search in tandem and goes to the first children of  $(t, s)$ , i.e.,  $(v_1, x_3)$ . Since  $x_3$  is a variable node,  $unify1(v_1, x_3)$  will do  $up(x_3) := v_1$  and backtrack. The second children of  $(t, s)$  are  $(x_3, v_2)$ . Since  $up(x_3) = v_1$ , we use  $v_1$  for  $x_3$  and call  $unify1(v_1, v_2)$ . Since  $v_1$  and  $v_2$  have the same label, the search goes to their children. The first children of  $(v_1, v_2)$  are  $(v_3, x_2)$ . Since  $x_2$  is a variable node, we do  $up(x_2) := v_3$  and backtrack. The second children of  $(v_1, v_2)$  are  $(x_2, v_4)$ . Using  $v_3$  for  $x_2$  and call  $unify1(v_3, v_4)$ . Since  $v_3$  and  $v_4$  have the same label, the search goes to their children. The first children of  $(v_3, v_4)$  are  $(v_5, x_1)$ . We do  $up(x_1) := v_5$  and backtrack. The second children of  $(v_3, v_4)$  are  $(x_1, v_6)$ . Using  $v_5$  for  $x_1$  and call  $unify1(v_5, v_6)$ . Since  $(v_5, v_6)$  have the same constant label, the search backtracks after doing  $up(v_6) := v_5$ . Now, all children of  $(v_3, v_4)$  are done. After  $up(v_4) := v_3$ , the search backtracks. Similarly, all children of  $(v_1, v_2)$  are done, and we backtrack after doing  $up(v_2) := v_1$ . Finally, all children of  $(t, s)$  are done and  $unify1(t, s)$  finishes with  $up(s) := t$ .

The content of  $up(v)$  after the termination of  $unify1(t, s)$  is shown below:

node	$s$	$t$	$v_1$	$v_2$	$v_3$	$v_4$	$v_5$	$v_6$	$x_1$	$x_2$	$x_3$
$up$	$t$	$t$	$v_1$	$v_1$	$v_3$	$v_3$	$v_5$	$v_5$	$v_5$	$v_3$	$v_1$

The equivalence relation  $E$  generated by  $(v, up(v))$  can be represented by the following partition of  $V$ :

$$\{s, t\}, \{v_1, v_2, x_3\}, \{v_3, v_4, x_2\}, \{v_5, v_6, x_1\}$$

It is clear from the above example that each node of  $V$  is in a singleton set initially. When the terms represented by the two nodes are unifiable, the sets

containing the two nodes are merged into one through the assignment  $up(v) := u$ , i.e., we obtain “union” of the sets containing  $u$  and  $v$ . To know which node is represented by which node, we need the “find” operation. Thus, we may employ the well-known union/find data structure to support the implementation of *unify1*.

In view of the union/find data structure, each disjoint set is represented by a tree:  $up(v) = u$  means the parent of  $v$  is  $u$  in the tree if  $v \neq u$ . If  $up(u) = u$ , then  $u$  is the root of the tree, which also serves as the name of the set. *Find*( $u$ ) will find the name of the set containing  $u$  and will compress the tree if necessary. *Union*( $u, v$ ) will merge the two sets whose names are  $u$  and  $v$ , respectively. A typical implementation of *Union* will let the root of shorter/smaller tree point to the root of a taller/larger tree, and can be done in constant time.

As stated earlier, the algorithm *unify1* takes the two nodes representing  $s$  and  $t$ , respectively, and visits  $G$  in tandem by depth-first search. *unify1* will first replace the current two nodes by their set names, which are also the nodes in  $G$ . If the two set names are the same, backtrack; If one of them is a variable node, make the other as its parent in the union/find tree, and backtrack. If both are function nodes and their labels are different, report *clash failure*; otherwise, recursively visit their children. If no failure is found from the children, merge the two sets named by the current two nodes and backtrack. During the depth-first search, we also use a marking mechanism to detect if a cycle is created because a variable is unified with a term; if yes, an occur-check failure happened. Only a portion of occur-check failures can be found this way. The pseudo-code of the *unify* algorithm is given below.

**Algorithm 6.1.19.** The algorithm *lunify*( $s, t$ ) will take terms  $s$  and  $t$ , as input, create the formula graph  $G = (V, E)$  from  $s$  and  $t$ , merging all the nodes of the same variable into one, and call *unify1*( $v_s, v_t$ ), where  $v_s$  and  $v_t$  are the nodes representing  $s$  and  $t$ , respectively. If *unify1*( $v_s, v_t$ ) return true, call *postOccurCheck*( $v_s, v_t$ ), which returns true iff there exist no cycles in  $G = (V, E)$  by the depth-first search.

```

proc lunify( $s, t$ ): Boolean
1    $G = createGraph(s, t)$ ; //  $up(v) = v$  for each  $v$  in  $G$ 
2   if unify1( $v_s, v_t$ ) //  $v_s, v_t$  represent  $s, t$  in  $G$ 
3     return postOccurCheck( $v_s, v_t$ )
4   else return false;

```

```

proc unify1( $v_1, v_2$ ) : Boolean
1    $s_1 := Find(v_1)$ ;  $s_2 := Find(v_2)$ ;
2   if  $s_1 = s_2$  return true; // already unified

```

```

3   if ( $label(s_1)$  and  $label(s_2)$  are variables)  $Union(s_1, s_2)$ ; return true;
4   if ( $label(s_1)$  is variable)  $up(s_1) := s_2$ ; return true;
5   if ( $label(s_2)$  is variable)  $up(s_2) := s_1$ ; return true;
6   if  $label(s_1) \neq label(s_2)$  return false // clash failure
   // use  $mark(s)$  to detect if  $G$  has a cycle
7   for  $i := 1$  to  $arity(label(s_1))$  do
8      $c_1 := child(s_1, i)$ ;  $c_2 := child(s_2, i)$ ;
9     if ( $marked(c_1)$  or  $marked(c_2)$ ) return false // occur-failure
10     $mark(c_1)$ ;  $mark(c_2)$ ;
11    if  $unify1(c_1, c_2) = \text{false}$  return false;
12     $unmark(c_1)$ ;  $unmark(c_2)$ ;
13   $Union(s_1, s_2)$ ;
14  return true;

```

```

proc  $Union(s_1, s_2)$ 
  // Initially, for any node  $s$ ,  $h(s) = 0$ , the height of a singleton tree.

```

```

1   if  $h(s_1) < h(s_2)$   $up(s_1) := s_2$ ; return ;
2    $up(s_2) := s_1$ ; // the parent of  $s_2$  is  $s_1$ 
3   if  $h(s_1) = h(s_2)$   $h(s_1) := h(s_1) + 1$ ;

```

```

proc  $Find(v)$ 

```

```

1    $r := up(v)$ ; // look for the root of the tree containing  $v$ 
2   while ( $r \neq up(r)$ )  $r := up(r)$ ; //  $r$  is root iff  $r = up(r)$ 
3   if ( $r \neq v$ ) // compress the path from  $v$  to  $r$ 
4     while ( $r \neq v$ )  $up(v) := r$ ;  $v := up(v)$ ;
5   return  $r$ ;

```

As said earlier, in the union-find data structure, each set of the nodes is represented by a tree (defined by  $up(v) = u$ , where the parent of  $v$  is  $u$ ) and the root of the tree is the name of the set. The union operation implements the “union by height”, making the root of the higher tree as the parent of the root of the shorter tree.  $Find(v)$  will return the name of the set, i.e., the root of the tree containing  $v$ ; this is done by the first while loop. For the efficiency of future calls to  $Find$ , the path from  $v$  to the root is compressed after the root is found; this is completed by the second while loop.

There are two parent-child relations used in the above algorithms. The procedure  $unify1$  uses the depth-first search based on the parent-child relation defined by the formula graph, not the parent-child relation defined by  $up(v) = u$ . The later induces an equivalence relation and is used by  $Find$  and  $Union$ .

**Example 6.1.20.** Let  $s = f(x, x)$  and  $t = f(y, a)$ . A graph  $G = (V, E)$  will be created with  $V = \{v_s, v_t, v_x, v_y, v_a\}$  and  $E = \{(v_s, v_x)_1, (v_s, v_x)_2, (v_t, v_y)_1, (v_t, v_a)_2\}$ , where the subscripts are the order of children.  $unify1(v_s, v_t)$  will call  $unify1(v_x, v_y)$ , which changes  $up(v_x) = v_y$ , and  $unify(v_x, v_y)$ , which changes  $up(v_y) = v_a$ . After the two recursive calls,  $unify1(v_s, v_t)$  will change  $up(v_s) = v_t$  and return true.  $V$  is partitioned into two sets:  $\{v_s, v_t\}$  and  $\{v_x, v_y, v_a\}$ .

In algorithm  $unify1$ , the graph  $G = (V, D)$  never changes; only  $up(v)$  changed. The value of  $up(v)$  defines the equivalence relation in the union/find data structure:  $v$  and  $up(v)$  are in the same disjoint set. The value of  $up(v)$  is changed indirectly either by  $Find$  (called at line 1 of  $unify1$ ;  $up(v)$  is changed at line 4 of  $Find$ ), or by  $Union$  (called at lines 3 and 13;  $up(v)$  is modified at lines 1 and 2 of  $Union$ ), or directly at lines 4 and 5 of  $unify1$ , which are special cases of  $Union$ : a variable node's parent is set to be a function node (not the other direction). Thus, the algorithm uses only linear space. The algorithm uses almost linear time because of the following observations:

- Let  $G = (V, E)$  be the term graph of  $s$  and  $t$ , and  $n = |V|$  is bound by the total size of  $s$  and  $t$ .
- The number of find operations performed in  $unify1$  is twice of the number of recursive calls of  $unify1$ .
- Each union operation reduces the number of node sets by one, thus the number of union operations is bound by  $n$ .
- Once two subterms are shown to be unifiable, the two nodes representing them are put into the same set, thus we never try to unify again the same pair of subterms and the number of recursive calls of  $unify1$  is bound  $n$ .
- The total cost of  $O(n)$  finds and unions is  $O(n\alpha(n))$ , where  $\alpha(n)$  is the inverse of Ackermann function and grows very, very slow. In fact,  $\alpha(n)$  is practically a constant function, e.g.,  $\alpha(10^9) \leq 4$ .
- The special unions performed at lines 4 and 5 of  $unify1$  do not affect the almost linear complexity, because the total cost of such unions between a variable node and a non-variable is bounded by the number of variable nodes.
- The procedure  $postOccurCheck$  will check if  $G = (V, E)$  contains a cycle and can be implemented in  $O(n)$  by the depth-first search.

**Theorem 6.1.21.** *The algorithm  $\text{lunify}(s, t)$  takes  $O(n\alpha(n))$  time, where  $n$  is the total sizes of  $s$  and  $t$ , and  $\alpha(n)$  is the inverse of Ackermann function;  $\text{lunify}(s, t)$  returns **true** iff  $s$  and  $t$  are unifiable.*

The algorithm *lunify* tells us if two terms expressions are unifiable or not. However, this is different from knowing what their unifiers actually are. We leave the problem of finding the mgu of  $s$  and  $t$  after *lunify*( $s, t$ ) returns true as an exercise.

Since *lunify* does not change the term graph of  $s$  and  $t$ , we may wonder how to get the term  $t\sigma$  from  $G$  if  $\sigma$  is the mgu of  $s$  and  $t$ . The term  $t\sigma$ , which is equal to  $s\sigma$ , can be obtained by  $\text{term}'(v_s)$  (or  $\text{term}'(v_t)$ ), where  $\text{term}'$  can be recursively defined as follows: For any node  $v$  of  $G$ :

1. If  $\text{label}(v)$  is a variable and  $\text{up}(v) \neq v$ ,  $\text{term}'(v) = \text{term}'(\text{up}(v))$ ;
2. If  $\text{label}(v)$  is a constant or a variable and  $\text{up}(v) = v$ ,  $\text{term}'(v) = \text{label}(v)$ ;
3. If  $\text{label}(v)$  is a function  $f$  other than a constant,  $\text{term}'(v) = f(s_1, s_2, \dots, s_k)$ , where  $s_i = \text{term}'(\text{child}(v, i))$  for  $1 \leq i \leq \text{arity}(f)$ .

It is easy to check that  $\text{term}'(v) = \text{term}(v)$  before the execution of *lunify* as at that time,  $\text{up}(v) = v$  for any  $v$ . We may view  $\{(v, \text{up}(v)) \mid \text{up}(v) \neq v\}$  as a new set of edges added to  $G = (V, E)$  during the execution of *lunify*. This set of edges defines not only the equivalence relation represented by the union/find data structure, but also the mgu  $\sigma$  when  $s$  and  $t$  are unifiable.

The contrived examples in Table 6.1.1 illustrate behavior of *lunify*. The size of a term is the total number of symbols in the term, excluding parentheses and commas. The column **visited** gives the number of nodes visited before a result is returned.

In P1, the two terms are unifiable, and the occur-check can be very expensive if we don't postpone the check. In P2, the first occur-check (on  $z \leftarrow g(z)$ ) will fail; though *lunify* won't know it until `postOccurCheck` is called. In P3, the last occur-check will fail. It is interesting to notice that *unify1* will report this failure through the marking mechanism in a constant number of steps, after visiting only 8 nodes of the term graph. In P4, a clash-failure takes place at the bottom of the depth-first-search. P5 is similar to P3, though an occur-check failure can be found only after a through search. P6 is similar to P1, though the last occur-check will locate the failure.

	pair of terms	size	failure	visited
P1	$f(f(\dots f(f(a, x_1), x_2), \dots), x_n)$ $f(x_n, f(x_{n-1}, f(\dots, f(x_1, a)\dots)))$	$4n + 2$	no	$4n + 2$
P2	$f(z, g^n(x))$ $f(g(z), g^n(y))$	$n + 3$	occur	$n + 3$
P3	$h(x, f(x, f(x, \dots, f(x, f(x, x))\dots)), x)$ $h(f(y, \dots, f(y, f(y, y))\dots), y, y)$	$8n + 3$	occur	8
P4	$f(f(\dots f(f(a, x_1), x_2), \dots), x_n)$ $f(x_n, f(x_{n-1}, f(\dots, f(x_1, b)\dots)))$	$4n + 2$	clash	$4n + 2$
P5	$h(x, g^n(x), x)$ $h(g^n(y), y, y)$	$2n + 8$	occur	$2n + 8$
P6	$f(f(f(\dots f(f(a, x_1), x_2), \dots), x_n), z)$ $f(f(x_n, f(x_{n-1}, f(\dots, f(x_1, a)\dots))), g(z))$	$4n + 7$	occur	$4n + 7$

Table 6.1.1: The last column gives the number of visited nodes by *unify1* in the term graph for each example.

## 6.2 Resolution

### 6.2.1 The Resolution and Factoring Rules

**Definition 6.2.1. (resolution)** Suppose clause  $C_1$  is  $(A \mid \alpha)$  and  $C_2$  is  $(\neg B \mid \beta)$ , where  $A$  and  $B$  are atoms,  $\alpha$  and  $\beta$  are the rest literals in  $C_1$  and  $C_2$ , respectively, resolution is the following inference rule:

$$\frac{(A \mid \alpha) \quad (\neg B \mid \beta)}{(\alpha \mid \beta)\sigma}$$

where  $\sigma$  is the mgu of  $A$  and  $B$ . The clause  $(\alpha \mid \beta)\sigma$  produced by the resolution rule is called *resolvent of the resolution*;  $C_1$  and  $C_2$  are the *parents of the resolvent*. We use  $\text{resolve}(C_1, C_2)$  to denote the clause  $(\alpha \mid \beta)\sigma$ .

This resolution rule is also called *binary resolution*, as it involves two clauses as premises.

**Example 6.2.2.** Given a pair of clauses  $(p(f(x_1)) \mid q(x_1))$  and  $(\overline{p(x_2)} \mid \overline{q(g(x_2))})$ . The resolution on  $p$  will generate the resolvent  $(q(x_1) \mid \overline{q(g(f(x_1)))})$  with the mgu  $[x_2 \leftarrow f(x_1)]$ ; the resolution on  $q$  will generate the resolvent  $(p(f(g(x_2))) \mid \overline{p(x_2)})$  with the mgu  $[x_1 \leftarrow g(x_2)]$ .

**Proposition 6.2.3.** *The resolution rule is sound, that is,  $C_1 \wedge C_2 \models \text{resolve}(C_1, C_2)$ .*

*Proof.* Suppose  $C_1$  is  $(A \mid \alpha)$  and  $C_2$  is  $(\neg B \mid \beta)$ ,  $A$  and  $B$  are unifiable and their mgu is  $\sigma$ . Then  $\text{resolve}(C_1, C_2) = (\alpha, \beta)\sigma$ . For any model  $I$  of  $C_1 \wedge C_2$ ,  $I$  is also a model of  $C_1\sigma = A\sigma \vee \alpha\sigma$  and  $C_2\sigma = \neg B\sigma \vee \beta\sigma$ , because free variables are treated as universally quantified variables. Since  $\sigma$  is a unifier of  $A$  and  $B$ ,  $\alpha\sigma = \beta\sigma$ . If  $A\sigma$  is false in  $I$ , then  $\alpha\sigma$  must be true in  $I$ , otherwise  $C_1\sigma$  will be false in  $I$ ; If  $A\sigma$  is true in  $I$ , then  $\neg B\sigma$  is false in  $I$ , so  $\beta\sigma$  must be true in  $I$ . In both cases, either  $\alpha\sigma$  or  $\beta\sigma$  is true in  $I$ , so  $(\alpha \vee \beta)\sigma$  is true in  $I$ .  $\square$

By the above proposition, every resolvent is a logical consequence of the input clauses, because the entailment relation is transitive. If the empty clause is generated, then we know that the input clauses are unsatisfiable. In other words, we may use the resolution rule to design a refutation prover. Suppose we have some axioms,  $A$ , and we want to see if some conjecture,  $B$ , is a logical consequence of  $A$ . By the refutational strategy, we want to show that  $A \wedge \neg B$  is unsatisfiable, by putting it into clausal form and finding a contradiction through resolution. We wish the following properties hold.

- $B$  is a logical consequence of  $A$ , i.e.,  $A \models B$ , iff
- $A \wedge \neg B$  is unsatisfiable, iff
- $S$ , the clause set derived from  $A \wedge \neg B$ , is unsatisfiable, iff
- The empty clause can be generated by resolution from  $S$ .

Unfortunately, the last “iff” does not hold in first-order logic. Indeed, when  $S$  is satisfiable, resolution will never generate the empty clause. If the empty clause can be generated, we can claim that  $S$  is unsatisfiable, because resolution is sound. However, when  $S$  is unsatisfiable, resolution may not generate the empty clause.

**Example 6.2.4.** Let  $S = \{(p(x_1) \mid p(y_1)), (\overline{p(x_2)} \mid \overline{p(y_2)})\}$ . From the first clause, we have an instance of  $(p(a) \mid p(a))$ , which can be simplified to  $(p(a))$ ; from the second clause, we have an instance  $(\overline{p(a)} \mid \overline{p(a)})$ , which is simplified to  $(\overline{p(a)})$ . A resolution between  $(p(a))$  and  $(\overline{p(a)})$  will generate the empty clause. This shows that  $S$  is unsatisfiable.

One resolvent of the two clauses in  $S$  is  $(p(x_3) \mid \overline{p(y_3)})$ , and we may add it into  $S$ . Any other resolvent from  $S$  will be a renaming of the three clauses. In other words, resolution alone cannot generate the empty clause from  $S$ .

In propositional logic, we defined a proof system which uses the resolution rule as the only single inference rule and showed that this proof system is a decision procedure for propositional logic. In first-order logic, we do not have such a luck.

To get rid of the problem illustrated in the above example, we need the following inference rule:

**Definition 6.2.5. (factoring)** *Let clause  $C$  be  $(A \mid B \mid \alpha)$ , where  $A$  and  $B$  are literals and  $\alpha$  are the rest literals in  $C$ . factoring is the following inference rule:*

$$\frac{(A \mid B \mid \alpha)}{(A \mid \alpha)\sigma}$$

where  $\sigma$  is the mgu of  $A$  and  $B$ .

For the clause  $(p(x_1) \mid p(y_1))$  in the previous example, factoring will generate  $(p(x_1))$ ; for  $(\overline{p(x_2)} \mid \overline{p(y_2)})$ , factoring will generate  $(\overline{p(x_2)})$ . Now resolution will generate the empty clause from  $(p(x_1))$  and  $(\overline{p(x_2)})$ .

**Proposition 6.2.6.** *The factoring rule is sound, that is,  $C \models C'$ , where  $C'$  is the clause generated by factoring from  $C$ .*

*Proof.*  $C'$  is an instance of  $C$ , i.e.,  $C' = C\sigma$ , where  $\sigma$  is the unifier used in factoring. Any instance of a formula is a logical consequence of the formula.  $\square$

Armed with resolution and factoring, now we can claim the following result without a proof.

**Theorem 6.2.7. (refutational completeness of resolution)** *A set  $S$  of clauses is unsatisfiable iff the empty clause can be generated by resolution and factoring.*

Though factoring is needed for the completeness of resolution, in practice, factoring is rarely needed.

### 6.2.2 A Refutational Proof Procedure

For propositional logic, we introduced three deletion strategies for deleting clauses without worrying about missing a proof:

- **pure literal deletion:** Clauses containing pure literals are discarded.
- **tautology deletion:** Tautology clauses are discarded.
- **subsumption deletion:** Subsumed clauses are discarded.

The pure literal deletion and tautology deletion can be used without modification. For the subsumption deletion, the rule is modified as follows:

**Definition 6.2.8. (subsumption)** *Clause  $C_1$  subsumes clause  $C_2$  if there exists a substitution  $\sigma$  such that every literal of  $C_1\sigma$  appears in  $C_2$ .*

For example,  $(p(x, y) \mid \overline{p(y, x)})$  subsumes  $(p(a, b) \mid \overline{p(b, a)} \mid q(a))$  with  $\sigma = [x \leftarrow a, y \leftarrow b]$ .

The above deletion strategies can be integrated into the algorithm *resolution* as follows.

**Procedure 6.2.9.** The procedure *resolution*( $C$ ) takes a set  $C$  of clauses and returns false iff  $C$  is unsatisfiable. It uses *preprocessing*( $C$ ) to simplify the input clauses  $C$ . The procedure *subsumedBy*( $A, S$ ) checks if clause  $A$  is subsumed by a clause in  $S$ . The procedure *factoring*( $A$ ) takes clause  $A$  as input and generates the set of results from the factoring rule on  $A$ , including  $A$  itself.

```

proc resolution( $C$ )
1    $G := preprocessing(C)$ ; //  $G$ : given clauses
2    $K := \emptyset$ ; //  $K$ : kept clauses
3   while  $G \neq \emptyset$  do
4      $A := pickClause(G)$ ; // heuristic for picking a clause
5      $G := G - \{A\}$ ;
6      $N := \emptyset$ ; // new clauses from  $A$  and  $K$  by resolution
7     for  $B \in K$  if resolvable( $A, B$ ) do
8        $D := resolve(A, B)$ ;
9       if  $D = ()$  return false; // the empty clause is found
10      if subsumedBy( $D, G \cup K$ ) continue;
11       $N := N \cup factoring(D)$ ;
12       $K := K \cup \{A\}$ ;
13      for  $A \in G$  if subsumedBy( $A, N$ ) do  $G := G - \{A\}$ ;
14      for  $B \in K$  if subsumedBy( $B, N$ ) do  $K := K - \{B\}$ ;
15       $G := G \cup N$ ;
16  return true; //  $K$  is saturated by resolution

```

The above procedure is almost identical to Algorithm 3.3.26 with two differences: (a) *factoring* is used (line 11); (b) the termination of the procedure is not guaranteed as it may generate an infinite number of new clauses. Thus, it is not an algorithm. Despite of this, this procedure is a base for many resolution

theorem provers, including McCune’s Otter and its successor Prover9, which will be introduced shortly. In Prover9, the subsumption check at line 10 is called *forward subsumption* and the subsumption checks at lines 13 and 14 are called *back subsumption*.

Prover9 also uses another simplification rule called *unit deletion*.

**Definition 6.2.10. (unit deletion)** *Unit deletion is the simplification rule that deletes literal  $B$  from clause  $C$  if there exists a unit clause  $(A)$  and a substitution  $\sigma$  such  $A\sigma = \neg B$  (or  $\neg A\sigma = B$ ).*

Unit deletion is sound because by the requirement, there is a resolution between  $(A)$  and  $C$  and the resolvent of this resolution will subsume  $C$ . For example, if we have a unit clause  $(\overline{p(x,b)})$ , then the literal  $p(a,b)$  in  $(p(a,b) \mid q(a))$  can be deleted to get a new clause  $(q(a))$ , which subsumes  $(p(a,b) \mid q(a))$ .

Unit deletion can be used at line 10, so that the new clause  $D$  is simplified by the unit clauses in  $G \cup K$ . This is called *forward unit deletion*. *Back unit deletion* can be inserted at lines 13 and 14 when the unit clauses in  $N$  are used to simplify clauses in  $G \cup K$ .

### 6.3 Ordered Resolution

In section 3.3.2, we listed some well-known resolution strategies which add restrictions to the use of the resolution rule: unit resolution, input resolution, ordered resolution, positive resolution, negative resolution, set of support, and linear resolution. They can be carried over to first-order logic without modification, with the exception of ordered resolution. The procedure *resolvable*( $A, B$ ) (line 7 of *resolution*) is the place where various restrictions are implemented so we can use these restricted resolution strategies.

For unit resolution and input resolution, they are still incomplete in general, but are equivalent in the sense that there is a resolution proof for one strategy, then there exists a resolution proof for the other strategy.

#### 6.3.1 Simplification Orders

To use ordered resolution in first-order logic, we need a partial order to compare atoms. We call this order as “term order”, although atoms, literals, and clauses are not terms, the term orders we present here apply to those objects as well.

There are several requirements on the term orders used for ordered resolution and we term it as “simplification order”.

**Definition 6.3.1. (stable and monotonic relation)** A binary relation  $R$  over  $T(F, X)$  is said to be *stable* if for any term  $s$  and  $t$ ,  $R(s, t)$  implies  $R(s\sigma, t\sigma)$  for any substitution  $\sigma : X \rightarrow T(F, X)$ .  $R$  is said to be *monotonic* if  $R(s, t)$  implies  $R(f(\dots, s, \dots), f(\dots, t, \dots))$  for any function or predicate symbol  $f$ .

For a stable order  $\succ$ , we cannot have  $f(y) \succ x$  because  $f(y)\sigma \succ x\sigma$  does not hold for  $\sigma = [x \leftarrow f(y)]$ .

**Definition 6.3.2. (simplification order)** The order  $\succ$  is called a *simplification order* if it is well-founded, stable and monotonic. We write  $s \prec t$  if  $t \succ s$ ;  $s \succeq t$  if  $s \succ t$  or  $s \approx t$  (i.e.,  $s$  and  $t$  are equivalent in the order).

The reason we stick to well-founded orders because well-founded orders allow us to use mathematical induction. Let  $S$  be a set with a well-founded order  $\succ$ . To show a property  $P(x)$  holds for every  $x \in S$ , we need to do two things:

- For every minimal element  $m \in S$ ,  $P(m)$  is true.
- Prove that  $P(x)$  is true, assuming  $P(x')$  is true for every  $x'$  and  $x \succ x'$ .

The above is called the *well-founded induction principle*.

There are many well-founded orders over different sets. For example,  $>$  over the set of natural numbers is well-founded;  $>$  on the set of integers is not well-founded. Any partial order over a finite set is well-founded.

There are many ways to obtain well-founded orders. A *multiset* or *bag* over a set  $S$  of elements is a modification of the concept of a subset of  $S$ , that, unlike a set, allows for multiple instances for each of its elements from  $S$ . The union, intersection, and subtraction operations can be extended over multisets. For example, let  $S_1 = \{a, b, c\}$  and  $S_2 = \{a, a, b, b\}$ , then  $S_1 \cup S_2 = \{a, a, a, b, b, b, c, c\}$ ,  $S_1 \cap S_2 = \{a, b, c\}$ , and  $S_1 - S_2 = \{a\}$ . If there exists an order over the elements of  $S$ , we may extend this order to compare the multisets of  $S$ .

**Proposition 6.3.3.** Let  $\succ$  be a well-founded order over a set  $S$ . The order  $\succ^{mul}$  over the multisets of  $S$  is well-founded, where, for any multisets  $S_1, S_2$  over  $S$ ,  $S_1 \succ^{mul} S_2$  iff for any  $y \in S_2 - S_1$ , there exists  $x \in S_1 - S_2$ ,  $x \succ y$ .

For example, if  $c \succ b \succ a$ , then  $S_1 \succ^{mul} S_2$  for  $S_1 = \{a, b, b, c\}$  and  $S_2 = \{a, a, b, c\}$ , because  $S_1 - S_2 = \{b\}$ ,  $S_2 - S_1 = \{a\}$ , and  $b \succ a$ .

One popular way is to use lexicographic combination from well-founded orders of simpler sets.

**Proposition 6.3.4.** *Let  $\succ_i$  be a well-founded order over the set  $S_i$ ,  $1 < i \leq n$ . The lexicographic order  $\succ^{lex}$  over  $S_1 \times S_2 \times \cdots \times S_n$  is well-founded, where for any  $\langle x_1, x_2, \dots, x_n \rangle, \langle y_1, y_2, \dots, y_n \rangle \in S_1 \times S_2 \times \cdots \times S_n$ ,  $\langle x_1, x_2, \dots, x_n \rangle \succ^{lex} \langle y_1, y_2, \dots, y_n \rangle$  iff there exists  $k$ ,  $1 \leq k \leq n$ , such that  $x_i \succeq_i y_i$  for  $1 \leq i < k$  and  $x_k \succ_k y_k$ .*

The lexicographic combination can be extended to compare two sequences of different lengths, assuming a non-empty sequence is greater than the empty sequence, like the dictionary order.

Given a first-order language  $L = (P, F, X, Op)$ , let  $\Sigma = P \cup F \cup Op$ . We say  $\succ$  is a *precedence* if  $\succ$  is a well-founded total order over  $\Sigma$ . Note that if  $\Sigma$  is finite,  $\succ$  is trivially well-founded. For any term  $t \in T(F, X)$ , let  $var(t)$  denote the set of variables appearing in  $t$ . We may extend any precedence to terms, as well as to atoms and literals, by several different ways.

- LPO (Lexicographic Path Order). The term order is determined entirely by the symbol precedence.
- RPO (Recursive Path Order). The term order is determined entirely by the symbol precedence.
- KBO (Knuth-Bendix Order). This order uses a weighting function on symbols as well as the symbol precedence. The weighting function is used first, and the symbol precedence breaks ties.

The formal definitions of LPO, RPO, and KBO are here:

**Definition 6.3.5.** ( $\succ_{lpo}$ ) *The lexicographical path order  $\succ_{lpo}$  over  $T(F, X)$  is defined recursively as follows:  $s \succ_{lpo} t$  if  $s \neq t$  and either  $t \in var(s)$  or  $s = f(s_1, \dots, s_m) \succ_{lpo} t = g(t_1, \dots, t_n)$  by one of the following conditions:*

1.  $f \succ g$  and  $s \succ_{lpo} t_i$  for  $1 \leq i \leq n$ , or
2.  $f \approx g$ ,  $[s_1, \dots, s_m] \succ_{lpo}^{lex} [t_1, \dots, t_n]$ , and  $s \succ_{lpo} t_i$  for  $1 \leq i \leq n$ , or
3.  $f \prec g$  and there exists  $j$ ,  $1 \leq j \leq m$ ,  $s_j \succeq_{lpo} t$ .

**Example 6.3.6.** Let  $* \succ +$ ,  $s$  be  $x * (y + z)$  and  $t$  be  $(x * y) + (x * z)$ . Then  $s \succ_{lpo} t$  by condition 1, because  $* \succ +$ ,  $s \succ_{lpo} x * y, x * z$ , both by condition 2. That is, from  $[x, (y + z)] \succ_{lpo}^{lex} [x, y], [x, z]$ , we have  $s \succ_{lpo} x * y, x * z$ .

**Proposition 6.3.7.**  $\succ_{lpo}$  is a simplification order and is total on ground terms, if the precedence  $\succ$  is a well-founded total order.

**Definition 6.3.8.** ( $\succ_{rpo}$ ) The recursive path order  $\succ_{rpo}$  over  $T(F, X)$  is defined recursively as follows:  $s \succ_{rpo} t$  if  $s \neq t$  and either  $t \in \text{var}(s)$  or  $s = f(s_1, \dots, s_m) \succ_{rpo} t = g(t_1, \dots, t_n)$  by one of the following conditions:

1.  $f \succ g$  and  $s \succ_{rpo} t_i$  for  $1 \leq i \leq n$ , or
2.  $f \approx g$  and  $\{s_1, \dots, s_m\} \succ_{rpo}^{mul} \{t_1, \dots, t_n\}$ , or
3.  $f \prec g$  and there exists  $j$ ,  $1 \leq j \leq m$ ,  $s_j \succeq_{rpo} t$ .

**Proposition 6.3.9.**  $\succ_{rpo}$  is a simplification order, if the precedence  $\succ$  is a well-founded total order.

**Example 6.3.10.** Let  $s = (x * y) * z$  and  $t = x * (y * z)$ , then  $s \succ_{lpo} t$  because of condition 2:  $[x * y, z] \succ_{lpo}^{lex} [x, y * z]$ . However,  $\succ_{rpo}$  cannot compare these two terms, because  $\{x * y, z\}$  and  $\{x, y * z\}$  are not comparable by  $\succ_{rpo}^{mul}$ . Moreover,  $\succ_{rpo}$  treats  $(a * a) * a$  and  $a * (a * a)$  as “equal terms”.

By definition, a partial order is an antisymmetry, transitive, and reflexive relation. Since  $\succ_{rpo}$  violates the antisymmetry condition, it is sometimes called “quasi-order”. We can easily modify  $\succ_{rpo}$  so that when  $\{s_1, \dots, s_m\} \approx_{rpo} \{t_1, \dots, t_n\}$ , we may compare them using  $\succ_{rpo}^{lex}$ , so that  $\succ_{rpo}$  can be a total order on ground terms.

Because  $\succ_{rpo}$  cannot compare  $s = x * (y * z)$  and  $t = (x * y) * z$ , and is neither a partial order nor a total order on ground terms in the strict sense, it is less popular than  $\succ_{lpo}$  in practice.

The *Knuth-Bendix order* ( $\succ_{kbo}$ ) uses the weights first to compare terms. When two terms have the same weight, the precedence relation is used to break ties. Using the weights gives us flexibility as well as complexity.

**Definition 6.3.11. (weight function)** A function  $w : X \cup \Sigma \rightarrow N$  is said to be a weight function if it satisfies (i)  $w(c) > 0$  for every constant  $c \in \Sigma$ , and (ii) there exists at most one unary function symbol  $f \in \Sigma$  such that  $w(f) = 0$  and  $f$  must be maximal in the precedence if  $w(f) = 0$ .

The weight function  $w$  can be extended to a function over terms or formulas as follows:

$$\begin{aligned} w(x) &= w_0 \text{ for any variable } x; \\ w(f(t_1, \dots, t_n)) &= w(f) + w(t_1) + \dots + w(t_n) \end{aligned}$$

where  $w_0$  denotes the minimal weight of all constants.

**Definition 6.3.12.** ( $\succ_{kbo}$ ) Let  $w$  be a weight function and  $mv(t)$  be the multiset of variables appearing in  $t$ . The Knuth–Bendix order  $\succ_{kbo}$  over  $T(F, X)$  defined recursively as follows:  $s \succeq_{kbo} t$  if  $w(s) \geq w(t)$  and  $mv(s) \supseteq mv(t)$ ;  $s \succ_{kbo} t$  if  $s \succeq_{kbo} t$  and (a)  $w(s) > w(t)$  or (b)  $s = f(s_1, \dots, s_m)$ ,  $t = g(t_1, \dots, t_n)$ , and one of the following conditions is true:

1.  $f \succ g$ , or
2.  $f \approx g$  and  $[s_1, \dots, s_m] \succ_{kbo}^{lex} [t_1, \dots, t_n]$ , or
3.  $f \prec g$  and there exists  $j$ ,  $1 \leq j \leq m$ ,  $s_j \succeq_{kbo} t$ .

We like to point out that the condition  $mv(s) \supseteq mv(t)$  is necessary. Let  $s = f(g(x), y)$  and  $t = f(y, y)$ . If  $w(g) > 0$ , then  $w(s) > w(t)$ , but  $mv(s) = \{x, y\} \not\supseteq \{y, y\} = mv(t)$  does not hold. If we ignore the condition  $mv(s) \supseteq mv(t)$ , then  $s \succ_{kbo} t$  would imply  $s\theta \succ_{kbo} t\theta$  for  $\theta = [y \leftarrow f(x)]$ , where  $s\theta = t\theta = f(g(x), g(x))$ .

**Example 6.3.13.** Let  $w(f) = w(a) = 1$ ,  $w(b) = 2$ , then  $f(x, b, a) \succ_{kbo} f(a, a, b)$  by condition (b.2), because  $x \succeq_{kbo} a$  and  $b \succ_{kbo} a$ . Note that neither  $\succ_{lpo}$  nor  $\succ_{rpo}$  can compare  $p(x, b, y)$  and  $p(a, a, y)$ .

**Proposition 6.3.14.**  $\succ_{kbo}$  is a simplification order and is total on ground terms, if the precedence  $\succ$  is a well-founded total order.

**Example 6.3.15.** Let the precedence relation be  $b/1 \succ a/1 \succ e/0$ , where  $a$  and  $b$  are unary function symbols and  $e$  is a constant. We will write  $a(a(b(a(e))))$  as  $aaba$  for brevity. Let  $w(e) = 1$  and the weights of  $a$  and  $b$ , as well as the least twelve terms (ignoring  $e$ ), with the exception of the last row, are listed (in the order of  $\prec_{kbo}$ ) in the following table:

$w(a)$	$w(b)$	the least 12 terms in the order of $\prec_{kbo}$
1	1	$a, b, aa, ab, ba, bb, aaa, aab, aba, abb, baa, bab, \dots$
2	1	$b, a, bb, ab, ba, bbb, aa, abb, bab, bba, bbbb, abb, \dots$
2	3	$a, b, aa, ab, ba, aaa, bb, aab, aba, baa, aaaa, abb, \dots$
2	0	$b, bb, \dots, b^i, \dots, a, ab, \dots, ab^i, \dots, ba, bab, babb, \dots$

On the other hand, the least terms under  $\succ_{lpo}$  or  $\succ_{rpo}$  are

$$a, \dots, a^i, \dots, b, ab, \dots, a^i b, \dots, ba, aba, \dots, a^i ba, \dots, baa, abaa, \dots, a^i baa, \dots$$

The above example shows the flexibility of  $\succ_{kbo}$ , which provides various ways of comparing two terms by changing weight functions as needed. However,  $\succ_{kbo}$  lacks the flexibility of comparing terms like  $x * (y + z)$  and  $x * y + x * z$ . We cannot

make  $x * (y + z) \succ_{kbo} x * y + x * z$ , while  $x * (y + z) \succ_{lpo} x * y + x * z$  if  $* \succ +$  and  $x * y + x * z \succ_{lpo} x * (y + z)$  if  $+ \succ *$ .

Now when we say to use a simplification order  $\succ$ , we could use either  $\succ_{lpo}$ ,  $\succ_{rpo}$ , or  $\succ_{kbo}$ .

**Example 6.3.16.** Let the operators  $\leftrightarrow, \oplus, \rightarrow, \neg, \vee, \wedge$  be listed in the descending order of the precedence  $\succ$ , then all the right sides of the following equivalence relations are less than their left sides by either  $\succ_{lpo}$  or  $\succ_{rpo}$ .

$$\begin{aligned} A \oplus B &\equiv (A \vee B) \wedge (\neg A \vee \neg B); \\ A \leftrightarrow B &\equiv (A \vee \neg B) \wedge (\neg A \vee B); \\ A \rightarrow B &\equiv \neg A \vee B; \\ \neg \neg A &\equiv A; \\ \neg(A \vee B) &\equiv \neg A \wedge \neg B; \\ \neg(A \wedge B) &\equiv \neg A \vee \neg B; \\ A \vee (B \wedge C) &\equiv (A \vee B) \wedge (A \vee C). \end{aligned}$$

These equivalence relations are used to translate a formula into CNF. Thus, the term orderings like  $\succ_{rpo}$  and  $\succ_{lpo}$  can be used to show the termination of the process of converting a formula into CNF.

If we want to convert formulas into DNF, then the precedence relation should be  $\wedge \succ \vee$ , so that the termination of converting a formula into DNF can be shown by  $\succ_{rpo}$  and  $\succ_{lpo}$ .

### 6.3.2 Completeness of Ordered Resolution

Let  $\succ$  be a simplification order, which is a well-founded, stable partial order over the set of atoms and total over the set of ground atoms, we will use  $\succ$  in ordered resolution. Given  $\succ$ , an atom  $A \in S$  is said to be *maximal* if there is no  $B \in S$  such that  $B \succ A$ .  $A$  is said to be *minimal* if there is no  $B \in S$  such that  $A \succ B$ .  $S$  may have several minimal or maximal atoms. In fact, an atom  $A$  can be both minimal and maximal in  $S$  if  $A$  cannot be compared by  $\succ$  with any other atom of  $S$ . Literals are compared by comparing their atoms. Since a clause is represented by a set of literals, we may say a literal is maximal or minimal in a clause. For example, the three atoms in  $(\overline{p(x, y)} \mid \overline{p(y, z)} \mid p(x, z))$  are not comparable under any stable  $\succ$ , so each atom is both minimal and maximal.

**Definition 6.3.17. (ordered resolution)** Given a simplification order  $\succ$  over the set of atoms, the ordered resolution is the resolution with the condition that  $A$  is a maximal atom in  $\alpha$  and  $B$  is a maximal atom in  $\beta$ :

$$\frac{(A \mid \alpha) \quad (\neg B \mid \beta)}{(\alpha \mid \beta)\sigma}$$

Given a set  $S$  of clauses, let  $S^*$  denote the set of clauses *saturated* by ordered resolution and factoring. That is, any resolvent from ordered resolution on any two clauses of  $S^*$  is in  $S^*$  and any clause from factoring on any clause in  $S^*$  is also in  $S^*$ . Let  $GC$  be the set of all ground instances of  $S^*$ .

