
Basics of Compiler Design

Anniversary edition

Torben Ægidius Mogensen



DEPARTMENT OF COMPUTER SCIENCE
UNIVERSITY OF COPENHAGEN

Published through lulu.com.

© Torben Ægidius Mogensen 2000 – 2010

torbenm@diku.dk

Department of Computer Science
University of Copenhagen
Universitetsparken 1
DK-2100 Copenhagen
DENMARK

Book homepage:

<http://www.diku.dk/~torbenm/Basics>

First published 2000

This edition: August 20, 2010

ISBN 978-87-993154-0-6

Contents

1	Introduction	1
1.1	What is a compiler?	1
1.2	The phases of a compiler	2
1.3	Interpreters	3
1.4	Why learn about compilers?	4
1.5	The structure of this book	5
1.6	To the lecturer	6
1.7	Acknowledgements	7
1.8	Permission to use	7
2	Lexical Analysis	9
2.1	Introduction	9
2.2	Regular expressions	10
2.2.1	Shorthands	13
2.2.2	Examples	14
2.3	Nondeterministic finite automata	15
2.4	Converting a regular expression to an NFA	18
2.4.1	Optimisations	20
2.5	Deterministic finite automata	22
2.6	Converting an NFA to a DFA	23
2.6.1	Solving set equations	23
2.6.2	The subset construction	26
2.7	Size versus speed	29
2.8	Minimisation of DFAs	30
2.8.1	Example	32
2.8.2	Dead states	34
2.9	Lexers and lexer generators	35
2.9.1	Lexer generators	41
2.10	Properties of regular languages	42
2.10.1	Relative expressive power	42
2.10.2	Limits to expressive power	44

2.10.3	Closure properties	45
2.11	Further reading	46
	Exercises	46
3	Syntax Analysis	53
3.1	Introduction	53
3.2	Context-free grammars	54
3.2.1	How to write context free grammars	56
3.3	Derivation	58
3.3.1	Syntax trees and ambiguity	60
3.4	Operator precedence	63
3.4.1	Rewriting ambiguous expression grammars	64
3.5	Other sources of ambiguity	66
3.6	Syntax analysis	68
3.7	Predictive parsing	68
3.8	<i>Nullable</i> and <i>FIRST</i>	69
3.9	Predictive parsing revisited	73
3.10	<i>FOLLOW</i>	74
3.11	A larger example	77
3.12	LL(1) parsing	79
3.12.1	Recursive descent	80
3.12.2	Table-driven LL(1) parsing	81
3.12.3	Conflicts	82
3.13	Rewriting a grammar for LL(1) parsing	84
3.13.1	Eliminating left-recursion	84
3.13.2	Left-factorisation	86
3.13.3	Construction of LL(1) parsers summarized	87
3.14	SLR parsing	88
3.15	Constructing SLR parse tables	90
3.15.1	Conflicts in SLR parse-tables	94
3.16	Using precedence rules in LR parse tables	95
3.17	Using LR-parser generators	98
3.17.1	Declarations and actions	99
3.17.2	Abstract syntax	99
3.17.3	Conflict handling in parser generators	102
3.18	Properties of context-free languages	104
3.19	Further reading	105
	Exercises	105

4	Scopes and Symbol Tables	113
4.1	Introduction	113
4.2	Symbol tables	114
4.2.1	Implementation of symbol tables	115
4.2.2	Simple persistent symbol tables	115
4.2.3	A simple imperative symbol table	117
4.2.4	Efficiency issues	117
4.2.5	Shared or separate name spaces	118
4.3	Further reading	118
	Exercises	118
5	Interpretation	121
5.1	Introduction	121
5.2	The structure of an interpreter	122
5.3	A small example language	122
5.4	An interpreter for the example language	124
5.4.1	Evaluating expressions	124
5.4.2	Interpreting function calls	126
5.4.3	Interpreting a program	128
5.5	Advantages and disadvantages of interpretation	128
5.6	Further reading	130
	Exercises	130
6	Type Checking	133
6.1	Introduction	133
6.2	The design space of types	133
6.3	Attributes	135
6.4	Environments for type checking	135
6.5	Type checking expressions	136
6.6	Type checking of function declarations	138
6.7	Type checking a program	139
6.8	Advanced type checking	140
6.9	Further reading	143
	Exercises	143
7	Intermediate-Code Generation	147
7.1	Introduction	147
7.2	Choosing an intermediate language	148
7.3	The intermediate language	150
7.4	Syntax-directed translation	151
7.5	Generating code from expressions	152
7.5.1	Examples of translation	155

7.6	Translating statements	156
7.7	Logical operators	159
7.7.1	Sequential logical operators	160
7.8	Advanced control statements	164
7.9	Translating structured data	165
7.9.1	Floating-point values	165
7.9.2	Arrays	165
7.9.3	Strings	171
7.9.4	Records/structs and unions	171
7.10	Translating declarations	172
7.10.1	Example: Simple local declarations	172
7.11	Further reading	172
	Exercises	173
8	Machine-Code Generation	179
8.1	Introduction	179
8.2	Conditional jumps	180
8.3	Constants	181
8.4	Exploiting complex instructions	181
8.4.1	Two-address instructions	186
8.5	Optimisations	186
8.6	Further reading	188
	Exercises	188
9	Register Allocation	191
9.1	Introduction	191
9.2	Liveness	192
9.3	Liveness analysis	193
9.4	Interference	196
9.5	Register allocation by graph colouring	199
9.6	Spilling	200
9.7	Heuristics	202
9.7.1	Removing redundant moves	205
9.7.2	Using explicit register numbers	205
9.8	Further reading	206
	Exercises	206
10	Function calls	209
10.1	Introduction	209
10.1.1	The call stack	209
10.2	Activation records	210
10.3	Prologues, epilogues and call-sequences	211

10.4	Caller-saves versus callee-saves	213
10.5	Using registers to pass parameters	215
10.6	Interaction with the register allocator	219
10.7	Accessing non-local variables	221
10.7.1	Global variables	221
10.7.2	Call-by-reference parameters	222
10.7.3	Nested scopes	223
10.8	Variants	226
10.8.1	Variable-sized frames	226
10.8.2	Variable number of parameters	227
10.8.3	Direction of stack-growth and position of FP	227
10.8.4	Register stacks	228
10.8.5	Functions as values	228
10.9	Further reading	229
	Exercises	229
11	Analysis and optimisation	231
11.1	Data-flow analysis	232
11.2	Common subexpression elimination	233
11.2.1	Available assignments	233
11.2.2	Example of available-assignments analysis	236
11.2.3	Using available assignment analysis for common subexpression elimination	237
11.3	Jump-to-jump elimination	240
11.4	Index-check elimination	241
11.5	Limitations of data-flow analyses	244
11.6	Loop optimisations	245
11.6.1	Code hoisting	245
11.6.2	Memory prefetching	246
11.7	Optimisations for function calls	248
11.7.1	Inlining	249
11.7.2	Tail-call optimisation	250
11.8	Specialisation	252
11.9	Further reading	254
	Exercises	254
12	Memory management	257
12.1	Introduction	257
12.2	Static allocation	257
12.2.1	Limitations	258
12.3	Stack allocation	258

12.4	Heap allocation	259
12.5	Manual memory management	259
12.5.1	A simple implementation of <code>malloc()</code> and <code>free()</code>	260
12.5.2	Joining freed blocks	263
12.5.3	Sorting by block size	264
12.5.4	Summary of manual memory management	265
12.6	Automatic memory management	266
12.7	Reference counting	266
12.8	Tracing garbage collectors	268
12.8.1	Scan-sweep collection	269
12.8.2	Two-space collection	271
12.8.3	Generational and concurrent collectors	273
12.9	Summary of automatic memory management	276
12.10	Further reading	277
	Exercises	277
13	Bootstrapping a compiler	281
13.1	Introduction	281
13.2	Notation	281
13.3	Compiling compilers	283
13.3.1	Full bootstrap	285
13.4	Further reading	288
	Exercises	288
A	Set notation and concepts	291
A.1	Basic concepts and notation	291
A.1.1	Operations and predicates	291
A.1.2	Properties of set operations	292
A.2	Set-builder notation	293
A.3	Sets of sets	294
A.4	Set equations	295
A.4.1	Monotonic set functions	295
A.4.2	Distributive functions	296
A.4.3	Simultaneous equations	297
	Exercises	297

List of Figures

2.1	Regular expressions	11
2.2	Some algebraic properties of regular expressions	14
2.3	Example of an NFA	17
2.4	Constructing NFA fragments from regular expressions	19
2.5	NFA for the regular expression $(a b)^*ac$	20
2.6	Optimised NFA construction for regular expression shorthands . .	21
2.7	Optimised NFA for $[0-9]^+$	21
2.8	Example of a DFA	22
2.9	DFA constructed from the NFA in figure 2.5	29
2.10	Non-minimal DFA	32
2.11	Minimal DFA	34
2.12	Combined NFA for several tokens	38
2.13	Combined DFA for several tokens	39
2.14	A 4-state NFA that gives 15 DFA states	44
3.1	From regular expressions to context free grammars	56
3.2	Simple expression grammar	57
3.3	Simple statement grammar	57
3.4	Example grammar	59
3.5	Derivation of the string aabbbcc using grammar 3.4	59
3.6	Leftmost derivation of the string aabbbcc using grammar 3.4 . .	59
3.7	Syntax tree for the string aabbbcc using grammar 3.4	61
3.8	Alternative syntax tree for the string aabbbcc using grammar 3.4 .	61
3.9	Unambiguous version of grammar 3.4	62
3.10	Preferred syntax tree for $2+3*4$ using grammar 3.2	63
3.11	Unambiguous expression grammar	66
3.12	Syntax tree for $2+3*4$ using grammar 3.11	67
3.13	Unambiguous grammar for statements	68
3.14	Fixed-point iteration for calculation of <i>Nullable</i>	71
3.15	Fixed-point iteration for calculation of <i>FIRST</i>	72
3.16	Recursive descent parser for grammar 3.9	81