**Lemma 6.3.18.**  *$GC$  is saturated by ordered resolution.*

*Proof.* Suppose  $(g \mid \alpha)$  and  $(\bar{g} \mid \beta)$  are in  $GC$ , where  $g$  is the maximal atom in both clauses. The resolvent of the two clauses is  $(\alpha \mid \beta)$ . We need to show that  $(\alpha \mid \beta) \in GC$ .

Let  $(g \mid \alpha)$  be an instance of  $(A \mid \alpha') \in S^*$ , i.e.,  $g = A\theta$  and  $\alpha = \alpha'\theta$  for some substitution  $\theta$ . We assume that  $A$  is the only literal in  $(A \mid \alpha')$  such that  $A\alpha = g$ ; otherwise, we use factoring to achieve this condition. Similarly, let  $(\bar{g} \mid \beta)$  be an instance of  $(\bar{B} \mid \beta') \in S^*$ , i.e.,  $\bar{g} = \bar{B}\gamma$  and  $\beta = \beta'\gamma$  for some  $\gamma$ , and  $B$  is the only literal in  $(\bar{B} \mid \beta')$  such that  $B\gamma = g$ .

Since  $g = A\theta = B\gamma$  is maximal,  $A$  and  $B$  must be unifiable and maximal in  $(A \mid \alpha')$  and  $(\bar{B} \mid \beta')$ , respectively. Let  $\lambda$  be the mgu of  $A$  and  $B$ , then the resolvent of  $(A \mid \alpha')$  and  $(\bar{B} \mid \beta')$  by ordered resolution is  $(\alpha' \mid \beta')\lambda$ , which must be in  $S^*$ . Since  $\theta = \lambda\theta'$  and  $\gamma = \lambda\gamma'$  for some  $\theta'$  and  $\gamma'$ , we must have  $(\alpha \mid \beta) = (\alpha' \mid \beta')\lambda\theta'\gamma'$ , which is an instance of  $(\alpha' \mid \beta')\lambda$ . Hence,  $(\alpha \mid \beta) \in GC$ .  $\square$

**Theorem 6.3.19. (refutational completeness of ordered resolution)**  *$S$  is unsatisfiable iff the empty clause is in  $S^*$ .*

*Proof.* If the empty clause is in  $S^*$ , since  $S \models \perp$ ,  $S$  must be unsatisfiable.

If the empty clause is not in  $S^*$ , we will construct a Herbrand model for  $S$  from  $GC$ , which is the set of all ground instances of  $S^*$ . Obviously, the empty clause is not in  $GC$ . Let  $GA$  be the set of ground atoms appearing in  $GC$ . By the assumption,  $\succ$  is a well-founded total order over  $GA$ .

We start with an empty model  $H$  in which no atoms have a truth value. We then add ground literals into  $H$  one by one, starting from the minimal atom, which exists because  $<$  is well-founded. Let  $g$  be any atom in  $GA$ , and

$$\begin{aligned} pc(g) &= \{(g \mid \alpha) \mid g \text{ is maximal in } (g \mid \alpha) \in G\} \\ nc(g) &= \{(\bar{g} \mid \beta) \mid g \text{ is maximal in } (\bar{g} \mid \beta) \in G\} \\ upto(g) &= \bigcup_{g' \in GA, g' \leq g} (pc(g') \cup nc(g')) \end{aligned}$$

That is,  $pc(g)$  contains all the clauses which contains  $g$  as its maximal literal;  $nc(g)$  contains all the clauses which contains  $\bar{g}$  as its maximal literal; and  $upto(g)$  contains all the clauses from  $GC$  whose maximal literals are less or equal to  $g$ .

We claim that after  $g$  or  $\bar{g}$  is added into  $H$ , the following the property is true:

**Claim:**  $H$  is a model of  $upto(p)$  after  $g$  or  $\bar{g}$  is added into  $H$ .

Let  $g_1$  be the minimal atom of  $GA$ . We cannot have both  $(g_1)$  and  $(\bar{g}_1) \in GC$ ; otherwise the empty clause will be generated from  $(g_1)$  and  $(\bar{g}_1)$ . By Lemma 6.3.18, the empty clause will appear in  $S^*$ , a contradiction. If  $(g_1) \in GC$ , we add  $g_1$  into  $H$  (it means  $H(g_1) = 1$ ); otherwise, we add  $\bar{g}_1$  into  $H$  (it means  $H(g_1) = 0$ ). It is trivial to check that  $H$  is a model of  $upto(g_1)$ , so the claim is true for  $g_1$ . This is the base case of the induction based on  $<$ .

Let  $g$  be the minimal atom in  $GA$  which has no truth value in  $H$ . Assume as the induction hypothesis that the claim is true for all  $g' < g$ . If there exists a clause  $(g \mid \alpha) \in pc(g)$  such that  $H(\alpha) = 0$ , we add  $g$  into  $H$ ; otherwise, add  $\bar{g}$  into  $H$ .

Suppose  $g$  is added into  $H$ , then every clause in  $pc(g)$  will be true under  $H$ . If there is exists a clause  $(\bar{g} \mid \beta) \in nc(g)$  such that  $H(\beta) = 0$ , then there exists an ordered resolution between  $(g \mid \alpha)$  and  $(\bar{g} \mid \beta)$ , and their resolvent is  $(\alpha \mid \beta)$ . By Lemma 6.3.18,  $(\alpha \mid \beta) \in GC$ . Since every literal  $g'$  of  $(\alpha \mid \beta)$  is less than  $g$ , by the induction hypothesis,  $(\alpha \mid \beta)$  is true in  $H$ . If  $H(\alpha) = 0$ , then  $H(\beta) = 1$ . That means no clause in  $nc(g)$  will be false in  $H$ .

Now suppose  $\bar{g}$  is added into  $H$ , the above analysis still holds by changing the role of  $p$  and  $\bar{p}$ . In both cases, all clauses in  $pc(g)$  and  $nc(g)$  are true in  $H$  after adding  $g$  or  $\bar{g}$  into  $H$ . Thus the claim is true for every  $g \in GA$ . Once every atom in  $GA$  are processed, we have found a model  $H$  of  $GC$  by the claim. Since  $GC$  is the set of all ground instances of  $S^*$ ,  $H$  is also a model of  $S^*$  or  $S$ . In other words,  $S$  is satisfiable.  $\square$

Note that the model construction in the above proof is not algorithmic, because  $GC$  is infinite in general. The induction is sound because it is based on the well-founded order  $\succ$ .

Now it is easy to prove Theorem 6.2.7: If the empty clause is generated from  $S$ , then  $S$  must be unsatisfiable; if  $S$  is unsatisfiable, then the empty clause will be generated by ordered resolution, which is a special case of resolution.

A traditional proof of Theorem 6.2.7 uses Herbrand's Theorem:

**Theorem 6.3.20. (Herbrand's Theorem)** *A set  $S$  of clauses is unsatisfiable iff there exists a finite, unsatisfiable, ground instances of  $S$ .*

A proof of Theorem 6.2.7 goes as follows:

1. If  $S$  is unsatisfiable, then by Herbrand's Theorem, there exists a finite set  $G$  of ground instances of  $S$ .
2. Resolution for propositional logic will find a proof from  $G$ , treating each ground atom in  $G$  as a propositional variable.

3. This ground resolution proof can be lifted to a general resolution proof in  $S$ , using a lemma similar to Lemma 6.3.18.

## 6.4 Prover9: A Resolution Theorem Prover

Prover9 is a full-fledged automated theorem prover for first-order logic based on resolution. Prover9 is the successor of the Otter prover; both Otter and Prover9 are created by William McCune (1953-2011).

### 6.4.1 Input Formulas to Prover9

Prover9 has a fully automatic mode in which the user simply gives it formulas representing the problem. A good way to learn about Prover9 is to browse and study the example input and output files that are available with the distribution of Prover9. Let us look at an example.

**Example 6.4.1.** Once we can specify the following puzzle in first-order logic, it is trivial to find a solution by Prover9:

Jack owns a dog. Every dog owner is an animal lover. No animal lover kills an animal. Either Jack or Curiosity killed the cat, who is named Tuna. Did Curiosity kill the cat?

The following first-order formulas come from the statements of the puzzle:

- (1)  $\exists x(\text{dog}(x) \wedge \text{owns}(\text{Jack}, x))$
- (2)  $\forall x(\exists y(\text{dog}(y) \wedge \text{owns}(x, y)) \rightarrow \text{animalLover}(x))$
- (3)  $\forall x(\text{animalLover}(x) \rightarrow (\forall y(\text{animal}(y) \rightarrow \neg \text{kills}(x, y))))$
- (4)  $\text{kills}(\text{Jack}, \text{Tuna}) \vee \text{kills}(\text{Curiosity}, \text{Tuna})$
- (5)  $\text{cat}(\text{Tuna})$
- (6)  $\forall x(\text{cat}(x) \rightarrow \text{animal}(x))$

The goal is to show that  $\text{kills}(\text{Curiosity}, \text{Tuna})$ . The above formulas can be written in Prover9's syntax as follows:

```
exists x (dog(x) & owns(Jack, x)).
all x (exists y (dog(y) & owns(x, y)) -> animalLover(x)).
all x (animalLover(x) -> (all y (animal(y) -> -kills(x, y)))).
kills(Jack, Tuna) | kills(Curiosity, Tuna).
cat(Tuna).
all x (cat(x) -> animal(x)).
```

From this example, we can see that each formula in Prover9 is ended with “.”, and we use “->” for  $\rightarrow$ , “|” for  $\vee$ , “&” for  $\wedge$ , “-” for  $\neg$ , “all” for  $\forall$ , and “exists” for  $\exists$ . In fact, the top most “all” is optional as free variables are assumed to universally quantified. Prover9 will convert formulas in the “assumptions” list and the negation of the formulas in the “goals” list into clauses before resolution is called.

Note that Prover9 uses a different precedence for the quantifiers “all” and “exists”; they are higher than other Boolean operators. For instance,

```
exists x dog(x) & owns(Jack, x)
```

represents the formula  $(\exists x \text{dog}(x)) \wedge \text{owns}(\text{Jack}, x)$ , not  $\exists x(\text{dog}(x) \wedge \text{owns}(\text{Jack}, x))$  as assumed in this book. To avoid confusions, insert a pair of parentheses when you are uncertain about the scope of a quantified variable.

Prover9 implements the procedure *resolution* (Procedure 6.2.9) with bells and whistles. The given list is called the “sos” (set of support) list and the kept list is called the “usable” list. If you already have a set of clauses ready, you may use the “usable” list and the “sos” list as follows:

```
formulas(usable). % the kept list
...
end_of_list.

formulas(sos). % the given list
...
end_of_list.
```

The resolution is always done between a clause from the sos list (the given list) and some clauses from the usable list (the kept list).

A basic Prover9 command on a linux machine will look like

```
prover9 -f Tuna.in > Tuna.out
```

or

```
prover9 < Tuna.in > Tuna.out
```

where the file named “Tuna.in” contains the input to Prover9 and the file named “Tuna.out” contains the output of Prover9. If the file “Tuna.in” contains the formulas in Example 6.4.1, then the file “Tuna.out” will contain the following proof:

```

===== PROOF =====

% ----- Comments from original proof -----
% Proof 1 at 0.03 (+ 0.05) seconds.
% Length of proof is 19.
% Level of proof is 6.
% Maximum clause weight is 6.
% Given clauses 0.

1 (exists x (dog(x) & owns(Jack,x))). [assumption].
2 (all x ((exists y (dog(y) & owns(x,y))) -> animalLover(x))).
   [assumption].
3 (all x (animalLover(x) -> (all y (animal(y) -> -kills(x,y))))))
   [assumption].
4 (all x (cat(x) -> animal(x))). [assumption].
5 kills(Curiosity,Tuna). [goal].
6 -dog(x) | -owns(y,x) | animalLover(y). [clausify(2)].
7 dog(c1). [clausify(1)].
8 -owns(x,c1) | animalLover(x). [resolve(6,a,7,a)].
9 owns(Jack,c1). [clausify(1)].
10 animalLover(Jack). [resolve(8,a,9,a)].
11 -animalLover(x) | -animal(y) | -kills(x,y). [clausify(3)].
12 -cat(x) | animal(x). [clausify(4)].
13 cat(Tuna). [assumption].
14 animal(Tuna). [resolve(12,a,13,a)].
15 -animal(x) | -kills(Jack,x). [resolve(10,a,11,a)].
16 kills(Jack,Tuna) | kills(Curiosity,Tuna). [assumption].
17 -kills(Curiosity,Tuna). [deny(5)].
18 -kills(Jack,Tuna). [resolve(14,a,15,a)].
19 $F. [back_unit_del(16),unit_del(a,18),unit_del(b,17)].

===== end of proof =====

```

In Prover9, the empty clause is denoted by \$F. In each clause, the literals are numbered by *a*, *b*, *c*, .... For example, to generate clause 18, `-kills(Jack, Tuna)`, from

```

14 animal(Tuna). [resolve(12,a,13,a)].
15 -animal(x) | -kills(Jack,x). [resolve(10,a,11,a)].

```

Prove9 performed a resolution on the first literal (a) of clause 14 and the first literal (a) of clause 15, with the mgu  $x \leftarrow \text{Tuna}$ , and the resolvent is

```
18 -kills(Jack,Tuna). [resolve(14,a,15,a)].
```

The message followed by each clause can be checked by either human or machine, to ensure the correctness of the proof.

#### 6.4.2 Inference Rules and Options

There are two types of parameters in Prover9: Boolean flags and numeric parameters. To change the default values of the former type, use `set(flag)` or `clear(flag)`. The latter can be changed by `assign(parameter, value)`. For example, the default automatic mode is set by

```
set(auto).
```

Turning “auto” on will cause a list of other flags to turn on, including setting up “hyper-resolution” as the main inference rule for resolution. Hyper-resolution reduces the number of intermediate resolvents by combining several resolution steps into a single inference step.

**Definition 6.4.2. (hyper-resolution)** Positive hyper-resolution *consists of a sequence of positive resolutions between one non-positive clause (called nucleus) and a set of positive clauses (called satellites), until a positive clause or the empty clause is produced.*

The number of positive resolutions in a hyper-resolution is equal to the number of negative literals in the nucleus clause.

**Example 6.4.3.** Let the nucleus clause be  $(\overline{p(x)} \mid \overline{q(x)} \mid r(x))$ , and the satellite clauses be  $(q(a) \mid r(b))$  and  $(p(a) \mid r(c))$ . The first resolution between the nucleus and the first satellite produces  $(\overline{p(a)} \mid r(a) \mid r(b))$ . The second resolution between the resolvent and the second satellite produces the second resolvent  $(r(a) \mid r(b) \mid r(c))$ , which is the result of positive hyper-resolution between the nucleus and the satellites.

Negative hyper-resolution can be defined similarly by replacing “positive” by “negative”. Hyper-resolution means both positive hyper-resolution and negative hyper-resolution.

**Example 6.4.4.** In section 3.2.2, we discussed the Hilbert system which consist of three axioms and one inference rule, i.e., *Modus Ponens*. To prove  $x \rightarrow x$  in the Hilbert system, we may use the following input:

```
op(400, infix_right, ["->"]). % infix operator
```

```
formulas(usable).
P((n(y) -> n(x)) -> (x -> y)).
-P(x) | -P(x -> y) | P(y).
end_of_list.
```

```
formulas(sos).
P(x -> (y -> x)).
P((x -> (y -> z)) -> ((x -> y) -> (x -> z))).
-P(a -> a).
end_of_list.
```

Prover9 will produce the following proof:

```
1 -P(x) | -P(x -> y) | P(y). [assumption].
2 P(x -> y -> x). [assumption].
3 P((x -> y -> z) -> (x -> y) -> x -> z). [assumption].
4 -P(a -> a). [assumption].
5 P(x -> y -> z -> y). [hyper(1,a,2,a,b,2,a)].
7 P((x -> y) -> x -> x). [hyper(1,a,2,a,b,3,a)].
15 P(x -> x). [hyper(1,a,5,a,b,7,a)].
16 $F. [resolve(15,a,4,a)].
```

The hyper-resolution is the default inference rule for the automatic mode. If you like to see a binary resolution proof, you have to turn off “auto” and turn on “binary\_resolution”:

```
clear(auto).
set(binary_resolution).
```

Turning off hyper-resolution, Prover9 will produce the following proof:

```
1 -P(x) | -P(x -> y) | P(y). [assumption].
2 P(x -> y -> x). [assumption].
3 P((x -> y -> z) -> (x -> y) -> x -> z). [assumption].
4 -P(a -> a). [assumption].
5 -P(x) | P(y -> x). [resolve(2,a,1,b)].
6 -P(x -> y -> z) | P((x -> y) -> x -> z). [resolve(3,a,1,b)].
10 P((x -> y) -> x -> x). [resolve(6,a,2,a)].
19 P(x -> (y -> z) -> y -> y). [resolve(10,a,5,a)].
```

```

20  $\neg P(x \rightarrow y) \mid P(x \rightarrow x)$ . [resolve(10,a,1,b)].
45  $P(x \rightarrow x)$ . [resolve(20,a,19,a)].
46 $F. [resolve(45,a,4,a)].

```

**Definition 6.4.5.** Unit-resulting resolution *consists of a sequence of unit resolutions between one non-unit clause (called nucleus) and a set of unit clauses (called electrons), until a unit clause or the empty clause is produced.*

Since positive resolution is complete and unit resolution is not complete, hyper-resolution is complete and unit-resulting resolution is not complete. To use “unit-resulting resolution”, use the command

```
set(ur_resolution).
```

To use “hyper-resolution” in Prover9, use the command

```
set(hyper_resolution).
```

This option will cause the flags `pos_hyper_resolution` and `neg_hyper_resolution` to be true.

The resolution inference rules provided by Prover9 include “binary resolution”, “hyper resolution”, “ordered resolution”, and “unit-resulting resolution”. To use “ordered resolution”, use the command

```
set(ordered_res).
```

This option puts restrictions on the binary and hyper-resolution inference rules. It says that resolved literals in one or more of the parents must be maximal in the clause. Continuing from the previous example, if we turn on the ordered resolution flag, the proof generated by Prover9 will be the following.

```

1  $\neg P(x) \mid \neg P(x \rightarrow y) \mid P(y)$ . [assumption].
2  $P(x \rightarrow y \rightarrow x)$ . [assumption].
3  $P((x \rightarrow y \rightarrow z) \rightarrow (x \rightarrow y) \rightarrow x \rightarrow z)$ . [assumption].
4  $\neg P(a \rightarrow a)$ . [assumption].
5  $\neg P(x) \mid P(y \rightarrow x)$ . [resolve(2,a,1,b)].
6  $\neg P(x \rightarrow y \rightarrow z) \mid P((x \rightarrow y) \rightarrow x \rightarrow z)$ . [resolve(3,a,1,b)].
8  $P(x \rightarrow y \rightarrow z \rightarrow y)$ . [resolve(5,a,2,a)].
9  $P(x \rightarrow y \rightarrow z \rightarrow u \rightarrow z)$ . [resolve(8,a,5,a)].
10  $P(x \rightarrow y \rightarrow z \rightarrow u \rightarrow w \rightarrow u)$ . [resolve(9,a,5,a)].
14  $P((x \rightarrow y) \rightarrow x \rightarrow x)$ . [resolve(6,a,2,a)].
19  $\neg P(x \rightarrow y) \mid P(x \rightarrow x)$ . [resolve(14,a,1,b)].
21  $P(x \rightarrow x)$ . [resolve(19,a,10,a)].
22 $F. [resolve(21,a,4,a)].

```

### 6.4.3 Simplification Orders in Prover9

Prover9 has available several methods for comparing terms. The term orders are partial (and sometimes total on ground terms), and they are used to decide which literals in clauses are admissible for application of ordered resolution. Several of the resolution rules require that some of the literals be maximal in their clause.

The symbol precedence is a total order on function and predicate symbols (including constants). The symbol weighting function maps symbols to nonnegative integers. Prover9 supports three term orders: LPO (Lexicographic Path Order), RPO (Recursive Path Order), and KBO (Knuth-Bendix Order).

Here is the command for choosing a simplification order.

```
assign(order, string). % default string=lpo, range [lpo,rpo,kbo]
```

This option is used to select the primary term order to be used for determining maximal literals in clauses. The choices are “lpo” (Lexicographic Path Order), “rpo” (Recursive Path Order), and “kbo” (Knuth-Bendix Order).

The default symbol precedence (for LPO, RPO, and KBO) is given by the following rules (in order).

- function symbols < non-equality predicate symbols;
- for function symbols:  $c/0 < f/2 < g/1 < h/3 < i/4 < \dots$ , where  $c$  is any constant,  $f$  is any function of arity 2,  $g$  is any function of arity 1, ... (note the position of  $g/1$ );
- for predicate symbols: lower arity < higher arity;
- non-Skolem symbols < Skolem symbols;
- for Skolem symbols, the lower index is the lesser;
- for non-Skolem symbols, more occurrences < fewer occurrences;
- the lexical ASCII order (UNIX strcmp() function).

The `function_order` and `predicate_order` commands can be used to change the default symbol precedence. They contain lists of symbols ordered by increasing precedence. For example,

```
predicate_order([=, <=, P, Q]).           % = < <= < P < Q
function_order([a, b, c, +, *, h, g]).     % a < b < c < + < * < h < g
```

We need two separate commands for defining the precedence, because we assume that predicate symbols are always greater than function symbols in the precedence. The used symbol precedence for a problem is always printed in the output file (in the section PROCESS INPUT).

### 6.5 Exercise Problems

1. Let  $\theta = [x \leftarrow u, y \leftarrow f(g(u)), z \leftarrow f(u)]$  and  $\sigma = [u \leftarrow x, x \leftarrow f(u), y \leftarrow a, z \leftarrow f(u)]$ . Compute  $\theta\sigma$  and  $\sigma\theta$ . Among  $\theta$ ,  $\sigma$ ,  $\theta\sigma$ , and  $\sigma\theta$ , which substitutions are idempotent?
2. Use the rule-based unification algorithm *unify* to decide if the following pairs of terms are unifiable or not, and provide each state in *unify* after applying a transformation rule.
  - (a)  $p(x, a) \doteq p(f(y), y)$
  - (b)  $p(x, f(x)) \doteq p(f(y), y)$
  - (c)  $p(f(a), g(x)) \doteq p(y, y)$
  - (d)  $q(x, y, f(y)) \doteq q(u, h(v, v), u)$
  - (e)  $q(a, x, f(g(y))) \doteq q(z, f(z), f(u))$
3. How to modify *unify*( $s, t$ ) so that it can be used to decide if a set  $S$  of terms are unifiable? That is, find an idempotent substitution  $\sigma$  and a term  $t$  such that  $s\sigma = t$  for every  $s \in S$ .
4. Find the mgu of  $s = f(f(\dots f(f(a, x_1), x_2), \dots), x_n)$  and  $t = f(x_n, f(x_{n-1}, f(\dots, f(x_1, a)\dots)))$  for  $n = 5$ .
5. Create the term graph  $G = (V, E)$  for each of the following pairs of terms:
  - (a)  $s = p(x, a)$  and  $t = p(f(y), y)$
  - (b)  $s = p(x, f(x))$  and  $t = p(f(y), y)$
  - (c)  $s = q(x, y, f(y))$  and  $t = q(u, h(v, v), u)$
  - (d)  $s = q(a, x, f(g(y)))$  and  $t = q(z, f(z), f(u))$
  - (e)  $s = f(f(f(f(a, x_1), x_2), x_3), x_4)$  and  $t = f(x_4, f(x_3, f(x_2, f(x_1, a))))$

and show the equivalence classes of  $V$  after calling *unify1* with a success.
6. Provide the pseudo code for the procedure *postOccurCheck* used in *lunify*. What is the complexity of your algorithm?

7. Design an efficient algorithm and provide its pseudo code for creating the mgu when the procedure *lunify* returns true. What is the complexity of your algorithm?
8. Show that all the left side  $s$  of each equivalence relation in Example 6.3.16 is greater than their right side  $t$  by  $\succ_{lpo}$ , and explain what precedence is used and what conditions of  $\succ_{lpo}$  are used for  $s \succ_{lpo} t$ .
9. Show that all the left side  $s$  of each following equation is greater than their right side  $t$  by  $\succ_{lpo}$ , and explain what precedence is used and what conditions of  $\succ_{lpo}$  are used for  $s \succ_{lpo} t$ .

$$\begin{aligned}x \times (y + z) &= (x \times y) + (x \times z) \\(y + z) \times x &= (y \times x) + (z \times x) \\(x \times y) \times z &= x \times (y \times z)\end{aligned}$$

Can you use  $\succ_{rpo}$  or  $\succ_{kbo}$  to achieve the same result? Why?

10. Choose with justification a simplification order from  $\succ_{lpo}$ ,  $\succ_{rpo}$  and  $\succ_{kbo}$  such that the left side  $s$  is greater than the right side  $t$  for each of the following equations (which defines the Ackermann's function):

$$\begin{aligned}A(0, y) &= s(y) \\A(s(x), 0) &= A(x, s(0)) \\A(s(x), s(0)) &= A(x, A(s(x), y))\end{aligned}$$

11. Given the following statements:

Tony, Tom and Liz belong to the Hoofers Club. Every member of the Hoofers Club is either a skier or a mountain climber or both. No mountain climber likes rain, and all skiers like snow. Liz dislikes whatever Tony likes and likes whatever Tony dislikes. Tony likes rain and snow.

The “Hoofers Club problem” asks for the solution of the following question: “Is there a member of the Hoofers Club who is a mountain climber but not a skier?” Please find the following resolution proofs for “Hoofers Club problem”: (a) unit, (b) input, (c) positive, (d) negative, and (e) linear resolutions.

12. Use Prover9 to answer the question in the “Hoofers Club problem” of the previous problem. Please prepare the input to Prover9 and turn in the output file of Prover9.

13. Use binary resolution, hyper-resolution and ur-resolution, respectively, of Prover9 to answer the question that “Is John happy?” from the following statements:

Anyone passing his logic exams and winning the lottery is happy.  
 But anyone who studies or is lucky can pass all his exams. John did not study but he is lucky. Anyone who is lucky wins the lottery.

14. Use hyper-resolution and ur-resolution, respectively, of Prover9 to show that “the sprinklers are on” from the following statements:

Someone living in your house is wet. If a person is wet, it is because of the rain, the sprinklers, or both. If a person is wet because of the sprinklers, the sprinklers must be on. If a person is wet because of rain, that person must not be carrying any umbrella. There is an umbrella in your house, which is not in the closet. An umbrella that is not in the closet must be carried by some person who lives in that house. Nobody are carrying any umbrella.

Please prepare the input to Prover9 and turn in the output file of Prover9.

15. Following the approach illustrated in Example 6.4.4, use Prover9 to prove that the following properties are true in the Hilbert system:

- (a)  $(x \rightarrow y) \rightarrow ((y \rightarrow z) \rightarrow (x \rightarrow z))$
- (b)  $(x \rightarrow (y \rightarrow z)) \rightarrow (y \rightarrow (x \rightarrow z))$
- (c)  $\neg x \rightarrow (x \rightarrow y)$
- (d)  $x \rightarrow (\neg x \rightarrow y)$
- (e)  $\neg\neg x \rightarrow x$
- (f)  $x \rightarrow \neg\neg x$
- (g)  $(x \rightarrow \neg x) \rightarrow \neg x$
- (h)  $(\neg x \rightarrow x) \rightarrow x$
- (i)  $(x \rightarrow y) \rightarrow (\neg y \rightarrow \neg x)$

## CHAPTER 7

### EQUATIONAL LOGIC

There are many equivalence relations in logic. At semantic level, we have logical equivalence and equal satisfiability, which are denoted by the symbols  $\equiv$  and  $\approx$ , respectively, for both propositional and first-order logics. At syntactic level, we have the logical operator  $\leftrightarrow$  to express the equality between two statements. Since a first-order language has two types of objects, i.e., formulas and terms,  $\leftrightarrow$  is used for the formulas, the conventional symbol for the equality of terms is “=”, which is a predicate symbol. Instead, we have been using “=” as a syntactical relation outside of the logical languages.

#### 7.1 Equality of Terms

The equality is an omnipresent and important relation in every field of the mathematics. For example, how would you specify that “for every  $x$ , there exists a unique  $y$  such that the relation  $p(x, y)$  holds”? For example, if the meaning of  $p(x, y)$  is “ $x$  is the mother of  $y$ ”, then how would you state in first-order logic that “everybody has a unique mother”? If the meaning of  $p(x, y)$  is  $f(x) = y$ , then how would you state that “the function  $f(x)$  has a unique value”? Using =, the answer is easy:

$$\forall x \exists y p(x, y) \wedge (\forall z p(x, z) \rightarrow y = z)$$

The above formula is often denoted by  $\forall x \exists! y p(x, y)$  in mathematics. In fact, when we translate a first-order formula into a propositional formula for searching a finite model, we used the formula

$$(f(x_1, \dots, x_k) \neq y_1 \mid f(x_1, \dots, x_k) \neq y_2 \mid y_1 = y_2),$$

which states that  $f(x_1, \dots, x_k)$  has a unique value. Using “=” gives us higher expressive power to specify statements.

- We can express “there are at least two elements  $x$  such that  $A(x)$  holds” as  $\exists x \exists y x \neq y \wedge A(x) \wedge A(y)$ .
- We can express “there are at most two elements  $x$  such that  $A(x)$  holds” as  $\forall x \forall y \forall z (A(x) \wedge A(y) \wedge A(z)) \rightarrow x = y \vee y = z \vee x = z$  This states that if we have three elements  $a$  for which  $A(a)$  holds, then two of them must be equal.

- We can express “there are exactly two elements  $x$  such that  $A(x)$  holds” as the conjunction of the above two statements.

### 7.1.1 Axioms of Equality

The axioms of the equality consist of five types of formulas, each is sufficiently famous to have earned itself a name.

- **reflexivity:**  
 $x = x$ . That is, an object is always equal to itself.
- **symmetry:**  
 $(x = y) \rightarrow (y = x)$ . That is, it does not matter how the parameters of  $=$  are ordered.
- **transitivity:**  
 $(x = y) \wedge (y = z) \rightarrow (x = z)$ . That is, the equality relation can be inherited.
- **function monotonicity:**  
 $x_i = y_i \rightarrow f(\dots, x_i, \dots) = f(\dots, y_i, \dots)$ . That is, if the arguments of a function are equal, the composed term will be equal.
- **predicate monotonicity:**  
 $x_i = y_i \wedge p(\dots, x_i, \dots) \rightarrow p(\dots, y_i, \dots)$ . That is, if the augments of a predicate are equal, then the resulting atoms are logically equivalent.

The first three axioms together say that “ $=$ ” is an equivalence relation. The fourth and fifth are very similar and share the name *monotonicity*. The first monotonicity rule governs terms and the second monotonicity rule governs formulas. Both of them assert that having equal arguments ensures equal results. In fact, these are axiom schemata as we need a version for each function  $f$  of arity  $k$  and for each value  $i$ ,  $1 \leq i \leq k$ . The similar requirement applies to each predicate  $p$ . For example, if the predicate  $p$  is the equality  $=$ , we obtain two axioms from the predicate monotonicity:

$$\begin{aligned} x_1 = y_1 \wedge (x_1 = x_2) &\rightarrow (y_1 = x_2), \\ x_2 = y_2 \wedge (x_1 = x_2) &\rightarrow (x_1 = y_2). \end{aligned}$$

Of course, these two axioms can be easily deduced from the symmetry and transitivity of  $=$ . In other words, the predicate monotonicity is needed for predicates other than  $=$ . We will need different versions of the fourth and fifth axioms in different

first-order languages. These monotonicity axioms are not to be confused with the substitution rule, which allows the monotonicity of terms for variables.

All the five axioms can be easily converted into clauses so that we can use resolution to prove theorems involving  $=$ .

**Example 7.1.1.** Let  $f(a) = b$  and  $f(b) = a$ , prove that  $f(f(a)) = a$ . Since resolution is a refutational prover, we need to add the negation of  $f(f(a)) = a$ ,  $f(f(a)) \neq a$ , into the clauses from the equality axioms and the premises, and show the clause set is unsatisfiable.

1	$(f(a) = b)$	premise
2	$(f(b) = a)$	premise
3	$(f(f(a)) \neq a)$	negation of the goal
4	$(x \neq y \mid y \neq z \mid x = z)$	transitivity
5	$(x \neq y \mid f(x) = f(y))$	monotonicity
6	$(f(f(a)) = f(b))$	resolvent from 1 and 5
7	$(f(b) \neq z \mid f(f(a)) = z)$	resolvent from 6 and 4
8	$(f(f(a)) = a)$	resolvent from 2 and 7
9	$()$	resolvent from 3 and 8

The above example shows clearly that without equality axioms, we cannot have a resolution proof of  $f(f(a)) = a$ .

### 7.1.2 Semantics of “=”

The equality symbol is meant to model the equality of objects in a domain. For example, when we say “ $2 + 0 = 2$ ”, or  $add(s(s(0)), 0) = s(s(0))$ , we meant that “ $2+0$ ” and “ $2$ ” represent the same object in a domain. That is,  $s = t$  is true if  $s$  and  $t$  represent “equal” or “identical” objects in any domain. We are asserting that two different descriptions refer to the same object. Because the notion of identity can be applied to virtually any domain of objects, the equality is often assumed to be omnipresent in every logic. However, talk of “equality” or “identity” raises messy philosophical questions. “Am I the same person I was three days ago?”

We need a simple and clear way to define the meaning of  $=$ . For a first-order language  $L = (P, F, X, Op)$ , a Herbrand model for a CNF  $A$  can be simply represented by a set  $H$  of ground atoms such that only those atoms in  $H$  are interpreted to be true. Because of the existence of the equality axioms, the set  $H$  must hold some properties.

**Definition 7.1.2. (congruence)** *A relation  $=_m$  over the set of terms  $T(F, X)$  is*

called a congruence if  $=_m$  is an equivalence relation satisfying the monotonicity:

$$f(s_1, \dots, s_k) =_m f(s_2, \dots, s_k) \text{ if for any } 1 \leq i \leq k, s_i =_m t_i$$

for every function  $f/k$ .

**Proposition 7.1.3.** *Let  $H$  be a Herbrand model of  $A$  which is in CNF and contains all the equality axioms. Define  $s =_H t$  iff  $(s = t) \in H$  for every  $s, t \in T(F)$ . Then  $=_H$  is a congruence over  $T(F)$ .*

*Proof.* The first three equality axioms ensure that  $=_H$  is an equivalence relation and the fourth axiom ensures that  $=_H$  is monotonic.  $\square$

For any congruence relation  $=_m$ , since  $=_m$  is an equivalence relation, we may partition  $T(F)$  into equivalence classes such that each class contains equivalent objects. For every  $s \in T(F)$ , let

$$[s] = \{t \mid s =_m t, t \in T(F)\},$$

the *congruence class* of  $s$ . Let

$$T(F)/=_m = \{[s] \mid s \in T(F)\},$$

the set of all congruence classes, which is also called the *Herbrand base modulo  $=_m$* .

**Example 7.1.4.** Let  $F = \{0, s^{(1)}\}$  and  $A$  contains  $s(s(s(x))) = x$  and the equality axioms. Let  $H$  be a Herbrand model of  $A$  in which  $(0 = s(0)) \notin H$ , then

$$\begin{aligned} T(F) &= \{0, s(0), s(s(0)), s(s(s(0))), s^4(0), \dots\}, \\ [0] &= \{0, s(s(s(0))), s^6(0), s^9(0), \dots\}, \\ [s(0)] &= \{s(0), s^4(0), s^7(0), s^{10}(0), \dots\}, \\ [s(s(0))] &= \{s(s(0)), s^5(0), s^8(0), s^{11}(0), \dots\}, \\ T(F)/=_H &= \{[0], [s(0)], [s(s(0))]\}. \end{aligned}$$

Now, it is easy to check if two ground terms  $s$  and  $t$  are equal or not under  $H$ :  $(s = t) \in H$  iff  $t \in [s]$  (or  $s \in [t]$ ), i.e.,  $s =_H t$ . In fact, when we define a Herbrand model for a formula with equality, at first we can define a congruence  $=_m$ , and then regard  $T(F)/=_m$  as the domain of the Herbrand model with equality:  $(s = t) \in H$  iff  $s =_m t$ . This way, we do not have to concern about the equality axioms explicitly.

### 7.1.3 Theory of Equations

A set of equations is simply a set of positive unit clauses where the only predicate symbol is “=”.

**Example 7.1.5.** In modern algebra, a *group* is  $G = (S, *)$ , where  $S$  is a set of elements and  $*$  is a binary operator  $*$  satisfying the properties that (a)  $*$  is associative and closed on  $S$ ; (b) there exists an identity element in  $S$ ; and (c) every element of  $S$  has an inverse in  $S$ . Using the equational approach, the identity is denoted by  $e$ , for each  $x \in S$ , the inverse of  $x$  is denoted by  $i(x)$ , with the following three equations as axioms:

$$\begin{array}{ll} 1 & x * e = x \quad e \text{ is the identity element} \\ 2 & x * i(x) = e \quad i(x) \text{ is the inverse of } x \\ 3 & (x * y) * z = x * (y * z) \quad * \text{ is associative} \end{array}$$

From these three equations, we may prove many interesting properties of the group theory, such as  $i(i(x)) = x$ .

**Definition 7.1.6.** Given a set  $E$  of equations, the relation  $=_E$  is recursively defined as follows:

$$\begin{array}{l} 1 \quad s =_E t \text{ if } s = t \in E; \\ 2 \quad t =_E t \text{ for any } t \in T(F, X); \\ 3 \quad s =_E t \text{ if } t =_E s; \\ 4 \quad s =_E t \text{ if } s =_E r, r =_E t; \\ 5 \quad f(\dots, s_i, \dots) =_E f(\dots, t_i, \dots) \text{ if } s_i =_E t_i \text{ for any } f \in F; \\ 6 \quad s\sigma =_E t\sigma \text{ if } s =_E t \text{ for any } \sigma : X \longrightarrow T(F, X). \end{array}$$

Recall that every binary relation  $R$ , including  $=_E$ , can be regarded as a set of pairs of items, i.e.,  $\{(s, t) \mid R(s, t)\}$ . Similarly,  $E$  can be also regarded as a set of pairs of terms.

**Definition 7.1.7.** Given a binary relation  $R$  over  $T(F, X)$ ,  $R$  is said to be closed under

- reflexivity if  $(t, t) \in R$  for any  $t \in T(F, X)$ ;
- symmetry if  $(t, s) \in R$  whenever  $(s, t) \in R$ ;
- transitivity if  $(r, t) \in R$  whenever  $(r, s), (s, t) \in R$ ;
- monotonicity if  $(f(\dots, s, \dots), f(\dots, t, \dots)) \in R$  whenever  $(s, t) \in R$ ;
- instantiation if  $(s\sigma, t\sigma) \in R$  whenever  $(s, t) \in R$  for any  $\sigma : X \longrightarrow T(F, X)$ .

The transitive closure is the minimum superset of  $R$  that is closed under transitive; the equivalence closure is the minimum superset of  $R$  that is closed under reflexivity, symmetry and transitivity; the congruence closure is the minimum equivalence closure of  $R$  that is closed under monotonicity.

**Proposition 7.1.8.** For any set  $E$  of equations,  $=_E$  is a minimum stable congruence containing  $E$  and closed under instantiation.

*Proof.* Let us check the conditions for  $s =_E t$  to be true: Condition 1 says  $E \subset =_E$ ; condition 2 says  $(t, t) \in =_E$ . Conditions 3-6 says  $=_E$  is closed under symmetry, transitivity, monotonicity, and instantiation.  $\square$

**Example 7.1.9.** Using  $F$  from Example 7.1.4, let  $E = \{s(s(s(x))) = x\}$  and  $X = \{x\}$ . Since  $=_E$  is a congruence, we can check that  $T(F)/=_E = T(F)/=_H = \{[0], [s(0)], [s(s(0))]\}$ ,  $[x] = \{x, s(s(s(x))), s^6(x), \dots\}$ , and

$$T(F, X)/=_E = \{[0], [s(0)], [s(s(0))], [x], [s(x)], [s(s(x))]\}.$$

**Proposition 7.1.10.** Given  $E$ , let  $E_{ax}$  be the equality axioms associated with  $E$ , then for every  $(s, t) \in =_E$ ,  $E \cup E_{ax} \models (s = t)$ .

*Proof.* (Sketch) Check all the six conditions for  $s =_E t$ : Condition 1 is sound because  $A \models A$ . Conditions 2-5 are sound because of  $E_{ax}$ . Condition 6 is sound because the free variables are assumed to universally quantified.  $\square$

Since  $=_E$  contains all the equations that are logical consequence of  $E$  and the equality axioms, we call  $=_E$  the *theory* of  $E$ .

**Example 7.1.11.** Let  $E$  be the three equations in Example 7.1.5, we may show that  $e * x =_E x$ ,  $i(x) * x =_E e$ ,  $i(e) =_E e$ ,  $i(i(x)) =_E x$ ,  $i(x * y) =_E i(y) * i(x)$ , etc. The proofs will be provided in the next section.

Given any two terms, say  $s$  and  $t$ , does  $s =_E t$ ? This is an important decision problem with many applications in mathematics. In computation theory, all computable functions can be constructed by equations. Unfortunately, it is undecidable in general to answer  $(s, t) \in =_E$ , in other words we do have an algorithm which takes  $E, s$ , and  $t$  as input and returns yes if  $s =_E t$ .

In computer science, the *congruence closure problem* refers to the problem of deciding  $s =_E t$  when  $E$  is a set of ground equations. This is a decidable problem and there exist efficient algorithms, some of them are based on the union-find data structures we used for the unification algorithm. There are a number of applications using congruence closures.

## 7.2 Rewrite Systems

In Example 7.1.1, we showed how to prove  $f(f(a)) = a$  from  $f(a) = b$  and  $f(b) = a$ , using resolution and the equality axioms. A straightforward proof is to replace  $f(a)$  in  $f(f(a)) = a$  by  $b$  to get  $f(b) = a$ , and then replace  $f(b)$  by  $a$  to get  $a = a$ . In this simple example, we treat  $f(a) = b$  and  $f(b) = a$  as two rewrite rules and each replacement is one application of a rewrite rule, called *rewriting*. To replace “equal” by “equal” is a common practice in mathematics and logics. For example, we used equivalence relations as rewrite rules to simplify formulas to normal forms. Rewrite rules are not only a common inference system in everyday mathematics, but also a very efficient computational technique, as they can involve very little search. We will devote this section to a brief study of them, including a discussion of some interesting theoretical results.