3.17	LL(1) table for grammar 3.9	82
3.18	Program for table-driven LL(1) parsing	83
3.19	Input and stack during table-driven LL(1) parsing	83
3.20	Removing left-recursion from grammar 3.11	85
3.21	Left-factorised grammar for conditionals	87
3.22	SLR table for grammar 3.9	90
3.23	Algorithm for SLR parsing	91
3.24	Example SLR parsing	91
3.25	Example grammar for SLR-table construction	92
3.26	NFAs for the productions in grammar 3.25	92
3.27	Epsilon-transitions added to figure 3.26	93
3.28	SLR DFA for grammar 3.9	94
3.29	Summary of SLR parse-table construction	95
3.30	Textual representation of NFA states	103
5.1	Example language for interpretation	123
5.2	Evaluating expressions	125
5.3	Evaluating a function call	127
5.4	Interpreting a program	128
6.1	The design space of types	134
6.2	Type checking of expressions	137
6.3	Type checking a function declaration	139
6.4	Type checking a program	141
7.1	The intermediate language	150
7.2	A simple expression language	152
7.3	Translating an expression	154
7.4	Statement language	156
7.5	Translation of statements	158
7.6	Translation of simple conditions	159
7.7	Example language with logical operators	161
7.8	Translation of sequential logical operators	162
7.9	Translation for one-dimensional arrays	166
7.10	A two-dimensional array	168
7.11	Translation of multi-dimensional arrays	169
7.12	Translation of simple declarations	173
8.1	Pattern/replacement pairs for a subset of the MIPS instruction set	185
9.1	Gen and kill sets	194
9.2	Example program for liveness analysis and register allocation	195

9.3	<i>succ</i> , <i>gen</i> and <i>kill</i> for the program in figure 9.2	196
9.4	Fixed-point iteration for liveness analysis	197
9.5	Interference graph for the program in figure 9.2	198
9.6	Algorithm 9.3 applied to the graph in figure 9.5	202
9.7	Program from figure 9.2 after spilling variable <i>a</i>	203
9.8	Interference graph for the program in figure 9.7	203
9.9	Colouring of the graph in figure 9.8	204
10.1	Simple activation record layout	211
10.2	Prologue and epilogue for the frame layout shown in figure 10.1	212
10.3	Call sequence for $x := \text{CALL } f(a_1, \dots, a_n)$ using the frame layout shown in figure 10.1	213
10.4	Activation record layout for callee-saves	214
10.5	Prologue and epilogue for callee-saves	214
10.6	Call sequence for $x := \text{CALL } f(a_1, \dots, a_n)$ for callee-saves	215
10.7	Possible division of registers for 16-register architecture	216
10.8	Activation record layout for the register division shown in figure 10.7	216
10.9	Prologue and epilogue for the register division shown in figure 10.7	217
10.10	Call sequence for $x := \text{CALL } f(a_1, \dots, a_n)$ for the register division shown in figure 10.7	218
10.11	Example of nested scopes in Pascal	223
10.12	Adding an explicit frame-pointer to the program from figure 10.11	224
10.13	Activation record with static link	225
10.14	Activation records for <i>f</i> and <i>g</i> from figure 10.11	225
11.1	Gen and kill sets for available assignments	235
11.2	Example program for available-assignments analysis	236
11.3	<i>pred</i> , <i>gen</i> and <i>kill</i> for the program in figure 11.2	237
11.4	Fixed-point iteration for available-assignment analysis	238
11.5	The program in figure 11.2 after common subexpression elimination	239
11.6	Equations for index-check elimination	242
11.7	Intermediate code for for-loop with index check	244
12.1	Operations on a free list	261

Chapter 1

Introduction

1.1 What is a compiler?

In order to reduce the complexity of designing and building computers, nearly all of these are made to execute relatively simple commands (but do so very quickly). A program for a computer must be built by combining these very simple commands into a program in what is called *machine language*. Since this is a tedious and error-prone process most programming is, instead, done using a high-level *programming language*. This language can be very different from the machine language that the computer can execute, so some means of bridging the gap is required. This is where the *compiler* comes in.

A compiler translates (or *compiles*) a program written in a high-level programming language that is suitable for human programmers into the low-level machine language that is required by computers. During this process, the compiler will also attempt to spot and report obvious programmer mistakes.

Using a high-level language for programming has a large impact on how fast programs can be developed. The main reasons for this are:

- Compared to machine language, the notation used by programming languages is closer to the way humans think about problems.
- The compiler can spot some obvious programming mistakes.
- Programs written in a high-level language tend to be shorter than equivalent programs written in machine language.

Another advantage of using a high-level level language is that the same program can be compiled to many different machine languages and, hence, be brought to run on many different machines.

On the other hand, programs that are written in a high-level language and automatically translated to machine language may run somewhat slower than programs that are hand-coded in machine language. Hence, some time-critical programs are still written partly in machine language. A good compiler will, however, be able to get very close to the speed of hand-written machine code when translating well-structured programs.

1.2 The phases of a compiler

Since writing a compiler is a nontrivial task, it is a good idea to structure the work. A typical way of doing this is to split the compilation into several phases with well-defined interfaces. Conceptually, these phases operate in sequence (though in practice, they are often interleaved), each phase (except the first) taking the output from the previous phase as its input. It is common to let each phase be handled by a separate module. Some of these modules are written by hand, while others may be generated from specifications. Often, some of the modules can be shared between several compilers.

A common division into phases is described below. In some compilers, the ordering of phases may differ slightly, some phases may be combined or split into several phases or some extra phases may be inserted between those mentioned below.

Lexical analysis This is the initial part of reading and analysing the program text: The text is read and divided into *tokens*, each of which corresponds to a symbol in the programming language, *e.g.*, a variable name, keyword or number.

Syntax analysis This phase takes the list of tokens produced by the lexical analysis and arranges these in a tree-structure (called the *syntax tree*) that reflects the structure of the program. This phase is often called *parsing*.

Type checking This phase analyses the syntax tree to determine if the program violates certain consistency requirements, *e.g.*, if a variable is used but not declared or if it is used in a context that does not make sense given the type of the variable, such as trying to use a boolean value as a function pointer.

Intermediate code generation The program is translated to a simple machine-independent intermediate language.

Register allocation The symbolic variable names used in the intermediate code are translated to numbers, each of which corresponds to a register in the target machine code.

Machine code generation The intermediate language is translated to assembly language (a textual representation of machine code) for a specific machine architecture.

Assembly and linking The assembly-language code is translated into binary representation and addresses of variables, functions, *etc.*, are determined.

The first three phases are collectively called *the frontend* of the compiler and the last three phases are collectively called *the backend*. The middle part of the compiler is in this context only the intermediate code generation, but this often includes various optimisations and transformations on the intermediate code.

Each phase, through checking and transformation, establishes stronger invariants on the things it passes on to the next, so that writing each subsequent phase is easier than if these have to take all the preceding into account. For example, the type checker can assume absence of syntax errors and the code generation can assume absence of type errors.

Assembly and linking are typically done by programs supplied by the machine or operating system vendor, and are hence not part of the compiler itself, so we will not further discuss these phases in this book.

1.3 Interpreters

An *interpreter* is another way of implementing a programming language. Interpretation shares many aspects with compiling. Lexing, parsing and type-checking are in an interpreter done just as in a compiler. But instead of generating code from the syntax tree, the syntax tree is processed directly to evaluate expressions and execute statements, and so on. An interpreter may need to process the same piece of the syntax tree (for example, the body of a loop) many times and, hence, interpretation is typically slower than executing a compiled program. But writing an interpreter is often simpler than writing a compiler and the interpreter is easier to move to a different machine (see chapter 13), so for applications where speed is not of essence, interpreters are often used.

Compilation and interpretation may be combined to implement a programming language: The compiler may produce intermediate-level code which is then interpreted rather than compiled to machine code. In some systems, there may even be parts of a program that are compiled to machine code, some parts that are compiled to intermediate code, which is interpreted at runtime while other parts may be kept as a syntax tree and interpreted directly. Each choice is a compromise between speed and space: Compiled code tends to be bigger than intermediate code, which tend to be bigger than syntax, but each step of translation improves running speed.

Using an interpreter is also useful during program development, where it is more important to be able to test a program modification quickly rather than run

the program efficiently. And since interpreters do less work on the program before execution starts, they are able to start running the program more quickly. Furthermore, since an interpreter works on a representation that is closer to the source code than is compiled code, error messages can be more precise and informative.

We will discuss interpreters briefly in chapters 5 and 13, but they are not the main focus of this book.

1.4 Why learn about compilers?

Few people will ever be required to write a compiler for a general-purpose language like C, Pascal or SML. So why do most computer science institutions offer compiler courses and often make these mandatory?

Some typical reasons are:

- a) It is considered a topic that you should know in order to be “well-cultured” in computer science.
- b) A good craftsman should know his tools, and compilers are important tools for programmers and computer scientists.
- c) The techniques used for constructing a compiler are useful for other purposes as well.
- d) There is a good chance that a programmer or computer scientist will need to write a compiler or interpreter for a domain-specific language.

The first of these reasons is somewhat dubious, though something can be said for “knowing your roots”, even in such a hastily changing field as computer science.