### 7.2.1 Rewrite Rules

**Definition 7.2.1.** A rewrite rule is an ordered pair of terms  $(l, r)$ , often denoted by  $l \rightarrow r$ , where  $l$  is the left side and  $r$  is the right side of the rule. A term  $t$  is rewritten to  $t'$  by  $l \rightarrow r$ , denoted by  $t \Rightarrow t'$ , if there exists a substitution  $\sigma$  and a position  $p$  of  $t$  such that  $l\sigma = t/p$  and  $t' = t[p \leftarrow r\sigma]$ , where  $t/p$  denotes the subterm of  $t$  at position  $p$ .

To rewrite  $t$  by  $l \rightarrow r$ , we have to find a subterm of  $t$  at position  $p$  and a substitution  $\sigma$  such that  $t/p = l\sigma$ . In this case, we say  $l$  matches  $t/p$  and  $\sigma$  is the *matching*. Matching is one-way unification and finding matching is easier than finding a unifier as the variables in  $t$  cannot be affected by  $\sigma$ . If we treat the variables of  $t$  as constants, then finding a matching is the same as finding a unifier of  $t/p$  and  $l$ .

Once the matching is found, we need to replace the matched subterm in  $t$  by the matched right side. Note that  $t[p \leftarrow s]$  denote the term where the subterm at position  $p$ ,  $t/p$ , (ref. Definition 5.1.4) is replaced by  $s$ . For example, if  $g(x) \rightarrow h(x)$  is a rewrite rule, then term  $t = f(a, g(b))$  can be rewritten to  $t' = f(a, h(b))$ , because  $t/1 = g(b) = g(x)\sigma$  and  $t'/1 = h(x)\sigma$ , where  $\sigma = [x \leftarrow b]$ .

**Definition 7.2.2.** A rewrite system  $R$  is a set of rewrite rules. A term  $t$  is written to  $t'$  by  $R$ , denoted by  $t \Rightarrow_R t'$ , or simply  $t \Rightarrow t'$ , if there exists a rewrite rule in  $R$  which rewrites  $t$  to  $t'$ .

As a convention, we will use  $\Rightarrow^+$  and  $\Rightarrow_*$  to denote the transitive closure and the reflexive and transitive closure of  $\Rightarrow$ , respectively. That is,  $s \Rightarrow_R^* t$  means that  $s$  is rewritten to  $t$  by  $R$  using any number of rewriting.

While our major use of  $R$  is to represent a stable congruence relation,  $R$  can be used for representing any transitive, stable and monotonic relation. This might be inequality (like  $>$  or  $\geq$ ) or implication.

**Example 7.2.3.** In Chapter 1 (section 1.3.4.1), we have seen that the functions can be constructed by equations. For instance, we used equations to define the predecessor, addition, subtraction, and multiplication functions over the natural numbers. We may treat these equations as rewrite rules.

$$\begin{aligned}
 pre(0) &\rightarrow 0; \\
 pre(s(x)) &\rightarrow x. \\
 add(0, y) &\rightarrow y; \\
 add(s(x), y) &\rightarrow s(add(x, y)). \\
 sub(x, 0) &\rightarrow x; \\
 sub(x, s(y)) &\rightarrow sub(pre(x), y). \\
 mul(0, y) &\rightarrow 0; \\
 mul(s(x), y) &\rightarrow add(mul(x, y), y).
 \end{aligned}$$

To compute  $2 \times 2 = 4$ , we rewrite  $mul(s(s(0)), s(s(0)))$  to  $s(s(s(s(0))))$ . We give below also the positions for each rewriting and you may find out which rule is used for each rewriting.

$$\begin{aligned}
 &mul(s(s(0)), s(s(0))) && \text{at } \epsilon \\
 \Rightarrow &add(mul(s(0), s(s(0))), s(s(0))) && \text{at } 1 \\
 \Rightarrow &add(add(mul(0, s(s(0))), s(s(0))), s(s(0))) && \text{at } 1.1 \\
 \Rightarrow &add(add(0, s(s(0))), s(s(0))) && \text{at } 1 \\
 \Rightarrow &add(s(s(0)), s(s(0))) && \text{at } \epsilon \\
 \Rightarrow &s(add(s(0), s(s(0)))) && \text{at } 1 \\
 \Rightarrow &s(s(add(0, s(s(0)))) && \text{at } 1.1 \\
 \Rightarrow &s(s(s(s(0)))) &&
 \end{aligned}$$

In this language, if a term is a ground, it is not hard to prove by induction that the rewriting process will eventually terminate with a specific number, i.e. a term made only from 0 and s. Using this rewrite system, we can do all the addition, subtraction, and multiplication over the natural numbers by rewriting, if we do not care about the speed.

## 7.2.2 Termination of Rewrite Systems

A nice property of many rewrite systems, including the previous example, is that the application of rules to terms, ground or not, cannot go on forever; it will eventually terminate.

**Definition 7.2.4.** A rewrite system  $R$  is said to be terminating if there exists no infinite sequences of terms  $t_0, t_1, \dots, t_i, \dots$  such that  $t_i \Rightarrow t_{i+1}$ .

To prove a rewrite system  $R$  is terminating, we may use a simplification order  $\succ$  (ref. section 6.3.1).

**Proposition 7.2.5.** Let  $\succ$  be a simplification order and  $R$  a rewrite system. If  $l \succ r$  for all every  $l \rightarrow r$  in  $R$ , then  $R$  is terminating.

*Proof.* By definition,  $\succ$  is a well-founded, stable and monotonic order on terms. If  $t_i \Rightarrow t_{i+1}$  by  $l \rightarrow r$ , then  $t_i \succ t_{i+1}$ , because  $l \succ r$ , and  $\succ$  is stable and monotonic. If there exists an infinite sequence,  $t_0, t_1, \dots, t_i, \dots$  such that  $t_i \succ t_{i+1}$ , then  $\succ$  cannot be well-founded.  $\square$

The above proposition provides a sufficient condition for the termination of a rewrite system. In section 6.3.1, we discussed lexicographic path order (LPO), recursive path order (RPO), and Knuth-Bendix order (KBO). These orders are simplification orders and can be used to prove the termination of  $R$ .

**Example 7.2.6.** To show the termination of  $R$  in Example 7.2.3, we may use the lexicographic path order  $\succ_{lpo}$  with the precedence  $s \prec pre \prec add \prec sub \prec mul$  and every binary operator has the left-to-right status, with the exception of  $sub$ :  $sub$  must have the right-to-left status, as we want  $[s(y), x] \succ^{lex} [y, pre(x)]$ .

In general, how to show that the rewriting process using  $R$  will terminate? Unfortunately, this is an undecidable problem. That is, we do not have an algorithm to answer this question for general  $R$ .

**Definition 7.2.7.** Given a rewrite system  $R$ , if  $s \Rightarrow_R^* t$  and there exists no  $t'$  such that  $t \Rightarrow_R t'$ , we say  $t$  is a normal form of  $s$ .

Why is termination an important property of a rewrite system? the termination of  $R$  ensures that there exists a normal form for every term. Without termination, we will have hard time to control rewriting. In computation theory, termination is a divider for telling whether a decision problem is decidable or not. Some equations, such as the commutativity of  $+$  or  $\vee$ , can never be made into a terminating rule. It is undecidable to tell if a rewrite system is terminating or not.

### 7.2.3 Confluence of Rewriting

Termination is one important property of a set of rewrite rules. Other important properties are being confluent and being canonical.

**Example 7.2.8.** In Example 7.1.5, we used three equations to specify a group  $(S, *)$ . These equations can be made as a rewrite system  $R$ :

$$\begin{array}{l} 1 \quad e * x \rightarrow x \\ 2 \quad i(x) * x \rightarrow e \\ 3 \quad (x * y) * z \rightarrow x * (y * z) \end{array}$$

Let  $t = (i(x) * x) * y$ , and there are two ways for  $R$  to rewrite  $t$ :

$$\underline{(i(x) * x)} * y \Rightarrow \left\langle \begin{array}{l} (3) \ i(x) * (x * y) \\ (2) \ \underline{e} * y \Rightarrow (1) \ y \end{array} \right.$$

That is, there are more than one way to rewrite a term. For  $t$ , we have two normal forms:  $i(x) * (x * y)$  and  $y$ .

**Definition 7.2.9.** A rewrite system  $R$  is said to be confluent if, for any term  $s$ , if  $s \Rightarrow_R^* t_1$  and  $s \Rightarrow_R^* t_2$ , then there exists a term  $t'$  such that  $t_1 \Rightarrow_R^* t'$  and  $t_2 \Rightarrow_R^* t'$ .  $R$  is said to be canonical if  $R$  is both terminating and confluent.

That is, if  $R$  is confluent, it does not matter how you go about applying the rules to a term, you will get the same result. There are several concepts of “confluence” proposed in the literature. They are equivalent for terminating rewrite systems. The termination of  $R$  ensures the existence of a normal form for every term; the confluence of  $R$  ensures the uniqueness of normal forms. When  $R$  is canonical, you do not need to worry about at which positions and choose which rule to rewrite a term. Any choices can be made and the end result will be the same. Let us denote the unique normal form of  $t$  in a canonical rewrite system by  $t \downarrow_R$ , the canonical form of  $t$ .

**Definition 7.2.10.** A canonical rewrite system  $R$  is a decision procedure for a set  $E$  of equations if for any two terms  $s$  and  $t$ ,  $s =_E t$  iff  $s \downarrow_R = t \downarrow_R$ .

We say  $R$  is a decision procedure for  $E$  because the computation of  $t \downarrow_R$  is terminating.

**Theorem 7.2.11.** Let  $R$  be a canonical rewrite system satisfying the following conditions:

1. For each rewrite rule  $l \rightarrow r$  of  $R$ ,  $l =_E r$ ;
2. For each equation  $s_1 = s_2$  of  $E$ ,  $s_1 \downarrow_R = s_2 \downarrow_R$ .

Then  $R$  is a decision procedure for  $E$ .

*Proof.* For each rule  $l \rightarrow r$  of  $R$ , since  $l =_E r$ , if  $s \Rightarrow t$  by  $l \rightarrow r$ , then  $s =_E t$  by applying the instantiation and monotonicity rules on  $l =_E r$ . Applying this result to each rewriting, we have  $s =_E s \downarrow_R$  and  $t =_E t \downarrow_R$ . If  $s \downarrow_R = t \downarrow_R$ , then  $s =_E t$ .

If  $s =_E t$ , check all the six conditions for obtaining  $s =_E t$ : If  $s =_E t$  because  $s = t \in E$ , then  $s \downarrow_R = t \downarrow_R$ . If  $s =_E t$  because  $t =_E s$ , then  $t \downarrow_R = s \downarrow_R$  implies  $s \downarrow_R = t \downarrow_R$ . If  $s =_E t$  because  $s =_E r$  and  $r =_E t$ , then  $s \downarrow_R = r \downarrow_R = t \downarrow_R$ . If  $s = f(\dots, s_i, \dots) =_E f(\dots, t_i, \dots) = t$  because  $s_i =_E t_i$ , then  $s_i \downarrow_R = t_i \downarrow_R$  implies  $s \downarrow_R = t \downarrow_R$ . If  $s = s'\sigma =_E t'\sigma = t$  because  $s' =_E t'$ , then  $s' \downarrow_R = t' \downarrow_R$  implies  $s \downarrow_R = t \downarrow_R$ . In all cases, we have  $s \downarrow_R = t \downarrow_R$  if  $s =_E t$ .  $\square$

Even when a rewrite system is canonical we may still want to exercise carefully our choices of rules and positions in a term to apply rules, in order to minimize the amount of computational effort required to find the canonical form. We may want to explore only the shortest rewriting steps. For instance, choosing a position for rewriting, we may use an outermost-position-first strategy or an innermost-position-first strategy. A careful inspection of the rules to be applied may reveal that which strategy works better.

#### 7.2.4 The Knuth-Bendix Completion Procedure

In Example 7.2.8, we have seen that the term  $t = (i(x) * x) * y$  has two normal forms, i.e.,  $i(x) * (x * y)$  and  $y$ . The rewrite system in the example is not confluent. Knuth and Bendix suggested a procedure which takes a set of equations and a simplification order, and each equation is oriented into rewrite rules according to the simplification order. So-called *critical pairs* are then computed from the rewrite rules. If all the critical pairs have the same normal form, we obtain a canonical rewrite system. If not, such critical pairs are added into the set of equations and continue.

**Definition 7.2.12.** *Given two rewrite rules which share no variables and may differ only by variable names, say  $l_1 \rightarrow r_1$  and  $l_2 \rightarrow r_2$ , a critical pair of the two rules is a pair of terms  $s = t$ , where  $s = r_1\sigma$ ,  $t = l_1[p \leftarrow r_2]\sigma$ ,  $p$  is a non-variable position of  $l_1$ , and  $\sigma$  is the mgu of  $l_1/p$  and  $l_2$ . We say  $l_2$  is superposed at  $p$  of  $l_1$  to produce the critical pair and this process is called superposition.*

In Example 7.2.8, we are given three rewrite rules. If we rename the variables in rule 3 as  $(x_3 * y_3) * z_3 \rightarrow x_3 * (y_3 * z_3)$ , we can compute a critical pair from rule 3 and rule 2, by superposing  $l_2 = i(x) * x$  at position 1 of  $l_3$ , that is, unifying  $l_2$  and  $l_3/1 = x_3 * y_3$  with the mgu  $\sigma = [x_3 \leftarrow i(x), y_3 \leftarrow x, z_3 \leftarrow y]$  and generate the critical pair  $(i(x) * (x * y), e * y)$  from  $l_3\sigma = (i(x) * x) * y$  as illustrated in Example 7.2.8.

During the computation of a critical pair, applying the mgu  $\sigma$  to  $l_1$ ,  $l_1\sigma$  can be rewritten to  $r_1\sigma$  by  $l_1 \rightarrow r_1$ . At the same time, since  $l_1/p\sigma = l_2\sigma$ ,  $l_1\sigma$  contains an instance of  $l_2$  at position  $p$ , so  $l_1\sigma$  can be also rewritten to  $l_1[p \leftarrow l_2]\sigma$  by  $l_2 \rightarrow r_2$ . In other words, the critical pair is obtained by two rewritings on  $l_1\sigma$  using the two rules, respectively:

$$l_1\sigma \Rightarrow \left\langle \begin{array}{ll} r_1\sigma & \text{by } l_1 \rightarrow r_1 \\ l_1[p \leftarrow r_2]\sigma & \text{by } l_2 \rightarrow r_2 \end{array} \right.$$

**Lemma 7.2.13.** *If  $l_1 =_E r_1$  and  $l_2 =_E r_2$ , and  $(s, t)$  is a critical pair from them, then  $s =_E t$ .*

*Proof.* From  $l_1 =_E r_1$ , we have  $l_1\sigma =_E r_1\sigma$ . From  $l_2 =_E r_2$  and  $l_1/p\sigma = l_2\sigma$ , we have  $l_1\sigma =_E l_1[p \leftarrow r_2]\sigma$ . Hence  $s = r_1\sigma =_E l_1[p \leftarrow r_2]\sigma = t$ .  $\square$

**Example 7.2.14.** We will show to obtain a canonical rewrite system from the three rewrite rules in Example 7.2.8:

$$\begin{array}{l} 1 \quad e * x \rightarrow x \\ 2 \quad i(x) * x \rightarrow e \\ 3 \quad (x * y) * z \rightarrow x * (y * z) \end{array}$$

The lexicographical path order is used with the precedence  $e \prec * \prec i$ . When computing a critical pair of the two rules, we assume that the variables in one of the rules are implicitly subscribed by the rule number. For example, the third rule is  $(x_3 * y_3) * z_3 \rightarrow x_3 * (y_3 * z_3)$ .

The first critical pair is  $i(x) * (x * y) = e * y$ , which, after rewriting  $e * y$  to  $y$ , will be oriented into the fourth rewrite rule:

$$4 \quad i(x) * (x * y) \rightarrow y$$

A critical pair comes from rule 4 and itself, and the mgu of  $l_4/2$  and  $l_4$  is  $\sigma = [x_4 \leftarrow i(x), y_4 \leftarrow x * y]$ :

$$i(i(x)) * \underline{i(x) * (x * y)} \Rightarrow \left\langle \begin{array}{l} (4) \quad x * y \\ (4) \quad i(i(x)) * \underline{y} \end{array} \right.$$

The critical pair is  $x * y = i(i(x)) * y$ , which produces the fifth rule:

$$5 \quad i(i(x)) * y \rightarrow x * y$$

A critical pair comes from rule 4 and rule 2, and the mgu of  $l_4/2$  and  $l_2$  is  $\sigma = [x_4 \leftarrow i(x), y_4 \leftarrow x]$ :

$$i(i(x)) * \underline{(i(x) * x)} \Rightarrow \left\langle \begin{array}{l} (4) \ x \\ (2) \ i(i(x)) * \underline{e} \rightarrow (5) \ x * e \end{array} \right.$$

The critical pair is  $x = x * e$ , which produces the sixth rule:

$$6 \quad x * e \rightarrow x$$

A critical pair comes from rule 6 and rule 2 with  $\sigma = [x_6 \leftarrow i(e)]$  and generates the seventh rule:

$$7 \quad i(e) \rightarrow e$$

A critical pair comes from rule 5 and rule 2 with  $\sigma = [x_2 \leftarrow i(x), y \leftarrow i(x)]$ :

$$i(i(x)) * i(x) \Rightarrow \left\langle \begin{array}{l} (5) \ x * i(x) \\ (2) \ e \end{array} \right.$$

and generates the eighth rule:

$$8 \quad x * i(x) \rightarrow e$$

A critical pair comes from rule 6 and rule 5 with  $\sigma = [x_6 \leftarrow i(i(x)), y \leftarrow e]$ :

$$i(i(x)) * e \Rightarrow \left\langle \begin{array}{l} (6) \ i(i(x)) \\ (5) \ x * e \rightarrow (6) \ x \end{array} \right.$$

and generates the ninth rule:

$$9 \quad i(i(x)) \rightarrow x$$

A critical pair comes from rule 4 and rule 3, and the mgu of  $l_4/2$  and  $l_3$  is  $\sigma = [x_4 \leftarrow x * y, y_4 \leftarrow z]$ :

$$i(x * y) * \underline{((x * y) * z)} \Rightarrow \left\langle \begin{array}{l} (4) \ z \\ (3) \ i(x * y) * (x * (y * z)) \end{array} \right.$$

and generates the tenth rule:

$$10 \quad i(x * y) * (x * (y * z)) \rightarrow z$$

A critical pair comes from rule 10 and rule 8, and the mgu of  $l_{10}/2.2$  and  $l_9$  is  $\sigma = [x_9 \leftarrow y, z \leftarrow i(y)]$ :

$$i(x * y) * (x * \underline{(y * i(y))}) \Rightarrow \left\langle \begin{array}{l} (10) \ i(y) \\ (8) \ i(x * y) * (x * \underline{e}) \Rightarrow (6) \ i(x * y) * x \end{array} \right.$$

and generates the eleventh rule:

$$11 \quad i(x * y) * x \rightarrow i(y)$$

A critical pair comes from rule 11 and rule 4, and the mgu of  $l_{11}/1.1$  and  $l_4$  is  $\sigma = [x_{11} \leftarrow i(x), y_{11} \leftarrow x * y]$ :

$$i(\underline{i(x) * (x * y)}) * i(x) \Rightarrow \left\langle \begin{array}{l} (11) \quad i(x * y) \\ (4) \quad i(\underline{y}) * i(x) \end{array} \right.$$

and generates the twelfth rule:

$$12 \quad i(x * y) \rightarrow i(y) * i(x)$$

The final critical pair comes from rule 4 and rule 5, and generates the thirteenth rule:

$$13 \quad x * (i(x) * y) \rightarrow y$$

Since rules 5, 10, and 11 can be rewritten to the identity i.e.,  $t = t$ , by rules 5 and 12, the final set of rewrite rules are

$$\begin{array}{ll} 1 & e * x \rightarrow x \\ 2 & i(x) * x \rightarrow e \\ 3 & (x * y) * z \rightarrow x * (y * z) \\ 4 & i(x) * (x * y) \rightarrow y \\ 6 & x * e \rightarrow x \\ 7 & i(e) \rightarrow e \\ 8 & x * i(x) \rightarrow e \\ 9 & i(i(x)) \rightarrow x \\ 12 & i(x * y) \rightarrow i(y) * i(x) \\ 13 & x * (i(x) * y) \rightarrow y \end{array}$$

All the critical pairs from this set of 10 rules can be reduced to the identity. These 10 rules are sufficient to solve the *word problem* for a free group with no other properties. The *word problem* is to find a decision procedure to decide whether any two terms are equal, according to the axioms of group theory. The above rewrite system can serve as a decision procedure by checking if the two terms have the same canonical form.

The above example illustrates the execution of the Knuth-Bendix completion procedure, which is described by the following pseudo-code.

**Procedure 7.2.15.** The procedure  $KnuthBendix(E, \succ)$  takes a set  $E$  of equations and a simplification order  $\succ$ , and generates a canonical rewrite system from  $E$  when it succeeds. It calls four procedures:

- $pickEquation(E)$  will pick an equation from  $E$ ;
- $NF(t, R)$  returns a normal form of term  $t$  by  $R$ .
- $rewritable(t, r)$  will check if rewrite rule  $r$  can rewrite term  $t$ ;
- $criticalPairs(r, R)$  will return the set of all critical pairs between rule  $r$  and every rule in  $R$ .

```

proc KnuthBendix( $E, \succ$ )
1    $R := \emptyset$ ;
2   while ( $E \neq \emptyset$ ) do
3        $(s, t) := pickEquation(E)$ ;  $E := E - \{(s, t)\}$ ;
4        $s := NF(s, R)$ ;  $t := NF(t, R)$ ; // normalize  $s$  and  $t$  by  $R$ 
5       if ( $s = t$ ) continue; // identity is discarded
6       else if ( $s \succ t$ )  $r := (s, t)$ ; //  $r$  is " $s \rightarrow t$ "
7       else if ( $t \succ s$ )  $r := (t, s)$ ; //  $r$  is " $t \rightarrow s$ "
8       else return "failure"; //  $s$  and  $t$  are not comparable by  $\succ$ 
9        $R := R \cup \{r\}$ ; // add new rule into  $R$ 
10      for  $(s, t) \in R$  do
11          if  $rewritable(s, r)$   $R := R - \{(s, t)\}$ ;  $E := E \cup \{(s, t)\}$ ;
12          if  $rewritable(t, r)$   $R := R - \{(s, t)\} \cup \{(s, NF(t, R))\}$ ;
13           $E := E \cup criticalPairs(r, R)$ ;
14  return  $R$ ;

```

The procedure is not an algorithm because it may go on forever due to ever-growing set of rules and critical pairs. Note that the procedure may fail if an equation cannot be oriented into a rewrite rule. There are many attempts to reduce the cases of failures. For instance, if an equation like  $f(x, g(x)) = h(y, g(y))$  cannot be oriented, we may introduce a new constant  $c$  to make two equations  $f(x, g(x)) = c$  and  $h(y, g(y)) = c$ . In general, if we have an equation  $s(x, z) = t(y, z)$ , where  $x$  does not appear in  $t$  and  $y$  does not appear in  $s$ , we may introduce a new function  $h(z)$  to create two new equations:  $s(x, z) = h(z)$  and  $t(y, z) = h(z)$ . For the axioms  $AC = \{(x * y) * z = x * (y * z), x * y = y * x\}$ , we may consider the rewriting over  $T(F, X)/=_{AC}$ , the so-called rewriting modulo  $AC$  (associativity and commutativity).

When *KnuthBendix* terminates successfully, we obtain a canonical rewrite system.

**Theorem 7.2.16.** *If  $KnuthBendix(E, \succ)$  returns  $R$ , then  $R$  is canonical.*

*Proof.* At first,  $R$  is terminating, because for every  $l \rightarrow r$  of  $R$ ,  $l \succ r$ , where  $\succ$  is a simplification order. Secondly, all the critical pairs of  $R$  are computed and each pair is rewritten to the same term by  $R$  (if not, a new rule will be added into  $R$ ).

To show that  $R$  is confluent, we do induction on  $t$  based on  $\succ$ . As a base case, if  $t$  is not rewritable by  $R$ , then  $R$  is confluent trivially at  $t$ . Suppose, for  $i = 1, 2$ ,  $t \Rightarrow t_i$  by  $l_i \rightarrow r_i \in R$  at position  $p_i$  of  $t$ , i.e.,  $t_i = t[p \leftarrow r_i \sigma_i]$ .

Case 1: If  $p_1$  and  $p_2$  are disjoint positions, then both  $t_1$  and  $t_2$  can be written to a common term  $t' = t[p_1 \leftarrow r_1 \sigma_1, p_2 \leftarrow r_2 \sigma_2]$ .

Case 2: If  $p_1$  and  $p_2$  are not disjoint positions, without loss of generality, assume  $p_1 = \epsilon$ , then  $t = l_1 \sigma_1$ , and  $t/p_2 = l_1 \sigma_1 / p_2 = l_2 \sigma_2$ . There are two subcases to consider:

Case 2.1: If  $p_2$  is a non-variable position of  $l_1$ , then  $l_1 \sigma_1 / p_2 = l_1 / p_2 \sigma_1 = l_2 \sigma_2$ , that is,  $l_1 / p_2$  and  $l_2$  are unifiable. Thus, there exists a critical pair  $(s_1, s_2)$  between the two rules such that  $t_1 = s_1 \theta$  and  $t_2 = s_2 \theta$  for some substitution  $\theta$ . Since  $s_1$  and  $s_2$  are rewritten to the same term by  $R$ , so are  $t_1$  and  $t_2$ .

Case 2.2: When  $p_2$  is not a non-variable position of  $l_1$ , finding a common successor of  $t_1$  and  $t_2$  is left as an exercise. Note that in this case there must exist a variable position  $p$  of  $l_1$  such that  $p$  is a prefix of  $p_2$ .

By induction hypotheses,  $R$  is confluent at  $t_1$  and  $t_2$ , respectively, because  $t_1 \succ t$  and  $t_2 \succ t$ . Now we have shown that  $t_1$  and  $t_2$  are rewritten to a common term in all cases,  $R$  must be confluent at  $t$ .  $\square$

**Example 7.2.17.** To illustrate case 2.2 in the above proof, let  $t = i(i(i(z))) * i(i(z))$  and  $R$  be the ten rules in Example 7.2.14. Then  $t$  can be rewritten to  $e$  by rule 2 and to  $i(z) * i(i(z))$  by rule 9 at position 1.1.  $i(z) * i(i(z))$  can be rewritten to  $i(z) * z$  by rule 9 at position 2, and to  $e$  by rule 2. Note that both 1.1 and 2 are variable positions of  $t$ .

**Theorem 7.2.18.** *If  $\text{KnuthBendix}(E, \succ)$  returns  $R$ , then  $R$  is a decision procedure for  $E$ .*

*Proof.* (sketch) The previous theorem tells us that  $R$  is canonical. In light of Theorem 7.2.11, the condition 1 is ensured by Lemma 7.2.13; the condition 2 is true because every equation is either made into a rule or simplified to identity.  $\square$

In conclusion, rewrite systems can be used to put terms in normal form and to prove terms equal, equivalent or some transitive monotonic relations. Their application represents a powerful method of mathematical reasoning, because they can sometimes overcome the combinatorial explosions caused by other proof procedures, e.g. resolution. The desirable properties of rewrite systems are: termination and

confluence. Rewrite systems can be shown to terminate with the aid of simplification orders. Rewrite systems can be shown to be confluent by looking at critical pairs formed from pairs of rules.

### 7.3 Inductive Theorem Proving

In Example 7.2.3, we converted a set  $E$  of equations into rewrite rules, where  $add$  is defined using two equations:  $E = \{add(0, y) = y, add(s(x), y) = s(add(x, y))\}$ . It is easy to prove by induction on  $i$  that  $add(s^i(0), s^j(0)) = s^{i+j}(0) = add(s^j(0), s^i(0))$ . Since  $add(s^i(0), s^j(0)) = add(s^j(0), s^i(0))$ , it is natural to state that  $add(x, y) = add(y, x)$  for all ground terms  $x$  and  $y$ . Obviously,  $add(x, y) = add(y, x)$  is not a logical consequence of  $E$  and  $E_{ax}$  (the equality axioms), because  $E \cup E_{ax}$  may have models in which  $add(x, y) = add(y, x)$  does hold. It is known that some properties of a first-order formula cannot be proved by deduction or contradiction; they can be proved only by induction. That is why we call these properties *inductive theorems*,  $add(x, y) = add(y, x)$  is just one of them.

Mathematical induction is a fundamental and powerful rule of inference. Inductive theorem proving has important applications in formal verification and has been used to prove the correctness of computer architecture or algorithms. Like the previous chapters, we are particularly interested in inductive proving methods which can be automated in a computer.

#### 7.3.1 Inductive Theorems

**Definition 7.3.1.** *Given a first-order language  $L = (P, F, X, Op)$ , a quantifier-free formula  $B$  is said to be an inductive theorem of  $A$ , if for all substitution  $\sigma : X \rightarrow T(F)$ ,  $A \models B\sigma$ .*

In other words, a quantifier-free formula is an inductive theorem of  $A$  if all its ground instances are logical consequence of  $A$ . We can show that  $add(x, y) = add(y, x)$  and  $add(add(x, y), z) = add(x, add(y, z))$  are inductive theorems of  $E = \{add(0, y) = y, add(s(x), y) = s(add(x, y))\}$ . Typically,  $A$  is also quantifier-free, with the understanding that all free variables act as universally quantified variables.

To prove an inductive theorem, we need an induction inference rule, which has many different versions.

**Example 7.3.2.** In Example 7.2.3, we convert a set  $E$  of equations into a rewrite

system so that the computation can be carried out by rewriting.

$$\begin{array}{ll}
 1 & pre(0) \rightarrow 0; \\
 2 & pre(s(x)) \rightarrow x. \\
 3 & add(0, y) \rightarrow y; \\
 4 & add(s(x), y) \rightarrow s(add(x, y)). \\
 5 & sub(x, 0) \rightarrow x; \\
 6 & sub(x, s(y)) \rightarrow sub(pre(x), y). \\
 7 & mul(0, y) \rightarrow 0; \\
 8 & mul(s(x), y) \rightarrow add(mul(x, y), y).
 \end{array}$$

In the above example, the function symbols are  $F = \{0, s, pre, add, sub, mul\}$ , which can be divided into two parts:  $C = \{0, s\}$  and  $D = \{pre, add, sub, mul\}$ . For each function in  $D$ , we have a complete definition over 0 and  $s$ :

- $pre(0) = 0$  and  $pre(s^i(0)) = s^{i-1}(0)$  for  $i > 0$ ;
- $add(s^i(0), s^j(0)) = s^{i+j}(0)$  for any  $i, j \geq 0$ ;
- $sub(s^i(0), s^j(0)) = 0$  if  $i \leq j$  and  $pre(s^i(0), s^j(0)) = s^{i-j}(0)$  for  $i > j$ ;
- $mul(s^i(0), s^j(0)) = s^{i*j}(0)$  for any  $i, j \geq 0$ ;

It is easy to check that  $T(F)/=E = \{[0], [s(0)], \dots, [s^i(0)], \dots\}$ , which is the domain of the Herbrand model with equality. We say  $T(F)/=E$  and  $T(C)$  are *isomorphic* because there is a one-to-one mapping between  $T(F)/E$  and  $T(C)$ .  $C$  is called *constructors* and  $D$  is called *defined functions* (or *interpreted functions*). Inductive theorems are about properties of defined functions.

### 7.3.2 Structural Induction

One basic induction rule is called the *structural induction rule*, which is based on the structure of terms built up by the constructors.

#### Structural Induction Rule

$$\frac{B(0) \quad B(y) \rightarrow B(s(y))}{B(x)}$$

where  $B(x)$  stands for an inductive theorem to be inducted on  $x$ .

The above rule states that to prove an inductive theorem  $B(x)$ , we need to prove two formulas  $B(0)$  and  $B(x) \rightarrow B(s(x))$ . Its soundness is based on the condition that the isomorphic relation of  $T(F)/=E$  and  $T(C)$ .

**Example 7.3.3.** Let  $B(x)$  be  $add(x, 0) = x$ . Applying the structural induction on this one, we need to prove (1)  $add(0, 0) = 0$  and (2)  $add(x, 0) = x \rightarrow add(s(x), 0) = s(x)$ . (1) is trivial; (2) is simplified by rule 4, the definition of  $add$ , to:  $add(x, 0) = x \rightarrow s(add(x, 0)) = s(x)$ , or

$$(2') \quad (add(x, 0) \neq x \mid s(add(x, 0)) = s(x)),$$

which is true by *equality crossing*, as described below.

### Equality Crossing

$$\frac{(s \neq t \mid A(s))}{(s \neq t \mid A(t))}$$

For instance, let  $s \neq t$  be  $add(x, 0) \neq x$  in (2') of the above example, then  $s(add(x, 0)) = s(x)$  is simplified to  $s(x) = s(x)$ . In other words, we used the induction hypothesis  $add(x, 0) = x$  to show that  $s(add(x, 0)) = s(x)$  is true.

Equality crossing, called *crossing fertilization* by Boyer and Moore, is a simplification rule in the presence of equality axioms, due to the following result.

**Proposition 7.3.4. (soundness of equality crossing)**  $E_{ax} \models (s \neq t \mid A(s)) \leftrightarrow (s \neq t \rightarrow A(t))$ , where  $E_{ax}$  is the equality axioms.

*Proof.* From the monotonicity of equality, i.e.,  $x_i = y_i \wedge p(\dots, x_i, \dots) \rightarrow p(\dots, y_i, \dots)$ , we can deduce  $(s = t) \wedge A(s) \rightarrow A(t)$ , whose clausal form is

$$(s \neq t \mid \overline{A(s)} \mid A(t))$$

The resolvent of this clause with  $(s \neq t \mid A(s))$  is  $(s \neq t \mid A(t))$ . Switching  $s$  and  $t$ , the other implication also holds.  $\square$

The fundamental difference between using an inference rule like resolution and using an induction rule is not that a well-founded order, e.g.  $s(x) \succ x$ , is needed when applying induction. When using resolution for theorem proving by refutation, the empty clause is the goal and the resolution rule is used to derive this goal from the input clauses. When using an induction rule for theorem proving, the goal is the formula to be proved. However, we never try (and it is almost impossible) to derive this goal from the input formulas by the induction rule. Instead, we try to generate subgoals from the goal by using the induction rule backward. That is, to prove  $B(x)$  by induction, we generate two subgoals  $B(0)$  and  $B(y) \rightarrow B(s(y))$  from  $B(x)$  and try to establish the validity of these subgoals. Remember that when

$T(F)/=E$  is isomorphic to  $T(C)$ , every defined function is completely defined on the constructors and these definition equations ensure the soundness of the induction rule from these equations.

**Example 7.3.5.** Continuing from Example 7.3.2, let us prove that  $sub(add(x, y), y) = x$ . If we do induction on  $x$ , the proof would not go through. On the other hand, we can do induction on  $y$ . By inspecting the definitions of  $add$  and  $sub$ , we see that  $x$  is the induction position for  $add$  and  $y$  is the induction position for  $sub$ .

The base case is  $sub(add(x, 0), 0) = x$ , which is simplified to true by rules 3 and 5. The inductive case is

$$sub(add(x, y), y) = x \rightarrow sub(add(x, s(y)), s(y)) = x$$

To simplify  $add(x, s(y))$ , we need  $add(x, s(y)) = s(add(x, y))$ , which is also an inductive theorem and its proof is left as an exercise. Suppose this theorem is available to us, then

$$\begin{aligned} & sub(add(x, s(y)), s(y)) \\ = & sub(s(add(x, y)), s(y)) && \text{by } add(x, s(y)) = s(add(x, y)), \\ = & sub(pre(s(add(x, y)), y) && \text{by rule 6,} \\ = & sub(add(x, y), y) && \text{by rule 2,} \\ = & x && \text{by equality crossing.} \end{aligned}$$

### 7.3.3 Induction on Two Variables

The previous example shows that we need to choose the right variable to do the induction. Sometimes, we need to do induction on multiple variables.

**Example 7.3.6.** Adding the following two predicates to Example 7.3.2:

$$\begin{aligned} 11 & \quad y < 0 \rightarrow \perp; \\ 12 & \quad 0 < s(x) \rightarrow \top; \\ 13 & \quad s(y) < s(x) \rightarrow y < x. \\ 14 & \quad 0 \leq y \rightarrow \top; \\ 15 & \quad s(x) \leq 0 \rightarrow \perp; \\ 16 & \quad s(x) \leq s(y) \rightarrow x \leq y. \end{aligned}$$

To prove  $(x \leq y) = \neg(y < x)$  by structural induction, neither induct on  $x$  alone nor on  $y$  alone will work. We need an induction on two variables.

### Structural Induction Rule on Two Variables

$$\frac{B(0, y) \quad B(s(x), 0) \quad B(x, y) \rightarrow B(s(x), s(y))}{B(x, y)}$$

where  $B(x, y)$  stands for an inductive theorem to be inducted on both  $x$  and  $y$ .

This is not the only induction rule for the two variables. Why should we use this rule for  $(x \leq y) = \neg(y < x)$ ? The clue is in the definitions of  $\leq$  and  $<$ : For  $(x \leq y)$ , we defined  $\leq$  on  $\langle 0, y \rangle, \langle s(x), 0 \rangle$  and  $\langle s(x), s(y) \rangle$  with a recursive call on  $\langle x, y \rangle$ . For  $(y < x)$ , we also defined  $<$  on  $\langle 0, y \rangle, \langle s(x), 0 \rangle, \langle s(x), s(y) \rangle$  with a recursive call on  $\langle x, y \rangle$  (the positions of  $x$  and  $y$  are switched here). In other words, the definition equations of these two predicates tell us what cases needed to consider: the base cases come from the definition equations without recursion; the inductive cases come from the definition equations with recursion. This is the major induction heuristic in Boyer and Moore's computational logic, where lisp functions are used instead of equational definitions.

Now back on the proof of  $(x \leq y) = \neg(y < x)$ , by the above induction rule, we need to prove three subgoals:

1  $(0 \leq y) = \neg(y < 0)$ , which can be easily proved by showing  $0 \leq y = \top$  by rule 14 and  $\neg(y < 0) = \neg(\perp) = \top$  by rule 11.

2  $(s(x) \leq 0) = \neg(0 < s(x))$ , which can be proved by showing  $(s(x) \leq 0) = \perp$  by rule 15 and  $\neg(0 < s(x)) = \neg(\top) = \perp$  by rules 12.

3  $(x \leq y) = \neg(y < x) \rightarrow (s(x) \leq s(y)) = \neg(s(y) < s(x))$ . Proof of 3 is also easy:  $s(x) \leq s(y)$  is simplified to  $x \leq y$  by rule 16 and  $\neg(s(y) < s(x))$  to  $\neg(y < x)$  by rule 13, then 3 becomes a trivial implication:  $(x \leq y) = \neg(y < x) \rightarrow (x \leq y) = \neg(y < x)$ .

#### 7.3.4 Multi-Sort Algebraic Specifications

By far, we have focused on the natural number functions based on the constructors  $0$  and  $s$ . In many applications, we need other data structures, such as lists, trees, etc., which have their own constructors. Multi-sort algebraic specification will provide a convenient tool for such applications. For example, to work on the lists of natural numbers, we need two types: natural numbers and lists, and their domains can be specified as follows:

$$\begin{aligned} \langle nat \rangle & ::= 0 \mid s(\langle nat \rangle) \\ \langle list \rangle & ::= nil \mid cons(\langle nat \rangle, \langle list \rangle) \end{aligned}$$

where the constructors for  $nat$  are  $0 : nat$  and  $s : nat \rightarrow nat$  and the constructors for  $list$  are  $nil : list$  and  $cons : nat, list \rightarrow list$ . We have seen the functions on

*nat*; the functions over *list* can be also defined using equations.

**Example 7.3.7.** The *append* and *reverse* operations of lists can be defined as follows:

$$\begin{array}{ll}
 & app : list, list \longrightarrow list \\
 1 & app(nil, y) \rightarrow y; \\
 2 & app(cons(x, y), z) \rightarrow cons(x, app(y, z)). \\
 & rev : list \longrightarrow list \\
 3 & rev(nil) \rightarrow nil; \\
 4 & rev(cons(x, y)) \rightarrow app(rev(y), cons(x, nil)).
 \end{array}$$

To prove  $rev(app(y, z)) = app(rev(z), rev(y))$  for all lists  $y$  and  $z$ , we need the following rule:

### Structural Induction Rule for List

$$\frac{B(nil) \quad B(y) \rightarrow B(cons(x, y))}{B(y)}$$

where  $B(y)$  stands for an inductive theorem to be inducted on  $y$ .

The above rule reduces  $rev(app(y, z)) = app(rev(z), rev(y))$  to two subgoals:

$$\begin{array}{l}
 G_1 \quad rev(app(nil, z)) = app(rev(z), rev(nil)) \\
 G_2 \quad rev(app(y, z)) = app(rev(z), rev(y)) \rightarrow \\
 \quad \quad \quad rev(app(cons(x, y), z)) = app(rev(z), rev(cons(x, y)))
 \end{array}$$

We can simplify  $G_1$  by rules 1 and 3 to  $rev(z) = app(rev(z), nil)$ , which suggests a more general property:  $app(y, nil) = y$ . This heuristic is called *generalization* and can be implemented automatically. Indeed,  $app(y, nil) = y$  is an inductive theorem and its proof is left as an exercise. Using this lemma,  $rev(z) = app(rev(z), nil)$  becomes true.

For  $G_2$ , we can simplify

$$rev(app(cons(x, y), z)) = app(rev(z), rev(cons(x, y)))$$

by rule 2 to  $rev(cons(x, app(y, z))) = app(rev(z), rev(cons(x, y)))$ , then to

$$app(rev(app(y, z)), cons(x, nil)) = app(rev(z), app(rev(y), cons(x, nil)))$$

by rule 4. Using equality crossing,  $rev(app(y, z))$  is rewritten to  $app(rev(z), rev(y))$ , and we are facing

$$app(app(rev(y), rev(z)), cons(x, nil)) = app(rev(z), app(rev(y), cons(x, nil)))$$

A generalization of the above formula is the associativity of *app*:

$$\text{app}(\text{app}(x, y), z) = \text{app}(x, \text{app}(y, z))$$

which can be shown to be an inductive theorem separately. Once this is done, the proof of  $G_2$  is complete.