Reason “b” is more convincing: Understanding how a compiler is built will allow programmers to get an intuition about what their high-level programs will look like when compiled and use this intuition to tune programs for better efficiency. Furthermore, the error reports that compilers provide are often easier to understand when one knows about and understands the different phases of compilation, such as knowing the difference between lexical errors, syntax errors, type errors and so on.

The third reason is also quite valid. In particular, the techniques used for reading (*lexing* and *parsing*) the text of a program and converting this into a form (*abstract syntax*) that is easily manipulated by a computer, can be used to read and manipulate any kind of structured text such as XML documents, address lists, *etc.*

Reason “d” is becoming more and more important as domain specific languages (DSLs) are gaining in popularity. A DSL is a (typically small) language designed for a narrow class of problems. Examples are data-base query languages, text-formatting languages, scene description languages for ray-tracers and languages

for setting up economic simulations. The target language for a compiler for a DSL may be traditional machine code, but it can also be another high-level language for which compilers already exist, a sequence of control signals for a machine, or formatted text and graphics in some printer-control language (*e.g.* PostScript). Even so, all DSL compilers will share similar front-ends for reading and analysing the program text.

Hence, the methods needed to make a compiler front-end are more widely applicable than the methods needed to make a compiler back-end, but the latter is more important for understanding how a program is executed on a machine.

1.5 The structure of this book

The first part of the book describes the methods and tools required to read program text and convert it into a form suitable for computer manipulation. This process is made in two stages: A lexical analysis stage that basically divides the input text into a list of “words”. This is followed by a syntax analysis (or *parsing*) stage that analyses the way these words form structures and converts the text into a data structure that reflects the textual structure. Lexical analysis is covered in chapter 2 and syntactical analysis in chapter 3.

The second part of the book (chapters 4 – 10) covers the middle part and back-end of interpreters and compilers. Chapter 4 covers how definitions and uses of names (*identifiers*) are connected through *symbol tables*. Chapter 5 shows how you can implement a simple programming language by writing an interpreter and notes that this gives a considerable overhead that can be reduced by doing more things before executing the program, which leads to the following chapters about static type checking (chapter 6) and compilation (chapters 7 – 10). In chapter 7, it is shown how expressions and statements can be compiled into an *intermediate language*, a language that is close to machine language but hides machine-specific details. In chapter 8, it is discussed how the intermediate language can be converted into “real” machine code. Doing this well requires that the registers in the processor are used to store the values of variables, which is achieved by a *register allocation* process, as described in chapter 9. Up to this point, a “program” has been what corresponds to the body of a single procedure. Procedure calls and nested procedure declarations add some issues, which are discussed in chapter 10. Chapter 11 deals with analysis and optimisation and chapter 12 is about allocating and freeing memory. Finally, chapter 13 will discuss the process of *bootstrapping* a compiler, *i.e.*, using a compiler to compile itself.

The book uses standard set notation and equations over sets. Appendix A contains a short summary of these, which may be helpful to those that need these concepts refreshed.

Chapter 11 (on analysis and optimisation) was added in 2008 and chapter 5

(about interpreters) was added in 2009, which is why editions after April 2008 are called “extended”. In the 2010 edition, further additions (including chapter 12 and appendix A) were made. Since ten years have passed since the first edition was printed as lecture notes, the 2010 edition is labeled “anniversary edition”.

1.6 To the lecturer

This book was written for use in the introductory compiler course at DIKU, the department of computer science at the University of Copenhagen, Denmark.

At DIKU, the compiler course was previously taught right after the introductory programming course, which is earlier than in most other universities. Hence, existing textbooks tended either to be too advanced for the level of the course or be too simplistic in their approach, for example only describing a single very simple compiler without bothering too much with the general picture.

This book was written as a response to this and aims at bridging the gap: It is intended to convey the general picture without going into extreme detail about such things as efficient implementation or the newest techniques. It should give the students an understanding of how compilers work and the ability to make simple (but not simplistic) compilers for simple languages. It will also lay a foundation that can be used for studying more advanced compilation techniques, as found *e.g.* in [35]. The compiler course at DIKU was later moved to the second year, so additions to the original text has been made.

At times, standard techniques from compiler construction have been simplified for presentation in this book. In such cases references are made to books or articles where the full version of the techniques can be found.

The book aims at being “language neutral”. This means two things:

- Little detail is given about how the methods in the book can be implemented in any specific language. Rather, the description of the methods is given in the form of algorithm sketches and textual suggestions of how these can be implemented in various types of languages, in particular imperative and functional languages.
- There is no single through-going example of a language to be compiled. Instead, different small (sub-)languages are used in various places to cover exactly the points that the text needs. This is done to avoid drowning in detail, hopefully allowing the readers to “see the wood for the trees”.

Each chapter (except this) has a section on further reading, which suggests additional reading material for interested students. All chapters (also except this) has a set of exercises. Few of these require access to a computer, but can be solved on paper or black-board. In fact, many of the exercises are based on exercises that

have been used in written exams at DIKU. After some of the sections in the book, a few easy exercises are listed. It is suggested that the student attempts to solve these exercises before continuing reading, as the exercises support understanding of the previous sections.

Teaching with this book can be supplemented with project work, where students write simple compilers. Since the book is language neutral, no specific project is given. Instead, the teacher must choose relevant tools and select a project that fits the level of the students and the time available. Depending on how much of the book is used and the amount of project work, the book can support course sizes ranging from 5 to 15 ECTS points.

1.7 Acknowledgements

The author wishes to thank all people who have been helpful in making this book a reality. This includes the students who have been exposed to draft versions of the book at the compiler courses “Dat 1E” and “Oversættelse” at DIKU, and who have found numerous typos and other errors in the earlier versions. I would also like to thank the instructors at Dat 1E and Oversættelse, who have pointed out places where things were not as clear as they could be. I am extremely grateful to the people who in 2000 read parts of or all of the first draft and made helpful suggestions.

1.8 Permission to use

Permission to copy and print for personal use is granted. If you, as a lecturer, want to print the book and sell it to your students, you can do so if you only charge the printing cost. If you want to print the book and sell it at profit, please contact the author at torbenm@diku.dk and we will find a suitable arrangement.

In all cases, if you find any misprints or other errors, please contact the author at torbenm@diku.dk.

See also the book homepage: <http://www.diku.dk/~torbenm/Basics>.

Chapter 2

Lexical Analysis

2.1 Introduction

The word “lexical” in the traditional sense means “pertaining to words”. In terms of programming languages, words are objects like variable names, numbers, keywords *etc.* Such words are traditionally called *tokens*.

A *lexical analyser*, or *lexer* for short, will as its input take a string of individual letters and divide this string into tokens. Additionally, it will filter out whatever separates the tokens (the so-called *white-space*), *i.e.*, lay-out characters (spaces, newlines *etc.*) and comments.

The main purpose of lexical analysis is to make life easier for the subsequent syntax analysis phase. In theory, the work that is done during lexical analysis can be made an integral part of syntax analysis, and in simple systems this is indeed often done. However, there are reasons for keeping the phases separate:

- **Efficiency:** A lexer may do the simple parts of the work faster than the more general parser can. Furthermore, the size of a system that is split in two may be smaller than a combined system. This may seem paradoxical but, as we shall see, there is a non-linear factor involved which may make a separated system smaller than a combined system.
- **Modularity:** The syntactical description of the language need not be cluttered with small lexical details such as white-space and comments.
- **Tradition:** Languages are often designed with separate lexical and syntactical phases in mind, and the standard documents of such languages typically separate lexical and syntactical elements of the languages.

It is usually not terribly difficult to write a lexer by hand: You first read past initial white-space, then you, in sequence, test to see if the next token is a keyword, a

number, a variable or whatnot. However, this is not a very good way of handling the problem: You may read the same part of the input repeatedly while testing each possible token and in some cases it may not be clear where the next token ends. Furthermore, a handwritten lexer may be complex and difficult to maintain. Hence, lexers are normally constructed by *lexer generators*, which transform human-readable specifications of tokens and white-space into efficient programs.

We will see the same general strategy in the chapter about syntax analysis: Specifications in a well-defined human-readable notation are transformed into efficient programs.

For lexical analysis, specifications are traditionally written using *regular expressions*: An algebraic notation for describing sets of strings. The generated lexers are in a class of extremely simple programs called *finite automata*.

This chapter will describe regular expressions and finite automata, their properties and how regular expressions can be converted to finite automata. Finally, we discuss some practical aspects of lexer generators.

2.2 Regular expressions

The set of all integer constants or the set of all variable names are sets of strings, where the individual letters are taken from a particular alphabet. Such a set of strings is called a *language*. For integers, the alphabet consists of the digits 0-9 and for variable names the alphabet contains both letters and digits (and perhaps a few other characters, such as underscore).

Given an alphabet, we will describe sets of strings by *regular expressions*, an algebraic notation that is compact and easy for humans to use and understand. The idea is that regular expressions that describe simple sets of strings can be combined to form regular expressions that describe more complex sets of strings.

When talking about regular expressions, we will use the letters (*r*, *s* and *t*) in italics to denote unspecified regular expressions. When letters stand for themselves (*i.e.*, in regular expressions that describe strings that use these letters) we will use typewriter font, *e.g.*, a or b. Hence, when we say, *e.g.*, “The regular expression *s*” we mean the regular expression that describes a single one-letter string “s”, but when we say “The regular expression *s*”, we mean a regular expression of any form which we just happen to call *s*. We use the notation $L(s)$ to denote the language (*i.e.*, set of strings) described by the regular expression *s*. For example, $L(a)$ is the set {“a”}.

Figure 2.1 shows the constructions used to build regular expressions and the languages they describe:

- A single letter describes the language that has the one-letter string consisting of that letter as its only element.