The above example illustrates an interesting feature of inductive theorem proving: In order to prove  $\text{rev}(\text{app}(y, z)) = \text{app}(\text{rev}(z), \text{rev}(y))$ , we need to prove extra lemmas:  $\text{app}(y, \text{nil}) = y$  and  $\text{app}(\text{app}(x, y), z) = \text{app}(x, \text{app}(y, z))$ . These lemmas can sometimes be generated automatically, and sometimes need to be provided by the user. For rigorousness, they are needed to be proved before they are used in a proof.

Let us prove another theorem  $\text{rev}(\text{rev}(y)) = y$  by the same induction rule. The induction rule reduces  $\text{rev}(\text{rev}(y)) = y$  to two subgoals:  $G_3 : \text{rev}(\text{rev}(\text{nil})) = \text{nil}$ , which is easy to prove to be true by rule 2, and

$$G_4 : \text{rev}(\text{rev}(y)) = y \rightarrow \text{rev}(\text{rev}(\text{cons}(x, y))) = \text{cons}(x, y).$$

Applying rule 4, we rewrite  $\text{rev}(\text{rev}(\text{cons}(x, y)))$  to  $\text{rev}(\text{app}(\text{rev}(y), \text{cons}(x, \text{nil})))$ , which is then simplified to  $\text{app}(\text{rev}(\text{cons}(x, \text{nil})), \text{rev}(\text{rev}(y)))$ . Since  $\text{rev}(\text{cons}(x, \text{nil}))$  is rewritten to  $\text{cons}(x, \text{nil})$  by rules 4, 3 and 1, so  $\text{app}(\text{rev}(\text{cons}(x, \text{nil})), \text{rev}(\text{rev}(y)))$  is rewritten  $\text{cons}(x, \text{rev}(\text{rev}(y)))$  by rules 2 and 1. Thus  $G_4$  becomes

$$G_4 : \text{rev}(\text{rev}(y)) = y \rightarrow \text{cons}(x, \text{rev}(\text{rev}(y))) = \text{cons}(x, y)$$

Applying equality crossing,  $\text{rev}(\text{rev}(y))$  is replaced by  $y$  and  $G_4$  becomes true.

We have shown how use the language of algebraic specification of abstract data types to write axioms (definitions of functions and predicates), theorems or lemmas (intermediate inferences), etc. for inductive reasoning. In particular, we require that the axioms as well as theorems in an algebraic specification be expressed as a set of (conditional) equations. In other words, we are interested in automatic methods which can prove equations from equations by mathematical induction.

**Example 7.3.8.** We may define the insertion sort algorithm using equations.

$$\begin{array}{ll}
& \textit{isort} : \textit{list} \longrightarrow \textit{list} \\
5 & \textit{isort}(\textit{nil}) \rightarrow \textit{nil}; \\
6 & \textit{isort}(\textit{cons}(x, y)) \rightarrow \textit{insert}(x, \textit{isort}(y)). \\
& \textit{insert} : \textit{nat}, \textit{list} \longrightarrow \textit{list} \\
7 & \textit{insert}(x, \textit{nil}) \rightarrow \textit{cons}(x, \textit{nil}); \\
8 & \textit{insert}(x, \textit{cons}(y, z)) \rightarrow \textit{cons}(x, \textit{cons}(y, z)) \textbf{ if } x \leq y; \\
9 & \textit{insert}(x, \textit{cons}(y, z)) \rightarrow \textit{cons}(y, \textit{insert}(x, z)) \textbf{ if } \neg(x \leq y); \\
& \textit{sorted} : \textit{list} \longrightarrow \textit{boolean} \\
10 & \textit{sorted}(\textit{nil}) \rightarrow \top; \\
11 & \textit{sorted}(\textit{cons}(x, \textit{nil})) \rightarrow \top; \\
12 & \textit{sorted}(\textit{cons}(x, \textit{cons}(y, z))) \rightarrow (x \leq y) \wedge \textit{sorted}(y, z);
\end{array}$$

where  $\leq$  is defined in Example 7.3.6.

We may prove that  $\textit{sorted}(\textit{isort}(x))$  is an inductive theorem (an exercise problem). If we define  $\textit{perm}(x, y)$  to be true iff  $x$  is a permutation of  $y$ , then we may also prove that  $\textit{perm}(x, \textit{isort}(x))$  is an inductive theorem. Both theorems ensure the correctness of the insertion sort algorithm. These inductive theorems can be easily proved by today's inductive theorem prover. Many algorithms can be formally verified automatically using this approach.

## 7.4 Resolution with Equality

Up to this point in this chapter, our focus was on equations, which are positive unit clauses with the equality as the only predicate. Equations are expressive enough to describe all the computations. For instance, all Turing machines can be expressed as a set of equations. On the other hand, it is more convenient to express axioms and theorems in general clauses. The Knuth-Bendix completion procedure suggests that we do not need the equality axioms explicitly if we have an inference like superposition which computes critical pairs from two rewrite rules. This is also true in resolution-based proving as Robinson and Wos suggested the paramodulation rule, shortly after the invention of the resolution rule.

### 7.4.1 Paramodulation

**Definition 7.4.1. (paramodulation)** Suppose clause  $C_1$  is  $(s = t \mid \alpha)$  and  $C_2$  is  $(A \mid \beta)$ , where  $A$  is a literal and  $p$  is a non-variable position of  $A$ ,  $\alpha$  and  $\beta$  are the rest literals in  $C_1$  and  $C_2$ , respectively, the paramodulation rule is defined by the

following schema:

$$\frac{(s = t \mid \alpha), \quad (A \mid \beta)}{(A[p \leftarrow t] \mid \alpha \mid \beta)\sigma}$$

where  $\sigma$  is the mgu of  $s$  and  $A/p$ . The clause  $(A[p \leftarrow t] \mid \alpha \mid \beta)\sigma$  produced by the paramodulation rule is called *paramodulant of the paramodulation*;  $C_1$  and  $C_2$  are the parents of the paramodulant. We use  $\text{paramod}(C_1, C_2)$  to denote the clause  $(A[p \leftarrow t] \mid \alpha \mid \beta)\sigma$ .

**Example 7.4.2.** Let  $C_1$  be  $(f(g(y)) = a \mid r(y))$  and  $C_2$  be  $(p(g(f(x))) \mid q(x))$ , then  $f(g(y))$  and  $p(g(f(x)))/1.1 = f(x)$  are unifiable with  $\sigma = [x \leftarrow g(y)]$ . The paramodulant of this paramodulation is  $(p(g(a)) \mid q(g(y)) \mid r(y))$ .

**Proposition 7.4.3. (soundness of paramodulation)** *For any clauses  $C_1$  and  $C_2$ , if  $\text{paramod}(C_1, C_2)$  exists, then  $E_{ax} \cup \{C_1, C_2\} \models \text{paramod}(C_1, C_2)$ , where  $E_{ax}$  is the equality axioms.*

*Proof.* Let  $\sigma$  be the mgu used in the paramodulation,  $C_1\sigma$  is  $(s' = t') \mid \alpha'$  and  $C_2\sigma$  is  $(A'[s'] \mid \beta')$ , then  $\text{paramod}(C_1, C_2)$  is  $(A'[t'] \mid \alpha' \mid \beta')$ .

From the monotonic axiom, we can deduce  $(s' = t') \wedge A'(s') \rightarrow A'(t')$ , whose clausal form is  $C_3 : (s' \neq t' \mid \overline{A'(s')} \mid A'(t'))$ .  $\text{resolve}(C_1\sigma, C_3) = (\overline{A'(s')} \mid A'(t') \mid \alpha')$ , named as  $C_4$ .  $\text{resolve}(C_2\sigma, C_4) = (A'(t') \mid \alpha' \mid \beta')$ , which is the same as  $\text{paramod}(C_1, C_2)$ .  $\square$

Using paramodulation, we do not need most of the equality axioms. It can be shown that any consequence from the equality axioms can be deduced from paramodulation, with the only exception of the reflexivity:  $x = x$ . The axiom  $x = x$  is indispensable in showing that  $A = \{f(x) \neq f(y)\}$  is unsatisfiable. We may introduce the following inference rule which makes the axiom  $x = x$  redundant.

**Definition 7.4.4. (reflexing)** *Suppose clause  $C$  is  $(s \neq t \mid \alpha)$ , reflexing is an inference rule defined by the following schema:*

$$\frac{(s \neq t \mid \alpha)}{(\alpha)\sigma}$$

where  $\sigma$  is the mgu of  $s$  and  $t$ .

Reflexing works in the same way as factoring: factoring removes one literal by unifying two literals; reflexing removes  $s \neq t$  by unifying  $s$  and  $t$ .

Paramodulation combined with resolution, factoring and reflexing is refutationally complete for theorem proving with equality. That is, if  $A \cup E_{ax}$  is unsatisfiable, where  $E_{ax}$  is the equality axioms, then the empty clause can be derived from

$A$  by paramodulation, resolution, factoring, and reflexing. It remains complete with any complete resolution strategies, such as positive resolution or ordered resolution. We may add restrictions on paramodulation as we put restrictions on resolution, without sacrificing the refutational completeness.

**Definition 7.4.5. (ordered paramodulation)** *Given a simplification order  $\succ$ , suppose  $s = t$  is maximal in clause  $C_1 (s = t \mid \alpha)$ ,  $s \not\prec t$ , and literal  $A$  is maximal in  $C_2 (A \mid \beta)$ ,  $p$  is a non-variable position of  $A$ ,  $\alpha$  and  $\beta$  are the rest literals in  $C_1$  and  $C_2$ , respectively. The ordered paramodulation rule will generate the paramodulant  $(A[p \leftarrow t] \mid \alpha \mid \beta)\sigma$ , where  $\sigma$  is the mgu of  $s$  and  $A/p$ .*

The paramodulation in Example 7.4.2 is an ordered paramodulation if  $f(g(y)) = a$  is maximal in  $C_1$  and  $p(g(f(x)))$  is maximal in  $C_2$ .

Using the same simplification order, ordered paramodulation can be combined with ordered resolution without losing the refutational completeness. When a set of clauses is just a set of equations (and regarded as a set of rewrite rules), ordered paramodulation is reduced to superposition. Robinson and Wos's work on paramodulation preceded, but apparently did not influence, the work of Knuth and Bendix, who independently formulated superposition and rewriting, called paramodulation and demodulation by Robinson and Wos, and then applied them to the notion of a canonical rewrite system in their 1970 paper.

#### 7.4.2 Simplification Rules

For a theorem prover based on resolution and paramodulation, huge number of new clauses are generated and we need to control redundancy in large search spaces. Restricted strategies for using resolution and paramodulation are effective to reduce redundant clauses, deleting unnecessary clauses in the search for the empty clause is also important to control the search space. As shown before, we may apply the following deletion strategies in a resolution-based theorem prover:

- **pure literal deletion:** Clauses containing pure literals are discarded.
- **tautology deletion:** Tautology clauses are discarded.
- **subsumption deletion:** Subsumed clauses are discarded.
- **unit deletion:** Those literals which are the negation of instances of some unit clauses are discarded from a clause.

Simplification strategies allow us to keep clauses in simplified forms so that multiple copies of equivalent clauses are not needed. For instance, if a clause contains a term  $s$ , say  $(p[s] \mid q)$ , and we have a rewrite rule  $s \rightarrow t$  or simply  $s = t$ , then this clause can be rewritten to  $(p[t] \mid q)$  by  $s = t$ . Rewriting is not the only way to obtain  $(p[t] \mid q)$ , because paramodulation can generate it from  $(p[s] \mid q)$  and  $s = t$ . However, there is no need to keep both  $(p[s] \mid q)$  and  $(p[t] \mid q)$ , because both are logically equivalent under equality. In other words, when  $(p[s] \mid q)$  is rewritten to  $(p[t] \mid q)$ , we keep the latter and throw the former away. That is why we call this kind of strategies a *simplification strategy*.

To work with ordered resolution and paramodulation, we may put a restriction on equality crossing as a simplification rule: **Ordered Equality Crossing**

$$\frac{(s \neq t) \mid A(s)}{(s \neq t) \mid A(t)}$$

if  $s \succ t$  for the simplification order  $\succ$  used in the ordered paramodulation.

Its soundness follows from Proposition 7.3.4.

**Definition 7.4.6. (contextual rewriting)** *Given a clause  $C_1 : (s = t \mid \alpha)$ , where  $s \succ t$ , clause  $C_2 : (A \mid \beta)$  can be rewritten by  $C_1$  if  $p$  is a position of  $A$ , and there exists a substitution  $\sigma$  such that  $A/p = s\sigma$ , and  $\alpha\sigma \subseteq \beta$ . The result of this rewriting is  $C'_2 : (A[p \leftarrow t\sigma] \mid \beta)$ .*

**Example 7.4.7.** Let  $C_1$  be  $(f(g(y)) = a \mid q(y))$  and  $C_2$  be  $(p(f(g(b))) \mid q(b))$ , then  $C_1$  can rewrite  $C_2$  to  $(p(a) \mid q(b))$  because  $p(f(g(b)))/1 = f(g(b)) = f(g(y))\sigma$ , where  $\sigma = [y \leftarrow b]$ , and  $q(y)\sigma$  appears in  $C_2$ .

Note that if  $\alpha$  is empty, then  $C_1$  is a positive unit clause and contextual rewriting is reduced ordinary rewriting.

**Proposition 7.4.8. (soundness of contextual rewriting)** *If  $C_1$  rewrites  $C_2$  to  $C'_2$  by contextual rewriting, then  $E_{ax} \cup \{C_1\} \models C_2 \leftrightarrow C'_2$ , where  $E_{ax}$  is the equality axioms.*

*Proof.* Suppose  $C_1$  is  $(s = t \mid \alpha)$ ,  $s \succ t$ ,  $C_2$  is  $(A[s] \mid \beta)$ , and  $C'_2$  is  $(A[s] \mid \beta)$ . Then  $C'_2$  can be obtained from  $C_1$  and  $C_2$  by paramodulation. Switching the role of  $s$  and  $t$ ,  $C'_2$  can also be obtained from  $C_1$  and  $C'_2$  by paramodulation. The soundness of paramodulation ensures the equivalence of  $C_2$  and  $C'_2$  under  $C_1$  and  $E_{ax}$ .  $\square$

There are two possible extensions to contextual rewriting. In contextual rewriting, clause  $C_1 (s = t \mid \alpha)$  serves as a conditional rewrite rule  $s \rightarrow t$  **if**  $\neg\alpha$ . If the second clause  $C_2$  contains a term  $s'$  such that  $s' = s\sigma$  for some substitution

$\sigma$ , assuming  $C_2$  is  $(A[s'] \mid \beta)$ , then  $s'$  is replaced by  $t\sigma$  if the condition  $\neg\alpha$  is true by enforcing  $\alpha\sigma \subseteq \beta$ . This condition can be replaced by the condition  $\neg\beta \models \neg\alpha\sigma$ , where  $\neg\beta$  is called the *context* of  $\neg\alpha\sigma$ .

The second extension of contextual rewriting is to treat non-equality literals as rewrite rules. That is, if  $C_1$  is  $(A \mid \alpha)$ , where  $A$  is a maximal non-equality literal in  $C_1$ , then we regard  $C_1$  as a rewrite rule  $A = \top$  **if**  $\neg\alpha$  if  $A$  is positive; or  $B = \perp$  **if**  $\neg\alpha$  if  $A = \neg B$ . Then  $C_1$  can rewrite atoms in other clauses to  $\top$  or  $\perp$ . If the literals in the other clauses are rewritten to  $\top$ , this is equivalent to subsumption. If they are rewritten to false, the other clauses will have one less literal. In particular, if  $\alpha$  is empty, this is equivalent to unit deletion (Definition 6.2.10).

The soundness of contextual rewriting remains true with the two extensions, thus contextual rewriting can serve as a good simplification rule which covers ordinary rewriting, subsumption, and unit deletion.

**Theorem 7.4.9.** *Ordered paramodulation combined with ordered resolution, factoring, relexing, contextual rewriting, and ordered equality crossing is refutationally complete for theorem proving with equality.*

The proof of this theorem follows the same proof line of Theorem 3.3.12: If the empty clause is generated, then the input clauses are unsatisfiable because all inference and simplification rules are sound. If the empty clause can never be generated and the clause set is saturated by the inference rules, then a Herbrand model can be constructed using the simplification order on  $T(F)$ .

### 7.4.3 Equality in Prover9

All the inference and simplification rules discussed in this section have been implemented in Prover9, with the exception of contextual rewriting. Rewriting using equations, called *demodulation*, is supported in Prover9. In fact, most examples provided in the distribution of Prover9 contain the equality.

**Example 7.4.10.** In Example 7.1.5, we are given three equations which specify a free group  $(S, *)$ . Can we reduce the number of equations to one? What is the minimal size of such a single equation? These are examples of questions interested by mathematicians. Here is an answer to the first question:

$$y * i(z * (((u * i(u)) * i(x * z)) * y)) = x$$

which involves four variables.

Prover9's input file to this problem is simply the single axiom and the goals will be the three equations in Example 7.1.5, with the only exception that we use here  $y * i(y)$  for the identity  $e$ .

```
formulas(sos).
y * i(z * ((u * i(u)) * i(x * z)) * y)) = x # label(oneAxiom).
end_of_list.
```

```
formulas(goals).
(x * y) * z = x * (y * z) # label(associativity).
x * i(x) = y * i(y) # label(inverse).
x * (y * i(y)) = x # label(identity).
end_of_list.
```

It took Prover9 less than a half second to find proofs for all the three equations after generating 1587 clauses. The proof of the associativity is the hardest, which consists of 48 steps after the proofs of the other equations.

The equality symbol “=” is built in; so is “!=”, which stands for the negation of equality. The major options involving paramodulation are the following:

```
clear(ordered_para). % default set
set(ordered_par).
```

```
assign(para_lit_limit, n). % default n=-1, range [-1 .. INT_MAX]
```

If  $n \neq -1$ , each parent in paramodulation can have at most  $n$  literals. This option may cause incompleteness of the inference system.

```
set(para_units_only).
clear(para_units_only). % default clear
```

This flag says that both parents for paramodulation must be unit clauses. The only effect of this flag is the same as to assign 1 to the parameter `para_lit_limit`.

The major options involving rewriting are the following:

```
assign(demod_step_limit, n). % default n=1000, range [-1 .. INT_MAX]
```

This parameter limits the number of rewrite steps that are applied to a clause during demodulation. If  $n = -1$ , there is no limit.

```
assign(demod_size_limit, n). % default n=1000, range [-1 .. INT_MAX]
```

This parameter limits the size (measured as symbol count) of terms as they are demodulated. If any term being demodulated has more than  $n$  symbols, demodulation of the clause stops. If  $n = -1$ , there is no limit.

```
set(back_demod).
clear(back_demod). % default clear
```

If this flag is set, a new rewrite rule will try to rewrite equations in the `usable` and `sos` lists. Prover9 also provides options for non-orientable equations to be rewrite rules under certain conditions and rewriting using these rules are also restricted.

## 7.5 Mace4: Finite Model Finding in FOL

Mace4 is a software tool for finite model finding in a first-order language. Mace4 is distributed together with Prover9 and was created by William McCune. Prover9 and Mace4 share a first-order language and there is a GUI interface for both of them. Mace4 searches for finite models satisfying the formulas in a first-order language based on the DPLL procedure. If the formula is the denial of some conjecture, any model found by Mace4 are counterexamples to the conjecture. Mace4 is a valuable complement to Prover9, looking for counterexamples before (or at the same time as) using Prover9 to search for a proof. It can also be used to help debugging input clauses and formulas for Prover9.

### 7.5.1 Use of Mace4

For the most part, Mace4 accepts the same input files as Prover9. If the input file contains commands that Mace4 does not understand, then the argument “-c” must be given to tell Mace4 to ignore those commands. For example, we can run the following two jobs in parallel, with Prover9 looking for a proof, and Mace4 looking for a counterexample.

```
prover9 -f x2.in > x2.prover9.out
mace4 -c -f x2.in > x2.mace4.out
```

Most of the options accepted by Mace4 can be given either on the command line or in the input file. The following command lists the command-line options accepted by Mace4.

```
mace4 -help
```

Mace4 searches for unsorted finite models only. That is, a model has one underlying finite set, called the domain (or the universe), and the members are always  $0, 1, \dots, n - 1$  for a set of size  $n$ . The models are the structures which define functions and relations over the domain, as an interpretation to the function and predicate symbols in the formula (using the same symbols). By default, Mace4 starts searching for a structure of domain size 2, and then it increments the size until it succeeds or reaches some limit. The size of the initial domain or the incremental size can be specified by the user.

If a formula contains constants that are natural-numbers,  $\{0, 1, \dots\}$ , Mace4 assumes they are members of the domain of some structure, that is, they are distinct objects; in effect, Mace4 operates under the assumptions  $0 \neq 1$ ,  $0 \neq 2$ , and so on. To Prover9, natural numbers are just ordinary constants. This is a subtle difference between Prover9 and Mace4. Because Mace4 assumes that natural numbers are members of the domain, if a formula contains a natural number that is out of range ( $\geq n$ , when searching for a structure of size  $n$ ), Mace4 will terminate with a fatal error.

Mace4 and Prover9 have the same restrictions on the goal formulas it accepts. Mace4 negates the goals and translates them to clauses in the same way as Prover9. The term “goal” is not particularly intuitive for Mace4 users, because Mace4 does not prove things. It makes more sense, however, when one thinks of Mace4 as searching for a counterexample to the goal.

Mace4 uses the following commands to specify the initial domain size, the maximal domain size, and the increment of the domain size.

```
assign(domain_size, n).    % default n=2, range [2 .. 200]
                          % command-line -n n
assign(iterate_up_to, n). % default n=10, range [-1 .. 200]
                          % command-line -N n
assign(increment, n).     % default n=1, range [1 .. 200]
                          % command-line -i n
```

These three parameter work together to determine the domain sizes to be searched. The search starts for structures of size `domain_size`; if that search fails, the size is incremented, and another search starts. This continues up through the value `iterate_up_to` (or until some other limit terminates the process).

**Example 7.5.1.** In section 4.3.3, we showed how to find a Latin square with special constraints by a SAT solver. We used the propositional variable  $p_{x,y,z}$  to represent the meaning that the cell at row  $x$  and column  $y$  contains value  $z$ . The same relation

between  $x$ ,  $y$  and  $z$  can be equally specified by the predicate  $p(x, y, z)$ . Thus, the five sets of propositional clauses can be specified by the five first-order formulas as follows:

1.  $\forall x \forall y \exists z p(x, y, z)$ .
2.  $\forall x \forall y \forall z \forall w ((z = w) \mid \overline{p(x, y, z)} \mid \overline{p(x, y, w)})$ .
3.  $\forall x \forall y \forall z \forall w ((y = w) \mid \overline{p(x, y, z)} \mid \overline{p(x, w, z)})$ .
4.  $\forall x \forall y \forall z \forall w ((x = w) \mid \overline{p(x, y, z)} \mid \overline{p(w, y, z)})$ .
5.  $\forall x \forall y \forall z \forall w (\overline{p(y, x, z)} \mid \overline{p(z, y, w)} \mid p(w, y, x))$ .

These formulas can be easily converted into clauses and can be specified in Mace4 as the following:

```
assign(domain_size, 5).
formulas(assumptions).
  p(x, y, f(x, y)).
  (z = w) | -p(x, y, z) | -p(x, y, w).
  (y = w) | -p(x, y, z) | -p(x, w, z).
  (x = w) | -p(x, y, z) | -p(w, y, z).
  -p(y, x, z) | -p(z, y, w) | p(w, y, x).
end_of_list.
```

Note that the function  $f$  in the first clause is a Skolem function from  $\exists z$ . Mace4 will find a model of size 5 for this input and it displays the value of  $p(x, y, z)$  as a Boolean matrix of 25 rows by 5 column, not a very readable format for a Latin square. From the meaning of  $p(x, y, z)$ , we may reconstruct the Latin square as follows.

```
0 2 1 4 3
4 1 3 2 0
3 4 2 0 1
1 0 4 3 2
2 3 0 1 4
```

If we specify the Latin square using  $x * y = z$  instead of  $p(x, y, z)$ , the input to Mace4 is the following.

```

assign(domain_size, 5).

formulas(assumptions).
  x * z != y * z | x = y.      % Latin square axioms
  x * y != x * z | y = z.
  x * x = x.                   % idempotent law
  ((y * x) * y) * y = x.      % special constraint
end_of_list.

```

Note that  $x * z! = y * z \mid x = y$  specifies that if rows  $x$  and  $y$  contain the same value at column  $z$ , then  $x = y$ ;  $x * z! = y * z \mid x = y$  specifies that if column  $y$  and  $z$  contain the same value at row  $x$ , then  $y = z$ . An alternative specification will use  $\backslash$  and  $/$  as the Skolem functions for  $u$  and  $v$ , respectively, in

$$\forall x, y \exists! u, v (x * u = y) \wedge (v * y = x)$$

where  $\exists!$  means “there exists uniquely”. Now Mace4’s input will look like

```

assign(domain_size, 5).

formulas(assumptions).
  % quasigroup axioms (equational)
  x * (x \ y) = y.
  x \ (x * y) = y.
  (x / y) * y = x.
  (x * y) / y = x.

  ((y * x) * y) * y = x. % special constraint
end_of_list.

```

Mace4 will find the following model:

```

interpretation( 5, [number=1, seconds=0], [

function(*(_,_), [      function(/(_,_), [      function(\(_,_), [
  0, 2, 1, 4, 3,        0, 4, 3, 1, 2,          0, 2, 1, 4, 3,
  3, 1, 4, 0, 2,        4, 1, 0, 2, 3,          3, 1, 4, 0, 2,
  4, 3, 2, 1, 0,        3, 0, 2, 4, 1,          4, 3, 2, 1, 0,
  2, 4, 0, 3, 1,        1, 2, 4, 3, 0,          2, 4, 0, 3, 1,
  1, 0, 3, 2, 4 ]),    2, 3, 1, 0, 4 ]),    1, 0, 3, 2, 4 ]))
]).

```

The Latin square defined by  $*$  is different from the one obtained by  $p(x, y, z)$ . If we use the command `assign(max_models, -1)`, Mace4 will find all the six models when the domain size is 5; it will find 120 models when the domain size is 7 (you need to turn off the `auto` command).

### 7.5.2 Finite Model Finding by SAT Solvers

The working principle of Mace4 is based on the fact that a set  $C$  of first-order clauses has a finite model if the propositional formula translated from  $C$  is satisfiable. Let  $D = \{0, 1, \dots, n-1\}$  be the domain of the finite model  $I$ , we will do the following to obtain a set of propositional clauses from  $C$ .

- For each predicate symbol  $p/k$ , define a set of propositional variables  $q_{x_1, \dots, x_k}^p$ ,  $x_i \in D$ , such that  $q_{x_1, \dots, x_k}^p$  is true iff  $p^I(x_1, \dots, x_k)$  is true.
- For each function symbol  $f/k$ , define a set of propositional variables  $q_{x_1, \dots, x_k, y}^f$ ,  $x_i, y \in D$ , such that  $q_{x_1, \dots, x_k, y}^f$  is true iff  $f^I(x_1, \dots, x_k) = y$ , and add the following clause into  $C$ :

$$(f(x_1, \dots, x_k) \neq y_1 \mid f(x_1, \dots, x_k) \neq y_2 \mid y_1 = y_2)$$

where  $s \neq t$  stands for  $\neg(s = t)$ .

- For each clause  $A$  of  $C$ , if  $f(t_1, \dots, t_k)$  appears in  $A$  and not in the format of  $f(x_1, \dots, x_k) = y$  (or  $y = f(x_1, \dots, x_k)$ ), then replace  $A$  by  $A[f(t_1, \dots, t_k) \leftarrow y] \mid f(t_1, \dots, t_k) \neq y$ , where  $y$  is a new variable. This step is called *flattening*. Once this step stops, every atom in  $C$  will be either  $p(x_1, \dots, x_k)$  or  $f(x_1, \dots, x_k) = y$ , where  $x_1, \dots, x_k, y$  are free variables.
- For each flattened clause  $A$  of  $C$ , for each variable  $x$  of  $A$ , for each value  $d$  of  $D$ , create all the ground instances of  $A$ . Let the set of all ground clauses be named as  $G$ .
- Finally, for each ground clause in  $G$ , replace  $p(d_1, \dots, d_k)$  by  $q_{d_1, \dots, d_k}^p$  and  $f(d_1, \dots, d_k) = d$  by  $q_{d_1, \dots, d_k, d}^f$ . Let the set of all propositional clause be named as  $propo(C)$ .

**Example 7.5.2.** For the equational specification of the Latin square problem, we used three function symbols:  $*$ ,  $\backslash$ , and  $/$ . The propositional variables will be  $p_{x,y,z}$  for  $x * y = z$ ,  $q_{x,y,z}$  for  $x \backslash y = z$ , and  $r_{x,y,z}$  for  $x / y = z$ ; if  $|D| = 5$ , there will be 375

propositional variables. The flattened clauses from the axioms plus the functional constraints are

$$\begin{array}{ll}
(x \setminus y \neq z \mid x * z = y) & // \ x * (x \setminus y) = y. \\
(x * y \neq z \mid x \setminus z = y) & // \ x \setminus (x * y) = y. \\
(x / y \neq z \mid z * y = x) & // \ (x / y) * y = x. \\
(x * y \neq z \mid z / y = x) & // \ (x * y) / y = x. \\
(y * z \neq z \mid z * y \neq w \mid w * y = x) & // \ ((y * x) * y) * y = x. \\
(x * y \neq z) \mid x * y \neq w \mid z = w) & // \text{“} * \text{” is a function} \\
(x \setminus y \neq z) \mid x \setminus y \neq w \mid z = w) & // \text{“} \setminus \text{” is a function} \\
(x / y \neq z) \mid x / y \neq w \mid z = w) & // \text{“} / \text{” is a function}
\end{array}$$

If  $|D| = 5$ , a clause with three free variables will generate  $5^3 = 125$  ground clauses; a clause with four free variable will generate  $5^4 = 625$  ground clauses. The ground instances of the above clauses are easy to convert to propositional clauses. For the literal  $(z = w)$ , its ground instance becomes  $\top$  if we have  $z\sigma = w\sigma = d$  and  $\perp$  if  $z\sigma = d_1 \neq d_2 = w\sigma$ .

While the conversion from a first-order formula to a set of propositional clauses uses the equality “=”, the equality axioms (introduced in the first section of this chapter) are not needed for the conversion because they produce only trivial propositional clauses. Only the function axiom is needed in the conversion for each function  $f$ :

$$(f(x_1, \dots, x_k) \neq y_1 \mid f(x_1, \dots, x_k) \neq y_2 \mid y_1 = y_2)$$

**Proposition 7.5.3.** *A first-order formula  $C$  in CNF has a finite model iff the propositional formula  $\text{propo}(C)$  is satisfiable.*

*Proof.* Let  $C$  be a CNF formula of  $(P, F, X, Op)$  and  $D = \{0, 1, \dots, n - 1\}$  be the finite domain.

If  $\text{propo}(C)$  is satisfiable, then  $\text{propo}(C)$  has a model  $\sigma$  and  $\sigma$  defines a relation  $R^p$  over  $D^k$  for each predicate symbol  $p/k$ :

$$R^p = \{\langle d_1, \dots, d_k \rangle \mid \sigma(q_{d_1, \dots, d_k}^p) = 1\}$$

and a function  $f^\sigma : D^k \rightarrow D$  for each function  $f/k$ :

$$f^\sigma(d_1, \dots, d_k) = d \text{ iff } \sigma(q_{d_1, \dots, d_k, d}^f) = 1.$$

Then it is ready to check that  $I = (D, \{R^p \mid p \in P\}, \{f^m \mid f \in F\})$  is a model of  $C$ .

On the other hand, if  $C$  has a finite model  $I = (D, R, G)$ , where  $D$  is finite, then we can create the interpretation  $\sigma_I$  for  $\text{propo}(C)$  as follows:

$$\sigma_I(q_{d_1, \dots, d_k}^p) = 1 \text{ iff } p^I(d_1, \dots, d_k) = 1$$

and

$$\sigma_I(q_{d_1, \dots, d_k, d}^f) = 1 \text{ iff } f^I(d_1, \dots, d_k) = d$$

It is ready to check that  $\sigma$  is a model of  $\text{propo}(C)$ .  $\square$

To use SAT provers to solve a problem of modest difficulty, we need to write an *encoder* which generates propositional clauses and then feed these clauses into a SAT solver. We also need a *decoder* which translates SAT solver's solution into a solution of the problem. Mace4 releases the burden of writing encoders/decoders by converting automatically the first-order clauses into propositional clauses with the free variables being substituted by all the values in the finite domain and then call a SAT solver to find a model; it is then convert the propositional model into a first-order model.

## 7.6 Exercise Problems

1. Show that  $(f(x_1, \dots, x_k) \neq y_1 \mid f(x_1, \dots, x_k) \neq y_2 \mid y_1 = y_2)$  from the equality axioms.
2. Write out in first order logic the statements that there are at least, at most, and exactly three elements  $x$  such that  $A(x)$  holds.
3. Complete the proof of case 2.2 in Theorem 7.2.16.
4. Given the definition of *add* in Example 7.3.2, prove by structural induction that  $\text{add}(x, s(y)) = s(\text{add}(x, y))$ .
5. Given the definition of *add* in Example 7.3.2, prove by structural induction that  $\text{add}(x, y) = \text{add}(y, x)$  and  $\text{mul}(x, y) = \text{mul}(y, x)$ .
6. Given the definitions of *add* and *mul* in Example 7.3.2, prove by structural induction that  $\text{add}(\text{add}(x, y), z) = \text{add}(x, \text{add}(y, z))$  and  $\text{mul}(\text{mul}(x, y), z) = \text{mul}(x, \text{mul}(y, z))$ .
7. Adding the following function to Example 7.3.2:

$$\begin{array}{l} 9 \quad \text{exp}(x, 0) \rightarrow s(0); \\ 10 \quad \text{exp}(x, s(y)) \rightarrow \text{mul}(x, \text{exp}(x, y)). \end{array}$$

where  $\text{exp}(x, y)$  computes  $x^y$ , the exponential function. Prove by structural induction that  $\text{exp}(x, \text{add}(y, z)) = \text{mul}(\text{exp}(x, y), \text{exp}(x, z))$ .

8. Adding the following functions to Example 7.3.2:

$$\begin{array}{ll}
 11 & fib(0) \rightarrow 0; \\
 12 & fib(s(0)) \rightarrow s(0); \\
 13 & fib(s(s(x))) \rightarrow add(fib(s(x)), fib(x)). \\
 14 & sum(0) \rightarrow 0; \\
 15 & sum(s(x)) \rightarrow add(sum(x), fib(s(x))).
 \end{array}$$

where  $fib(x)$  stands for the Fibonacci number and  $sum(x)$  stands for the summation of the first  $x$  Fibonacci numbers. Prove by the structural induction that  $sum(x) = pre(fib(s(s(x))))$ .

9. Adding the following function to Example 7.3.2:

$$\begin{array}{ll}
 16 & A(0, y) \rightarrow s(y); \\
 17 & A(s(x), 0) \rightarrow A(x, s(0)); \\
 18 & A(s(x), s(y)) \rightarrow A(x, A(s(x), y)).
 \end{array}$$

where  $A(x, y)$  stands for the Ackermann function. Prove by the structural induction that  $A(s(0), y) = s(s(y))$  and  $A(s(s(0)), y) = s(s(s(add(y, y))))$ .

10. Given the definition of  $rev$  in Example 7.3.8, prove that  $sorted(isort(y)) = \top$  is an inductive theorem by structural induction.
11. Use Prover9 to obtain the set of 10 rewrite rules from the three equations in Example 7.2.8.
12. Use Mace4 to decide if there exists Latin squares of size between 4 and 9 for each of the following constraints (called “short conjugate-orthogonal identities”), respectively.

Code	Constraint	Name
QG3	$(x * y) * (y * x) = x$	<i>Schröder quasigroup</i>
QG4	$(y * x) * (x * y) = x$	<i>Stein's third law</i>
QG5	$((y * x) * y) * y = x$	
QG6	$(x * y) * y = x * (x * y)$	<i>Schröder's first law</i>
QG7	$(y * x) * y = x * (y * x)$	<i>Stein's second law</i>
QG8	$x * (x * y) = y * x$	<i>Stein's first law</i>
QG9	$((x * y) * y) * y = x$	<i>C<sub>3</sub>-quasigroup</i>

13. Repeat the previous problem when an additional constraint,  $x * x = x$ , is added into the axiom set.

## CHAPTER 8

### PROLOG: PROGRAMMING IN LOGIC

Prolog is a programming language whose name comes from programming in logic. Unlike many other programming languages, Prolog comes from first-order logic and is intended primarily as a declarative programming language: the program is expressed in terms of relations, represented as facts and rules. A computation is initiated by running a query over these relations.

Based upon Robinson's resolution principle, Prolog was first developed and implemented in 1972 by Alain Colmerauer with Philippe Roussel, called "Marseille Prolog". The Japanese Fifth-Generation Computer Project, announced in 1981, adopted Prolog as a development language, and thereby focused considerable attention on the language and its capabilities. Prolog is the major example of a fourth generation programming language supporting the declarative programming paradigm, and is widely regarded as excellent languages for "exploratory" and "prototype programming".

There are several free implementations of Prolog available today, such as GNU's gprolog and SWI's SWI-prolog, and they are the descendants of Marseille Prolog. The language has been used for theorem proving, expert systems, automated planning, natural language processing, and is a popular tool in artificial intelligence and computational linguistics. Prolog is well-suited for specific tasks that benefit from rule-based logical queries such as searching databases, voice control systems, and filling templates.

#### 8.1 Prolog's Working Principle

In Prolog, program logic is expressed in terms of relations, and a computation is initiated by running a query over these relations.

##### 8.1.1 Horn Clauses in Prolog

Typically, relations and queries are constructed by Horn clauses. In section 3.4, we have seen that Horn clauses are clauses with at most one positive literal. In Prolog, Horn clauses are further divided into three types.

- **fact**: a unit positive clause;

- **rule**: a clause with one positive literal and one or more negative literals;
- **query**: a negative clause.

A Prolog program consists of facts and rules, which are collectively called *definite clause*. For a rule clause  $(A \mid \overline{B_1} \mid \dots \mid \overline{B_n})$ , where  $A, B_1, \dots, B_n, n \geq 1$ , are atomic formulas, it is written in Prolog as  $A : -B_1, \dots, B_n$ . This format looks closely to the equivalent formula

$$B_1 \wedge \dots \wedge B_n \rightarrow A.$$

In other words, the symbol  $:-$  is the inverse of  $\rightarrow$  and the commas are  $\wedge$  in Prolog. For a rule  $H : -B$ ,  $H$  is the *head* and  $B$  is the *body* of the rule. It is read as “H is true if B is true”. Clauses with an empty body are facts, which do not need the symbol “:-” or we may write it as “H :- true” in Prolog.

Negative clauses like  $(\overline{B_1} \mid \dots \mid \overline{B_n})$  are Prolog queries and denoted by  $?-B_1, \dots, B_n$ , which resembles the equivalent formula  $\neg(B_1 \wedge \dots \wedge B_n)$ .

**Example 8.1.1.** Suppose we want to study the family tree of a big family, we may state the father-child and mother-child relations by facts, and define other relations by rules. Here is a Prolog program.

```

mom(rose, john).    % Rose is the mother of John
mom(rose, bob).
mom(rose, ted).
papa(joe, john).    % Joe is the father of John
papa(joe, bob).
papa(joe, ted).
papa(ted, patrick).
papa(john, caroline).
papa(bob, kathleen).
papa(bob, joseph).
papa(bob, mary).
parent(X, Y) :- papa(X, Y). % X is parent of Y if X is papa of Y.
parent(X, Y) :- mom(X, Y). % X is parent of Y if X is mom of Y.
grandpa(X, Y) :- papa(X, Z), parent(Z, Y).
grandmo(X, Y) :- mom(X, Z), parent(Z, Y).

```

From this example, we can see that variables in Prolog are identifies whose first character is a capital letter. Once the above program is loaded into Prolog, we may issue queries after the symbol  $?-$ .

To query who are the parents of `ted`, we type `parent(X, ted)`. after the symbol `?-`. Prolog will answer with `X = joe`. If you type `“;”`, Prolog will give you another answer: `X = rose`. If you hit `“return”`, you will terminate the query.

To query the sibling relation, we may define the sibling relation as a rule,

```
sibling(X, Y) :- parent(Z, X), parent(Z, Y), X \= Y.
                % X and Y share parent Z.
```

where `\=` stands for “not equal”, and issue the query `?- sibling(X, Y)`. On the other hand, we may use the query

```
?- parent(Z, X), parent(Z, Y), X \= Y.
```

This query also prints out the common parent  $Z$  of  $X$  and  $Y$ . For instance the first answer found by Prolog is

```
X = john, Y = bob, Z = joe.
```

### 8.1.2 Resolution Proofs in Prolog

In Prolog, associated with every successful query, there exists a resolution proof. For instance, for the query `?- parent(X, ted)`, we have an answer `X = joe`. The corresponding resolution proof is given below.

Prolog format	clausal format
1 <code>papa(joe, ted)</code> .	$(f(joe, ted))$ // input
2 <code>parent(X, Y) :- papa(X, Y)</code> .	$(p(X, Y) \mid \overline{f(X, Y)})$ // input
3 <code>?- parent(X, ted)</code> .	$(\overline{p(X, ted)})$ // query
4 <code>?- papa(X, ted)</code> .	$(f(X, ted))$ // resolvent of 2 and 3
5 $\perp$ .	$\perp$ // resolvent of 1 and 4.

In the last resolution, the mgu is  $X \mapsto \text{joe}$ , which is the answer of the query. The query comes from the negation of the formula  $\exists X \text{parent}(X, \text{ted})$ , which is a logical consequence of the input clauses.