Regular expression	Language (set of strings)	Informal description
a	$\{“a”\}$	The set consisting of the one-letter string “a”.
ϵ	$\{“”\}$	The set containing the empty string.
$s t$	$L(s) \cup L(t)$	Strings from both languages
st	$\{vw \mid v \in L(s), w \in L(t)\}$	Strings constructed by concatenating a string from the first language with a string from the second language. Note: In set-formulas, “ ” is not a part of a regular expression, but part of the set-builder notation and reads as “where”.
s^*	$\{“”\} \cup \{vw \mid v \in L(s), w \in L(s^*)\}$	Each string in the language is a concatenation of any number of strings in the language of s .

Figure 2.1: Regular expressions

- The symbol ϵ (the Greek letter *epsilon*) describes the language that consists solely of the empty string. Note that this is not the empty set of strings (see exercise 2.10).
- $s|t$ (pronounced “ s or t ”) describes the union of the languages described by s and t .
- st (pronounced “ s t ”) describes the concatenation of the languages $L(s)$ and $L(t)$, *i.e.*, the sets of strings obtained by taking a string from $L(s)$ and putting this in front of a string from $L(t)$. For example, if $L(s)$ is {“a”, “b”} and $L(t)$ is {“c”, “d”}, then $L(st)$ is the set {“ac”, “ad”, “bc”, “bd”}.
- The language for s^* (pronounced “ s star”) is described recursively: It consists of the empty string plus whatever can be obtained by concatenating a string from $L(s)$ to a string from $L(s^*)$. This is equivalent to saying that $L(s^*)$ consists of strings that can be obtained by concatenating zero or more (possibly different) strings from $L(s)$. If, for example, $L(s)$ is {“a”, “b”} then $L(s^*)$ is {“”, “a”, “b”, “aa”, “ab”, “ba”, “bb”, “aaa”, ...}, *i.e.*, any string (including the empty) that consists entirely of as and bs.

Note that while we use the same notation for concrete strings and regular expressions denoting one-string languages, the context will make it clear which is meant. We will often show strings and sets of strings without using quotation marks, *e.g.*, write {a, bb} instead of {“a”, “bb”}. When doing so, we will use ϵ to denote the empty string, so the example from $L(s^*)$ above is written as { ϵ , a, b, aa, ab, ba, bb, aaa, ...}. The letters u , v and w in italics will be used to denote unspecified single strings, *i.e.*, members of some language. As an example, abw denotes any string starting with ab.

Precedence rules

When we combine different constructor symbols, *e.g.*, in the regular expression $a|ab^*$, it is not *a priori* clear how the different subexpressions are grouped. We can use parentheses to make the grouping of symbols explicit such as in $(a|(ab))^*$. Additionally, we use precedence rules, similar to the algebraic convention that $3 + 4 * 5$ means 3 added to the product of 4 and 5 and not multiplying the sum of 3 and 4 by 5. For regular expressions, we use the following conventions: $*$ binds tighter than concatenation, which binds tighter than alternative ($|$). The example $a|ab^*$ from above, hence, is equivalent to $a|(a(b^*))$.

The $|$ operator is associative and commutative (as it corresponds to set union, which has these properties). Concatenation is associative (but obviously not commutative) and distributes over $|$. Figure 2.2 shows these and other algebraic prop-

erties of regular expressions, including definitions of some of the shorthands introduced below.

2.2.1 Shorthands

While the constructions in figure 2.1 suffice to describe *e.g.*, number strings and variable names, we will often use extra shorthands for convenience. For example, if we want to describe non-negative integer constants, we can do so by saying that it is one or more digits, which is expressed by the regular expression

$$(0|1|2|3|4|5|6|7|8|9)(0|1|2|3|4|5|6|7|8|9)^*$$

The large number of different digits makes this expression rather verbose. It gets even worse when we get to variable names, where we must enumerate all alphabetic letters (in both upper and lower case).

Hence, we introduce a shorthand for sets of letters. Sequences of letters within square brackets represent the set of these letters. For example, we use `[ab01]` as a shorthand for `a|b|0|1`. Additionally, we can use interval notation to abbreviate `[0123456789]` to `[0-9]`. We can combine several intervals within one bracket and for example write `[a-zA-Z]` to denote all alphabetic letters in both lower and upper case.

When using intervals, we must be aware of the ordering for the symbols involved. For the digits and letters used above, there is usually no confusion. However, if we write, *e.g.*, `[0-z]` it is not immediately clear what is meant. When using such notation in lexer generators, standard ASCII or ISO 8859-1 character sets are usually used, with the hereby implied ordering of symbols. To avoid confusion, we will use the interval notation only for intervals of digits or alphabetic letters.

Getting back to the example of integer constants above, we can now write this much shorter as `[0-9][0-9]*`.

Since s^* denotes *zero or more* occurrences of s , we needed to write the set of digits twice to describe that *one or more* digits are allowed. Such non-zero repetition is quite common, so we introduce another shorthand, s^+ , to denote one or more occurrences of s . With this notation, we can abbreviate our description of integers to `[0-9]+`. On a similar note, it is common that we can have zero or one occurrence of something (*e.g.*, an optional sign to a number). Hence we introduce the shorthand $s?$ for $s|\epsilon$. $^+$ and $?$ bind with the same precedence as * .

We must stress that these shorthands are just that. They do not add anything to the set of languages we can describe, they just make it possible to describe a language more compactly. In the case of s^+ , it can even make an exponential difference: If $^+$ is nested n deep, recursive expansion of s^+ to ss^* yields $2^n - 1$ occurrences of * in the expanded regular expression.

$(r s) t = r s t = r (s t)$	is associative
$s t = t s$	is commutative
$s s = s$	is idempotent
$s? = s \epsilon$	by definition
$(rs)t = rst = r(st)$	concatenation is associative
$s\epsilon = s = \epsilon s$	ϵ is a neutral element for concatenation
$r(s t) = rs rt$	concatenation distributes over
$(r s)t = rt st$	concatenation distributes over
$(s^*)^* = s^*$	* is idempotent
$s^*s^* = s^*$	0 or more twice is still 0 or more
$ss^* = s^+ = s^*s$	by definition

Figure 2.2: Some algebraic properties of regular expressions

2.2.2 Examples

We have already seen how we can describe non-negative integer constants using regular expressions. Here are a few examples of other typical programming language elements:

Keywords. A keyword like `if` is described by a regular expression that looks exactly like that keyword, *e.g.*, the regular expression `if` (which is the concatenation of the two regular expressions `i` and `f`).

Variable names. In the programming language C, a variable name consists of letters, digits and the underscore symbol and it must begin with a letter or underscore. This can be described by the regular expression $[a-zA-Z_][a-zA-Z_0-9]^*$.

Integers. An integer constant is an optional sign followed by a non-empty sequence of digits: $[+]?[0-9]^+$. In some languages, the sign is a separate symbol and not part of the constant itself. This will allow whitespace between the sign and the number, which is not possible with the above.

Floats. A floating-point constant can have an optional sign. After this, the mantissa part is described as a sequence of digits followed by a decimal point and then

another sequence of digits. Either one (but not both) of the digit sequences can be empty. Finally, there is an optional exponent part, which is the letter e (in upper or lower case) followed by an (optionally signed) integer constant. If there is an exponent part to the constant, the mantissa part can be written as an integer constant (*i.e.*, without the decimal point). Some examples:

3.14 -3. .23 3e+4 11.22e-3.

This rather involved format can be described by the following regular expression:

$$[+-]?((([0-9]^+ \cdot [0-9]^* \cdot [0-9]^+) ([eE] [+-]? [0-9]^+) ?) | [0-9]^+ [eE] [+-]? [0-9]^+)$$

This regular expression is complicated by the fact that the exponent is optional if the mantissa contains a decimal point, but not if it does not (as that would make the number an integer constant). We can make the description simpler if we make the regular expression for floats also include integers, and instead use other means of distinguishing integers from floats (see section 2.9 for details). If we do this, the regular expression can be simplified to

$$[+-]?((([0-9]^+ (\cdot [0-9]^*) ? \cdot [0-9]^+) ([eE] [+-]? [0-9]^+) ?)$$

String constants. A string constant starts with a quotation mark followed by a sequence of symbols and finally another quotation mark. There are usually some restrictions on the symbols allowed between the quotation marks. For example, line-feed characters or quotes are typically not allowed, though these may be represented by special “escape” sequences of other characters, such as “\n\n” for a string containing two line-feeds. As a (much simplified) example, we can by the following regular expression describe string constants where the allowed symbols are alphanumeric characters and sequences consisting of the backslash symbol followed by a letter (where each such pair is intended to represent a non-alphanumeric symbol):

$$"([a-zA-Z0-9] | [a-zA-Z])^*"$$

Suggested exercises: 2.1, 2.10(a).

2.3 Nondeterministic finite automata

In our quest to transform regular expressions into efficient programs, we use a stepping stone: Nondeterministic finite automata. By their nondeterministic nature, these are not quite as close to “real machines” as we would like, so we will later see how these can be transformed into *deterministic* finite automata, which are easily and efficiently executable on normal hardware.

A finite automaton is, in the abstract sense, a machine that has a finite number of *states* and a finite number of *transitions* between these. A transition between states is usually labelled by a character from the input alphabet, but we will also use transitions marked with ϵ , the so-called *epsilon transitions*.

A finite automaton can be used to decide if an input string is a member in some particular set of strings. To do this, we select one of the states of the automaton as the *starting state*. We start in this state and in each step, we can do one of the following:

- Follow an epsilon transition to another state, or
- Read a character from the input and follow a transition labelled by that character.