As another example, for the query `parent(Z, X), parent(Z, Y), X \= Y`, the corresponding resolution proof is the following.

```
1 papa(joe, john). // input
2 papa(joe, bob). // input
3 parent(X, Y) :- papa(X, Y). // input
4 ?- parent(Z, X), parent(Z, Y), X \= Y. // query
```

```

5 ?- papa(Z, X), parent(Z, Y), X \= Y.    // resolvent of 3 and 4
6 ?- parent(joe, Y), john \= Y.           // resolvent of 1 and 5
7 ?- papa(joe, Y), john \= Y.             // resolvent of 3 and 6
8 ?- john \= bob.                         // resolvent of 2 and 7
9 ⊥.                                       // built-in \ = returns ⊥.

```

The mgus used in the resolutions provide the solution of the query: Resolution between 1 and 5 provides  $X = \text{john}$ ,  $Z = \text{joe}$  and resolution between 2 and 7 provides  $Y = \text{bob}$ .

Both above resolution proofs are a negative, input and linear resolution proof. That is, every resolution involves a negative clause and an input clause as parents, and the latest resolvent (which is negative) is used in the next resolution. This is true for every resolution proof found by Prolog. Corollary 3.4.11 claims that, for propositional Horn clauses, if there exists only one negative clause, then the negative, input and linear resolution is complete for Horn clauses. This result can be extended to first-order clauses. In other words, the negative, input and linear resolution strategy is complete for first-order Horn clauses.

Thus, given a query, the Prolog engine attempts to find a resolution refutation of the query. If the empty clause can be derived, i.e., an instantiation for all free variables is found that makes the input clauses plus the query false, it follows that the negation of the query, with the found instantiation applied, is a logical consequence of the program. This makes Prolog particularly useful for symbolic computation and language parsing applications.

**Example 8.1.2.** Suppose we are asked to write a program to convert a propositional formula into negation normal form. This would be a week-long exercise using a conventional programming language. Using Prolog, it is a hour-long exercise. Suppose we use only three propositional variables, say  $p$ ,  $q$ ,  $r$ , and the logical operators are “n” for  $\neg$ , “a” for  $\wedge$ , “o” for  $\vee$ , “i” for  $\rightarrow$ , and “e” for  $\leftrightarrow$ , the Prolog program is given below:

```

pvar(p).    % definition of propositional variables
pvar(q).
pvar(r).

literal(n(L)) :- pvar(L). % definition of literals
literal(L) :- pvar(L).

nnf(n(n(X)), NNF) :- nnf(X, NNF). % double negation

```

```

nnf(a(A, B), a(N1, N2)) :- nnf(A, N1), nnf(B, N2). % and
nnf(n(a(A, B)), o(N1, N2)) :- nnf(n(A), N1), nnf(n(B), N2).

nnf(o(A, B), o(N1, N2)) :- nnf(A, N1), nnf(B, N2). % or
nnf(n(o(A, B)), a(N1, N2)) :- nnf(n(A), N1), nnf(n(B), N2).

nnf(i(A, B), o(N1, N2)) :- nnf(n(A), N1), nnf(B, N2). % implication
nnf(n(i(A, B)), a(N1, N2)) :- nnf(A, N1), nnf(n(B), N2).

nnf(e(A, B), o(a(NA, NB), a(NnA, NnB))) :- % equivalence
    nnf(A, NA), nnf(n(A), NnA), nnf(B, NB), nnf(n(B), NnB).
nnf(n(e(A, B)), o(a(NA, NnB), a(NnA, NB))) :-
    nnf(A, NA), nnf(n(A), NnA), nnf(B, NB), nnf(n(B), NnB).

nnf(X, X) :- literal(X). % NNF of literal is itself

```

The following is a query to the above program:

```

?- nnf(n(o(n(a(p,q)),e(p,q))), X).
X = a(a(p,q),o(a(p,n(q)),a(n(p),q))) ?

```

### 8.1.3 A Goal-Reduction Procedure

The Prolog engine can be regarded as a goal-reduction process: For each  $G$  in the original query, the engine will search for a rule in the Prolog program whose head unifies with  $G$  and recursively apply the engine to the subgoals in the body of the rule. When all subgoals are solved, or when the body is empty, goal  $G$  is considered as *solved* and a substitution for the variables of  $G$  is returned; otherwise, the label “fail” is returned.

**Procedure 8.1.3.** Assuming  $P$  is a Prolog program, procedure *engine1* takes a single goal  $G$  as input and returns a substitution as output if one rule of  $P$  can solve  $G$ ; otherwise, *fail* is returned. Procedure *engine* takes a list  $Gs$  of goals as input and returns a substitution as output if every goal of  $Gs$  is solved by  $P$ .

```

proc engine1( $G$ )
1   for ( $H : - B$ )  $\in P$  do
2       if ( $H$  and  $G$  are unifiable with mgu  $\sigma$ )
3            $\theta := engine(B\sigma)$ ;

```

```

4         if  $\theta \neq fail$  return  $\sigma\theta$ ;
5     return fail;

proc engine(Gs)
1      $\sigma = \theta := \emptyset$ ;
2     for  $A \in Gs$  do
3          $\sigma := engine1(A\theta)$ ;
4         if  $\sigma = fail$  break; else  $\theta := \theta\sigma$ ;
5     if  $\sigma = fail$  return fail; else return  $\theta$ ;

```

We assume that  $B$  is empty if  $H : -B$  is a fact. The first call is  $engine(Qs)$ , where  $Qs$  are the goals in a query.

Obviously,  $engine1$  and  $engine$  are mutually recursive procedures. If we examine the recursion tree of  $engine1$  and  $engine$ , the links going out of the  $engine1$  node are labeled by rules of  $P$ ; and the links going out of the  $engine$  node are labeled by a goal in  $Gs$ . The  $engine1$  node is an  $\vee$ -node: if one of its children succeeds, it will succeed. The  $engine$  node is an  $\wedge$ -node: if one of its children fails, it will fail. In other words, the recursion tree of  $engine$  and  $engine1$  is an  $\wedge$ - $\vee$  tree.

Procedure  $engine1$  as given intends to find one answer; if more answers are desired, then line 4 of  $engine1$  will remember which rule is used to get  $\sigma$  and will use the next rule when more solutions are desired. In this case, the  $engine1$  node may have more than one successful child.

A leaf node in the  $\wedge$ - $\vee$  tree is said to be *ok* if (a) it is an  $engine$  node and  $Gs$  is empty or (b) it is an  $engine1$  node and  $G$  is a built-in atom evaluated to be true. An internal node in the  $\wedge$ - $\vee$  tree is said to be *ok* if (c) it is an  $engine$  node and every child node is *ok*, or (d) it is an  $engine1$  node and one of its children is *ok*.

In this case, the  $\wedge$ - $\vee$  tree has a solution iff its root node is *ok* and it represents all the solutions dictated by the  $\vee$  nodes.

**Example 8.1.4.** The following program implements a depth-first search procedure for a directed graph:

```

% s(u, v): directed edge (u, v)
1: s(a,b).
2: s(a,c).
3: s(b,d).
4: s(c,g).
5: s(d,g).
6: goal(e). % goal(x): x is a goal

```

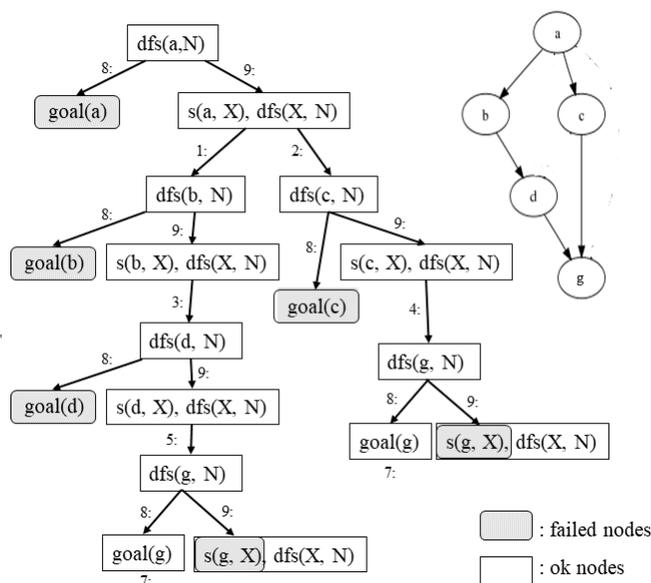


Figure 8.1.1:  $\vee$ -nodes in the  $\wedge$ - $\vee$  tree for the query `dfs(a, N)`.

```

7: goal(g).
8: dfs(N, N) :- goal(N). % dfs(startNode, solution)
9: dfs(N, l(N, R)) :- s(N, X), dfs(X, R).

```

There are two solutions to the query `?- dfs(a, N)`. That is,  $N = l(a, l(b, l(d, g)))$  and  $N = l(a, l(c, g))$ , where a term like  $l(a, l(c, g))$  represents a path (a, c, g) in the graph. Later, after introducing the list, the same path can be represented by `[a, c, g]`. Figure 8.1.1 shows the  $\vee$ -nodes in the search tree for the query `dfs(a, N)`; the  $\wedge$ -nodes are implicit. The shaded nodes are the *failed* nodes; the other nodes are *ok* nodes.

Note that the two recursive procedures *engine* and *engine1* just illustrate the goal-reduction process; the actual Prolog engine is an interactive procedure and uses a stack to save all the information in a search path of the  $\wedge$ - $\vee$  tree.

Procedure *engine1* shows that the order of rules in  $P$  is important. Procedure *engine* shows that the order of atomic formulas in the body of a rule is important. In other words, the links in the  $\wedge$ - $\vee$  tree are ordered and it implements a special case of the negative, input and linear resolution strategy: the literals in a rule are handled from left to right, and the commas, which represent  $\wedge$  in logic, are not commutative. The input clauses are tried to unify with a given goal from first to last. From the view point of search strategies, this is the depth-first search in the  $\wedge$ - $\vee$  tree.

**Example 8.1.5.** Adding the following rules to the Prolog program in Example 8.1.1:

```
descen(X, Y) :- parent(Y, X). % X is a descendant of Y
descen(X, Z) :- parent(Y, X), descen(Y, Z).
```

The last rule is different from the previous rules in that `descen` is recursively defined. If we ask “who are the descendants of Joe”, the query is `?- descen(X, joe)`. Prolog will find that `X = john, bob, ted, patrick, caroline, kathleen, joseph,` and `mary`, in the given order. For the first answer, i.e., `X = john`, the recursive calls of `engine` and `engine1` go as follows:

```
engine([descen(X, joe)])
  engine1(descen(X, joe))
    engine([parent(joe, X)])
      engine1(parent(joe, X))
        engine([papa(joe, X)])
          engine1(papa(joe, X)) return X = john
...
engine([descen(X, joe)]) return X = john
```

For the last answer, i.e., `X = mary`, the recursive calls go as follows (the calls to `engine` are ignored if the list contains a single goal).

```
engine1(descen(X, joe))
  engine([parent(Y, X), descen(Y, joe)])
    engine1(parent(Y, X))
      engine1(papa(Y, X)) return {X = mary, Y = bob}
    engine1(descen(bob, joe))
      engine1(parent(bob, joe))
        engine1(papa(bob, joe)) return { }
...
  engine([parent(Y, X), descen(Y, joe)]) return {X = mary, Y = bob}
engine1(descen(X, joe)) return {X = mary, Y = bob}
```

However, if we change the last rule to a logically equivalent one:

```
descen(X, Z) :- descen(Y, Z), parent(Y, X).
```

The same query, `?- descen(X, joe)`, will produce `X = john, bob, ted, caroline, kathleen, joseph, mary, and patrick`, the same set of eight answers in different orders, and will cause Prolog to crash if you ask for more answers. If we switch the order of the last two clauses, i.e., the last two clauses are

```
descen(X, Z) :- descen(Y, Z), parent(Y, X).
descen(X, Y) :- parent(Y, X).
```

Then the same query will cause Prolog to crash without generating any answer. The first-order logic cannot explain the difference because these two clauses are logically equivalent. The procedures *engine* and *engine1* can explain the difference perfectly:

The initial goal list is `[descen(X, joe)]` in all the three cases. In the last case, the rule `descen(X, Z) :- descen(Y, Z), parent(Y, X)` is applied first, and the recursive calls of *engine* and *engine1* go as follows.

```
engine1(descen(X, joe))
  engine([descen(Y, joe), parent(Y, X)])
    engine1(descen(Y, joe))
      engine([descen(Y1, joe), parent(Y1, Y)])
        engine1(descen(Y1, joe))
          engine([descen(Y2, joe), parent(Y2, Y1)])
            ...
```

The unsolved goal list contains a long sequence of goals like

```
..., parent(Y3, Y2), parent(Y2, Y1), parent(Y1, Y), parent(Y, X)
```

until Prolog runs out of memory. Here we changed the variable names because each clause is assumed to have different variables from others.

The above example shows that Prolog uses a depth-first strategy, which is not a fair strategy, to implement resolution, and becomes an incomplete theorem prover for Horn clauses. As a Prolog programmer, we need to be careful to avoid the pitfall of infinite loops associated with the depth-first strategy. We will address this issue in the next section.

**Example 8.1.6.** Hanoi's tower asks to move  $n$  disks from peg  $A$  to peg  $B$  using peg  $C$  as an auxiliary holding peg. At no time can a larger disk be placed upon a smaller disk. The Prolog code contains two predicate definitions: `move` writes out the move information; `hanoi(N,A,B,C)` checks if  $N = 1$ , then moves disk 1 from  $A$  to  $B$ ; otherwise, it tries to move the top  $N-1$  disks from  $X$  to  $C$ ; move bottom disk  $N$  from  $X$  to  $B$ ; and then  $N-1$  disks from  $C$  to  $B$ .

```

move(N,A,B) :-          % print a move info
    write('Move disk '), write(N), write(' from '),
    write(A), write(' to '), write(B), nl.

hanoi(1,A,B,C) :-      % disk 1 can move freely from A to B
    move(1,A,B).
hanoi(N,A,B,C) :-      % To move N disks from A to B,
    N>1, M is N-1,     % where N > 1,
    hanoi(M,A,C,B),    % we move the top N-1 disks from A to C
    move(N,A,B),       % then move disk N, the largest, from A to B
    hanoi(M,C,B,A).    % then move the N-1 disks from C to B

| ?- hanoi(3, a, b, c).
Move disk 1 from a to b
Move disk 2 from a to c
Move disk 1 from b to c
Move disk 3 from a to b
Move disk 1 from c to a
Move disk 2 from c to b
Move disk 1 from a to b

```

. The code comes directly from the description of the solution. It involves two recursive calls of `hanoi`. The recursive calls will terminate because the first argument of `hanoi` is reduced by one in each successive call.

## 8.2 Prolog's Data Types

As a programming language, Prolog has several built-in data types, such as numbers (floats or integers), strings, and lists, and associated functions. These data types are subtypes of a supertype, called *terms*. In this view, *term* is the only data type language in Prolog. Terms are either atoms, numbers, variables or compound terms. Atomic formulas in Prolog are called *callable terms*.

### 8.2.1 Atoms, Numbers, and Variables

We have been using “atom” for “atomic formula”; in this section, we will use “atomic formula” explicitly because “atom” has special meaning in Prolog: an *atom* is either an identifier other than variables or a single-quoted string of characters.

Examples of atoms include `x`, `red`, `'Taco'`, and `'some atom'`. An atom is a general-purpose name with no inherent meaning.

Prolog provides the predicate `atom/1`, where `atomX` is true iff `X` is an atom. For example, `atom('some atom')` is true; `atom(123)` is false.

Numbers can be floats or integers. The range of numbers that can be used is likely to be dependent on the number of bits (amount of computer memory) used to represent the number. Most of the major Prolog systems support arbitrary length integer numbers. The treatment of float numbers in Prolog is likely to vary from implementation to implementation; float numbers are not all that heavily used in Prolog programs, as the emphasis in Prolog is on symbol manipulation.

Prolog provides all the usually predicates on numbers, such as `>/2`, `</2`, `=</2`, `>=/2`, etc., and arithmetic operations, such as `-/1`, `+/1`, `+/2`, `-/2`, `*/2`, `//2`, `mod/2`, `rem/2`, `gcd/2`, `lcm/2`, `abs/1`, `sign/1`, `max/2`, `min/2`, `truncate/1`, `floor/1`, `ceiling/1`, `sqrt/1`, `sin/1`, `cos/1`, etc. You use `:=/2` to check if two numbers are equal; and `=~/2` for not equal. The symbol `:=/2` is reserved for unification: `s = t` returns true if `s` and `t` are unifiable.

Prolog also provides predicate `number/1` to check if a term is a number or not; `integer/1` for being an integer; `float/1` for being a float number. For example, `float(1.0)` returns true; `float(1)` returns false. Note that `integer(1+1)` returns false as `1+1` is a compound term. If we check the value of `1+1<3`, Prolog will return true because `1 + 1` will be evaluated to `2` before `<` is called on `(2, 3)`.

Variables are denoted by an identifier consisting of letters, numbers and underscore characters, and beginning with an upper-case letter or underscore. The underscore alone, `"_"`, denotes a nameless distinct variable. A nameless variable is used when we do not care about the value of the variable. For example, if we define the multiplication relation `mul(X, Y, Z)`, where `Z = X*Y`, the first rule can be `mul(0, _, 0)`, where the second position is a nameless variable.

Prolog variables closely resemble variables in logic in that it can be instantiated by arbitrary terms through unification. If a variable is instantiated, we say "the variable is bounded". Some operations require the involved variables are bounded. For instance, to show

```
descen(X, Z) :- parent(Y, X), descen(Y, Z).
```

is a terminating rule in Example 8.1.5, we need to check that the bindings of `X`, `Y` and `Z`. The predicate `parent(Y, X)` ensures that `X` and `Y` will be bounded to atoms when we reach `descen(Y, Z)`. When these variables are bounded by atoms, the atom bounded to `Y` is one step closer than that of `X` to `Z` in the descendant relation, this meaning defines a well-founded ordering `>` such that `X > Y` when we compare `descen(X, Z)` and `descen(Y, Z)`.

In Prolog, there is a special infix predicate called `is/2`, which was used in the Hanoi example. It resembles the assignment statement of a variable by a numeric expression in a conventional programming language. For example, the value of `X is 1+1` is true with `X` instantiated with 2. That is, Prolog will evaluate the right side of `is`, if it is a number, that number is bounded to the variable on the left side of `is`; otherwise it will abort with error. In contrast, `X = 1+1` will also return true, but it will bound the variable `X` to the compound term `1+1`. As another example, “`X = a`” will return true and “`X is a`” will be aborted, because `a` is not a number.

### 8.2.2 Compound Terms and Lists

A *compound term* is composed of a predicate or function symbol, called a “functor” and a number of “arguments”, which are again terms. Compound terms are ordinarily written as a functor followed by a comma-separated list of argument terms, which is contained in parentheses. The number of arguments is called the term’s arity. An identifier can be regarded as a compound term with arity zero.

There is a built-in predicate called `functor/3` which allows us to extract the functor and its arity from a compound term:

```
?- functor(f(a, b), F, N).
F = f, N = 2
```

Compound terms are indispensable for representing various data structures in conventional languages. For example, to represent a binary tree where each non-empty tree node contains a key, we may use the compound term `node(key, leftChild, rightChild)`. We can define easily predicate `tree/1`, which checks if its argument is a tree, and predicate `dfs/2`, which tries to find a key by depth-first search in a tree.

```
tree(nil).                % nil denotes the empty tree
tree(node(_, Left, Right) :- tree(Left), tree(Right).

dfs(node(K, _, _), K) :- write('found '), write(K).
dfs(node(_, L, _), K) :- dfs(L, K).                % go to Left branch
dfs(node(_, _, R), K) :- dfs(R, K).                % go to Right branch
```

In the definition of `dfs`, we used “`_`” twice in each rule; each “`_`” represents a unique variable.

One special member of compound terms is *list*, which is an ordered collection of terms. It is denoted by square brackets with the terms separated by commas, or

in the case of the empty list, by `[]`. Here are some lists of three elements, `[1,2,3]`, `[red,green,blue]`, or `[34,tom,[2,3]]`. Lists are the replacement of arrays in conventional programming languages.

For a non-empty list, the first member of a list is separated from the rest members by the symbol `|`. For example, the result of `[X | Y] = [1,2,3]` is `X = 1` and `Y = [2, 3]` and the result of `[X | Y] = [1]` is `X = 1` and `Y = []`.

Two lists unify if they are the same length and all their elements unify pairwise.

<code> ?-[a,B,c,D]=[A,b,C,d].</code>	<code> ?-[(a+X),(Y+b)]=[W+c),(d+b)].</code>
<code>A = a,</code>	<code>W = a,</code>
<code>B = b,</code>	<code>X = c,</code>
<code>C = c,</code>	<code>Y = d?</code>
<code>D = d ?</code>	<code>yes</code>
<code>yes</code>	

<code> ?-[a,b,c,d]=[Head Tail].</code>	<code> ?-[a,b,c,d]=[X,Y Z].</code>
<code>Head = a,</code>	<code>X = a,</code>
<code>Tail = [b,c,d]?</code>	<code>Y = b,</code>
<code>yes</code>	<code>Z = [c,d];</code>
	<code>yes</code>

### 8.2.3 Popular Predicates over Lists

Most Prolog implementations provide efficient implementations of the following predicates over lists. We also are providing hypothetical Prolog codes of these predicates, even though they are not identical to actual implementations.

- `list/1`: Check if its argument is a list.

```
list([]).
list([_ | X]) :- list(X).
```

- `member/2`: Check if the first argument is a member of the second argument.

```
member(X, [X | _]).
member(X, [_ | Y]) :- member(X, Y).
```

- `append/3`: Check if the third list is the concatenation of the first two lists.

```
append([], Y, Y).
append([F|X], Y, [F|Z]) :- append(X, Y, Z).
```

- **select/3:** `select(X, List1, List2)` is true when `List1`, with `X` removed, results in `List2`.

```
select(_, [], []).
select(X, [X|Y], Y).
select(X, [F|Y], [F|Z]) :- select(X, Y, Z).
```

For instance, `select(X, [a,b], Z)` has two solutions: (1) `X = a` and `Z = [b]`; (2) `X = b` and `Z = [a]`. An alternative implementation of `select` is to use `append`:

```
select(E, L, R) :- append(X, [E|Y], L), append(X, Y, R).
```

- **delete/3:** `delete(X, List1, List2)` is true when `List1`, with all members matching to `X` are removed, results in `List2`.

```
delete(_, [], []).
delete(X, [X|Y], Z) :- delete(X, Y, Z).
delete(X, [F|Y], [F|Z]) :- X \= F, delete(X, Y, Z).
```

where `\=` means “not unifiable”. For instance, `delete(X, [a,b], Z)` has only one solution: `Z = []`. This shows the difference of `select` and `delete`.

- **reverse/2:** `reverse(List1, List2)` is true when the elements of `List2` are in reverse order compared to `List1`.

```
reverse([], []).
reverse([F|X], Z) :- reverse(X, Y), append(Y, [F], Z).
```

- **permutation/2:** `permutation(List1, List2)` is true when `List2` is a permutation of `List1`.

```
permutation([], []).
permutation([F|X], Y) :- permutation(X, Z), select(X, Y, Z).
```

- **length/2:** `length(List, N)` is true when the length of `List` is `N`.

```
length([], 0).
length(_|X, N) :- length(X, M), N is M+1.
```

One common feature of these predicates is that they are all defined recursively. The termination of these definitions is easy to establish as one of the arguments to the recursive calls goes smaller, or we may show that the head of each rule is strictly greater than its body under a simplification ordering.

Due to the relational nature of many built-in predicates, they can typically be used in several directions. For example, `append/3` can be used both to append two lists (`append(ListA, ListB, X)` given lists `ListA` and `ListB`) as well as to split a given list into parts (`append(X, Y, List)`, given `List`). For instance, the query `append(X, Y, [1,2,3])` will produce the following four answers:

```
X = [], Y = [1,2,3]
X = [1], Y = [2,3]
X = [1,2], Y = [3]
X = [1,2,3], Y = []
```

For this reason, a comparatively small set of library predicates suffices for many Prolog programs. Note that the suggested implementation for `length` can compute the answer `N = 3` from `length([a,b,c], N)`, but cannot generate a list of 3 placeholders, `L = [_1, _2, _3]`, from `length(L, 3)`. Some version of Prolog, e.g., SWI Prolog, can do so.

For some applications, we need to obtain a set  $S = \{x \mid A(x)\}$ , that is,  $S$  is a set of elements satisfying the condition  $A(x)$ . It is not easy to write a Prolog program which generates such a set. Fortunately, in most implementations of Prolog, such a set can be obtained by calling `bagof(X, Goal(X), S)`. For example, in Example 8.1.1, we may define `children(F) = { C | papa(F, C) }`. A query like the following will give us a satisfactory answer:

```
?- bagof(C, papa(Father, C), Children).
Father = bob, Children = [kathleen,joseph,mary] ? ;
Father = joe, Children = [john,bob,ted] ? ;
Father = john, Children = [caroline] ? ;
Father = ted, Children = [patrick]
```

Some Prolog also provides the set operations, such as `intersection/3`, `union/3`, `subtract/3`, `subset/2`, when a list is viewed as a set.

In Prolog, a string is a sequence of characters surrounded by double quotes and is equivalent to a list of numeric character codes. One common mistake is to

use strings in the `write` predicate instead of a single-quoted sequence of characters (an atom).

#### 8.2.4 Sorting Algorithms in Prolog

To illustrate how an algorithm is implemented in Prolog, we show below the Prolog code of insertion sort and quick-sort.

**Example 8.2.1.** The following program defines four predicates for the implementation of insertion sort and quick-sort: `isort/2`, `insert/3`, `split/4`, and `qsort/2`.

```
% isort(L, R): R is the output of L by insert sort
isort([], []).
isort([X|L], R) :-
    isort(L, T),          % sort L into T
    insert(X, T, R).      % insert X into T to obtain R.

% insert(X, L, R): insert X into list L and R is the result.
insert(X, [], [X]).
insert(X, [Y|L], [X, Y|L]) :- X =< Y.
insert(X, [Y|L], [Y|R]) :- X > Y, insert(X, L, R).

% split(P, I, S, B): split the input list I into S and B by pivot P.
split(_, [], [], []).
split(P, [H|I], [H|S], B) :- H =< P, split(P, I, S, B).
split(P, [H|I], S, [H|B]) :- H > P, split(P, I, S, B).

% qsort(L, R): L is the input list and R is the output by quicksort
qsort([], []).
qsort([X|Y], Z) :- split(X, Y, S, B),
    qsort(S, S1), qsort(B, B1), append(S1, [X|B1], Z).
```

The above code is lucid and neat, easy to understand and maintain. This is typical feature of a declarative programming language. Of course, its execution is not as efficient as conventional programming languages. That is why Prolog is used mostly for prototype programming, and for popular functions, we have efficient built-in implementations.

Another feature of declarative programming is that the correctness of the algorithm can be verified. For sorting, the correctness is ensured by the two properties: (1) the output is sorted; (2) the output is a permutation of the input list. In Prolog, (2) is a built-in function and we can define (1) with easy.

**Example 8.2.2.** The following Prolog program defines the predicate `sorted(L)`, which returns true iff  $L$  is a sorted list, and some utility predicates for the verification purpose.

```
sorted([]).
sorted([_]).
sorted([X,Y|L]) :- X<Y, sorted([Y|L]).

verifyIsort(X, Y) :- isort(X, Y), permutation(X, Y), sorted(Y).
verifyQsort(X, Y) :- qsort(X, Y), permutation(X, Y), sorted(Y).
slowsort(X, Y) :- permutation(X, Y), sorted(Y).
```

In the above code, the first three clauses define `sorted/1`. The next two clauses illustrate the verification of `isort` and `qsort`. The last clause shows how the requirement can be the code itself. The code of `permutation/2` can be found in section 8.2.3.

Most Prolog implementations provide two versions of predicates for sorting: `sort/2`, which removes duplicate elements, and `msort/2`, which does not remove duplicates. For instance, the query `sort([b,a,c,a], X)` will produce  $X = [a,b,c]$ , while `msort([b,a,c,a], X)` produces  $X = [a,a,b,c]$ . These built-in predicates run faster than the user defined sorting algorithms.

### 8.3 Recursion in Prolog

As is commonly the case in most programming tasks, we often wish to repeatedly perform some operation either over a whole collection of data, or until a certain point is reached. Conventional programming languages provide various loop controls for performing such tasks. The only way we can do this in Prolog is by recursion. In section 8.2.3, we have seen many Prolog programs which contain recursive definitions. Recursion allows us to write clear and elegant code. Data such formulas or lists are often recursively and best processed recursively. Recursion allows Prolog to perform complex search of a problem space without any dedicated algorithms. Without recursion, Prolog goes nowhere. To master Prolog, we have to learn how to use recursion correctly and efficiently.

Any recursive definition, whether in Prolog or some other language needs two things: *base cases* and *recursive cases*. The base cases specify a condition of when the recursion terminates. Without this the recursion would never stop! For example, the base case of `append(X, Y, Z)` is `append([], Y, Y)` when  $X=[]$ . Base cases almost always come before recursive cases. Recursive cases ask a program to

handle a similar problem of smaller size. For example, `append([U|X], Y, [U|Z]) :- append(X, Y, Z).`

### 8.3.1 Program Termination

The danger of recursion is that the program may loop forever and never terminate. Termination problem is more complicated in Prolog than in conventional programs because the input to a Prolog program may contain variables and we may ask for multiple solutions from a single query.

**Definition 8.3.1.** *A Prolog program  $P$  is said to have first termination on an input  $x$  if  $P$  terminates with either no solution or the first solution on  $x$ .  $P$  is said to have last termination on  $x$  if  $P$  terminates after all solutions of  $x$  are generated. If  $P$  has first (last) termination on any ground input  $x$ , then  $P$  is said to have first (last) ground termination.*

**Example 8.3.2.** Given the following program  $P$ ,

```
ancestor(tom, jack).
ancestor(X, Y) :- ancestor(X, Z), ancestor(Z, Y).
```

$P$  has first termination on the query `?- ancestor(A, B)` with the solution `A=tom, B=jack`, but not on `?- ancestor(jerry, tom)`, as the second clause will be used forever during the search. Since the second input is ground,  $P$  has no first ground termination. Since last termination implies first termination,  $P$  has no last ground termination.

Prolog is a relational programming language, and relations can hold between multiple entities that are reported on backtracking. Therefore, first termination does not fully characterize the procedural behavior of a Prolog query and we need the concept of last termination. We may use the termination of the query `?- Q, false.` to check if  $Q$  has last termination, because all the solutions of  $Q$  must be generated before the query terminates.

Consider the `member/2` program in section 8.2.3, the program has ground termination when the second argument of `member` is ground. However, the query `?- member(X, Y), false.` will fail with a stack overflow error, thus showing that `member(X, Y)` has no last termination.

**Example 8.3.3.** Given a propositional formula  $A$  in negation normal form (NNF), we may use the following program to decide if  $A$  is satisfiable, assuming the formula uses Prolog variables for propositional variables, and `and`, `or`, `not` for  $\wedge$ ,  $\vee$ ,  $\neg$ , respectively.

```

sat(true).                % positive literal
sat(not(false)).         % negative literal
sat(or(X, _)) :- sat(X).
sat(or(_, Y)) :- sat(Y).
sat(and(X, Y)) :- sat(X), sat(Y).

test1(X, Y) :- sat(and(not(X), X)).    % test examples
test2(X, Y) :- sat(and(X, not(Y))).
test3(X, Y) :- sat(and(X, or(not(X), Y))).

```

The *NNF* formulas can be defined recursively by the BNF as follows:

$$\langle NNF \rangle ::= \langle Literal \rangle \mid \langle NNF \rangle \text{ or } \langle NNF \rangle \mid \langle NNF \rangle \text{ and } \langle NNF \rangle$$

The first two clauses are the base cases of `sat` which handle literals, and the next three clauses are the recursive cases which handle `or` and `and`.

It is easy to see that the argument to the recursive call is strictly less than the argument in the head of each rule. Thus, we expect that the program will terminate with the test examples. For the three test examples, the query `?- test1(X, Y)` will say `no`, thus `sat` has last termination on `and(not(X), X)`. For the second query, `?- test2(X, Y)`, Prolog will give the first solution `X=true, Y=false`. The third query `?- test3(X, Y)` will produce the first solution `X=true, Y=true`. Thus, `sat` has first termination on both inputs.

However, if we look for more answers from `?- test2(X, Y)`, we will see the following answers:

```

X = not(false), Y = false
X = or(true,_), Y = false
X = or(not(false),_), Y = false
X = or(or(true,_),_), Y = false
X = or(or(not(false),_),_), Y = false
X = or(or(or(true,_),_),_), Y = false
X = or(or(or(not(false),_),_),_), Y = false
...

```

These answers are correct in Prolog, but not correct for propositional satisfiability (which asks for the models of a formula). If you ask for more answers from `?- test3(X, Y)`, you will see an infinite number of answers, too. Thus, `sat` has last termination neither on `and(X, not(Y))` nor on `and(X, or(not(X), Y))`. Ground termination does not make sense for `sat` because propositional variables are represented by variables.

The problem illustrated by the above example is that variables appearing in the arguments of a recursive program are not bounded and they can be substituted by terms containing other unbound variables, thus creating an infinite number of possibilities. This is the same reason why `member(X, Y)` has no last termination.

### 8.3.2 Focused Recursion

A poorly definition of a recursive predicate may cause the program to loop. In the next chapter, it is shown that checking if a Prolog program terminates on a given input is a undecidable problem. Thus, people are interested in sufficient conditions which ensure the termination of a Prolog program. Focused recursion is one of such conditions.

**Definition 8.3.4.** *A recursive predicate is said to be focused on an input  $t_0$  if there exists a well-founded order  $\succ$  such that  $t_1 \succ t_2 \succ \dots \succ t_n$ , where  $t_i$  are the values of the same argument appearing in the predicate's consecutive recursive calls.*

The argument in the above definition is called the *focused position* of the predicate. Given the definition of `append/3` in section 8.2.3, the predicate `append` is focused on the first argument when the first argument is bound to a list of items (which can be variables) of finite length. Each recursive call reduces the number of items in the list by one. On the other hand, if the first argument is bounded to a variable, the call may go on forever. The same can be said about `member/2` or `sat` in Example 8.3.3, where one sequence of calls of `sat` in the query `test2(X, Y)` may contain the following values as the argument of `sat`:

`and(X, not(Y)), X, or(X1, _), X1, or(X2, _), X2, or(X3, _), X3, ...`

**Example 8.3.5.** In Example 8.3.2, the transitivity of the `ancestor` relation is expressed by the rule `ancestor(X, Y) :- ancestor(X, Z), ancestor(Z, Y)`. This rule cannot be used in Prolog as it causes an infinite loop. To define `ancestor`, we may use the `parent/2` relation, which is defined by a set of ground facts, and can be represented by a directed acyclic graph (DAG) among individuals. Then we can show that `ancestor/2` defined by the following rules is a focused recursion on the first argument.

```
ancestor(X, Y) :- parent(X, Y).
ancestor(X, Y) :- parent(X, Z), ancestor(Z, Y).
```

Let  $X_0, X_1, \dots, X_n$  be the values of the first argument in the recursive calls of `ancestor`. Under the assumption that `parent` defines a DAG among the individuals, we can see that the distance from  $X_i$  to  $Y$  in the DAG is reduced by one

when  $i$  increases by one. Thus, the length of longest paths of the DAG provides an upper bound for the number of recursive calls of `ancestor`. It is easy to see that `ancestor` is the transitive closure of `parent`.

**Proposition 8.3.6.** *If a recursive predicate  $P$  is focused on an input  $t_0$ , then  $P$  has last termination on  $t_0$ .*

### 8.3.3 Tail Recursions

In Prolog, due to the power of logic variables, many predicates can be naturally written in a tail recursive way, where the recursion happens in the last position of the body. For example, the predicates `list`, `member`, `append`, `select`, and `delete`, are tail recursive as given in section 8.2.3. On the other hand, `reverse`, `permutation` and `length` are not tail recursive.

In many cases, tail recursion is good for performance. For conventional programs, a smart compiler can change tail recursion to an interactive one, thus saving the time and space of using a stack to store the environments of a recursive call. In Prolog, the same technique can be used so that a tail call means that the Prolog system can automatically reuse the allocated space of the environment on the local stack. In typical cases, this measurably reduces memory consumption of your programs, from  $O(N)$  in the number of recursive calls to  $O(1)$ . Since decreased memory consumption also reduces the stress on memory allocation and garbage collection, writing tail recursive predicates often improves both space and time efficiency of your Prolog programs.

**Example 8.3.7.** The `reverse/2` predicate defined in section 8.2.3 is copied here:

```
reverse([], []).
reverse([F|X], Z) :- reverse(X, Y), append(Y, [F], Z).
```

Suppose the query is `?- reverse([a, b, c, d], L)`. The subgoal list can become `?- reverse([], []), append([], [d], L3), append(L3, [c], L2), append(L2, [b], L1), append(L1, [a], L)`, before it is shrinking. To avoid this long list of subgoals, we may introduce a tail recursive `rev/3`:

```
reverse(X, Y) :- rev(X, [], Y).
rev([], R, R).
rev([U|X], Y, R) :- rev(X, [U|Y], R).
```

Now, the subgoal list for `?- reverse([a, b, c, d], L)` never contains more than one item at any time:

```

?- reverse([a, b, c, d], L)
?- rev([a, b, c, d], [], L)
?- rev([b, c, d], [a], L)
?- rev([c, d], [b, a], L)
?- rev([d], [c, b, a], L)
?- rev([], [d, c, b, a], L)

```

This example illustrates that tail recursion can save space and time. However, you should not overemphasize it, as for beginners, it is more important to understand termination, and to focus on clear declarative descriptions.

### 8.3.4 A Prolog program for $N$ -Queen Puzzle

**Example 8.3.8.** Below is a simple Prolog program for generating all legal positions for the  $N$ -queens problem. The predicate `queen( $N, L$ )` takes  $N$  as the number of queens and generate a list  $L$  for the row numbers of  $N$  queens in the  $N$  columns.

```

% queen(N, R) generates at first a list R of N variables.
% Assuming one variable per column, then generate values for R,
% which gives the row number of the queen in each column.
queen( N, Res) :-
    length( Res, N),           % Res = N variables for columns
    gen_list( 1, N, L),       % L = [1, 2, ..., N]
    solution( N, Res, L).    % Res is a permutation of L

% gen_list( X, N, R) is true iff R = [X, X+1, X+2, ..., N].
gen_list( X, N, []) :- X>N, !.
gen_list( X, N, [X | Res]) :- Y is X+1, gen_list( Y, N, Res).

% solution( N, RowValues, List) if
%   RowValues represent a list of N non-attacking queens
solution( 0, [], _).        % run out of columns
solution( C, [R|Others], L) :- % 1st queen at position [R,C]
    C1 is C-1,
    select( R, L, L2),      % pick row R for column C
    solution( C1, Others, L2), % find a partial solution
    noattack( R/C, C1, Others). % 1st queen doesn't attack Others

noattack( _, _, []).      % Nothing to attack

```

```

noattack( R/C, C1, [R1 | Others]) :-
    C1-C =\= R1-R,           % Different diagonals
    C1-C =\= R-R1,
    C2 is C1-1,
    noattack( R/C, C2, Others). % no attacks for Others

```

Given the query is `?- queen(4, R).` the first answer is `L = [2,4,1,3]`, which says that the queens are placed at row 2 in the first column, row 4 in the second column, etc., and the second answer is `L = [3,1,4,2]`. In the above program, each position is generated and checked against the positions generated previously, so that no two positions can be attacked each other by the move of queens. Using Prolog's depth-first search, the program is simple to write but highly inefficient. When  $N > 12$ , it will take quite a while to find one solution, must slower than a SAT solver. Sudoku puzzles can also be solved by Prolog using the same logic. However, Sudoku cannot be solved efficiently using this generate-and-test approach, as the search space is huge. SAT solvers are a much better tool for the Sudoku puzzles.

#### 8.4 Beyond Clauses and Logic

As a general purpose language, Prolog provides various built-in predicates to perform routine activities like input/output, using graphics and otherwise communicating with the operating system. These predicates are not given a relational meaning and are only useful for the side-effects they exhibit on the system. For example, the predicate `write/1` displays a term on the screen, and returns true when it succeeds. Besides this kind of extensions, various kinds of extensions are proposed, some intended to improve the expressive power, some speed up the execution. Some of these extensions follow the first-order logic; some go beyond the first-order logic.

For instance, since the occur-check in the unification algorithm is expensive, most prolog implementations choose to ignore the occur-check. This does not cause any problem in most applications. If you run into this problem and need a sound unification algorithm, you may call the built-in function `unify_with_occurs_check`.

```
?- X = f(Y,Y).
```

```
X = f(Y,Y)
```

```
yes
```

```
?- X = f(X,Y).
```

```
cannot display cyclic term for X
```

```
?- unify_with_occurs_check(X, f(X,Y)).
no
```

#### 8.4.1 The Cut Operator !

The cut, in Prolog, is a predicate of arity 0, written as `!`, which always succeeds, but cannot be backtracked. It is best used to prevent unwanted backtracking, including the finding of extra solutions by Prolog and to avoid unnecessary computations.

**Example 8.4.1.** The figure shows how the cut operator affects the search procedure *engine* and *engine1* by a contrived example, where `p` is reduced to a list of subgoals by rule 0, which contains `!`. When `!` is reached by *engine1*, alternative rules, such as rule 1 for `p`, are eliminated; alternative rules for all the left siblings of `!`, such as rule 4 for `a`, are also eliminated. The right siblings of `!` are not affected.

Some programmers call the cut a controversial control facility because it was added for efficiency reasons only and is a function of logic.

A *green cut* is a use of cut which only improves efficiency. Green cuts are used to make programs more efficient without changing program output. For example, in the definition of `split` predicate of `qsort`, we may insert a cut after `H =< P` of the second rule, that is, the three clauses become the following.

```
split(_, [], [], []).
split(P, [H|I], [H|S], B) :- H =< P, !, split(P, I, S, B).
split(P, [H|I], S, [H|B]) :- H > P, split(P, I, S, B).
```

If the current goal  $G$  is `split(3, [2, 4, 1], L, R)`, then  $G$  unifies with the heads of the last two rules. Once the subgoal `H =< P`, i.e., `2 =< 3`, succeeds, the last clause will never be tried due to the cut operator, even if  $G$  unifies with the head of this clause and we want to see more answers. This cut is a green cut, because we know the subgoal `H>P`, i.e., `2 > 3`, will fail and the last clause cannot be applied. The cut operator saves us the time for trying unifying  $G$  with the head of the last clause. With or without the cut, the program performs the same way as desired.

A *red cut* is a use of cut which changes the meaning of the program if the cut is missing. For the above `split` example, if `H =< P` fails, then `H > P` will succeed. To save the cost of `H > P`, we can write the `split` program as the following:

```

split(_, [], [], []).
split(P, [H|I], [H|S], B) :- H =< P, !, split(P, I, S, B).
split(P, [H|I], S, [H|B]) :- split(P, I, S, B).

```

This program will give us the same result as the previous one for any input. However, this program will perform differently from the following program, where the cut is gone:

```

split(_, [], [], []).
split(P, [H|I], [H|S], B) :- H =< P, split(P, I, S, B).
split(P, [H|I], S, [H|B]) :- split(P, I, S, B).

```

So the above cut is a red cut. Red cuts are potential pitfalls for bugs. For example, if the order of the two rules is reversed, the green cut will work correctly and the red cut will produce wrong results. If the saving is significant when using a red cut, a detailed documentation should be provided along the code. The cut should be used sparingly. Proper placement of the cut operator and the order of the rules is required to determine their logical meaning.

#### 8.4.2 Negation as Failure

A Non-Horn clause like  $A \vee B \vee \neg C$  can be represented as  $\neg B \wedge C \rightarrow A$ . Can we write it in Prolog as  $A : -\neg B, C$ ? The answer lies on how to solve the subgoal  $\neg B$  by the goal-reduction process. One approach is called “negation as failure”: Try to solve  $B$ . If  $B$  succeeds, then  $\neg B$  fails; if  $B$  fails, then  $\neg B$  succeeds.

Using cut together with the built-in predicate `fail`, we may define this kind of negation, so that properties that cannot be specified by Horn clauses can be specified.

For example, to express “Mary likes any animals except reptiles” in Prolog,

```

likes(mary, X) :- reptile(X), !, fail.
likes(mary, X) :- animal(X).

```

We can use the same idea of “cut fail” to define the predicate `not`, which takes a term as an argument, and uses the built-in predicate “call”. `not(G)` will “call” the term `G`, evaluates `G` as though `G` as a goal. If `G` succeeds, so is `call(G)`, and `not(G)` will fail. Otherwise, `not(G)` succeeds.

```

not(G) :- call(G), !, fail.
not(_).

```

Most Prolog systems have a built-in predicate like `not`. For instance, SWI-Prolog calls it `\+`. Remember, “not” is not the same as  $\neg$ , because it is based on the success/failure of goals. It can, however, be useful:

```
likes(mary, X) :- not(reptile(X)).
different(X, Y) :- not(X = Y).
```

Negation as failure can be misleading. Suppose the database held the names of members of the public, marked by whether they are innocent or guilty of some offense, and expressed in Prolog:

```
innocent(peter_pan).
innocent(winnie_the_pooh).
innocent(X) :- occupation(X, nun).
guilty(joe_bloggs).
guilty(X) :- occupation(X, thief).
```

```
?- innocent(einstein).
no.
```

If we add one more rule into the above program:

```
guilty(X) :- not(innocent(X)).
```

```
?- guilty(einstein).
yes.
```

It is one thing to show that “Einstein cannot be demonstrated to be innocent”; but it is quite another thing to incorrectly show that he is guilty. To justify the behavior of “negation as failure”, people suggest the *closed world assumption*: Every ground atomic formula is assumed to be false unless it can be shown to be true by the program. Under this assumption, since we cannot show that `innocent(einstein)` is true in the program, so `innocent(einstein)` is assumed to be false, thus it justifies the answer `guilty(einstein)`.

Some disturbing behavior is even more subtle than the innocent/guilty problem, and can lead to some extremely obscure programming errors. Here is a Prolog program about restaurants:

```
good_standard(goedels).
good_standard(hilberts).
```

```

expensive(goedels).
reasonable(R) :- not(expensive(R)).

?- good_standard(X), reasonable(X).
X = hilberts
yes

?- reasonable(X), good_standard(X).
no.

```

Logically, the query `good_standard(X), reasonable(X)` is equivalent to the query `reasonable(X), good_standard(X)`. Why do we get different answers for what seem to be logically equivalent queries? The difference between the questions is as follows. In the first query, the variable `X` is always instantiated when `reasonable(X)` is executed. In the second question, `X` is not instantiated when `reasonable(X)` is executed. The semantics of `reasonable(X)` differ depending on whether its argument is instantiated. To avoid this kind of problems, some implementations of Prolog ask that  $G$  is ground in  $not(G)$ .

Because of the problems associated with “negation as failure”, extra care is needed when using it.

### 8.4.3 Beyond Clauses

Typically, a rule’s body consists of atomic formulas, separated by the commas, and denotes the conjunction of subgoals. If the commas are replaced by the semicolons (`;/2`), then it will denote the disjunction of subgoals. A rule like `H :- (C; D), E` is no longer a clause: It is equivalent to two clauses: `H :- C, E` and `H :- D, E`, because “`;`” is considered to be  $\vee$ . Conjunctions “`,`” and disjunctions “`;`” can only appear in the body, not in the head of a rule.

Prolog also supports the if-then relation by using the infix symbol “`-;/2`”. A rule like `H :- C -> G` has the same affect as `H :- C, G`: if `C` returns true, `G` will be called; if `C` returns false, `C -> G` returns false. When `->` is combined with “`,`”, we then have the if-then-else relation implemented as `H :- C -> G1; G2`: If `C` returns true, `G1` will be called; otherwise, `G2` will be called. Logically, `H :- C -> G1; G2`, which is not a clause, is equivalent to two clauses: `H :- C, G1` and `H :-  $\neg$  C, G2` (the latter is not a Horn clause).

Once a predicate is declared as *dynamic*, Prolog allows the clauses of that predicate can be inserted or removed during the execution. For Example 8.1.1,

if we have a file which contains the pair of father-children names, we can declare `papa` as dynamic: `dynamic(papa)`. During the execution, we read names into `Father` and `Child`, and call `assertz(papa(Father, Child))`, which adds the clause `papa(Father, Child)` at the end of all existing clauses of the predicate `papa`. If we want to add the new clause at the beginning of the existing clauses, the command is `asserta/1`. To remove a clause during the execution, the command is `extract/1`. This set of commands allows the user to dynamically change the program that occurs as the result of executing `assertz/1` or `retract/1`. The change does not affect any activation that is currently being executed. Thus the database is frozen during the execution of a goal, and the list of clauses defining a predicate is fixed at the moment of its execution.

## 8.5 Exercise Problems

1. Given the following Prolog program,

```
append([], Y, Y).
append([H | X], Y, [H | Z]) :- append(X, Y, Z).
```

Find the four solutions for the query `?- append(X, Y, [a, b, c])`. For each of the four solutions, find a resolution proof.

2. For the program in Example 8.3.3, place the clause for handling `and` before those for `or`, and run the query `?- test2(X, Y)` for three answers. Given the resolution proofs for each of the three answers.
3. Implement in Prolog the selection sort, `ssort(X, Y)`, where  $X$  is a list of integers and  $Y$  is the result of sorting  $X$ .
4. Implement in Prolog the merge sort, `msort(X, Y)`, where  $X$  is a list of integers and  $Y$  is the result of sorting  $X$ .
5. The predicate `slowsort` is defined in Example 8.2.2. Construct the resolution proof for the query `?- slowsort([3,2,1], Y)`.
6. Assume `xor/2` is the logical operation  $\oplus$  (exclusive-or). Add the rules for converting the formula containing `xor` into NNF in the Prolog program in Example 8.1.2.
7. Add a predicate `cnf/2` into the Prolog program in Example 8.1.2, where `cnf(F, X)` takes a propositional formula  $F$  as input and converts  $F$  into an equivalent formula  $X$  in CNF.

8. Write a Prolog program which converts a first-order formula  $A$  into NNF (negation normal form), with the following assumptions: The predicate symbols are  $p/2$ ,  $q/1$ , and  $r/0$ . The quantifiers are  $\text{all}(X, A)$  and  $\text{ex}(X, A)$ , where  $X$  takes a value in  $\{x, y, z\}$ .
9. Decide with evidence if the following Prolog programs are terminating or not. Do they have first termination? last termination? ground termination? on what inputs? under what conditions?
- (a) `goHome(no_12).`  
`goHome(X) :- get_next_house(X,Y), goHome(Y).`
- (b) `hasFlu(rebecca).`  
`hasFlu(john).`  
`hasFlu(X) :- kisses(X, Y), hasFlu(Y).`  
`kisses(janet, john).`
- (c) `search(end).`  
`search(X) :- path(X,Y), search(Y).`
10. Introduce a tail recursive predicate for each of the following predicate: (a) `length/1` and (b) `permutation/2`, so that the same result can be obtained through the tail recursive calls.
11. Given the following Prolog program which gives reviews on restaurants:

```
not(P) :- call(P), !, fail.
not(_).
good_standard(goedels).
good_standard(hilberts).
expensive(goedels).
reasonable(R) :- not(expensive(R)).
```

Do the following two queries provide the same result? Why?

```
?- good_standard(X), reasonable(X).
?- reasonable(X), good_standard(X).
```

## CHAPTER 9

### UNDECIDABLE PROBLEMS IN FIRST-ORDER LOGIC

Undecidability concerns about the problems which are beyond the power of computers. To talk about “undecidable problems”, we have to start with Turing machines. Turing machines are important in computability theory due to the *Church-Turing thesis*, which states that every computing model, invented in the past or to be invented in the future, can be simulated by a Turing machine.

Let us quote an account of historic events from the Wikipedia on the “Church-Turing thesis”: In the 1930s, several independent attempts were made to formalize the notion of computability:

- In 1933, Kurt Gödel, with Jacques Herbrand, created a formal definition of a class called *general recursive functions*. The class of general recursive functions is the smallest class of functions (possibly with more than one argument) which includes all constant functions, projections, the successor function, and which is closed under function composition, recursion, and minimization.
- In 1936, Alonzo Church created a method for defining functions called  $\lambda$ -calculus. Within  $\lambda$ -calculus, he defined an encoding of the natural numbers called *Church numerals*. A function on the natural numbers is called  $\lambda$ -computable if the corresponding function on the Church numerals can be represented by a term of the  $\lambda$ -calculus.
- Also in 1936, before learning of Church’s work, Alan Turing created a theoretical model for machines, now called *Turing machines*, that could carry out calculations from inputs by manipulating symbols on a tape. Given a suitable encoding of the natural numbers as sequences of symbols, a function on the natural numbers is called *Turing computable* if some Turing machine computes the corresponding function on the encoded natural numbers.

Church and Turing proved that these three formally defined classes of computable functions coincide: a function is  $\lambda$ -computable iff it is Turing computable, and iff it is general recursive. This has led mathematicians and computer scientists to believe that the concept of computability is accurately characterized by these three equivalent computing models. Other formal attempts to characterize computability have subsequently strengthened this belief. Since we cannot exhaust

all possible computing models, the Church–Turing thesis, although it has near-universal acceptance, cannot be formally proven.

In computability theory, the halting problem is the problem of determining, from a description of an arbitrary computer program and an input, whether the program will finish running, or continue to run forever. Alan Turing proved in 1936 that a general algorithm to solve the halting problem for all possible pairs  $\langle M, w \rangle$ , where  $M$  is a Turing machine and  $w$  is an input to  $M$ , cannot exist. Turing’s proof is one of the first cases of undecidable decision problems to be concluded. The theoretical conclusion that it is not undecidable is significant to practical computing efforts, defining a class of problems which no computing devices can possibly perform perfectly.

The halting problem is not just a theoretic problem. This is a practical problem for many researchers. For example, the unsolved Collatz conjecture is related to whether the program below terminates for all positive integer  $x$ :

```
while  $x > 1$  do
    if  $even(x)$   $x := x/2$  else  $x := 3 * x + 1$ 
```

The same code can be expressed as a recursive function

```
proc  $f(x)$ 
    if  $x = 1$  return 1
    else if  $even(x)$  return  $f(x/2)$ 
    else return  $f(3 * x + 1)$ 
```

or as a rewrite system:

$$\begin{aligned} f(1) &\rightarrow 1 \\ f(2x) &\rightarrow f(x) \\ f(2x + 1) &\rightarrow f(3x + 2) \end{aligned}$$

Does  $f(x) = 1$  for any positive integer  $x$ ? Or will the above rewrite system terminate? The renowned mathematician Paul Erdos (1913-1996) once said about the Collatz conjecture: “Mathematics is not yet ready for such problems.” He offered \$500 for its solution.

## 9.1 Turing Machines

A *Turing machine* (TM) is a mathematical model of computation first proposed by Alan Turing in 1936. A Turing machine manipulates symbols on a tape of symbols according to a set of rules. According to the Church-Turing Thesis, any function on the natural numbers can be calculated by an effective method iff

it is computable by a TM. In other words, despite the model's simplicity, Turing machines are capable of simulating any algorithm on the natural numbers.

A Turing machine operates on a tape containing infinite cells and each cell contains one symbol. The machine positions its "head" over a cell and "reads" the symbol there. Then, as per the symbol and the machine's own present state in a user-specified set of instructions, the machine (*i*) writes a symbol in the cell, (*ii*) either moves the tape head one cell left or right, (*iii*) possibly changes the current state, and (*iv*) goes to the next round, or halts the computation with either "success" or "failure".

Initially, the beginning of the tape is filled with an input string of symbols and is blank everywhere else. The machine positions its head to the first symbol with the designated state called the "initial state". The machine moves as described above through the steps (*i*) – (*iv*). There are three possible outcomes: (1) the machine stops with "success"; (2) it stops with "failure"; or (3) it loops forever.

Turing was able to answer a fundamental question in the negative: Does an algorithm exist that can determine whether any arbitrary TM stops on a given input string? This is so-called the "halting problem" of Turing machines.

A Turing machine can do everything that a real computer can do. However, their minimalist design makes them unsuitable for computation in practice: real-world computers are based on different designs that, unlike Turing machines, use random-access memory.

While a Turing machine can express arbitrary computations, nonetheless, even a Turing machine cannot solve certain problems. In a very real sense, these problems are beyond the theoretical limits of computation. Turing machines proved the existence of fundamental limitations on the power of mechanical computation.

### 9.1.1 Formal Definition of Turing Machines

The TM model uses an infinite tape as its unlimited memory. The symbols appearing in the tape is denoted by  $\Gamma$ , called the *tape alphabet*. By default,  $\Gamma$  is a super set of  $\Sigma$ , which are the symbols used in all the input strings. By default, the blank symbol,  $\sqcup$ , is in  $\Gamma - \Sigma$ .  $\Sigma^*$  denotes the set of all input strings (see Example 1.3.12) and  $\Gamma^*$  denotes the set of all tape contents excluding the trailing blanks.

For any string  $w \in \Sigma^*$ , let  $w^0 = \epsilon$  if  $i = 0$  and  $w^{i+1} = ww^i$ , where  $i \geq 0$ . For example,  $a^3 = aaa$  and  $(ab)^2 = abab$ .

**Example 9.1.1.** Suppose  $\Sigma = \{a, b\}$  and we would like to design a TM which checks whether input string  $w \in L_1 = \{a^j b^k \mid 0 \leq j \leq k\}$ . To design a Turing

machine is the same as to design an algorithm working on the array of symbols. If the machine reads 'a', we may cross it off, which corresponds to  $j := j - 1$ , and then look for the symbol 'b' and cross it off, which corresponds to  $k := k - 1$ . The pseudo code of the algorithm may read as:

```

while ( $j > 0 \wedge k > 0$ ) do
     $j := j - 1; k := k - 1;$ 
if ( $j = 0$ ) return "success" else return "failure";

```

The machine has a tape head that can read and write symbols and move around on the tape. Initially the tape contains only the input string in the beginning portion of the tape and is blank everywhere else. If the machine needs to store information, it may write this information on the tape. To read the information that it has written, the machine can move its head back over it. The machine continues computing until it decides to produce an output. The outputs "success" (or "accept") and "failure" (or "reject") are obtained by entering designated accepting and rejecting states. If it doesn't enter an accepting or a rejecting state, it will go on forever, never halting.

To report "success", the Turing machine uses the state  $q_a$  (or  $q_{accept}$ ); to report "failure", the machine either uses the state  $q_r$  (or  $q_{reject}$ ) or provides no move information in that situation ("reject" by blocking the next move). Since each move is decided by the current symbol and the current state, we may define each move as a function  $\delta$  which takes a state and a symbol as input and outputs (i) the next state, (ii) the symbol to be written on the tape, and (iii) the move direction of the tape head.

**Definition 9.1.2.** A Turing machine is a 6-tuple  $M = (Q, \Sigma, \Gamma, \delta, q_0, q_a)$ , where  $Q$  is a finite set of states;  $\Sigma$  is the alphabet of input strings,  $\Gamma$  is the set of tape symbols,  $\Sigma \subset \Gamma$ ,  $\sqcup \in \Gamma - \Sigma$ , the initial state  $q_0 \in Q$ , and the final state (or the accepting state)  $q_a \in Q$  and the (partial) function  $\delta$ ,

$$\delta : Q \times \Gamma \longrightarrow Q \times \Gamma \times \{L, R\}$$

defines each move of  $M$ , where  $L$  means the tape head moves left and  $R$  means the tape head moves right. If the tape head points to the first symbol on the tape, it cannot move left ("move off the tape" is blocked).

**Example 9.1.3.** Continuing from Example 9.1.1, We would like to construct  $M_1 = (Q, \Sigma, \Gamma, \delta, q_0, q_a)$  to accept the strings in  $L_1 = \{a^j b^k \mid 0 \leq j \leq k\}$ , where  $\Sigma = \{a, b\}$ ,  $\Gamma = \{a, b, x, y, \sqcup\}$ ,  $Q = \{q_0, q_1, q_2, q_3, q_a\}$  and the functions of each state are:

- $q_0$ : at the beginning of the while loop for checking  $j > 0$  or  $k > 0$ .
- $q_1$ : done  $j := j - 1$ , look for  $b$ .
- $q_2$ : done  $k := k - 1$ , rewind to the last  $x$ .
- $q_3$ : found  $j = 0$ , check  $j \geq 0$  and no  $a$ 's after  $b$ .

and  $\delta$  is defined as follows.

1	$\delta(q_0, \sqcup)$	$= (q_a, \sqcup, R)$	$// j = k = 0, \epsilon \in L_1$
2	$\delta(q_0, a)$	$= (q_1, x, R)$	$// j > 0, j := j - 1, \text{goto } q_1$
3	$\delta(q_0, b)$	$= (q_3, b, R)$	$// j = 0, k > 0, \text{goto } q_3$
4	$\delta(q_0, y)$	$= (q_3, y, R)$	$// \text{skip } y, \text{goto } q_3$
5	$\delta(q_1, a)$	$= (q_1, a, R)$	$// \text{skip } a, \text{continue right}$
6	$\delta(q_1, y)$	$= (q_1, y, R)$	$// \text{skip } y, \text{continue right}$
7	$\delta(q_1, b)$	$= (q_2, y, L)$	$// k := k - 1, \text{goto } q_2$
8	$\delta(q_2, y)$	$= (q_2, y, L)$	$// \text{skip } y, \text{continue left}$
9	$\delta(q_2, a)$	$= (q_2, a, L)$	$// \text{skip } a, \text{continue left}$
10	$\delta(q_2, x)$	$= (q_0, x, R)$	$// \text{found } x, \text{goto } q_0$
11	$\delta(q_3, y)$	$= (q_3, y, R)$	$// \text{skip } y, \text{continue right}$
12	$\delta(q_3, b)$	$= (q_3, b, R)$	$// \text{skip } b, \text{continue right}$
13	$\delta(q_3, \sqcup)$	$= (q_a, \sqcup, L)$	$// k \geq 0, b^k \in L_1$

**Definition 9.1.4.** Given an input string  $w$  to a TM  $M = (Q, \Sigma, \Gamma, \delta, q_0, q_a)$ , let  $q \in Q$  be the current state,  $\alpha$  and  $\beta$  be the strings of tape symbols before and after the tape head, respectively. The triple  $\langle \alpha, q, \beta \rangle$ , or simply  $\alpha q \beta$ , is called a configuration of  $M$ . The initial configuration is  $q_0 w$ . A legal move of  $M$  is a pair  $(C_1, C_2)$  of configurations, written as  $C_1 \vdash C_2$ , such that either  $C_1 = \alpha q a \beta$ ,  $\delta(q, a) = (p, b, R)$ , and  $C_2 = \alpha b p \beta$ , or  $C_1 = \alpha c q a \beta$ ,  $\delta(q, a) = (p, b, L)$ , and  $C_2 = \alpha p c b \beta$ .

**Example 9.1.5.** Continuing from the previous example,  $q_0 a a b b \vdash x q_1 a b b \vdash x a q_1 b b \vdash x q_2 a y b \vdash q_2 x a y b \vdash x q_0 a y b \vdash x x q_1 y b \vdash x x y q_1 b \vdash x x q_2 y y \vdash x q_2 x y y \vdash x x q_0 y y \vdash x x y q_3 y \vdash x x y y q_3 \sqcup \vdash x x y q_a y$ . So  $a a b b$  is accepted by  $M_1$ .

On the other hand,  $q_0 a b a \vdash x q_1 b a \vdash q_2 x y a \vdash x q_0 y a \vdash x y q_3 a$ . Since  $\delta(q_3, a)$  is not defined, the move is blocked and the input  $a b a$  is rejected by  $M_1$ .

**Definition 9.1.6.** Let  $M = (Q, \Sigma, \Gamma, \delta, q_0, q_a)$  be a TM, and  $w \in \Sigma^*$ ,  $w$  is said to be accepted by  $M$  if

$$q_0 w \vdash^* \alpha q_a \beta$$

where  $\vdash^*$  denotes the reflexive and transitive closure of  $\vdash$ , and  $\alpha, \beta \in \Gamma^*$ .

The language recognized by  $M$  is the set  $L(M)$ , where

$$L(M) = \{w \mid w \in \Sigma^*, w \text{ is accepted by } M\}.$$

It is easy to see that  $a^2b^2 \in L(M_1)$  and  $aba \notin L(M_1)$ , as shown by Example 9.1.5. It is not an easy task to show that  $M_1$  recognizes  $L_1 = \{a^j b^k \mid 0 \leq j \leq k\}$ , i.e.,  $L(M_1) = L_1$ . We would like to point out that  $M_1$  can be also regarded performing  $k - j$ , as the number of  $b$ 's left on the tape when  $M_1$  accepts is exactly  $k - j$ .

### 9.1.2 High-level Description of Turing Machines

In general, a detailed description of a TM is often tedious. We often prefer to provide a high level description such that every line in the description can be implemented by a finite number of moves of the machine.

**Example 9.1.7.** If we are asked to design a TM  $M_2$  to recognize  $L_2 = \{b^i a^j b^k \mid i * j = k, i, j, k \geq 0\}$ , the basic algorithm is the following.

```

Check if  $i = 0$  or  $j = 0$ , then  $k$  must be 0;
while ( $i > 0$ ) do
     $i := i - 1; k := k - j;$ 
if ( $k = 0$ ) return "success" else return "failure";

```

From Example 9.1.1, we know how to check  $i = 0$ , how to do  $i := i - 1$  and  $k := k - j$ . There is no technical difficulty to design  $M_2$  for  $L_2$ .

$M_2 =$  "On input  $w \in \{a, b\}^*$

1. Scan the tape to check if  $i = 0$  or  $j = 0$ . If yes, *accept* if  $k = 0$ , else *reject*;
2. Check if  $w$  is of form  $b^i a^j b^k$ ; if not, *reject*, else go to the first symbol;
3. Decrease  $i$  by one by crossing off the first  $b$ ;
4. Zigzag between  $a^j$  and  $b^k$  to check if  $k < j$ , if yes, "reject", else  $k := k - j$ ;
5. Go back to the first  $b$  and restore  $a$ ;
6. If  $i > 0$ , goto 3; otherwise check if  $k = 0$ ; if yes, *accept*, else *reject*."

Once we know how to check  $i * j = k$ , we can modify the machine code to compute  $i * j$ ,  $k/i$  (an exercise), or  $j^k$ . We can also compare two numbers, or two strings are identical, or copy one string from one place to another. This kind of operations can be assumed in a high-level description of a TM. In many cases, a high-level description is the only way to design a TM. However, a basic algorithm or procedure, which can be expressed by a high-level description, is a must in the design of some Turing machines.

### 9.1.3 Recognizable vs Decidable

Given any TM  $M = (Q, \Sigma, \Gamma, \delta, q_0, q_a)$ , if we feed  $w \in \Sigma^*$  to  $M$ , there will be three possible outcomes: (1)  $M$  halts with “accept”; (2)  $M$  halts with “reject”; and (3)  $M$  loops on  $w$ . In other words,  $M$  partitions  $\Sigma^*$  into three disjoint sets:

- $L(M) = \{w \mid w \in \Sigma^*, M \text{ halts and accepts } w\}$ ,
- $reject(M) = \{w \mid w \in \Sigma^*, M \text{ halts and rejects } w\}$ , and
- $loop(M) = \{w \mid w \in \Sigma^*, M \text{ loops on } w\}$ .

Obviously,  $L(M) \cup reject(M) \cup loop(M) = \Sigma^*$ . By convention,  $\overline{L(M)} = \Sigma^* - L(M) = reject(M) \cup loop(M)$ .

In computability theory, any subset  $L$  of  $\Sigma^*$  is said to be a *formal language*.  $L$  can be  $\emptyset$  or  $\Sigma^*$ . Turing machines can be used to classify formal languages.

**Definition 9.1.8. (recognizable and decidable)** *A formal language  $L \subseteq \Sigma^*$  is said to be recognizable, if there exists a TM  $M$  such that  $L = L(M)$ , and we say  $M$  is a recognizer of  $L$ .*

*$L$  is said to be decidable if there exists a TM  $M$  such that  $L = L(M)$  and  $loop(M) = \emptyset$ , and we say  $M$  is a decider of  $L$ .*

If  $L$  is decidable and  $L(M) = L$  for some TM  $M$ , it does not mean that  $loop(M) = \emptyset$ . The definition just ensures that there exists a decider  $M'$  such that  $L = L(M')$  and  $loop(M') = \emptyset$ .

## 9.2 Decidability of Problems

In computability theory, a problem is typically a decision problem which has a simple yes/no answer. Optimization problems like “finding the longest simple cycle in graph  $G$ ” is always converted to a question like “does  $G$  have a simple cycle of length at least  $m$ ?” by introducing a number  $m$ , so that it becomes a decision problem.

### 9.2.1 Encoding of Decision Problems

Another common practice in computability theory is to encode every decision problem by a set of strings.

For instance, the problem that “does  $G$  have a simple cycle of length at least  $m$ ?” can be encoded as a set  $L_1$ :

$$L_1 = \{\langle G, m \rangle \mid G \text{ has a simple cycle of length at least } m\},$$

where  $\langle G, m \rangle$  is a string which represents  $G$  and  $m$ ; we may think  $\langle G, m \rangle = \langle G \rangle \langle m \rangle$ . If  $G$  contains  $n$  vertices and  $m \leq n$ , then the alphabet for  $\langle G, m \rangle$  can be  $\{1, 2, \dots, n, (, ), \cdot\}$ , where each vertex is represented by a number  $i$ ,  $1 \leq i \leq n$ , and each edge is represented by  $(i \cdot j)$ ,  $1 \leq i, j \leq n$ . The alphabet can also be  $\{0, 1\}$ , if every symbol in  $\{1, 2, \dots, n, (, ), \cdot\}$  is encoded by a binary string of equal length, as we use the ASCII code in a computer. Now, a graph  $G$  has a simple cycle of length at least  $m$  iff  $\langle G, m \rangle \in L_1$ .  $L_1$  catches the essence of this decision problem.

“Does there exist an algorithm  $A$  which finds a simple cycle in  $G$  of length at least  $m$ ?” This is also a decision problem and can be encoded as

$$L_2 = \{\langle A, G, m \rangle \mid \text{Algorithm } A \text{ finds in } G \text{ a simple cycle of length at least } m\}$$

where  $\langle A \rangle$  can be the code file for implementing  $A$  in a programming language.

Turing machines can also be encoded as a string. For instance, the definition of  $M_1$  in Example 9.1.1 takes less than 20 lines of this book and each line has less than 100 characters, so a string of 2000 characters will be more than enough to encode  $M_1$ .

There are many decision problems about Turing machines and the following are some of them.

- $A_{TM} = \{\langle M, w \rangle \mid \text{TM } M \text{ accepts input } w\}$ .
- $H_{TM} = \{\langle M, w \rangle \mid \text{TM } M \text{ halts on input } w\}$ .
- $E_{TM} = \{\langle M \rangle \mid L(M) \text{ is empty}\}$ .
- $One_{TM} = \{\langle M \rangle \mid |L(M)| = 1\}$ .
- $Fin_{TM} = \{\langle M \rangle \mid L(M) \text{ is finite}\}$ .
- $S10_{TM} = \{\langle M \rangle \mid M \text{ has at least 10 states}\}$ .
- $A10_{TM} = \{\langle M \rangle \mid M \text{ accepts at least 10 strings}\}$ .
- $B10_{TM} = \{\langle M \rangle \mid M \text{ moves at least 10 steps on any input}\}$ .

The language  $A_{TM}$  encodes the *acceptance problem* of Turing machines: Given any TM  $M$  and any input  $w$ , does  $M$  accept  $w$ ? This is one of the first problems shown to be undecidable, that is, there exists no algorithm to solve the acceptance problem.

The language  $H_{TM}$  encodes the well-known halting problem of Turing machines. Alan Turing proved in 1936 that a general algorithm to solve the halting

problem for all possible  $\langle M, w \rangle$  pairs cannot exist. Later, after giving formally the definition of “undecidable”, we will show that both  $A_{TM}$  and  $H_{TM}$  are Turing recognizable, but not decidable.

The language  $E_{TM}$  encodes the *emptiness problem* of Turing machines. While both  $A_{TM}$  and  $H_{TM}$  are undecidable but Turing recognizable,  $E_{TM}$  is even not Turing recognizable.

The rest languages in the above list encode some problems less popular. You may create your own languages for any decision problem.

Let the encodings of all TMs use a common alphabet, say  $\Sigma$ , and  $L_{TM} = \{\langle M \rangle \mid M \text{ is a TM}\}$ .

**Proposition 9.2.1.**  *$L_{TM}$  is countable.*

*Proof.* Since  $\langle M \rangle \in \Sigma^*$  for every TM  $M$ ,  $L_{TM} \subseteq \Sigma^*$ . According to Lemma 1.3.4,  $\Sigma^*$  is countable and any subset of a countable set is also countable.  $\square$

The above proposition claims that the set of possible TMs is countable.

### 9.2.2 Decidable Problems

Since every decision problem can be encoded as a formal language  $L \subseteq \Sigma^*$ , we wish to design a TM  $M$  which can answer the question of  $w \in L$  or not, for any  $w \in \Sigma^*$ . This is only possible when  $L$  is decidable; if  $L$  is not recognizable, then no such TMs can exist.

**Definition 9.2.2. (decidable and computable problem)** *A decision problem is said to be decidable if the language of its encoding is decidable. A decision problem is said to be computable if its encoding is recognizable.*

All the decision problems we have met in the courses on data structures and algorithms, where their complexity is obtained by a big  $O$  function, belong to the class of decidable problems. In computability theory, a decider for a decision problem is also call a *decision procedure*, or an *algorithm*, which coincides with the conventional definition of *algorithm*: a set of instructions which will be executed in a finite number of steps for any input.

Some problems are computable but undecidable; some are not computable. Any Boolean function  $f : N \rightarrow \{0, 1\}$ , over the natural numbers is a decision problem about the natural numbers: Does  $f(n) = 1$  for any  $n \in N$ ? According to Lemma 1.3.5, the set of all Boolean functions over the natural numbers is not countable. That is, there are more Boolean functions than the available TMs, which are countable. For any function  $f : N \rightarrow \{0, 1\}$ , it is not always possible to have a

TM  $M_f$  such that  $M$  accepts  $i$  iff  $f(i) = 1$  for any natural number  $i$ . There exist many, many functions which cannot be computed by any computing device in view of the Church-Turing thesis.

Here, we are in particular interested in decision problems about TMs.

**Example 9.2.3.** Let  $S10_{TM} = \{\langle M \rangle \mid M \text{ has at least 10 states}\}$ . Then  $S10_{TM}$  is decidable because we can design a TM which counts how many states in  $\langle M \rangle$ . It takes no more than the linear time to count and if the count goes beyond 10, stop with “accept”, otherwise, stop with “reject”.

**Example 9.2.4.** Let  $B10_{TM} = \{\langle M \rangle \mid M \text{ moves at least 10 steps on any input}\}$ . Let  $M = (Q, \Sigma, \Gamma, \delta, q_0, q_a)$  and  $M_B$  be the decider that we will construct for  $B10_{TM}$ .

The basic algorithm of  $M_B$  works as follows: The tape of  $M_B$  is divided into three sections separated by the special symbol  $\#$  and  $\$$ . Section 1 stores the input string  $\langle M \rangle$  and will never change.

Section 2 has 11 cells and takes turns to store all the input strings  $w \in \Sigma^*$  for  $M$ ,  $|w| \leq 10$ . The strings in  $\Sigma^*$  are generated by the following order: (i) shorter strings first; (ii) for the strings of the same length, smaller string under the lexicographic order first. Thus, if  $\Sigma = \{0, 1\}$ , then the strings are in this order:  $\epsilon, 0, 1, 00, 01, 10, 11, 000, \dots$  We say the *successor* of  $\epsilon$  is 0, the *successor* of 0 is 1, and so on.

Section 3 will be the space for simulating the moves of  $M$ .

$M_B =$  ”On input  $\langle M \rangle$ ,  $M$  is a TM

1. Place  $\#$  at the end of  $\langle M \rangle$ ; move right 11 steps and place  $\$$ .  
Treat the string between  $\#$  and  $\$$  as  $w$ , which is empty initially.
2. Write  $q_0$  of  $M$  after  $\$$  in section 3.
3. Copy string  $w$  in section 2 to section 3, after  $q_0$ .  
Now section 3 contains  $q_0w$ , an initial configuration of  $M$ .
4. Zigzag between section 1 and section 3, do the following:
  - 4.1 Using the states of  $M_B$  to remember the pair  $(q, a)$  in section 3;
  - 4.2 Find the delta definition for  $(q, a)$  in section 1;
  - 4.3 If no  $\delta(q, a)$  is found, stop with ”reject”;
  - 4.4 If found, use the states to remember  $\delta(q, a) = (p, b, D)$ ,  $D \in \{L, R\}$ ;
  - 4.5 Move to section 3, and update the configuration of  $M$  according to  $\delta(q, a) = (p, b, D)$ ;
  - 4.6 If 4.5 has been done less than 10 times, go to 4.1;
  - 4.7 else erase the content of section 3 and go to 5;
5. Go to section 2 and generate the successor of  $w$ , replacing  $w$ ;
6. If  $w$  exceeds 10 symbols, stop with “accept”, else go to 2. ”

Lines 4.1-4.7 simulates the moves of  $M$  for 10 steps. If this is possible,  $M_B$  continues for the successor of  $w$  (4.7); if the move of  $M$  is blocked before moving 10 steps, it stops with “reject” (4.3), because we have found an input string for  $M$  to move less than 10 steps. In line 5, we generate the successor of the current  $w$  and replace  $w$  by it. If all the strings of length 10 or less are generated and  $M$  moves over 10 steps on each of them, we stop with “accept”, because to read the 11<sup>th</sup> symbol,  $M$  has to move at least 10 steps. Thus, we do not need to consider the strings of length more than 10 for  $w$ .

Lines 4.1-4.5 describe how a move of  $M$  is simulated by  $M_B$ . Later, when we say “ $M'$  simulates  $M$ ”, we refer to this example for the details of the simulation. Note that the states of a TM can remember only a finite amount of information. For  $M_B$ , one definition of  $\delta(q, a) = (p, b, D)$  contains only five symbols. You cannot use the states to restore the information like  $a^i$ , where  $i$  is any number.

### 9.2.3 Undecidable Problems

**Theorem 9.2.5.** *The acceptance problem of Turing machines is undecidable. That is,  $A_{TM} = \{\langle M, w \rangle \mid \text{TM } M \text{ accepts input } w\}$  is undecidable.*

*Proof.* Assume that  $A_{TM}$  is decidable, then there exists a hypothetical decider  $H$  for  $A_{TM}$ . Since  $H$  is a decider,  $H$  either accepts  $x$  or rejects  $x$  for any input  $x$ , and  $H$  accepts  $\langle M, w \rangle$  iff  $M$  accepts  $w$ .

Using  $H$ , we construct another decider  $D$ , which checks if a TM  $M$  does not accept its own code:

$D =$  ”On input  $\langle M \rangle$ ,  $M$  is a TM,  
 1. Simulate  $H$  on  $\langle M, \langle M \rangle \rangle$  until  $H$  halts.  
 2. If  $H$  accepts, *reject*; if  $H$  rejects, *accept*.”

Obviously,  $D$  is a TM and  $L(D) = \{\langle M \rangle \mid M \text{ is a TM, } M \text{ rejects } \langle M \rangle\}$  is decidable, because  $D$  halts on any input.

Let us check if  $\langle D \rangle \in L(D)$  and there are only two cases:

- $\langle D \rangle \in L(D)$ : Every member  $\langle M \rangle$  of  $L(D)$ , including  $\langle D \rangle$ , satisfies the definition of  $L(D)$ : “ $M$  rejects  $\langle M \rangle$ ”. Replace  $M$  by  $D$ , we have “ $D$  rejects  $\langle D \rangle$ ”, which means  $\langle D \rangle \notin L(D)$ .
- $\langle D \rangle \notin L(D)$ : That means  $D$  rejects  $\langle D \rangle$ . By the definition of  $L(D)$ ,  $\langle D \rangle \in L(D)$  should be true.

In both cases, we have a contradiction, so  $D$  cannot exist.  $D$  is constructed from  $H$ , hence,  $H$  cannot exist.  $H$  comes from the assumption that  $A_{TM}$  is decidable. Thus  $A_{TM}$  is undecidable.  $\square$

The above proof is analogous to that of “Barber’s paradox” (see Paradox 1.3.8). The TM  $D$ , which accepts the encodings of those TMs that do not accept their own encoding, corresponds to “the Barber who shaves all those who do not shave themselves”. Both proofs originated from Cantor’s diagonalization method. To see this connection, the reader is recommended to Michael Sipser’s textbook on theory of computation.

### 9.2.3.1 Reduction

Once we have found one undecidable language, other undecidable languages can be found through reduction, as illustrated by the proof of the halting problem of Turing machines.

**Theorem 9.2.6.** *The halting problem of Turing machines is undecidable. That is,  $H_{TM} = \{\langle M, w \rangle \mid \text{TM } M \text{ halts on input } w\}$  is undecidable.*

*Proof.* Assume  $H_{TM}$  is decidable, then there exists a hypothetical decider  $H$  for  $H_{TM}$ . Since  $H$  is a decider,  $H$  will halt on any input, and  $H$  accepts  $\langle M, w \rangle$  iff  $M$  halts on  $w$ .

Using  $H$ , we can construct another decider  $S$  for  $A_{TM}$  as follows:

$S =$  ”On input  $\langle M, w \rangle$ ,  $M$  is a TM and  $w$  is an input to  $M$ ,

1. Simulate  $H$  on  $\langle M, w \rangle$  until  $H$  halts.
2. If  $H$  rejects, *reject*.
3. Simulate  $M$  on  $w$  until  $M$  halts.
4. If  $M$  accepts, *accept*; if  $M$  rejects, *reject*.”

At line 2, when  $H$  rejects  $\langle M, w \rangle$ , it means  $M$  loops on  $w$ , so  $w$  cannot be accepted by  $M$ . This justifies the rejection of  $w$  by  $S$ . At line 3, we know  $M$  will halt  $w$ . When  $M$  halts on  $w$ ,  $S$  returns the output of  $M$  at line 4. All four lines of  $S$  will be finished in a finite number of steps, so  $S$  is a decider for  $A_{TM}$ . Since  $S$  cannot exist,  $H$  cannot exist, either. Thus,  $H_{TM}$  must be undecidable.  $\square$

The above two theorems are what we needed to show the undecidability of some logical problems.

Reduction is a powerful tool to show that a problem  $X$  is undecidable: Assume  $X$  is decidable and  $A$  is a decider for  $X$ . We use  $A$  to construct another decider  $B$

for a known undecidable problem  $Y$ . Since  $B$  cannot exist, so  $A$  cannot exist and  $X$  is thus undecidable. In the proof of the above theorem,  $X$  is  $H_{TM}$ ,  $A$  is  $H$ ,  $Y$  is  $A_{TM}$  and  $B$  is  $S$ . We showed that if there exists a decider for  $H_{TM}$ , then there exists a decider for  $A_{TM}$ .

We may switch the role of  $H_{TM}$  and  $A_{TM}$  by showing that if there exists a decider for  $A_{TM}$ , then there exists a decider for  $H_{TM}$ . Suppose  $A_{TM}$  is decidable and  $A$  is its decider. We may construct a decider  $B$  for  $H_{TM}$  as follows:

$B =$  "On input  $\langle M, w \rangle$ ,  $M$  is a TM and  $w$  is an input to  $M$ ,

1. Let  $\langle M' \rangle$  be a copy of  $\langle M \rangle$  and we modify  $\langle M' \rangle$  as follows:  
     For any state  $q$  and tape symbol  $x$  of  $M$ , if  $q$  is not  $q_a$   
     (the final state), and  $\delta(q, x)$  is not defined, add in  $M'$ :  
      $\delta'(q, x) = (q_a, x, R)$ .
2. Simulate  $A$  on  $\langle M', w \rangle$  until  $A$  halts.
3. If  $A$  accepts, *accept*; if  $A$  rejects, *reject*."

It is easy to check when  $M$  rejects  $w$ ,  $M'$  will accept  $w$ ; when  $M$  accepts  $w$  or loops on  $w$ , so does  $M'$ . Thus,  $loop(M) = loop(M')$ ,  $reject(M') = \emptyset$ , and  $L(M') = L(M) \cup reject(M)$ .  $B$  accepts  $\langle M, w \rangle$  iff  $A$  accepts  $\langle M', w \rangle$ , iff  $M'$  accepts  $w$ , iff  $M$  accepts  $w$  or  $M$  rejects  $w$ . In other words,  $B$  accepts  $\langle M, w \rangle$  iff  $M$  halts on  $w$ ;  $B$  rejects  $\langle M', w \rangle$  iff  $M$  loops on  $w$ . This justifies the output of  $B$  at line 3. Since all the three lines of  $B$  can be done in a finite number of steps,  $B$  is a decider for  $H_{TM}$ .