When all characters from the input are read, we see if the current state is marked as being *accepting*. If so, the string we have read from the input is in the language defined by the automaton.

We may have a choice of several actions at each step: We can choose between either an epsilon transition or a transition on an alphabet character, and if there are several transitions with the same symbol, we can choose between these. This makes the automaton *nondeterministic*, as the choice of action is not determined solely by looking at the current state and input. It may be that some choices lead to an accepting state while others do not. This does, however, not mean that the string is sometimes in the language and sometimes not: We will include a string in the language if it is *possible* to make a sequence of choices that makes the string lead to an accepting state.

You can think of it as solving a maze with symbols written in the corridors. If you can find the exit while walking over the letters of the string in the correct order, the string is recognized by the maze.

We can formally define a nondeterministic finite automaton by:

Definition 2.1 A nondeterministic finite automaton consists of a set S of states. One of these states, $s_0 \in S$, is called the starting state of the automaton and a subset $F \subseteq S$ of the states are accepting states. Additionally, we have a set T of transitions. Each transition t connects a pair of states s_1 and s_2 and is labelled with a symbol, which is either a character c from the alphabet Σ , or the symbol ϵ , which indicates an epsilon-transition. A transition from state s to state t on the symbol c is written as $s^c t$.

Starting states are sometimes called *initial states* and accepting states can also be called *final states* (which is why we use the letter F to denote the set of accepting states). We use the abbreviations FA for finite automaton, NFA for nondeterministic finite automaton and (later in this chapter) DFA for deterministic finite automaton.

We will mostly use a graphical notation to describe finite automata. States are denoted by circles, possibly containing a number or name that identifies the state. This name or number has, however, no operational significance, it is solely used for identification purposes. Accepting states are denoted by using a double circle instead of a single circle. The initial state is marked by an arrow pointing to it from outside the automaton.

A transition is denoted by an arrow connecting two states. Near its midpoint, the arrow is labelled by the symbol (possibly ϵ) that triggers the transition. Note that the arrow that marks the initial state is *not* a transition and is, hence, not marked by a symbol.

Repeating the maze analogue, the circles (states) are rooms and the arrows (transitions) are one-way corridors. The double circles (accepting states) are exits, while the unmarked arrow to the starting state is the entrance to the maze.

Figure 2.3 shows an example of a nondeterministic finite automaton having three states. State 1 is the starting state and state 3 is accepting. There is an epsilon-transition from state 1 to state 2, transitions on the symbol *a* from state 2 to states 1 and 3 and a transition on the symbol *b* from state 1 to state 3. This NFA recognises the language described by the regular expression $a^*(a|b)$. As an example, the string *aab* is recognised by the following sequence of transitions:

from	to	by
1	2	ϵ
2	1	<i>a</i>
1	2	ϵ
2	1	<i>a</i>
1	3	<i>b</i>

At the end of the input we are in state 3, which is accepting. Hence, the string is accepted by the NFA. You can check this by placing a coin at the starting state and follow the transitions by moving the coin.

Note that we sometimes have a choice of several transitions. If we are in state

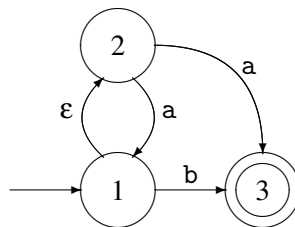


Figure 2.3: Example of an NFA

2 and the next symbol is an a, we can, when reading this, either go to state 1 or to state 3. Likewise, if we are in state 1 and the next symbol is a b, we can either read this and go to state 3 or we can use the epsilon transition to go directly to state 2 without reading anything. If we in the example above had chosen to follow the a-transition to state 3 instead of state 1, we would have been stuck: We would have no legal transition and yet we would not be at the end of the input. But, as previously stated, it is enough that there *exists* a path leading to acceptance, so the string aab is still accepted.

A program that decides if a string is accepted by a given NFA will have to check all possible paths to see if *any* of these accepts the string. This requires either backtracking until a successful path found or simultaneously following all possible paths, both of which are too time-consuming to make NFAs suitable for efficient recognisers. We will, hence, use NFAs only as a stepping stone between regular expressions and the more efficient DFAs. We use this stepping stone because it makes the construction simpler than direct construction of a DFA from a regular expression.

2.4 Converting a regular expression to an NFA

We will construct an NFA *compositionally* from a regular expression, *i.e.*, we will construct the NFA for a composite regular expression from the NFAs constructed from its subexpressions.

To be precise, we will from each subexpression construct an *NFA fragment* and then combine these fragments into bigger fragments. A fragment is not a complete NFA, so we complete the construction by adding the necessary components to make a complete NFA.

An NFA fragment consists of a number of states with transitions between these and additionally two incomplete transitions: One pointing into the fragment and one pointing out of the fragment. The incoming half-transition is not labelled by a symbol, but the outgoing half-transition is labelled by either ϵ or an alphabet symbol. These half-transitions are the entry and exit to the fragment and are used to connect it to other fragments or additional “glue” states.

Construction of NFA fragments for regular expressions is shown in figure 2.4. The construction follows the structure of the regular expression by first making NFA fragments for the subexpressions and then joining these to form an NFA fragment for the whole regular expression. The NFA fragments for the subexpressions are shown as dotted ovals with the incoming half-transition on the left and the outgoing half-transition on the right.

When an NFA fragment has been constructed for the whole regular expression, the construction is completed by connecting the outgoing half-transition to an accepting state. The incoming half-transition serves to identify the starting state of

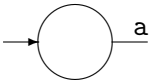
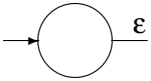
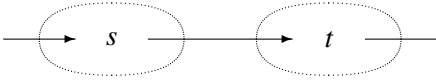
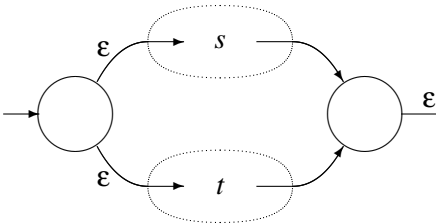
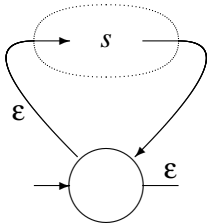
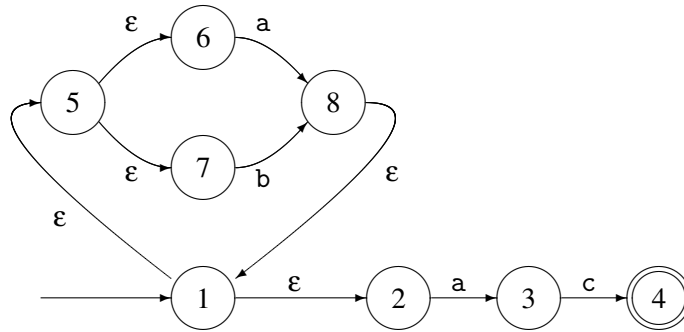
Regular expression	NFA fragment
a	
ϵ	
st	
$s t$	
s^*	

Figure 2.4: Constructing NFA fragments from regular expressions

Figure 2.5: NFA for the regular expression $(a|b)^*ac$

the completed NFA. Note that even though we allow an NFA to have several accepting states, an NFA constructed using this method will have only one: the one added at the end of the construction.

An NFA constructed this way for the regular expression $(a|b)^*ac$ is shown in figure 2.5. We have numbered the states for future reference.

2.4.1 Optimisations

We can use the construction in figure 2.4 for any regular expression by expanding out all shorthand, *e.g.* converting s^+ to ss^* , $[0-9]$ to $0|1|2|\dots|9$ and $s?$ to $s|\epsilon$, *etc.* However, this will result in very large NFAs for some expressions, so we use a few optimised constructions for the shorthands. Additionally, we show an alternative construction for the regular expression ϵ . This construction does not quite follow the formula used in figure 2.4, as it does not have two half-transitions. Rather, the line-segment notation is intended to indicate that the NFA fragment for ϵ just connects the half-transitions of the NFA fragments that it is combined with. In the construction for $[0-9]$, the vertical ellipsis is meant to indicate that there is a transition for each of the digits in $[0-9]$. This construction generalises in the obvious way to other sets of characters, *e.g.*, $[a-zA-Z0-9]$. We have not shown a special construction for $s?$ as $s|\epsilon$ will do fine if we use the optimised construction for ϵ .

The optimised constructions are shown in figure 2.6. As an example, an NFA for $[0-9]^+$ is shown in figure 2.7. Note that while this is *optimised*, it is not *optimal*. You can make an NFA for this language using only two states.

Suggested exercises: 2.2(a), 2.10(b).

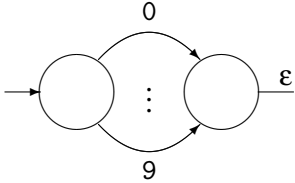
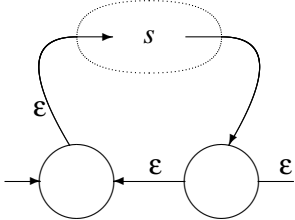
Regular expression	NFA fragment
ϵ	—
$[0-9]$	
s^+	

Figure 2.6: Optimised NFA construction for regular expression shorthands

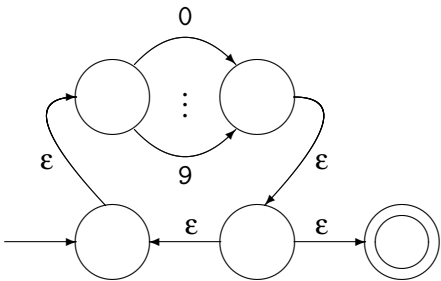


Figure 2.7: Optimised NFA for $[0-9]^+$

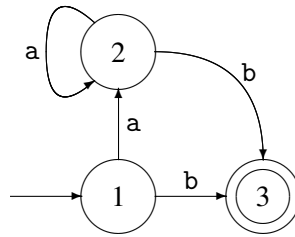


Figure 2.8: Example of a DFA

2.5 Deterministic finite automata

Nondeterministic automata are, as mentioned earlier, not quite as close to “the machine” as we would like. Hence, we now introduce a more restricted form of finite automaton: The deterministic finite automaton, or DFA for short. DFAs are NFAs, but obey a number of additional restrictions:

- There are no epsilon-transitions.
- There may not be two identically labelled transitions out of the same state.