We have just shown that both  $A_{TM}$  and  $H_{TM}$  are undecidable. Are they recognizable? The answer is yes for both of them.

$U =$  "On input  $\langle M, w \rangle$ ,  $M$  is a TM and  $w$  is an input to  $M$ ,

1. Simulate  $M$  on  $w$ .
2. When  $M$  halts on  $w$ , *accept* if  $M$  accepts; *reject* if  $M$  rejects."

It is clear that  $U$  will accept  $\langle M, w \rangle$  if  $M$  accepts  $w$ ;  $U$  will loop if  $M$  loops on  $w$ . Thus  $L(U) = A_{TM}$ ;  $U$  is called the *universal Turing machine*, because it simulates any Turing machine.

Similarly,  $H_{TM}$  can be recognized by a Turing machine.

$T =$  "On input  $\langle M, w \rangle$ ,  $M$  is a TM and  $w$  is an input to  $M$ ,

1. Simulate  $M$  on  $w$ .
2. When  $M$  halts on  $w$ , *accept*."

$T$  accepts iff  $M$  halts on  $w$ , thus  $L(T) = H_{TM}$ .

**Proposition 9.2.7.** *If a formal language  $L$  is recognizable but undecidable, then  $\overline{L}$  is not recognizable.*

*Proof.* (sketch) Let  $M_1$  be the recognizer of  $L$ . If  $\overline{L}$  is also recognizable, let  $M_2$  be its recognizer. We may use  $M_1$  and  $M_2$  to create a decider for  $L$ :

$T =$  “On input  $\langle M, w \rangle$ ,  $M$  is a TM and  $w$  is an input to  $M$ ,

1. Simulate  $M_1$  on  $w$  for one move;
2. If  $M_1$  accepts, *accept*;
3. Simulate  $M_2$  on  $w$  for one move;
4. If  $M_2$  accepts, *reject*;
5. go to 1. and continue the simulations. ”

For any  $w \in \Sigma^*$ , either  $w \in L$  or  $w \in \overline{L}$ . So either  $M_1$  accepts  $w$  or  $M_2$  accepts  $w$ .  $T$  will halt when one of them accepts. If  $w \in L$ ,  $T$  will accept; otherwise,  $T$  will reject. Thus,  $T$  is a decider for  $L$ .  $\square$

Applying the above proposition to  $A_{TM}$  and  $H_{TM}$ , we now know that  $\overline{A_{TM}}$  and  $\overline{H_{TM}}$  are not recognizable.

### 9.3 Turing Completeness

By the Church-Turing thesis, any computing model can be simulated by a Turing machine. On the other hand, it becomes easy to design a new computing model which has the same power as a Turing machine, as long as the new computing model is shown to do what a Turing machine can do. The concept of *Turing completeness* just expresses this property: If a computing model has the ability to simulate a Turing machine, it is *Turing complete*. Both Gödel’s recursive functions and Church’s  $\lambda$ -calculus are Turing complete.

**Example 9.3.1.** A two-stack machine is a variation of Turing machine in that the tape is replaced by two stacks: The first stack stores all the symbols from the tape head (inclusive) to right; the second stack stores all the symbols from the tape head (exclusive) to left. The state information remains the same. It is easy to show that the two-stack is Turing complete by simulating all the Turing moves by the operations on stacks. Initially, the input string is stored in the first stack with the first symbol on the top of the stack. The simulation of a Turing move goes as follows: A symbol  $a$  is popped off from the first stack; if the stack is empty, use the blank symbol. If  $\delta(q, a) = (p, b, R)$ , then push  $b$  into the second stack; if  $\delta(q, a) = (p, b, L)$ , then push  $b$  into the first stack, and also pop the top symbol from the second stack and push it into the first stack.

The significance of being Turing complete is that all the undecidable properties of Turing machines will be inherited by the new computing model. Thus, the halting problem of a two-stack machine is undecidable, because if we have an algorithm to tell the termination of a two-stack machine, this algorithm can be used to show the termination of a Turing machine. Similarly, if we can show the Turing completeness of a logic, we can derive the same result.

### 9.3.1 Turing Completeness of Prolog

**Proposition 9.3.2.** *Prolog is Turing complete.*

*Proof.* The proposition can be shown by using Prolog to simulate a Turing machine, using the idea of two-stack machines. The content of the Turing machine tape is split into two lists,  $L$  and  $R$ , which correspond to the two stacks in a two-stack machine. Initially,  $L = []$  and  $R$  contains the input string. We transform  $L$  to  $L'$  and  $R$  to  $R'$  according to the Turing machine, until we reach the final state. The transformation is illustrated by a Prolog program for the Turing machine in Example 9.1.1

```
% M1 checks  $a^j b^k$  for  $j \leq k$ .
delta(q0, b, qa, bl, right).
delta(q0, a, q1, x, right).
delta(q0, b, q3, b, right).
delta(q0, y, q3, y, right).
delta(q1, a, q1, a, right).
delta(q1, y, q1, y, right).
delta(q1, b, q2, y, left).
delta(q2, y, q2, y, left).
delta(q2, a, q2, a, left).
delta(q2, x, q0, x, right).
delta(q3, b, q3, b, right).
delta(q3, y, q3, y, right).
delta(q3, bl, qa, bl, left).

% General code for all Turing machines
turing(W, Tape) :-
    % W: input string, Tape: at the end
    write(' '), write(q0), write(W), nl, % initial configuration
    move([], L, q0, W, R), % ([], q0, W)  $\vdash^*$  (rev(L), qa, R)
    reverse(L, Lr),
    append(Lr, R, Tape).
```

```

move(L, L, qa, R, R) :- !.                % qa is the final state
move(L0, L, Q, R0, R) :-
    current(R0, A, Rest),                 % load the current symbol in A
    delta(Q, A, P, B, Direct),           % delta(Q, A) = (P, B, Direct)
    result(Direct, L0, L1, B, Rest, R1),  % update L and R
    write('|- '), write(rev(L1)), write(P), write(R1), nl,
    move(L1, L, P, R1, R).                % continue to move

current([], b1, []).                      % current symbol is a blank b1
current([Sym|Rest], Sym, Rest).           % current symbol is Sym

result(left, [C|L], L, B, R, [C,B|R]).    % move top symbol of L to R
result(right, L, [B|L], B, R, R).         % push B at the top of L

?- turing([a,a,b,b], Tape).
    q0[a,a,b,b]
|- rev([x])q1[a,b,b]
|- rev([a,x])q1[b,b]
|- rev([x])q2[a,y,b]
|- rev([])q2[x,a,y,b]
|- rev([x])q0[a,y,b]
|- rev([x,x])q1[y,b]
|- rev([y,x,x])q1[b]
|- rev([x,x])q2[y,y]
|- rev([x])q2[x,y,y]
|- rev([x,x])q0[y,y]
|- rev([y,x,x])q3[y]
|- rev([y,y,x,x])q3[]
|- rev([y,x,x])qa[y,b1]

```

This illustrates how any computation can be expressed declaratively as a sequence of Turing machine moves, implemented in Prolog as a relation between successive states of interest. The Prolog code may be used as a debugging tool for designing Turing machines.  $\square$

**Theorem 9.3.3.** *The termination of a Prolog program is undecidable.*

*Proof.* (sketch) If we have an algorithm to decide the termination of a Prolog program on any input, then we can use the algorithm to decide the termination of a

Turing machine on an input, as the moves of the Turing machine can be simulated by a Prolog program. The termination of the Prolog program in the simulation ensures the termination of the Turing machine.  $\square$

**Theorem 9.3.4.** *The satisfiability of a set of Horn clauses is undecidable.*

*Proof.* (sketch) If we have an algorithm to decide the satisfiability of a Horn clauses, then the algorithm can be used to decide the clauses in the Prolog program for simulating a Turing machine, together with the query derived from  $q_0w$ , is satisfiable or not. If it is unsatisfiable, then there exists a resolution proof derived from the query as this is the only negative clause (the clauses without this one is satisfiable). The resolution proof ensures that  $w$  is accepted by the Turing machine. Thus, the algorithm can tell if the Turing machine can accept  $w$  or not. This is a contradiction to the fact that the acceptance problem of Turing machines is undecidable.  $\square$

**Theorem 9.3.5.** *The satisfiability of a first-order formula is undecidable.*

The proof is left as an exercise.

### 9.3.2 Turing Completeness of Rewrite Systems

In theoretical computer science and mathematical logic, a *string rewrite system* (SRS), historically called a *semi-Thue system*, is a rewrite system over strings from an alphabet. Given a binary relation  $\rightarrow$  between fixed strings over the alphabet, called *rewrite rules*, denoted by  $s \rightarrow t$ , an SRS extends the rewriting relation to all strings in which the left- and right-hand side of the rules appear as substrings, that is  $usv \rightarrow utv$ , where  $s, t, u$ , and  $v$  are strings.

**Proposition 9.3.6.** *SRS is Turing complete.*

*Proof.* Any configuration of a Turing machine is a string of tape and state symbols, assuming the states and the tape do not share symbols. We may use an SRS to describe all the moves of a Turing machine as follows: If  $\delta(q, a) = (p, b, R)$ , we have a single rewrite rule  $qa \rightarrow bp$ ; if  $\delta(q, a) = (p, b, L)$ , we have a set of rewrite rules:  $cqa \rightarrow pcb$ , where  $c$  is any tape symbol. Then it is trivial to show that  $q_0w \vdash^* \alpha q_a \beta$  iff  $q_0w \Rightarrow^* \alpha q_a \beta$ .  $\square$

A term rewrite system (TRS) is a binary relation over terms. TRS are more expressive than SRS as an SRS is easily converted into a TRS, where each symbol  $f$  in the SRS is replaced by a unary function  $f/1$  in the TRS. For example, a rewrite rule  $abc \rightarrow cd$  in SRS is replaced by  $a(b(c(x))) \rightarrow c(d(x))$ . Since SRS is Turing complete, so are the term rewrite systems.

**Theorem 9.3.7.** *The termination of a rewrite system is undecidable.*

*Proof.* (sketch) If we have an algorithm to decide the termination of a rewrite system (either SRS or TRS), then we can use the algorithm to decide the termination of a Turing machine as the moves of the Turing machine can be expressed as a set of rewrite rules. The termination of the rewrite system ensures the termination of the Turing machine.  $\square$

#### 9.4 Exercise Problems

1. Some Turing machines allow the tape head stays after a move (only changing state or the current symbol). How to simulate this action in the Prolog program of section 9.3.1?
2. Some Turing machines allow the tape has infinite symbols on the left side of the head tape (so there is no “move off the tape”). How to simulate this action in the Prolog program of section 9.3.1?
3. Provide a complete design of Turing machine  $M_2$  for recognizing  $L_2 = \{b^i a^j b^k \mid i * j = k, i, j \geq 0\}$ . Test your code in the Prolog program on the input string  $w = b^2 a^1 b^2$ .
4. Provide a high-level description of a Turing machine  $M_3$  for recognizing  $L_3 = \{a^i b^j c^k \mid k = \lfloor i/j \rfloor, i, j \geq 0\}$ .
5. Prove that the satisfiability of a first-order formula is undecidable.

## CHAPTER 10

### HOARE LOGIC

One of the most important research objectives in mathematics and computer science is to obtain means of reasoning in and about complex systems. Proving properties of such systems is extremely important. In safety-critical systems even small mistakes can provoke disasters. Since the amount of data which has to be handled in all the application domains is usually huge, computer support is necessary. Many approaches have been proposed to formally verify the properties of software and hardware systems. For software formal verification, Hoare logic is one of the most influential approaches for proving the correctness of software designs in various programming languages.

*Hoare logic* (also known as *Floyd–Hoare logic* or *Hoare rules*) is a formal system with a set of logical rules for reasoning rigorously about the correctness of imperative programs, where a program describes steps that change the state of the computer. It was proposed in 1969 by the British computer scientist and logician Tony Hoare, and subsequently refined by Hoare and other researchers. The original ideas were seeded by the work of Robert Floyd, who had published a similar system for flowcharts.

We will use the word “state” to denote a valuation of all variables in an imperative program. A command in the program may change the value of variables, thus causing a transition of states. For example, suppose  $x$  holds the value 1 before the command  $x := x + 1$ , the value of  $x$  will be 2 after the execution of  $x := x + 1$ . A state can be easily specified by a first-order formula, called *assertions*, such as  $x = 1$ ,  $A[0] = 2$ , etc.

Executing an imperative program has the effect of changing the state. To use such a program, one first establishes an initial state by specifying the initial values of variables. One then executes the program. This transforms the initial state into a final one through a sequence of state transitions. One then inspects the values of variables in the final state to get the desired results. For example, to compute  $z := exp(x, y)$ , the initial state specifies the value of  $x$  and  $y$ , and the final state specifies the value of  $z$  satisfying  $z = x^y$  to judge that the program  $exp(x, y)$  is correct.

Hoare logic is an axiomatic approach to formal verification, which consists of (1) a first-order logic language for making assertions about programs; and (2) rules

for establishing state transitions, that is, what assertions hold after a transition of states. From the requirement on the input, we create a precondition of a program, and from the requirement on the output, we create a postcondition of the program. Both the precondition and the postcondition are assertions, i.e., formulas of the first-order logic. We then use a theorem prover to prove that the postcondition is a logical consequence of a formula derived from the precondition and the program code by the rules of Hoare logic. Other approaches to formal verification include model checking and abstract interpretation, which are directed at finding specific errors, such as the possible violation of a safety property.

### 10.1 Hoare Triples

The central idea of Hoare logic is *Hoare triples*, a notation proposed first by Hoare, which specify exactly what precondition that a program asks before the execution and what postcondition, i.e., the desired behaviour of the program, that the program will achieve. A Hoare triple describes how the execution of a piece of code changes the state of the computation. A *Hoare triple* is of the form

$$\{P\} C \{Q\}$$

where  $P$  and  $Q$  are assertions, which are formulas in the first-order logic, and  $C$  is a program code.  $P$  is named the *precondition* and  $Q$  the *postcondition*: when the precondition is met, executing the code  $C$  establishes the postcondition. For instance,  $\{x = 1\}x := x + 1\{x = 2\}$  is true, because the value of  $x$  changes from 1 to 2 after the execution of  $x := x + 1$ . On the other hand,  $\{x = 1\}x := x + 1\{x = 0\}$  is false. Since a triple may be true or false, we call it a *statement* in the general sense.

**Example 10.1.1.** The following code computes the sum of all numbers in an integer array.

```

{ P: A is an array of n integers, and n ≥ 0 }
s := 0;
for i := 0 to n - 1 do
    s := s + A[i];
{ Q: s = ∑i=0n-1 A[i] }

```

At first glance,  $P$  and  $Q$  are not the formulas we saw in the previous chapters. An array is not different from a list in logic. In Prolog, we have seen that a list is just a compound term. Thus, it is not strange to represent an array by a term in

first-order logic, and we may simply treat an array as a term in the logic. Similarly, we can also treat  $\sum_{i=0}^{n-1} A[i]$  as another term. In fact, we may define a function  $sum : N, Array \rightarrow N$  recursively as follows:

$$\begin{aligned} sum(0, A) &= 0; \\ sum(i + 1, A) &= sum(i, A) + A[i] \end{aligned}$$

Then  $\sum_{i=0}^{n-1} A[i]$  is equivalent to the term  $sum(n, A)$ .

Now suppose we can specify preconditions and postconditions in first-order logic, how do we prove that a postcondition is a theorem from the precondition and the program? This is handled by the Hoare rules for each type of commands in a small programming language. In this language, we consider only three types of commands:

1. assignment, i.e.,  $x := E$ ;
2. conditional, i.e., **if**  $B$  **then**  $S_1$  **else**  $S_2$  **fi**; and
3. while loop, i.e., **while**  $B$  **do**  $S$  **od**.

Commands are separated by semicolons. This language is small but covers the essential part of most imperative programming languages. There are many variations of the loop construct in imperative programming languages: the “while” forms, which use a continuation condition; the “do-until” forms, which always execute the loop body at least once, testing for the condition at the end rather than on entry; and the “for” or “do” forms, which iterate over an integer interval or a data structure. They can all be expressed in a straightforward way as a while loop, on which we will rely throughout this chapter. For instance, a *for-loop* can be expressed by a *while-loop*. The for-loop in Example 10.1.1 can be replaced by the equivalent while loop.

<pre> s := 0; <b>for</b> i := 0 <b>to</b> n - 1 <b>do</b>     s := s + A[i] <b>od</b> </pre>	<pre> s := 0; i := 0; <b>while</b> (i &lt; n) <b>do</b>     s := s + A[i]; i := i + 1 <b>od</b> </pre>
--	--

### 10.1.1 Hoare Rules

Hoare rules cover the three basic commands in the small programming language, plus one rule for the composition of commands and one rule for logical implications.

### 10.1.1.1 Assignment Rule

The assignment rule states that, after the assignment, any predicate that was previously true for the right-hand side of the assignment now holds for the variable. Formally, let  $P(x)$  be an assertion in which the variable  $x$  is free. Then:

$$\frac{\top}{\{P[x \leftarrow E]\} x := E \{P(x)\}}$$

where  $P[x \leftarrow E]$  denotes the assertion obtained from  $P(x)$  in which each free occurrence of  $x$  is replaced by the expression  $E$ .

The assignment rule means that the truth of  $P[x \leftarrow E]$  is equivalent to the after-assignment truth of  $P(x)$ . Thus if  $P[x \leftarrow E]$  is true prior to the assignment, by the assignment rule, then  $P(x)$  would be true after the assignment. All preconditions that are not modified by the assignment can be carried over to the postcondition.

Examples of valid triples include:  $\{x = 4\} y := x + 1 \{y = 5\}$  and  $\{x \leq N - 1\} x := x + 1 \{x \leq N\}$ . In the first example, to show that  $y = 5$  is true, we need the precondition  $x + 1 = 5$ , which is equivalent to  $x = 4$ . Since assigning  $y := x + 1$  does not change the fact that  $x = 4$ , so both  $x = 4$  and  $y = 5$  may appear in the postcondition. Formally, this result is obtained by applying the assignment rule with  $P(y)$  being  $x = 4 \wedge y = 5$  and  $P[y \leftarrow x + 1]$  is  $x = 4 \wedge x + 1 = 5$ , which can in turn be simplified to the given precondition  $x = 4$ .

Since we have not given a formal definition of “state”, we cannot prove rigorously the correctness of the assignment rule. However, we may argue informally. To show that  $\{P[x \leftarrow E]\} x := E \{P\}$  is sound, let  $s$  be the state before  $x := E$  and  $s'$  the state after. So,  $s' = s[x \mapsto E]$  (assuming  $E$  has no side-effect, i.e., no variables changed values by  $E$ ).  $P[x \leftarrow E]$  holds in  $s$  iff  $P$  holds in  $s' = s[x \mapsto E]$ , because

1. Every variable, except  $x$ , has the same value in  $s$  and  $s'$ ,
2.  $P[x \leftarrow E]$  has every  $x$  in  $P$  replaced by  $E$ , and
3.  $P$  has every  $x$  evaluated to  $E$  in  $s' = s[x \mapsto E]$ .

The assignment rule is equivalent to saying that to find the precondition, first take the postcondition and replace all occurrences of the left side of the assignment with the right side of the assignment. Be careful not to try to do this backwards by following this incorrect way of thinking:  $\{P(x)\} x := E \{P[x \leftarrow E]\}$ ; this rule leads to examples like:  $\{x = 4\} x := 5 \{5 = 4\}$ . Another incorrect rule looking tempting at first glance is  $\{P(x)\} x := E \{P \wedge x = E\}$ ; it leads to illogical examples like:  $\{x = 5\} x := x + 1 \{x = 5 \wedge x = x + 1\}$ .

While a given postcondition  $P(x)$  uniquely determines the precondition  $P[x \leftarrow E]$  of an assignment, the converse is not true. For example:

- $\{0 \leq y \wedge y \leq 9\} x := y \{0 \leq x \wedge x \leq 9\}$ ,
- $\{0 \leq y \wedge y \leq 9\} x := y \{0 \leq y \wedge x \leq 9\}$ ,
- $\{0 \leq y \wedge y \leq 9\} x := y \{0 \leq y \wedge x \leq 9\}$ , and
- $\{0 \leq y \wedge y \leq 9\} x := y \{0 \leq y \wedge y \leq 9\}$

are valid instances of the assignment rule.

The assignment rule proposed by Hoare can apply to multiple mutual assignments. For example, to swap two elements of an array, some programming languages support the command  $A[i], A[j] := A[j], A[i]$ . The assignment rule for this type of commands is

$$\frac{\top}{\{P[x \leftarrow E_1, y \leftarrow E_2]\} x, y := E_1, E_2 \{P(x)\}}$$

#### 10.1.1.2 Implication Rule

The implication rule, sometimes called the consequence rule, is purely logical as it does not involve any command of a programming language, but concerns about the implication relation of logical formulas.

$$\frac{P_1 \rightarrow P_2 \quad \{P_2\} S \{Q_2\} \quad Q_2 \rightarrow Q_1}{\{P_1\} S \{Q_1\}}$$

This rule can be split into two rules:

$$\frac{P_1 \rightarrow P_2 \quad \{P_2\} S \{Q\}}{\{P_1\} S \{Q\}} \quad \frac{\{P\} S \{Q_2\} \quad Q_2 \rightarrow Q_1}{\{P\} S \{Q_1\}}$$

We say formula  $A$  is weaker than formula  $B$  if  $B \models A$ , but not  $A \models B$ . If we regard the entailment relation  $\models$  as an order  $\succ$  on formulas, then  $\perp$  is the strongest (or the largest) and  $\top$  is the weakest (or the least). We say  $P$  is the *weakest precondition* for program  $S$  and the postcondition  $Q$  if (a)  $\{P\}S\{Q\}$  is true; and (b) for any precondition  $P'$ , if  $\{P'\}S\{Q\}$  is true, then  $P' \models P$ .

The implication rule allows us to strengthen the precondition  $P_2$  to  $P_1$ , i.e., assume more than we need, and to weaken the postcondition  $Q_2$  to  $Q_1$ , i.e., conclude less than we are allowed to.

For instance, if we have shown that  $\{n > 0\} \text{sort}(A) \{n > 0 \wedge \text{sorted}(A)\}$ , and the target assertion is  $\text{sorted}(A)$ , since  $n > 0 \wedge \text{sorted}(A) \rightarrow \text{sorted}(A)$ , the implication rule allows us to get  $\{n > 0\} \text{sort}(A) \{\text{sorted}(A)\}$ .

### 10.1.1.3 Conditional Rule

$$\frac{\{B \wedge P\} S \{Q\} \quad \{\neg B \wedge P\} T \{Q\}}{\{P\} \mathbf{if} B \mathbf{then} S \mathbf{else} T \mathbf{fi} \{Q\}}$$

The conditional rule states that a postcondition  $Q$  common to the two branches of an if-then-else statement is also a postcondition of the whole if-then-else statement, i.e.,  $\mathbf{if} B \mathbf{then} S \mathbf{else} T \mathbf{fi} \{Q\}$ .

In the **then** part and the **else** part, the unnegated and negated condition  $B$  is added to the precondition  $P$ , respectively. The condition  $B$  must not have side-effect, i.e., it does not change the variables.

**Example 10.1.2.** To show

$$\{P : \top\} \mathbf{if} x < 10 \mathbf{then} x := 1 \mathbf{else} x := 10 \mathbf{fi} \{Q : x = 1 \vee x = 10\},$$

by the conditional rule, we need to show

- (1)  $\{P_1 : \top \wedge x < 10\} x := 1 \{Q : x = 1 \vee x = 10\}$ ,
- (2)  $\{P_2 : \top \wedge x \geq 10\} x := 10 \{Q : x = 1 \vee x = 10\}$ .

To prove (1) and (2), we obtain the following by the assignment rule:

- (1')  $\{1 = 1 \vee 1 = 10\} x := 1 \{Q : x = 1 \vee x = 10\}$ ,
- (2')  $\{10 = 1 \vee 10 = 10\} x := 10 \{Q : x = 1 \vee x = 10\}$ .

Since  $1 = 1 \vee 1 = 10 \equiv \top$  and  $P_1 \models \top$ , by the implication rule, (1) is true. (2) is true by the same reasoning.

If the else-part is missing or  $T$  is *skip* or  $x := x$ , the conditional rule becomes:

$$\frac{\{B \wedge P\} S \{Q\} \quad \neg B \wedge P \rightarrow Q}{\{P\} \mathbf{if} B \mathbf{then} S \mathbf{fi} \{Q\}}$$

**Example 10.1.3.** To show the following statement

$$\{P : x = a \wedge y = b\} \mathbf{if} x > y \mathbf{then} y := x \mathbf{fi} \{Q : y = \max(a, b)\}$$

is true, by the conditional rule, we need to show

- (1)  $\{P_1 : x = a \wedge y = b \wedge x > y\} y := x \{Q : y = \max(a, b)\}$ ;
- (2)  $(x = a \wedge y = b \wedge \neg(x > y)) \rightarrow y = \max(a, b)$ .

From the assignment rule, we have

$$\{x = \max(a, b)\} y := x \{Q : y = \max(a, b)\}$$

To prove (1), by the implication rule, we need to show  $P_1 \rightarrow x = \max(a, b)$ , which can be simplified by equality crossing as  $x > y \rightarrow x = \max(x, y)$ , which holds for all  $x$  and  $y$ . (2) can also be simplified by equality crossing to  $\neg(x > y) \rightarrow y = \max(x, y)$ , which is true in arithmetic, too.

#### 10.1.1.4 Iteration Rule

$$\frac{\{I \wedge B\} S \{I\}}{\{I\} \mathbf{while} B \mathbf{do} S \mathbf{od} \{\neg B \wedge I\}}$$

Here  $I$  is the so-called *loop invariant*, which is true at the beginning of the loop and is preserved by the loop body  $S$ . That is, after the loop is finished, this invariant  $I$  still holds. If the loop condition  $B$  becomes false,  $\neg B$ , which is called the *exit condition* of the loop, is added to  $I$  as the postcondition of the loop. As in the conditional rule,  $B$  must not have side effects.

**Example 10.1.4.** To show that the following statement is true,

$$\{I : x \leq 5\} \mathbf{while} x < 5 \mathbf{do} x := x + 1 \mathbf{od} \{Q : x = 5\}$$

we apply the iteration rule and show that

$$\{x \leq 5 \wedge x < 5\} x := x + 1 \{I : x \leq 5\}$$

By the assignment rule, we have  $\{x + 1 \leq 5\} x := x + 1 \{I : x \leq 5\}$ , so we just need to prove  $(x \leq 5 \wedge x < 5) \rightarrow (x + 1 \leq 5)$  by the implication rule. When  $x$  is an integer,  $(x < 5) \rightarrow (x + 1 \leq 5)$  is true. So  $(x \leq 5 \wedge x < 5) \rightarrow (x + 1 \leq 5)$  is true. Finally, we need to show  $(I \wedge \neg B) \rightarrow Q$  by the implication rule. Since  $x \leq 5 \wedge \neg(x < 5) \rightarrow x = 5$  is true for all number  $x$ , the program is correct when  $x$  is an integer.

#### 10.1.1.5 Sequence Rule

The assignment, conditional and iteration rules cover the three basic commands of the small programming language. The sequence rule applies to sequentially executed programs  $S$  and  $T$ , where  $S$  executes prior to  $T$  and is written as  $S;T$ .

$$\frac{\{P\} S \{Q\} \quad \{Q\} T \{R\}}{\{P\} S; T \{R\}}$$

where  $Q$  is called the *middle assertion*. For example, consider the following two instances of the assignment axiom:  $\{x = 4\} y := x + 1 \{y = 5\}$  and  $\{y = 5\} z := y \{z = 5\}$ . By the sequence rule, one concludes:  $\{x = 4\} y := x + 1; z := y \{z = 5\}$ .

### 10.1.2 Examples of Formal Verification

Now, we have all the rules needed to show the correctness of any program in the small programming language.

**Example 10.1.5.** Consider the program in Example 10.1.1:

$$\begin{array}{ll} \{n \geq 0\} & sum(0, A) = 0 \\ s := 0; i := 0; & sum(i + 1, A) = sum(i, A) + A[i] \\ \mathbf{while} (i < n) \mathbf{do} & \\ \quad s := s + A[i]; i := i + 1 & \\ \mathbf{od} & \\ \{s = sum(n, A)\} & \end{array}$$

Since the code contains a while loop, a crucial step in applying the Hoare logic is to find a loop invariant. Examining the body of the loop, we can see that the values of  $s$  and  $i$  are changed inside the body. A good loop invariant must reflect these changes. Since the variable  $s$  accumulates the sum of the first  $i$  elements of  $A$ , the first guess for a suitable loop invariant would be  $I : s = sum(i, A)$ .

By the sequence rule, we need to prove the following statements:

- (1)  $\{n > 0\} s := 0; i := 0 \{I\}$ ;
- (2)  $\{I\} \mathbf{while} i < n \mathbf{do} s := s + A[i]; i := i + 1 \mathbf{od} \{s = sum(n, A)\}$

The proof of (1) is easy as we just need to prove  $n > 0 \rightarrow I[i \leftarrow 0, s \leftarrow 0]$ , which is from  $\{I[i \leftarrow 0, s \leftarrow 0]\} s := 0 \{I[i \leftarrow 0]\}$  and  $\{I[i \leftarrow 0]\} i := 0 \{I\}$ . Since  $I[i \leftarrow 0, s \leftarrow 0]$  is  $0 = sum(0, A)$ , which is true by the definition of  $sum$ , hence (1) is true.

The proof of (2) is reduced to the following two statements by the iteration rule and the implication rule:

- (2.1)  $\{I \wedge i < n\} s := s + A[i]; i := i + 1 \{I\}$ , and

- (2.2)  $I \wedge \neg(i < n) \rightarrow \{s = \text{sum}(n, A)\}$ .

The way to prove (2.1) is the same as that of (1), and we need to prove the following formula:

$$(2.1.1) \quad (I \wedge i < n) \rightarrow I[i \leftarrow i + 1, s \leftarrow s + A[i]]$$

Opening  $I$  in (2.2.1), we have

$$(s = \text{sum}(i, A) \wedge i < n) \rightarrow (s + A[i] = \text{sum}(i + 1, A))$$

By the definition of  $\text{sum}$ ,  $\text{sum}(i + 1, A) = \text{sum}(i, A) + A[i] = s + A[i]$ , so (2.1.1) is true, thus (2.1) is true.

Opening  $I$  in (2.2), we have  $(s = \text{sum}(i, A) \wedge i \geq n) \rightarrow (s = \text{sum}(n, A))$ . Since the condition  $i \geq n$  cannot give us  $i = n$ , the proof cannot go through. From this failure, we can see that we need to strengthen the loop invariant by adding  $i \leq n$  to it: the new loop invariant  $I'$  should be  $s = \text{sum}(i, A) \wedge i \leq n$ . Using  $I'$ , the proofs of (1) and (2.1) will go through as before and the proof of (2.2) becomes easy as  $i \leq n \wedge i \geq n \rightarrow i = n$ .

The above example shows that the loop invariant should contain all the information about how the loop works:

- It reflects what has been done so far together with what remains to be done and should contain all the changing variables;
- It holds at both the beginning and the end of each iteration;
- Together with the negation of the loop condition, it gives the desired result when the loop terminates;
- If a proof cannot go through, we may need to strengthen the loop invariant by adding a condition which is needed in the proof.

**Example 10.1.6.** The following program will compute the integer quotient.

```

{x, y > 0}
r := x; q := 0;
while y ≤ r do
    r := r - y;
    q := q + 1
od
{q = ⌊x/y⌋}

```

Since both  $r$  and  $q$  are changed inside the loop, let the loop invariant  $I$  be  $y > 0 \wedge r \geq 0 \wedge x = r + y \cdot q$ . By the sequence rule, we need to show the following statements:

- (1)  $\{x, y > 0\} r := x; q := 0 \{I\}$ ;
- (2)  $\{I\} \mathbf{while} \ y \leq r \ \mathbf{do} \ r := r - y; q := q + 1 \ \mathbf{od} \ \{q = \lfloor x/y \rfloor\}$

(1) is easy to prove by using the assignment rules twice, the sequence rule once, and the implication rule once:  $x, y > 0 \rightarrow I[r \leftarrow x, q \leftarrow 0]$ , that is,  $x, y > 0 \rightarrow (y > 0 \wedge x \geq 0 \wedge x = x + y \cdot 0)$  is true.

(2) needs the iteration rule and the implication rule: We need to prove

$$(2.1) \ \{I \wedge y \leq r\} r := r - y; q := q + 1 \{I\};$$

$$(2.2) \ (I \wedge \neg(y \leq r)) \rightarrow (q = \lfloor x/y \rfloor).$$

The way to prove (2.1) is the same as to prove (1), we need to show

$$(I \wedge y \leq r) \rightarrow I[r \leftarrow r - y, q \leftarrow q + 1]$$

which, when  $I$  is open up, is

$$(y > 0 \wedge r \geq 0 \wedge x = r + y \cdot q \wedge y \leq r) \rightarrow (y > 0 \wedge r - y \geq 0 \wedge x = (r - y) + y \cdot (q + 1))$$

Since  $(r - y) + y \cdot (q + 1) = r + y \cdot q$  and  $y > 0 \wedge y \leq r$  imply  $r - y \geq 0$ , (2.1) is true.

The proof of (2.2) is easy as  $x = r + y \cdot q$  iff  $x/y = r/y + q$  when  $y > 0$ . Since  $(y > r \wedge r \geq 0)$  implies  $0 \leq r/y < 1$ , so  $\lfloor x/y \rfloor = \lfloor r/y + q \rfloor = \lfloor r/y \rfloor + q = 0 + q = q$ .

### 10.1.3 Hoare Rule for Function Calls

In addition to the rules for the simple language in Hoare's original paper, rules for other language constructs have been developed since then by Hoare and many other researchers. There are rules for functions, procedures, jumps, and pointers.

### 10.1.4 Partial and Total Correctness

Hoare rules presented in the previous section concern about the *partial correctness* of programs: If the program stops, the postcondition will hold. If a program is both partially correct and terminating, we say the program is *totally correct*. The relationship between partial and total correctness can be informally expressed by the equation:

$$\text{Total correctness} = \text{Termination} + \text{Partial correctness}.$$

As shown in the previous chapter, the termination of programs is an undecidable problem and there are no decision procedures to solve this problem in general.

Using the Hoare triples, only partial correctness can be proven, while termination needs to be proved separately. Thus the intuitive reading of a Hoare triple is: Whenever  $P$  holds of the state before the execution of  $C$ , then  $Q$  will hold afterwards, or  $C$  does not terminate. In the latter case, there is no “after”, so  $Q$  can be any statement at all. Indeed, one can choose  $Q$  to be false to express that  $C$  does not terminate.

**Example 10.1.7.** Consider the statement  $\{T\} \mathbf{while} \ x \neq a \ \mathbf{do} \ \mathbf{od} \ \{x = a\}$ . Let the loop invariant  $I$  be  $\top$ , then both  $I \wedge x \neq a \rightarrow I$  and  $I \wedge \neg(x \neq a) \rightarrow x = a$  are trivially true. Thus, the program is partially correct: If the program stops, we must have  $x = a$ . The program will loop forever unless it starts with  $x = a$ , therefore it is not totally correct.

Total correctness can be proved with an extended version of the iteration rule, where the value of a function on the changing variables of the loop decreases under a well-founded order. Total correctness is what we are ultimately interested in, but it is usually easier to prove it by establishing partial correctness and termination separately. Termination is often straightforward to establish, and the complexity analysis of algorithms often provides finer results than the termination. For example, the program in Example 10.1.5 is terminating because the value of  $f(i) = n - i$  decreases by one after each iteration and the minimum value of  $f(i)$  is 0 when the looping condition  $i < n$  becomes false. For the program in Example 10.1.6, let  $f(q) = \lfloor x/y \rfloor - q$ , then  $f(q)$  decreases by one after each iteration and the minimum value of  $f(q)$  is 0 when the looping condition is false. The complexity of both algorithms is  $O(n)$  ( $n = \lfloor x/y \rfloor$  in the second example).

## 10.2 Automated Generation of Assertions

After seeing only a few examples, the following two things is painfully clear:

- Proofs are typically long and boring, even if the program being verified is quite simple.
- There are lots of fiddly little details to get right, many of which are trivial, e.g. proving  $r = x \wedge q = 0 \rightarrow (x = r + y \cdot q)$ .

Many attempts have been made (and are still being made) to automate correctness proofs by designing systems to do the boring and tricky bits of generating

formal proofs in Hoare logic. Unfortunately logicians have shown that it is impossible in principle to design a decision procedure to decide automatically the truth or falsehood of an arbitrary mathematical statement. However, this does not mean that one cannot have procedures that will prove many useful theorems. The non-existence of a general decision procedure merely shows that one cannot hope to prove everything automatically. In practice, it is quite possible to build a system that will mechanize many of routine aspects of verification. This section describes one commonly taken approach to doing this.

In the previous section it was shown how to prove  $\{P\} S \{Q\}$  by reducing this goal to several subgoals and then putting these together using the Hoare rules to get the desired property of  $S$  itself. For example, to prove  $\{P\} S; T \{Q\}$ , first prove  $\{R\} T \{Q\}$  and  $\{P\} S \{R\}$  for a middle assertion  $R$ , and then deduce  $\{P\} S; T \{Q\}$  by the sequence rule. This process is called *goal reduction*, or *backward reasoning*, because one works backwards: Starting from the goal of showing  $\{P\} S \{Q\}$ , one generates subgoals, subsubgoals, etc., until the problem is solved. Prolog also uses this process for computation.

**Example 10.2.1.** Suppose one wants to show:

$$\{x = x_0 \wedge y = y_0\} r := x; x := y; y := r \{y = x_0 \wedge x = y_0\}$$

then by the assignment and sequence rules, the above statement is reduced to the subgoal

$$\{x = x_0 \wedge y = y_0\} r := x; x := y \{r = x_0 \wedge x = y_0\}$$

because  $\{r = x_0 \wedge x = y_0\} y := r \{y = x_0 \wedge x = y_0\}$ . By a similar argument this subgoal can be reduced to

$$\{x = x_0 \wedge y = y_0\} r := x \{r = x_0 \wedge y = y_0\}$$

which is reduced further to  $(x = x_0 \wedge y = y_0) \rightarrow (x = x_0 \wedge y = y_0)$ . Every step of the above reduction clearly follows from the assignment rule. The middle assertions for the sequence rule are automatically generated.

We wish that the process illustrated by the above example can be extended to arbitrary programs so that all middle assertions can be generated automatically, releasing programmers from the burden of providing all assertions. A programmer just needs to provide the pre- and post-conditions of the program derived from the specification of the problem and all the needed assertions can be generated automatically. Coupling with a powerful theorem prover, all the formulas derived from Hoare rules can be proved automatically, and we arrive at the automation of software verification.

### 10.2.1 Verification Conditions

The formal verification of a program can be described succinctly as follows:

1. The program  $S$  is annotated by inserting into it the assertions expressing conditions that are meant to hold at various intermediate points. This step needs a good understanding of how the program works. The user needs to provide at least the pre- and post-conditions and loop invariants; other assertions can be generated automatically.
2. A set of formulas called *verification conditions* is then generated from the annotated specification. This process is purely mechanical and easily done by a software tool. We will describe how this done in this section.
3. The verification conditions are proved by a theorem prover. Automating this remains to be a big challenge.

**Definition 10.2.2.** *Given a Hoare triple  $\{P\} S \{Q\}$ , let  $vc(S, Q)$  denote the formula  $P'$  such that both (1)  $P \rightarrow P'$  and (2)  $\{P'\} S \{Q\}$  are true, then  $P'$  is called a verification condition of  $S$  with respect to  $Q$ .*

By the implication rule, if a verification condition  $vc(S, Q)$  exists, then the triple  $\{P\} S \{Q\}$  is true. Moreover, if  $vc(S, Q)$  can be generated automatically from  $S$  and  $Q$ , then the process of generating all assertions becomes automated. Besides being the basis for automated verification systems, verification conditions are a useful way of working with Hoare logic by hand.

We will view  $vc(S, Q)$  as a recursive procedure on  $S$ , coinciding with the goal reduction process. The definition of  $vc(S, Q)$  follows closely with Hoare rules on the three constructs and the sequence rule of our small programming language. We add a *skip* command, which does nothing, so that it is easy to use the if-then-else command for the if-then command: **if**  $B$  **then**  $S$  **fi** becomes **if**  $B$  **then**  $S$  **else** *skip* **fi**. In summary, there are only five cases to consider in the definition of  $vc(S, Q)$ .

In fact, the verification conditions are the *weakest preconditions* in most cases, with the exception of while loops, but we will neither make this more precise nor go into details here. An in-depth study of weakest preconditions can be found in Dijkstra's book.

For while loops, we do not have an algorithm which can generate  $vc(S, Q)$  for every loop; only some heuristic methods are available. To help the verification tool to generate  $vc(S, Q)$ , a programmer needs to suggest a loop invariant. That is, instead of the code **while**  $B$  **do**  $S$  **od**, we write **while**  $B$  **do**  $\{I\} S$  **od**, where  $I$  is the suggested loop invariant. The pre- and post-conditions can be also inserted

into a program code. This type of program codes with assertions is called “annotated program”, where assertions are provided by the programmer to assist the verification.

Another exception with while loops is that, in addition to returning a verification condition for the loop, it will create another formula whose validity is needed to ensure the validity of the Hoare triple. This formula will be pushed into a stack denoted by  $s$  during the execution of  $vc$ . The stack  $s$  is a global variable and is initially empty; it will hold as many formulas as the number of while loops in the program.

**Algorithm 10.2.3.** The algorithm  $vc$  will take an annotated program  $C$ , a formula  $Q$  and returns a formula as a precondition  $P$  for  $\{P\} C \{Q\}$  being true.

```

proc  $vc(C, Q)$ 
1   if  $C$  is “ $x := E$ ” return  $Q[x \leftarrow E]$ 
2   else if  $C$  is “skip” return  $Q$ 
3   else if  $C$  is “if  $B$  then  $S$  else  $T$  fi”
4     return  $(B \rightarrow vc(S, Q)) \wedge (\neg B \rightarrow vc(T, Q))$ 
5   else if  $C$  is “while  $B$  do  $\{I\}$   $S$  od”
6      $push(I \rightarrow ((B \rightarrow vc(S, I)) \wedge (\neg B \rightarrow Q)), s)$ ;
7     return  $((\neg B \wedge Q) \vee I)$ 
8   else  $C$  is “ $S; T$ ” return  $vc(S, vc(T, Q))$ 

```

In the above algorithm, lines 1-2 handles the base cases, when the command is an assignment or skip command; Lines 3-4 handle the conditional (for simplicity, we ignore the conditional case when the else part is missing); Lines 5-7 handle the while loop; and line 8 handles the sequence of commands. The last three cases involve the recursive calls of  $vc$ .

**Theorem 10.2.4.** *The output of  $vc(S, Q)$  is a formula  $P'$  in first-order logic such that (a)  $P \models P'$  and (b) all the formulas generated by  $vc(S, Q)$  in the stack  $s$  are valid, then  $\{P\} S \{Q\}$  is true.*

The following subsection will be devoted the proof of the above theorem.

### 10.2.2 Proof of Theorem 10.2.4

At first, we need to show that  $vc$  will terminate. This is an easy task because the first argument of  $vc$ , which is a piece of code, becomes strictly smaller for each recursive call, and it will come down to an assignment or a skip command, there cannot exist an infinite chain of recursive calls. So the algorithm  $vc$  must terminate.

The proof of the above theorem will be by induction on the structure of  $C$ . Such inductive arguments have two parts. First, it is shown that the result holds for the assignment and skip commands. Second, it is shown that when  $C$  is not an assignment command, then if the result holds for the constituent commands of  $C$  (this is called the induction hypothesis), then it holds also for  $C$ . The first of these parts is called the base cases of the induction and the second is called the inductive cases. From the base and the inductive cases, it follows that the result holds for all commands. For instance, to show that  $vc$  returns a formula, since all the base cases return a formula, assuming all the recursive calls will return a formula, then  $vc$  will return a formula by examining all the inductive cases.

The case of the skip command is trivial; the other cases will be discussed separately.

#### 10.2.2.1 Assignments

The only verification condition for  $x := E$  and  $Q$  is  $Q[x \leftarrow E]$ , which is the weakest precondition for  $x := E$  and  $Q$ . Thus, we have  $vc(x := E, Q) = Q[x \leftarrow E]$ .

#### 10.2.2.2 Sequences

The only verification condition for  $S; T$  and  $Q$  is  $vc((S; T), Q) = vc(S, vc(T, Q))$ . If both  $vc(S, Q')$  and  $vc(T, Q)$  are the weakest preconditions for  $S$  and  $T$ , respectively, then  $vc(S, vc(T, Q))$  is the weakest precondition for  $S; T$ .

**Example 10.2.5.** The verification condition for the code in Example 10.2.1 is easily obtained as it involves only two cases: assignment and sequence.

$$\begin{aligned}
 & vc((r := x; x := y; y := r), y = x_0 \wedge x = y_0) \\
 = & vc((r := x; x := y), vc(y := r, y = x_0 \wedge x = y_0)) \\
 = & vc((r := x; x := y), r = x_0 \wedge x = y_0) \\
 = & vc(r := x, vc(x := y, r = x_0 \wedge x = y_0)) \\
 = & vc(r := x, r = x_0 \wedge y = y_0) \\
 = & x = x_0 \wedge y = y_0
 \end{aligned}$$

#### 10.2.2.3 Conditionals

Let  $P' = vc(\mathbf{if } B \mathbf{ then } S \mathbf{ else } T \mathbf{ fi}, Q)$ , by the definition of  $vc$ ,  $P' = (B \rightarrow vc(S, Q)) \wedge (\neg B \rightarrow vc(T, Q))$ . Note that  $P'$  is the weakest precondition for the conditional command, if both  $vc(S, Q)$  and  $vc(T, Q)$  are the weakest preconditions for  $S$  and  $T$ , respectively. By the induction hypothesis, both  $\{vc(S, Q)\}S\{Q\}$  and

$\{vc(T, Q)\}T\{Q\}$  are true. If  $P \rightarrow P'$  is true, then both  $P \rightarrow (B \rightarrow vc(S, Q))$  and  $P \rightarrow (\neg B \rightarrow vc(T, Q))$  are true. Since  $P \rightarrow (B \rightarrow vc(S, Q)) \equiv (P \wedge B) \rightarrow vc(S, Q)$ , by the implication rule,  $\{P \wedge B\}S\{Q\}$ . Following the same reasoning, we have  $\{P \wedge \neg B\}T\{Q\}$ . By the conditional rule,

$$\{P\} \mathbf{if } B \mathbf{ then } S \mathbf{ else } T \mathbf{ fi } \{Q\}$$

is true.

### Example 10.2.6.

$$\begin{aligned} & vc(\mathbf{if } (x < 10) \mathbf{ then } x := 0 \mathbf{ else } x := 10 \mathbf{ fi}, x = 0 \vee x = 10) \\ = & (x < 10 \rightarrow vc(x := 0, x = 0 \vee x = 10)) \wedge (x \geq 10 \rightarrow vc(x := 10, x = 0 \vee x = 10)) \\ = & (x < 10 \rightarrow 0 = 0 \vee 0 = 10) \wedge (x \geq 10 \rightarrow 10 = 0 \vee 10 = 10) \\ = & (x < 10 \rightarrow \top) \wedge (x \geq 10 \rightarrow \top) \\ = & \top \end{aligned}$$

#### 10.2.2.4 Iterations

Let  $P' = vc(\mathbf{while } B \mathbf{ do } \{I\} S \mathbf{ od}, Q)$ , by the definition of  $vc$ ,  $P'$  is  $(\neg B \wedge Q) \vee I$ , which is a disjunction of two formulas: (a)  $\neg B \wedge Q$ , and (b)  $I$ . (a) covers the case when  $B$  is false in the beginning and the while loop is the same as a skip command. (a) can be omitted if we let  $vc$  to return a stronger precondition. (b) ensures that  $I$  holds before entering the loop.

Moreover, the formula  $I \rightarrow ((B \rightarrow vc(S, I)) \wedge (\neg B \rightarrow Q))$ , which is equivalent to  $(I \wedge B \rightarrow vc(S, I)) \wedge (I \wedge \neg B \rightarrow Q)$ , is created and stored in the stack  $s$ . The validity of  $I \wedge B \rightarrow vc(S, I)$  ensures that the iteration rule can be applied. By the induction hypothesis,  $\{vc(S, I)\}S\{I\}$  is true. Since  $I \wedge B \rightarrow vc(S, I)$ , by the implication rule,  $\{I \wedge B\}S\{I\}$  is true, that is,  $I$  is indeed a loop invariant. Thus, by the iteration rule,  $\{I\} \mathbf{while } B \mathbf{ do } S \mathbf{ od } \{I \wedge \neg B\}$  is true. Finally, the validity of  $(I \wedge \neg B) \rightarrow Q$  ensures that  $Q$  holds when the loop terminates. By the implication rule,  $\{I\} \mathbf{while } B \mathbf{ do } S \mathbf{ od } \{Q\}$  holds. This completes the proof of Theorem 10.2.4.  $\square$

The algorithm  $vc(S, Q)$  shows how assertions needed for Hoare rules are derived from the postcondition and loop invariants automatically in a goal-reduction or backward reasoning process. Backward reasoning can be useful for nested loops, where the inner loop's postcondition can be derived from the outer loop's invariant.

**Example 10.2.7.** The precondition for the following program is  $P : x = n \wedge y = 1$  and the loop invariant is  $I : x! \cdot y = n!$ , where  $x!$  is the factorial function of  $x$ .

$$\begin{aligned} & vc(\mathbf{while} (x \neq 0) \mathbf{do} \{I : x!y = n!\} y := x * y; x := x - 1 \mathbf{od}, y = n!) \\ &= (x = 0 \wedge y = n!) \vee (x!y = n!) \end{aligned}$$

The formula pushed into the stack and needs to be proved is

$$\begin{aligned} & x!y = n! \rightarrow ((x! = 1 \rightarrow vc([y := y * x; x := x - 1], x!y = n!)) \wedge (x = 1 \rightarrow y = n!)) \\ &= x!y = n! \rightarrow ((x! = 1 \rightarrow vc(y := y * x, vc(x := x - 1, x!y = n!))) \wedge (x = 1 \rightarrow y = n!)) \\ &= x!y = n! \rightarrow ((x! = 1 \rightarrow vc(y := y * x, (x - 1)!y = n!)) \wedge (x = 1 \rightarrow y = n!)) \\ &= x!y = n! \rightarrow ((x! = 1 \rightarrow (x - 1)!(yx) = n!) \wedge (x = 1 \rightarrow y = n!)) \\ &= x!y = n! \rightarrow ((x! = 1 \rightarrow x!y = n!) \wedge (x = 1 \rightarrow y = n!)) \\ &= x!y = n! \rightarrow (x = 1 \rightarrow y = n!) \end{aligned}$$

The last step used the equivalence  $X \rightarrow ((Y \rightarrow X) \wedge Z) \equiv X \rightarrow Z$ , where  $X$  is  $x!y = n!$ ,  $Y$  is  $x! = 1$ , and  $Z$  is  $x = 1 \rightarrow y = n!$ .

Can we get rid the stack used in  $vc$ ? That is, instead of storing  $I \rightarrow ((B \rightarrow vc(S, I)) \wedge (\neg B \rightarrow Q))$ , can we let  $vc(\mathbf{while} B \mathbf{do} \{I\} S \mathbf{od}, Q)$  returns  $I \rightarrow ((B \rightarrow vc(S, I)) \wedge (\neg B \rightarrow Q))$  as part of its output? The following counterexample shows that this approach is unsound.

**Example 10.2.8.** Let the Hoare triple be

$$\{x = 0\} \mathbf{while} (x < 5) \mathbf{do} x := x + 1 \mathbf{od} \{x = 5\},$$

which can be proved to be true by the Hoare rules with the loop invariant  $x \leq 5$ .

If we modify  $vc$  such that  $vc(\mathbf{while} B \mathbf{do} \{I\} S \mathbf{od}, Q)$  returns

$$(\neg B \wedge Q) \vee I \wedge (I \rightarrow ((B \rightarrow vc(S, I)) \wedge (\neg B \rightarrow Q))),$$

which is equivalent to

$$(\neg B \wedge Q) \vee I \wedge (B \rightarrow vc(S, I)) \wedge (\neg B \rightarrow Q).$$

Let  $I$  be  $x = 0 \vee x = 1$  for the above program, then  $vc(\mathbf{while} (x < 5) \mathbf{do} \{I\} x := x + 1 \mathbf{od}, x = 5)$  returns

$$(x \geq 5 \wedge x = 5) \vee I \wedge (x < 5 \rightarrow vc(x := x + 1, I)) \wedge (x \geq 5 \rightarrow x = 5).$$

which can be simplified to

$$P' : x = 5 \vee (x = 0 \vee x = 1) \wedge (x < 5 \rightarrow (x = -1 \vee x = 0)) \wedge (x \geq 5 \rightarrow x = 5),$$

because  $vc(x := x + 1, I)$  returns  $(x + 1 = 0 \vee x + 1 = 1)$ , which is equivalent to  $x = -1 \vee x = 0$ . Now it is easy to show (by equality crossing) that  $x = 0 \rightarrow P'$  is true (since  $P'[x \leftarrow 0]$  is true). By the original Hoare rule, we cannot show that  $x = 0 \vee x = 1$  is a valid loop invariant. Thus, the modification of  $vc$  is unsound.

### 10.2.3 Implementing *vc* in Prolog

Since Prolog is good at pattern matching, it is easy to implement the algorithm  $vc(S, Q)$  in Prolog. The first job is to let the Prolog program parse codes easily. For this purpose, we change the syntax of our programming language as follows:

<b>command</b>	<b>before</b>	<b>now</b>
assignment	$x := E$	$\text{let}(x, E)$
conditional	<b>if</b> $B$ <b>then</b> $S$ <b>else</b> $T$ <b>fi</b>	$\text{ite}(B, S, T)$
iteration	<b>while</b> $B$ <b>do</b> $\{I\}$ $S$ <b>od</b>	$\text{while}(B, I, S)$
skip	skip	skip
sequence	$S_1; S_2; \dots; S_n$	$[S_1, S_2, \dots, S_n]$

For the iteration, instead of using a stack to store formulas, we simply write the formulas on the screen. Below is the Prolog code for treating assignments, conditionals, iterations, and skip:

```
% vc1(S1, Post, VC): VC = vc(S1, Post), S1 is a single command

% assignment
vc1(let(X, E), Post, VC) :- substit(X, E, Post, VC).

% conditional
vc1(ite(B, S, T), Post, VC) :-
    vc(S, Post, VC1), vc(T, Post, VC2),
    VC = and((B -> VC1), (neg(B) -> VC2)).

% iteration
% use vc1(while(B, Inv, S), Post, Inv) as an alternative.
vc1(while(B, Inv, S), Post, or(and(neg(B), Post), Inv)) :-
    vc(S, Inv, VC),                % VC is vc(S, Inv)
    write((and(Inv, B) -> VC)), nl, % part 1 of a conjunction
    write((and(Inv, neg(B)) -> Post)), nl. % part 2 of a conjunction
```

The handling of assignment uses a utility function `substit`, which performs the substitution:

```
% substit(X, V, E, R) succeeds when every occurrence of X in
% expression E is replaced by V and the modified E is R.
substit(X, V, X, V):- !.
```

```

substit(X, V, E, R) :-
    E =.. [F|Args0],
    maplist(substit(X, V), Args0, Args),
    R =.. [F|Args].

```

For a sequence  $S$  of commands, we first reverse  $S$  into  $S_2$  and work on  $S_2$  from the beginning to the end, as illustrated by the following code.

```

% vc(S, Post, VC) succeeds when a VC is generated from
% from the program S and the post condition Post.

vc(P, Post, VC) :- is_list(P), !, reverse(P, P2), vcR(P2, Post, VC).
vc(P, Post, VC) :- vc1(P, Post, VC).      % if S is not a list

% vcR(SR, Post, VC) generates a VC from SR and Post, where
% SR is a list of instructions in reversed order.
vcR([], Post, Post).
vcR([Last | Rest], Post, VC) :-
    vc1(Last, Post, Post1), vcR(Rest, Post1, VC).

```

In the following examples, we define a triple  $tri_i(P, S, Q)$ , where  $i$  is a number. We then define  $q_i$  for the query of generating and displaying a verification condition for the triple  $tri_i(P, S, Q)$ . The validity of all the displayed formulas ensures that the triple is true.

**Example 10.2.9.** The triple and the query are specified as follows:

```

tri1(true,
    [ ite(a>b, let(x, b), let(x, a)), ift(x > c, let(x, c))],
    x=min(a,min(b,c))).

q1 :- tri1(P, S, Q), vc(S, Q, VC), write((P -> VC)), nl.

```

The query `?- q1.` produces one formula:

```

true->and((a>b->and((b>c->c=min(a,min(b,c))),
    (neg(b>c)->b=min(a,min(b,c))))),
    (neg(a>b)->and((a>c->c=min(a,min(b,c))),
    (neg(a>c)->a=min(a,min(b,c))))))

```

which represents the formula

$$\top \rightarrow ((a > b \rightarrow ((b > c \rightarrow c = m) \wedge (\neg(b > c) \rightarrow b = m))) \wedge (\neg(a > b) \rightarrow ((a > c \rightarrow c = m) \wedge (\neg(a >$$

where  $m = \min(a, \min(b, c))$ .

**Example 10.2.10.** The triple and the query are specified as follows:

```
tri2(and(n>=0, and(k=0, result=0)),
     while(k<n,
           and(0=<k, and(k=<n, result=sum(0, k, a))),
           [let(result, result+a(k)),
            let(k, k+1)]),
     result=sum(0, n, a)).
```

```
q2 :- tri2(P, S, Q), vc(S, Q, VC), write((P -> VC)), nl.
```

The query `?- q2` produces three formulas:

```
and(and(0=<k, and(k=<n, result=sum(0, k, a))), k<n)->
  and(0=<k+1, and(k+1=<n, result+a(k)=sum(0, k+1, a)))
and(and(0=<k, and(k=<n, result=sum(0, k, a))), neg(k<n))->
  result=sum(0, n, a)
and(n>=0, and(k=0, result=0))->
  or(and(neg(k<n), result=sum(0, n, a)),
     and(0=<k, and(k=<n, result=sum(0, k, a))))
```

where  $sum(b, e, A) = 0$  if  $(b \geq e)$ ; otherwise  $sum(b, e, A) = sum(b, e-1, A) + A[e-1]$ .

The three formulas can be rewritten as

$$\begin{aligned} (0 \leq k \wedge k \leq n \wedge result = sum(0, k, a) \wedge k < n) &\rightarrow (0 \leq k + 1 \wedge k + 1 \leq n \wedge result + a(k) = sum(0, k+1, a)) \\ (0 \leq k \wedge k \leq n \wedge result = sum(0, k, a) \wedge \neg(k < n)) &\rightarrow result = sum(0, n, a) \\ (n \geq 0 \wedge k = 0 \wedge result = 0) &\rightarrow ((\neg(k < n) \wedge result = sum(0, n, a)) \vee (0 \leq k \wedge k \leq n \wedge result = sum(0, k, a))) \end{aligned}$$

**Example 10.2.11.** The triple and the query are specified as follows:

```
tri3(and(n >= 0, and(x=n, y=1)),
     while(neg(x = 0),
           \{ fac(x)*y = fac(n)\},
           [let(y, x*y),
            let(x, x-1)]),
     y=fac(n)).
```

```
q3 :- tri3(P, S, Q), vc(S, Q, VC), write((P -> VC)), nl.
```

The query ?- q3 produces three formulas:

```
and(fac(x)*y=fac(n),neg(x=0))->fac(x-1)*(x*y)=fac(n)
and(fac(x)*y=fac(n),neg(neg(x=0)))->y=fac(n)
and(n>=0,and(x=n,y=1))->
  or(and(neg(neg(x=0)),y=fac(n)),fac(x)*y=fac(n))
```

where  $fac(n)$  is the factorial function:  $fac(0) = 1$ ,  $fac(i + 1) = fac(i) * (i + 1)$ . The three formulas can be rewritten as

$$\begin{aligned} (fac(x) * y = fac(n) \wedge \neg(x = 0)) &\rightarrow fac(x - 1) * (x * y) = fac(n) \\ (fac(x) * y = fac(n) \wedge \neg\neg(x = 0)) &\rightarrow y = fac(n) \\ (n \geq 0 \wedge x = n \wedge y = 1) &\rightarrow ((\neg\neg(x = 0) \wedge y = fac(n)) \vee fac(x) * y = fac(n)) \end{aligned}$$

The above three examples show how formulas can be generated from post-conditions and loop invariants. All middle assertions, with the exception of loop invariants, are generated automatically. These formulas are then fed to an automated theorem prover. Once they are proved to be valid, the Hoare triple is shown to be true and the program is partially correct.

More than often, a theorem prover cannot prove these formulas automatically; if it fails, advice is sought from the user. The analysis of the failure may suggest that we need to add a conjunct to a loop invariant, to the precondition, or modify the postcondition. We then run *vc* with new annotated code, obtain new formulas and feed them to the theorem prover. This process repeats until we find the right pre- and post-conditions and loop invariants.

The aim of much current research is to build systems which reduce the role of the slow and expensive human expert to a minimum. This can be achieved by:

- reducing the number and complexity of the annotations required, and
- increasing the power of the theorem prover.

From the implementation of *vc*, it is now clear that the success of Hoare logic to formal verification relies on (i) obtaining a good loop invariant for each while loop in a program; (ii) having a powerful theorem prover which can prove the validity of the formulas derived by *vc*. The first issue will be addressed in the next section. The second issue has been addressed in the previous chapters of this book and will be discussed in the next chapter. Since proofs of correctness of programs are typically very long and boring, it is thus useful to check proofs mechanically or automatically, even if they can only be carried out with human assistance.

### 10.3 Obtaining Good Loop Invariants

Finding an invariant for a loop is traditionally the responsibility of a human: either the person performing the verification or the programmer writing the loop in the first place. A better solution, when applicable, is *program synthesis*, a constructive approach to programming advocated by Dijkstra and others. Automated generation of loop invariants is still an ongoing research topic.

For the triple  $\{I\} \mathbf{while} B \mathbf{do} S \mathbf{od} \{Q\}$ ,  $I$  is a correct invariant for the loop if it satisfies the following conditions:

1.  $I$  is true in the state preceding the loop execution.
2. Every execution of the loop body  $S$ , started in any state in which both  $I$  and  $B$  are true, will yield a state in which  $I$  holds again.

If these properties hold, then any terminating execution of the loop will yield a state in which both  $I$  and  $\neg B$  hold; the latter is called the *exit condition*.

We may look at the notion of loop invariant from the constructive perspective of a programmer directing his or her program to reach a state satisfying a certain desired property, the postcondition. In this view, program construction is a form of problem solving, and the various control structures are problem-solving techniques; a loop solves a problem through successive approximation. This idea has been strongly advocated in Furia *et al.*'s survey article on loop invariants. In the following, our presentation of the idea is based on their survey.

- Generalize the postcondition (the characterization of possible solutions) into a broader condition: the invariant.
- As a result, the postcondition can be defined as the conjunction of the invariant and another condition, i.e., the exit condition.
- Find a way to reach the invariant from the previous state of the computation: the initialization.
- Find a way, given a state that satisfies the invariant, to get to another state, still satisfying the invariant but closer, in some appropriate sense, to the exit condition: the body.

The importance of this presentation of the loop process is that it highlights the nature of the invariant: it is a generalized form of the desired postcondition, which in a special case (represented by the exit condition) will give us that postcondition. This view of the invariant, as a particular way of generalizing the desired goal of the

loop computation, explains why the loop invariant is such an important property of loops; one can argue that understanding a loop means understanding its invariant.

### 10.3.1 Invariants from Generalizing Postconditions

To illustrate the idea of generalizing the postcondition to obtain an invariant, we will use the thousands-year-old Euclid's algorithm for computing the greatest common divisor (gcd) of two positive integers. The postcondition of the algorithm is

$$y = \text{gcd}(a, b),$$

where the precondition is  $a > 0$  and  $b > 0$ . The first step of the generalization is to replace this condition by

$$x = 0 \wedge \text{gcd}(x, y) = \text{gcd}(a, b)$$

with a new variable  $x$ , taking advantage of the mathematical property that, for every  $y$ ,  $\text{gcd}(0, y) = y$ . The second conjunct, i.e.,  $\text{gcd}(x, y) = \text{gcd}(a, b)$ , a generalization of the postcondition, will serve as the invariant; the first conjunct will serve as the exit condition. To obtain the loop body, we take advantage of another mathematical property: for every  $x$  and  $y$ ,  $\text{gcd}(x, y) = \text{gcd}(y, x)$  and  $\text{gcd}(x, y) = \text{gcd}(x - y, y)$  if  $x > y$ , yielding the well-known algorithm:

```

x := a; y := b;
while x ≠ 0 do
  {I : x ≥ 0 ∧ y > 0 ∧ gcd(x, y) = gcd(a, b)}
  if x < y then z := x; x := y; y := z fi;
  x := x - y
od

```

The proof of correctness follows directly from the mathematical property of  $\text{gcd}$ . This form of Euclid's algorithm uses subtraction; we may replace the subtraction by the remainder function, i.e., replace  $x - y$  by  $x \% y$ , resulting a more efficient algorithm. In this case, we need the property that  $\text{gcd}(x, y) = \text{gcd}(x \% y, y)$ .

Using the classification criteria proposed in Furia *et al.*'s survey, the invariant for the  $\text{gcd}$  algorithm is divided into two categories:

- The last conjunct of the invariant is an *essential* invariant, representing a weakening of the postcondition. The first two conjuncts are *bounding* invariants, indicating that the state remains within certain general boundaries and ensuring that the “essential” part is well defined.

- The strategy that leads to the essential invariant is *uncoupling*, which replaces a property of one variable ( $y$  in this case), used in the postcondition, by a property of two variables ( $y$  and  $x$ ), used in the invariant.

As illustrated by the above example, the essential invariant is a mutation (often, a weakening) of the loop's postcondition. The following mutation techniques are particularly common.

- **Constant Relaxation**

Constant relaxation replaces a constant  $n$  (more generally, an expression which does not change during the execution of the algorithm) by a variable  $i$ , and use  $i = n$  as part or all of the exit condition. In Example 10.1.5, constant relaxation is used to obtain the loop invariant, i.e.,  $s = \text{sum}(i, A)$ , from the postcondition  $s = \text{sum}(n, A)$ . This condition is trivial to establish initially for  $i = 0$  and  $s = 0$ , and is easy to extend to an incremented  $i$ , and yields the postcondition when  $i$  reaches  $n$ .

- **Uncoupling** Uncoupling is to replace a variable by two, using their equality as part or all of the exit condition. Uncoupling is used in the gcd algorithm where we replace  $y = \text{gcd}(a, b)$  by  $x = 0 \wedge \text{gcd}(x, y) = \text{gcd}(a, b)$ . Uncoupling is also used in Example 10.2.7 where we replace  $y = n!$  by  $x = 0 \wedge x!y = n!$ .
- **Term dropping** Term dropping removes a subformula (typically a conjunct), which gives a straightforward weakening. Term dropping often applies after uncoupling.  $x = 0$  is dropped from  $x = 0 \wedge \text{gcd}(x, y) = \text{gcd}(a, b)$  in the gcd algorithm, and  $x = 0 \wedge x!y = n!$  in the factorial algorithm to obtain the essential invariant. The dropped term often goes into the exit condition of the loop.
- **Aging** Aging is to replace a variable (more generally, an expression) by an expression that represents the value the variable had at previous iterations of the loop. Aging typically accommodates “off-by-one” discrepancies between when a variable is evaluated in the invariant and when it is updated in the loop body. Typically, aging is used for generating bounding invariants. For example, if  $i > 0$  appears in the loop condition and  $i := i - 1$  appears in the loop body, aging will generate  $i + 1 > 0$  (equivalently  $i \geq 0$ ) as a bounding invariant. On the other hand, if  $j < n$  appears in the loop condition and  $j := j + 1$  appears in the body, then aging will generate  $j - 1 < n$  (equivalently  $j \leq n$ ) as a bounding invariant.

The above generalization techniques can be implemented by heuristic methods and thus generating invariants automatically. It is pointed out in Furia *et al.*'s survey that the amount of research work on the automated generation of invariants is substantial and spread over more than three decades; this reflects the cardinal role that invariants play in the formal analysis and verification of programs. The methods for automatic generation of invariants can be classified as *static* and *dynamic*. Static methods use only the program text, whereas dynamic methods summarize the properties of many program executions with different inputs.

Historically, the earliest methods for invariant inference were static. Abstract interpretation and the constraint-based approach are the two most widespread frameworks for static invariant inference. Abstract interpretation is a symbolic execution of programs over abstract domains that overapproximates the semantics of loop iteration. Constraint-based techniques rely on sophisticated decision procedures over nontrivial mathematical domains (such as polynomials or convex polyhedra) to represent concisely the semantics of loops with respect to certain template properties. Static methods are sound and often complete with respect to the class of invariants that they can infer. Soundness and completeness are achieved by leveraging the decidability of the underlying mathematical domains they represent; this implies that the extension of these techniques to new classes of properties is often limited by undecidability.

Only in the last decade have dynamic techniques been applied to invariant generation. In a nutshell, the dynamic approach consists of testing a large number of candidate properties against several program runs; the properties that are not violated in any of the runs are retained as “likely” invariants. This implies that the method is not sound but only an “educated guess”: dynamic invariant generation is to static generation what testing is to program proofs. Nonetheless, just like testing is quite effective and useful in practice, dynamic invariant generation can work well if properly implemented. To date, dynamic invariant inference has been mostly used to infer pre- and postconditions or intermediate assertions, whereas it has been only rarely applied to loop invariant generation.

### 10.3.2 Program Synthesis from Invariants

Program synthesis is the task to construct a program that provably satisfies the postcondition. In contrast to program verification, the program is to be constructed rather than given. Both fields make use of formal proof techniques, and both comprise approaches of different degrees of automatization. In contrast to automatic programming techniques, postconditions in program synthesis are usually

non-algorithmic statements in an appropriate logical language.

If we are capable of generating loop invariants from the postcondition automatically, then it is relatively easy to construct the code that satisfies the postcondition. Since different algorithms for solving the same problem require different invariants, generating different invariants lead us to different algorithms.

Suppose we like to sort an array  $A$  of  $n$  elements. Let  $A_0$  denote the original array  $A$ . The postcondition of a sorting algorithm is  $perm(A, A_0) \wedge sorted(A, 0, n)$ , where  $perm(A, A_0)$  is true iff  $A$  is a permutation of  $A_0$  and  $sorted(A, 0, n)$  is true iff the elements of  $A$  from 0 to  $n - 1$  are sorted. Suppose the only way to change  $A$  is by  $swap(A, i, j)$ , which swaps the elements  $A[i]$  and  $A[j]$ , then  $perm(A, A_0)$  will remain true through the algorithm. That is, the triple

$$\{perm(A, A_0)\} swap(A, i, j) \{perm(A, A_0)\}$$

is always true ( $swap(A, i, j)$  represents the code  $z := A[i]; A[i] := A[j]; A[j] := z$ ). Thus, for simplicity, we drop  $perm(A, A_0)$  from our discussion and focus only on  $sorted(A, 0, n)$ .

Intuitively,  $sorted(A, i, j)$  is true iff  $A[i] \leq A[i+1] \leq \dots \leq A[j-1]$ . For verification purpose, the program library does not need an implementation of  $sorted(A, i, j)$ . On the other hand, the programmer and theorem prover needs to know the meaning of  $sorted(A, i, j)$ , which can be defined recursively as follows:

$$\begin{aligned} sorted(A, i, j) &= \text{true if } i \geq j; \\ sorted(A, i, j + 1) &= sorted(A, i, j) \wedge A[j - 1] \leq A[j] \text{ if } i < j + 1. \end{aligned}$$

**Example 10.3.1.** If we wish to generate an invariant for the insertion sort algorithm, the first technique to apply is constant relaxation: Replace  $sorted(A, 0, n)$  by  $k = n \wedge sorted(A, 0, k)$ . Dropping  $k = n$  as the exit condition, the outer loop looks as follows:

```

1  k := 1;
2  while k ≠ n do
3      {I : g > 0 ∧ sorted(A, 0, k)}
      ...
10     k := k + 1
11  od
```

It is easy to see that the exit condition of the outer loop, i.e.,  $k = n$ , will give us the postcondition  $sorted(A, 0, n)$ .

Now apply uncoupling and aging to  $sorted(A, 0, k)$  by introducing a new variable  $j$ , which breaks the first  $k+1$  elements of  $A$  into two sorted lists:  $A[0], \dots, A[j-1]$ , and  $A[j], \dots, A[k]$ . We obtain invariant  $I'$  for the inner loop:

$$I' : 0 \leq j \leq k \wedge sorted(A, 0, j) \wedge sorted(A, j, k+1)$$

which has three conjuncts. The first conjunct is the bounding invariant. When the exit condition of the inner loop, i.e.,  $\neg(j > 0 \wedge A[j-1] > A[j])$  becomes true, we have  $j = 0 \vee A[j-1] \leq A[j]$ . Together with the last two conjuncts of  $I'$ , we can show that  $sorted(A, 0, k+1)$  is true.

```

1  k := 1;
2  while k ≠ n do
3      {I : k > 0 ∧ sorted(A, 0, k)}
4      j := k;
5      while (j > 0 ∧ A[j-1] > A[j]) do
6          {I' : 0 ≤ j ≤ k ∧ sorted(A, 0, j) ∧ sorted(A, j, k+1)}
7          swap(A, j, j-1);
8          j := j-1
9      od
10     k := k+1;
11 od

```

It is easy to see that  $sorted(A, 0, j)$  and  $sorted(A, j, k+1)$  are true when  $j = k$ , because  $sorted(A, 0, j)$  comes from  $I$  and  $sorted(A, j, k+1)$  from the definition of  $sorted$ . To show that  $0 \leq j \leq k$  is true when  $j = k$ , we need to know  $k \geq 0$ . To have this, we may add  $k > 0$  as a bounding invariant into  $I$  of the outer loop. Once this is done,  $I'$  can be shown to be true before entering the inner loop.

To see that  $I'$  is an invariant for the inner loop, when the loop condition is true,  $0 \leq j-1$  is true because  $j > 0$ .  $sorted(A, 0, j-1)$  is true because of  $sorted(A, 0, j)$ ,  $sorted(A, j-1, k+1)$  is true because  $sorted(A, j, k+1)$  and  $A[j-1] < A[j]$  (after swapping). Thus,  $I'$  is indeed a loop invariant.

The exit condition of the inner loop is  $j \leq 0 \vee A[j-1] \leq A[j]$ . If  $j \leq 0$  is true, from  $0 \leq j$ , we know  $j = 0$ , hence  $sorted(A, 0, k+1)$  is true from  $sorted(A, j, k+1)$ . If  $A[j-1] \leq A[j]$  is true, from  $sorted(A, 0, j)$ ,  $A[j-1] < A[j]$ , and  $sorted(A, j, k+1)$ , we conclude also that  $sorted(A, 0, k+1)$  is true. After the execution of line 10,  $k := k+1$ , we arrive at the invariant of the outer loop, i.e.,  $k > 0 \wedge sorted(A, k)$ .

The insertion algorithm in the above example can be regarded as an example of “program synthesis”, if the invariants can be generated automatically.

**Example 10.3.2.** If we are asked to write a program that “store in  $y$  the maximum of  $x$  and  $y$ ”. The postcondition is naturally  $y = \max(x, y)$ . Assume the precondition is  $\top$ , what will be the code  $C$  in  $\{\top\} C \{y = \max(x, y)\}$ . In fact,  $C$  is not unique. For example, some choices are:

- **if**  $x > y$  **then**  $y := x$  **fi**;
- **if**  $x > y$  **then**  $x := y$  **fi**;
- $y := x$ ;
- $x := y$ ;

What will be the “correct” code? With respect to the postcondition, the above four codes are all correct. If the intended meaning is that  $y$  contains the maximum of the original values of  $x$  and  $y$ , in this case, we need auxiliary variables, such as  $x_0$  and  $y_0$ , to denote the original values of  $x$  and  $y$ . The correct postcondition is  $y = \max(x_0, y_0)$  and the first of the above codes is correct with respect to this postcondition. If we ask further that  $y = \max(x, y)$ ,  $x$  and  $y$  hold the original values of  $x$  and  $y$ , then the correct postcondition is  $y = \max(x_0, y_0) \wedge (x = x_0 \wedge y = y_0 \vee x = y_0 \wedge y = x_0)$ , and none of the above codes is correct.

The above simple example illustrates that there exist many codes to achieve the same postcondition; postconditions should be carefully chosen to reflect precisely the requirement. When loops invariants are needed, the complexity of program synthesis grows exponentially. Since there are too many choices to create proper loop invariants, thus leading to various algorithms, program synthesis remains a very challenging task. Here is a quotation from Wikipedia on an annual competition on program synthesis.

The early 21st century has seen a surge of practical interest in the idea of program synthesis in the formal verification community and related fields. In 2013, a unified framework for program synthesis problems was proposed by researchers at UPenn, UC Berkeley, and MIT. Since 2014 there has been a yearly program synthesis competition comparing the different algorithms for program synthesis in a competitive event, the Syntax-Guided Synthesis Competition or SyGuS-Comp. Still, the available algorithms are only able to synthesize small programs.

### 10.3.3 Choosing A Good Language for Assertions

Assertions are written in a first-order language. However, there are many ways to specify assertions, because first-order languages are so expressive and you can use various predicates and functions to specify the same property. There are two possibilities:

- Allow assertions, in particular postconditions and loop invariants, to use functions and predicates defined using some appropriate mechanism (often, the programming language's function declaration construct) to express high-level properties based on a domain theory covering specifics of the application area. We call this approach *domain theory*.
- Disallow the previous possibility, requiring assertions always to be expressed in terms of the constructs of the assertion language, without functions. We call this approach *atomic assertions*.

## 10.4 Exercise Problems

1. Prove formally the following triple is true by choosing a proper loop invariant.

$$\{x = 0\} \mathbf{while} (x < 100) \mathbf{do} x := x + 1 \mathbf{od} \{x = 100\}$$

2. Given the following annotated program, provide the loop invariant and prove its total correctness by Hoare logic. For every simplification step used in the proof, please state clearly which arithmetic/algebraic/logic law is used for simplification.

```

{A : b ≥ 0}
int x = a, y = b, z = 1;
while y ≠ 0 do
  if (y%0 = 1)
    y := y - 1; z := x * z
  else
    x := x * x; y := y/2
{B : z = ab}

```

## Index

- alphabet, 173
- argument, 5
- arity, 174
- assertion, 325
  - middle, 332
- atom
  - atomic formula, 175
- axiom, 23
  - consistent, 23
- Backus–Naur form, 26
- BCP
  - boolean constraint propagation, 115
- BDD
  - binary decision diagram, 56
- Boolean algebra, 22
- Boolean functions, 32
- Boolean variables, 4
- Church–Turing thesis, 28
- clause, 49
  - binary clause, 100
  - definite, 117, 279
  - empty clause, 100
  - fact, 117, 278
  - hard, 160
  - Horn, 117
  - negative, 100
  - positive, 100
  - query, 279
  - rule, 117, 279
  - soft, 160
  - unit clause, 100
- closure
  - congruence, 246
  - equivalence, 246
  - transitive, 246
- CNF
  - 2CNF, 127
  - conjunctive normal form, 49
- completeness, 23
- congruence, 244
  - class, 244
- consistency, 23
- constructor, 258
- correctness
  - partial, 334
  - total, 334
- countable, 16
- decision procedure, 23
- decision procedure
  - semi-decision procedure, 24
- deletion
  - pure, 113
  - subsumption, 113
  - tautology, 113
  - unit, 115
- DNF
  - disjunctive normal form, 51
- domain, 178
- entailment, 43
- equality crossing, 259
- exit condition, 331
- fair strategy, 25, 100
- fallacy, 5

- deductive, 5
- logical, 5
- flattening, 274
- form
  - canonical, 46
  - normal, 46
- formula
  - closed, 177
  - dual, 49
  - first-order, 175
  - instance, 35
  - model, 38
  - satisfiable, 38
  - valid, 39
- function
  - defined, 258
  - domain, 15
  - interpreted, 258
  - range, 15
- ground, 175
- group, 245
- Herbrand
  - base, 201
  - interpretations, 201
  - model, 202
  - universe, 201
- Hilbert system, 94
- Hoare triple, 326
- hyper-resolution, 234
  - negative, 234
- implication
  - contrapositive, 41
  - converse, 41
- independence, 23
- independent, 23
- induction
  - structural rule, 258
- INF
  - ite normal form, 54
- inference rule
  - sound, 24
- inference graph, 99
- inference rule, 24
  - axiom rule, 92
  - sound, 92
- inference system, 24
- injective, 15
- interpretation, 36
  - full, 120
  - partial, 120
- invariant
  - bounding, 347
  - essential, 347
  - loop, 331
- literal, 48
  - decision, 130
  - first-order, 193
  - implied, 130
  - negative, 48
  - positive, 48
- local search, 154
- logic, 1
  - equivalence, 41
  - formal logic, 2
  - Hoare, 325
  - logical connectives, 3
  - logical operator, 3
  - mathematical logic, 1
  - philosophical logic, 2
  - philosophy of logic, 2
  - propositional logic, 4
- logical consequence, 43
- loop, 327
  - invariant, 331

- MaxSAT
  - hybrid, 160
  - partial and weighted, 160
  - solution, 160
  - weighted, 160
- maxterm, 49
- minterm, 51
- model, 22, 184
- modus ponens, 24
- multiset, 224
- NNF
  - negation normal form, 48
- operator
  - minimum, 59
  - sufficient, 59
- order
  - minimal, 15
  - partial order, 15
  - simplification, 224
  - total order, 15
  - well-founded, 15
- PNF
  - prenex normal form, 193
  - PrenexNF, 196
- postcondition, 326
- precedence, 225
- precondition, 326
  - weakest, 329
- predicate, 172
- predicate calculus, 172
- probability logic, 4
- product
  - term, 51
- Prolog
  - fact, 278
  - query, 279
  - rule, 279
- proof, 24, 92
- proof procedure, 23
- proposition, 2
- propositional variable, 4
- prover
  - refutation prover, 85
  - tautology prover, 85
  - theorem prover, 85
- quantifier
  - existential, 173
  - universal, 173
- reason, 120
- reflexing, 265
- relation, 14
- resolution
  - binary, 101
  - core, 102
  - cut, 98
  - hyper-resolution, 234
  - input resolution, 103
  - linear resolution, 103
  - negative resolution, 103
  - ordered resolution, 103
  - positive resolution, 103
  - proof, 102
  - saturation, 110
  - set-of-support, 103
  - unit resolution, 103
- resolvent, 101
- SAT
  - decoder, 149
  - encoder, 149
  - MAX-SAT, 154
- satisfiable, 184
  - equally, 105
- self denotation, 201
- sentence, 177

- sequent, 96
- set, 13
- Skolem
  - function, 195
- soundness, 23
  - weighted clauses, 163
- stable, 224
- statement
  - composed, 3
  - simple, 3
- style
  - deduction, 86
  - enumeration, 86
  - reduction, 86
  - saturation, 109
- substitution, 206
  - affected variable, 206
  - idempotent, 207
  - renaming, 206
- subsumption
  - back, 223
  - forward, 223
- surjective, 15
- symbolic logic, 1
  
- tautology, 39
- term, 174
  - graph, 213
- theorem, 23, 44
  - inductive, 257
- theorem proving, 39
- theory, 44
  - equality, 246
- truth table, 32
- truth values, 3
- Tseitin transformation, 106
  
- unifiable, 207
- unification
  - clash, 209
  - occur check, 209
- unifier, 207
  - mgu, 209
  - more general, 209
- unit deletion, 223
- unit propagation, 115
- universal set, 13
- universe, 178
  
- variable
  - bounded, 177
  - free, 177
  - scope, 177
- verification condition, 337
  
- WCNF
  - weighted CNF, 159
- well-founded
  - induction, 224