This means that we never have a choice of several next-states: The state and the next input symbol uniquely determine the transition (or lack of same). This is why these automata are called *deterministic*. Figure 2.8 shows a DFA equivalent to the NFA in figure 2.3.

The transition relation of a DFA is a (partial) function, and we often write it as such: $move(s, c)$ is the state (if any) that is reached from state s by a transition on the symbol c . If there is no such transition, $move(s, c)$ is undefined.

It is very easy to implement a DFA: A two-dimensional table can be cross-indexed by state and symbol to yield the next state (or an indication that there is no transition), essentially implementing the *move* function by table lookup. Another (one-dimensional) table can indicate which states are accepting.

DFAs have the same expressive power as NFAs: A DFA is a special case of NFA and any NFA can (as we shall shortly see) be converted to an equivalent DFA. However, this comes at a cost: The resulting DFA can be exponentially larger than the NFA (see section 2.10). In practice (*i.e.*, when describing tokens for a programming language) the increase in size is usually modest, which is why most lexical analysers are based on DFAs.

Suggested exercises: 2.7(a,b), 2.8.

2.6 Converting an NFA to a DFA

As promised, we will show how NFAs can be converted to DFAs such that we, by combining this with the conversion of regular expressions to NFAs shown in section 2.4, can convert any regular expression to a DFA.

The conversion is done by simulating all possible paths in an NFA at once. This means that we operate with sets of NFA states: When we have several choices of a next state, we take all of the choices simultaneously and form a set of the possible next-states. The idea is that such a set of NFA states will become a single DFA state. For any given symbol we form the set of all possible next-states in the NFA, so we get a single transition (labelled by that symbol) going from one set of NFA states to another set. Hence, the transition becomes deterministic in the DFA that is formed from the sets of NFA states.

Epsilon-transitions complicate the construction a bit: Whenever we are in an NFA state we can always choose to follow an epsilon-transition without reading any symbol. Hence, given a symbol, a next-state can be found by either following a transition with that symbol or by first doing any number of epsilon-transitions and then a transition with the symbol. We handle this in the construction by first extending the set of NFA states with those you can reach from these using only epsilon-transitions. Then, for each possible input symbol, we follow transitions with this symbol to form a new set of NFA states. We define the *epsilon-closure* of a set of states as the set extended with all states that can be reached from these using any number of epsilon-transitions. More formally:

Definition 2.2 *Given a set M of NFA states, we define $\varepsilon\text{-closure}(M)$ to be the least (in terms of the subset relation) solution to the set equation*

$$\begin{aligned}\varepsilon\text{-closure}(M) \\ = M \cup \{t \mid s \in \varepsilon\text{-closure}(M) \text{ and } s^{\varepsilon}t \in T\}\end{aligned}$$

Where T is the set of transitions in the NFA.

We will later on see several examples of *set equations* like the one above, so we use some time to discuss how such equations can be solved.

2.6.1 Solving set equations

The following is a very brief description of how to solve set equations like the above. If you find it confusing, you can read appendix A and in particular section A.4 first.

In general, a set equation over a single set-valued variable X has the form

$$X = F(X)$$

where F is a function from sets to sets. Not all such equations are solvable, so we will restrict ourselves to special cases, which we will describe below. We will use calculation of epsilon-closure as the driving example.

In definition 2.2, $\varepsilon\text{-closure}(M)$ is the value we have to find, so we make an equation such that the value of X that solves the equation will be $\varepsilon\text{-closure}(M)$:

$$X = M \cup \{t \mid s \in X \text{ and } s^\varepsilon t \in T\}$$

So, if we define F_M to be

$$F_M(X) = M \cup \{t \mid s \in X \text{ and } s^\varepsilon t \in T\}$$

then a solution to the equation $X = F_M(X)$ will be $\varepsilon\text{-closure}(M)$.

F_M has a property that is essential to our solution method: If $X \subseteq Y$ then $F_M(X) \subseteq F_M(Y)$. We say that F_M is *monotonic*.

There may be several solutions to the equation $X = F_M(X)$. For example, if the NFA has a pair of states that connect to each other by epsilon transitions, adding this pair to a solution that does not already include the pair will create a new solution. The epsilon-closure of M is the *least* solution to the equation (*i.e.*, the smallest X that satisfies the equation).

When we have an equation of the form $X = F(X)$ and F is monotonic, we can find the least solution to the equation in the following way: We first guess that the solution is the empty set and check to see if we are right: We compare \emptyset with $F(\emptyset)$. If these are equal, we are done and \emptyset is the solution. If not, we use the following properties:

- The least solution S to the equation satisfies $S = F(S)$.
- $\emptyset \subseteq S$ implies that $F(\emptyset) \subseteq F(S)$.

to conclude that $F(\emptyset) \subseteq S$. Hence, $F(\emptyset)$ is a new guess at S . We now form the chain

$$\emptyset \subseteq F(\emptyset) \subseteq F(F(\emptyset)) \subseteq \dots$$

If at any point an element in the sequence is identical to the previous, we have a fixed-point, *i.e.*, a set S such that $S = F(S)$. This fixed-point of the sequence will be the least (in terms of set inclusion) solution to the equation. This is not difficult to verify, but we will omit the details. Since we are iterating a function until we reach a fixed-point, we call this process *fixed-point iteration*.

If we are working with sets over a finite domain (*e.g.*, sets of NFA states), we *will* eventually reach a fixed-point, as there can be no infinite chain of strictly increasing sets.

We can use this method for calculating the epsilon-closure of the set $\{1\}$ with respect to the NFA shown in figure 2.5. Since we want to find $\varepsilon\text{-closure}(\{1\})$, $M = \{1\}$, so $F_M = F_{\{1\}}$. We start by guessing the empty set:

$$\begin{aligned} F_{\{1\}}(\emptyset) &= \{1\} \cup \{t \mid s \in \emptyset \text{ and } s^\varepsilon t \in T\} \\ &= \{1\} \end{aligned}$$

As $\emptyset \neq \{1\}$, we continue.

$$\begin{aligned} F_{\{1\}}(\{1\}) &= \{1\} \cup \{t \mid s \in \{1\} \text{ and } s^\varepsilon t \in T\} \\ &= \{1\} \cup \{2, 5\} = \{1, 2, 5\} \end{aligned}$$

$$\begin{aligned} F_{\{1\}}(\{1, 2, 5\}) &= \{1\} \cup \{t \mid s \in \{1, 2, 5\} \text{ and } s^\varepsilon t \in T\} \\ &= \{1\} \cup \{2, 5, 6, 7\} = \{1, 2, 5, 6, 7\} \end{aligned}$$

$$\begin{aligned} F_{\{1\}}(\{1, 2, 5, 6, 7\}) &= \{1\} \cup \{t \mid s \in \{1, 2, 5, 6, 7\} \text{ and } s^\varepsilon t \in T\} \\ &= \{1\} \cup \{2, 5, 6, 7\} = \{1, 2, 5, 6, 7\} \end{aligned}$$

We have now reached a fixed-point and found our solution. Hence, we conclude that $\varepsilon\text{-closure}(\{1\}) = \{1, 2, 5, 6, 7\}$.

We have done a good deal of repeated calculation in the iteration above: We have calculated the epsilon-transitions from state 1 three times and those from state 2 and 5 twice each. We can make an optimised fixed-point iteration by exploiting that the function is not only monotonic, but also *distributive*: $F(X \cup Y) = F(X) \cup F(Y)$. This means that, when we during the iteration add elements to our set, we in the next iteration need only calculate F for the new elements and add the result to the set. In the example above, we get

$$\begin{aligned} F_{\{1\}}(\emptyset) &= \{1\} \cup \{t \mid s \in \emptyset \text{ and } s^\varepsilon t \in T\} \\ &= \{1\} \end{aligned}$$

$$\begin{aligned} F_{\{1\}}(\{1\}) &= \{1\} \cup \{t \mid s \in \{1\} \text{ and } s^\varepsilon t \in T\} \\ &= \{1\} \cup \{2, 5\} = \{1, 2, 5\} \end{aligned}$$

$$\begin{aligned} F_{\{1\}}(\{1, 2, 5\}) &= F(\{1\}) \cup F(\{2, 5\}) \\ &= \{1, 2, 5\} \cup (\{1\} \cup \{t \mid s \in \{2, 5\} \text{ and } s^\varepsilon t \in T\}) \\ &= \{1, 2, 5\} \cup (\{1\} \cup \{6, 7\}) = \{1, 2, 5, 6, 7\} \end{aligned}$$

$$\begin{aligned} F_{\{1\}}(\{1, 2, 5, 6, 7\}) &= F(\{1, 2, 5\}) \cup F_{\{1\}}(\{6, 7\}) \\ &= \{1, 2, 5, 6, 7\} \cup (\{1\} \cup \{t \mid s \in \{6, 7\} \text{ and } s^\varepsilon t \in T\}) \\ &= \{1, 2, 5, 6, 7\} \cup (\{1\} \cup \emptyset) = \{1, 2, 5, 6, 7\} \end{aligned}$$

A little explanation:

- The starting state of the DFA is the epsilon-closure of the set containing just the starting state of the NFA, *i.e.*, the states that are reachable from the starting state by epsilon-transitions.
- A transition in the DFA is done by finding the set of NFA states that comprise the DFA state, following all transitions (on the same symbol) in the NFA from all these NFA states and finally combining the resulting sets of states and closing this under epsilon transitions.
- The set S' of states in the DFA is the set of DFA states that can be reached from s'_0 using the *move* function. S' is defined as a set equation which can be solved as described in section 2.6.1.
- A state in the DFA is an accepting state if at least one of the NFA states it contains is accepting.

As an example, we will convert the NFA in figure 2.5 to a DFA.

The initial state in the DFA is $\epsilon\text{-closure}(\{1\})$, which we have already calculated to be $s'_0 = \{1, 2, 5, 6, 7\}$. This is now entered into the set S' of DFA states as unmarked (following the work-list algorithm from section 2.6.1).

We now pick an unmarked element from the uncompleted S' . We have only one choice: s'_0 . We now mark this and calculate the transitions for it. We get

$$\begin{aligned} \text{move}(s'_0, a) &= \epsilon\text{-closure}(\{t \mid s \in \{1, 2, 5, 6, 7\} \text{ and } s^a t \in T\}) \\ &= \epsilon\text{-closure}(\{3, 8\}) \\ &= \{3, 8, 1, 2, 5, 6, 7\} \\ &= s'_1 \end{aligned}$$

$$\begin{aligned} \text{move}(s'_0, b) &= \epsilon\text{-closure}(\{t \mid s \in \{1, 2, 5, 6, 7\} \text{ and } s^b t \in T\}) \\ &= \epsilon\text{-closure}(\{8\}) \\ &= \{8, 1, 2, 5, 6, 7\} \\ &= s'_2 \end{aligned}$$

$$\begin{aligned} \text{move}(s'_0, c) &= \epsilon\text{-closure}(\{t \mid s \in \{1, 2, 5, 6, 7\} \text{ and } s^c t \in T\}) \\ &= \epsilon\text{-closure}(\{\}) \\ &= \{\} \end{aligned}$$

Note that the empty set of NFA states is not an DFA state, so there will be no transition from s'_0 on c .

We now add s'_1 and s'_2 to our incomplete S' , which now is $\{s'_0, s'_1, s'_2\}$. We now pick s'_1 , mark it and calculate its transitions:

$$\begin{aligned}
\text{move}(s'_1, \mathbf{a}) &= \varepsilon\text{-closure}(\{t \mid s \in \{3, 8, 1, 2, 5, 6, 7\} \text{ and } s^{\mathbf{a}}t \in T\}) \\
&= \varepsilon\text{-closure}(\{3, 8\}) \\
&= \{3, 8, 1, 2, 5, 6, 7\} \\
&= s'_1
\end{aligned}$$

$$\begin{aligned}
\text{move}(s'_1, \mathbf{b}) &= \varepsilon\text{-closure}(\{t \mid s \in \{3, 8, 1, 2, 5, 6, 7\} \text{ and } s^{\mathbf{b}}t \in T\}) \\
&= \varepsilon\text{-closure}(\{8\}) \\
&= \{8, 1, 2, 5, 6, 7\} \\
&= s'_2
\end{aligned}$$

$$\begin{aligned}
\text{move}(s'_1, \mathbf{c}) &= \varepsilon\text{-closure}(\{t \mid s \in \{3, 8, 1, 2, 5, 6, 7\} \text{ and } s^{\mathbf{c}}t \in T\}) \\
&= \varepsilon\text{-closure}(\{4\}) \\
&= \{4\} \\
&= s'_3
\end{aligned}$$

We have seen s'_1 and s'_2 before, so only s'_3 is added: $\{s'_0, s'_1, s'_2, s'_3\}$. We next pick s'_2 :

$$\begin{aligned}
\text{move}(s'_2, \mathbf{a}) &= \varepsilon\text{-closure}(\{t \mid s \in \{8, 1, 2, 5, 6, 7\} \text{ and } s^{\mathbf{a}}t \in T\}) \\
&= \varepsilon\text{-closure}(\{3, 8\}) \\
&= \{3, 8, 1, 2, 5, 6, 7\} \\
&= s'_1
\end{aligned}$$

$$\begin{aligned}
\text{move}(s'_2, \mathbf{b}) &= \varepsilon\text{-closure}(\{t \mid s \in \{8, 1, 2, 5, 6, 7\} \text{ and } s^{\mathbf{b}}t \in T\}) \\
&= \varepsilon\text{-closure}(\{8\}) \\
&= \{8, 1, 2, 5, 6, 7\} \\
&= s'_2
\end{aligned}$$

$$\begin{aligned}
\text{move}(s'_2, \mathbf{c}) &= \varepsilon\text{-closure}(\{t \mid s \in \{8, 1, 2, 5, 6, 7\} \text{ and } s^{\mathbf{c}}t \in T\}) \\
&= \varepsilon\text{-closure}(\{\}) \\
&= \{\}
\end{aligned}$$

No new elements are added, so we pick the remaining unmarked element s'_3 :

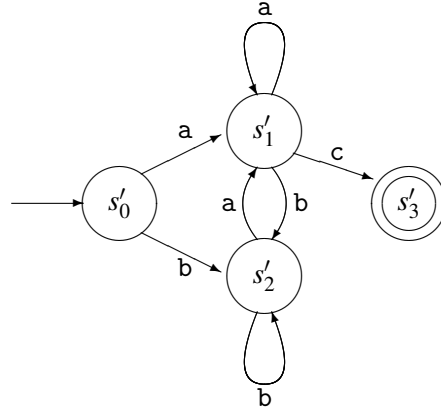


Figure 2.9: DFA constructed from the NFA in figure 2.5

$$\begin{aligned}
 \text{move}(s'_3, a) &= \varepsilon\text{-closure}(\{t \mid s \in \{4\} \text{ and } s^a t \in T\}) \\
 &= \varepsilon\text{-closure}(\{\}) \\
 &= \{\}
 \end{aligned}$$

$$\begin{aligned}
 \text{move}(s'_3, b) &= \varepsilon\text{-closure}(\{t \mid s \in \{4\} \text{ and } s^b t \in T\}) \\
 &= \varepsilon\text{-closure}(\{\}) \\
 &= \{\}
 \end{aligned}$$

$$\begin{aligned}
 \text{move}(s'_3, c) &= \varepsilon\text{-closure}(\{t \mid s \in \{4\} \text{ and } s^c t \in T\}) \\
 &= \varepsilon\text{-closure}(\{\}) \\
 &= \{\}
 \end{aligned}$$

Which now completes the construction of $S' = \{s'_0, s'_1, s'_2, s'_3\}$. Only s'_3 contains the accepting NFA state 4, so this is the only accepting state of our DFA. Figure 2.9 shows the completed DFA.

Suggested exercises: 2.2(b), 2.4.

2.7 Size versus speed

In the above example, we get a DFA with 4 states from an NFA with 8 states. However, as the states in the constructed DFA are (nonempty) sets of states from the NFA there may potentially be $2^n - 1$ states in a DFA constructed from an n -state

NFA. It is not too difficult to construct classes of NFAs that expand exponentially in this way when converted to DFAs, as we shall see in section 2.10.1. Since we are mainly interested in NFAs that are constructed from regular expressions as in section 2.4, we might ask ourselves if these might not be in a suitably simple class that do not risk exponential-sized DFAs. Alas, this is not the case. Just as we can construct a class of NFAs that expand exponentially, we can construct a class of regular expressions where the smallest equivalent DFAs are exponentially larger. This happens rarely when we use regular expressions or NFAs to describe tokens in programming languages, though.

It is possible to avoid the blow-up in size by operating directly on regular expressions or NFAs when testing strings for inclusion in the languages these define. However, there is a speed penalty for doing so. A DFA can be run in time $k * |v|$, where $|v|$ is the length of the input string v and k is a small constant that is independent of the size of the DFA¹. Regular expressions and NFAs can be run in time close to $c * |N| * |v|$, where $|N|$ is the size of the NFA (or regular expression) and the constant c typically is larger than k . All in all, DFAs are a lot faster to use than NFAs or regular expressions, so it is only when the size of the DFA is a real problem that one should consider using NFAs or regular expressions directly.

2.8 Minimisation of DFAs

Even though the DFA in figure 2.9 has only four states, it is not minimal. It is easy to see that states s'_0 and s'_2 are equivalent: Neither are accepting and they have identical transitions. We can hence collapse these states into a single state and get a three-state DFA.

DFAs constructed from regular expressions through NFAs are often non-minimal, though they are rarely very far from being minimal. Nevertheless, minimising a DFA is not terribly difficult and can be done fairly fast, so many lexer generators perform minimisation.

An interesting property of DFAs is that any regular language (a language that can be expressed by a regular expression, NFA or DFA) has a unique minimal DFA. Hence, we can decide equivalence of regular expressions (or NFAs or DFAs) by converting both to minimal DFAs and compare the results.

As hinted above, minimisation of DFAs is done by collapsing equivalent states. However, deciding whether two states are equivalent is not just done by testing if their immediate transitions are identical, since transitions to different states may be equivalent if the target states turn out to be equivalent. Hence, we use a strategy where we first assume all states to be equivalent and then distinguish them only if we can prove them different. We use the following rules for this:

¹If we do not consider the effects of cache-misses *etc.*

- An accepting state is *not* equivalent to a non-accepting state.
- If two states s_1 and s_2 have transitions on the same symbol c to states t_1 and t_2 that we have already proven to be different, then s_1 and s_2 are different. This also applies if only one of s_1 or s_2 have a defined transition on c .

This leads to the following algorithm.

Algorithm 2.4 (DFA minimisation) *Given a DFA D over the alphabet Σ with states S where $F \subseteq S$ is the set of the accepting states, we construct a minimal DFA D_{min} where each state is a group of states from D . The groups in the minimal DFA are consistent: For any pair of states s_1, s_2 in the same group G_1 and any symbol c , $move(s_1, c)$ is in the same group G_2 as $move(s_2, c)$ or both are undefined. In other words, we can not tell s_1 and s_2 apart by looking at their transitions.*

We minimize the DFA D in the following way:

- 1) *We start with two groups: the set of accepting states F and the set of non-accepting states $S \setminus F$. These are unmarked.*
- 2) *We pick any unmarked group G and check if it is consistent. If it is, we mark it. If G is not consistent, we split it into maximal consistent subgroups and replace G by these. All groups are then unmarked. A consistent subgroup is maximal if adding any other state to it will make it inconsistent.*
- 3) *If there are no unmarked groups left, we are done and the remaining groups are the states of the minimal DFA. Otherwise, we go back to step 2.*

The starting state of the minimal DFA is the group that contains the original starting state and any group of accepting states is an accepting state in the minimal DFA.

The time needed for minimisation using algorithm 2.4 depends on the strategy used for picking groups in step 2. With random choices, the worst case is quadratic in the size of the DFA, but there exist strategies for choosing groups and data structures for representing these that guarantee a worst-case time that is $O(n * \log(n))$, where n is the number of states in the (non-minimal) DFA. In other words, the method can be implemented so it uses little more than linear time to do minimisation. We will not here go into further detail but just refer to [3] for the optimal algorithm.

We will, however, note that we can make a slight optimisation to algorithm 2.4: A group that consists of a single state needs never be split, so we need never select such in step 2, and we can stop when all unmarked groups are singletons.

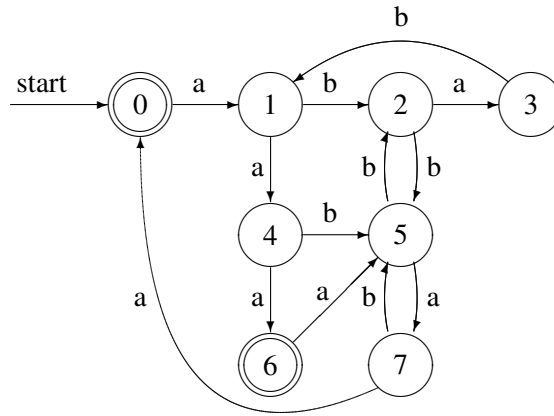


Figure 2.10: Non-minimal DFA

2.8.1 Example

As an example of minimisation, take the DFA in figure 2.10.

We now make the initial division into two groups: The accepting and the non-accepting states.

$$\begin{aligned} G_1 &= \{0, 6\} \\ G_2 &= \{1, 2, 3, 4, 5, 7\} \end{aligned}$$

These are both unmarked. We next pick any unmarked group, say G_1 . To check if this is consistent, we make a table of its transitions:

G_1	a	b
0	G_2	—
6	G_2	—

This is consistent, so we just mark it and select the remaining unmarked group G_2 and make a table for this

G_2	a	b
1	G_2	G_2
2	G_2	G_2
3	—	G_2
4	G_1	G_2
5	G_2	G_2
7	G_1	G_2

G_2 is evidently *not* consistent, so we split it into maximal consistent subgroups and erase all marks (including the one on G_1):

$$\begin{aligned} G_1 &= \{0, 6\} \\ G_3 &= \{1, 2, 5\} \\ G_4 &= \{3\} \\ G_5 &= \{4, 7\} \end{aligned}$$

We now pick G_3 for consideration:

G_3	a	b
1	G_5	G_3
2	G_4	G_3
5	G_5	G_3

This is not consistent either, so we split again and get

$$\begin{aligned} G_1 &= \{0, 6\} \\ G_4 &= \{3\} \\ G_5 &= \{4, 7\} \\ G_6 &= \{1, 5\} \\ G_7 &= \{2\} \end{aligned}$$

We now pick G_5 and check this:

G_5	a	b
4	G_1	G_6
7	G_1	G_6

This is consistent, so we mark it and pick another group, say, G_6 :

G_6	a	b
1	G_5	G_7
5	G_5	G_7

This, also, is consistent, so we have only one unmarked non-singleton group left: G_1 .

G_1	a	b
0	G_6	—
6	G_6	—

As we mark this, we see that there are no unmarked groups left (except the singletons). Hence, the groups form a minimal DFA equivalent to the one in figure 2.10. The minimised DFA is shown in figure 2.11.

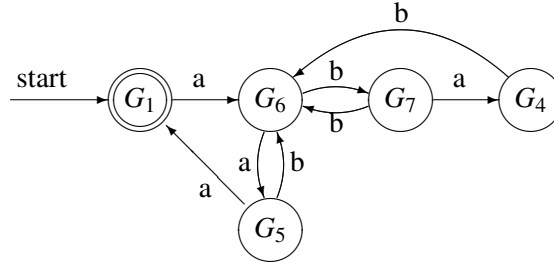


Figure 2.11: Minimal DFA

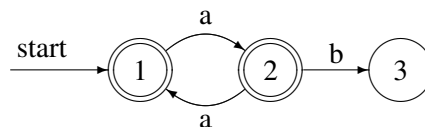
2.8.2 Dead states

Algorithm 2.4 works under some, as yet, unstated assumptions:

- The *move* function is total, *i.e.*, there are transitions on all symbols from all states, *or*
- There are no *dead states* in the DFA.

A dead state is a state from which no accepting state can be reached. Such do not occur in DFAs constructed from NFAs without dead states, and NFAs with dead states can not be constructed from regular expressions by the method shown in section 2.4. Hence, as long as we use minimisation only on DFAs constructed by this process, we are safe. However, if we get a DFA of unknown origin, we risk that it may contain both dead states and undefined transitions.

A transition to a dead state should rightly be equivalent to an undefined transition, as neither can yield future acceptance. The only difference is that we discover this earlier on an undefined transition than when we make a transition to a dead state. However, algorithm 2.4 will treat these differently and may hence decree a group to be inconsistent even though it is not. This will make the algorithm split a group that does not need to be split, hence producing a non-minimal DFA. Consider, for example, the following DFA:



States 1 and 2 are, in fact, equivalent, as starting from either one, any sequence of a's (and no other sequences) will lead to an accepting state. A minimal equivalent DFA has only one accepting state with a transition to itself on a.

But algorithm 2.4 will see a transition on b out of state 2 but no transition on b out of state 1, so it will not keep states 1 and 2 in the same group. As a result, no reduction in the DFA is made.

There are two solutions to this problem:

- 1) Make sure there are no dead states. This can be ensured by invariant, as is the case for DFAs constructed from regular expressions by the methods shown in this chapter, or by explicitly removing dead states before minimisation. Dead states can be found by a simple reachability analysis for directed graphs (if you can't reach an accepting state from state s , s is a dead state). In the above example, state 3 is dead and can be removed (including the transition to it). This makes states 1 and 2 stay in the same group during minimisation.
- 2) Make sure there are no undefined transitions. This can be achieved by adding a new dead state (which has transitions to itself on all symbols) and replacing all undefined transitions by transitions to this dead state. After minimisation, the group that contains the added dead state will contain all dead states from the original DFA. This group can now be removed from the minimal DFA (which will once more have undefined transitions). In the above example, a new (non-accepting) state 4 has to be added. State 1 has a transition to state 4 on b, state 3 has a transition to state 4 on both a and b, and state 4 has transitions to itself on both a and b. After minimisation, state 1 and 2 will be joined, as will state 3 and 4. Since state 4 is dead, all states joined with it are also dead, so we can remove the combined state 3 and 4 from the resulting minimised automaton.

Suggested exercises: 2.5, 2.10(c).

2.9 Lexers and lexer generators

We have, in the previous sections, seen how we can convert a language description written as a regular expression into an efficiently executable representation (a DFA). What we want is something more: A program that does lexical analysis, *i.e.*, a *lexer*:

- A lexer has to distinguish between several different types of tokens, *e.g.*, numbers, variables and keywords. Each of these are described by its own regular expression.

- A lexer does not check if its entire input is included in the languages defined by the regular expressions. Instead, it has to cut the input into pieces (tokens), each of which is included in one of the languages.
- If there are several ways to split the input into legal tokens, the lexer has to decide which of these it should use.

A program that takes a set of token definitions (each consisting of a regular expression and a token name) and generates a lexer is called a *lexer generator*.

The simplest approach would be to generate a DFA for each token definition and apply the DFAs one at a time to the input. This can, however, be quite slow, so we will instead from the set of token definitions generate a single DFA that tests for all the tokens simultaneously. This is not difficult to do: If the tokens are defined by regular expressions r_1, r_2, \dots, r_n , then the regular expression $r_1 \mid r_2 \mid \dots \mid r_n$ describes the union of the languages r_1, r_2, \dots, r_n and the DFA constructed from this combined regular expression will scan for all token types at the same time.

However, we also wish to distinguish between different token types, so we must be able to know *which* of the many tokens was recognised by the DFA. The easiest way to do this is:

- 1) Construct NFAs N_1, N_2, \dots, N_n for each of r_1, r_2, \dots, r_n .
- 2) Mark the accepting states of the NFAs by the name of the tokens they accept.
- 3) Combine the NFAs to a single NFA by adding a new starting state which has epsilon-transitions to each of the starting states of the NFAs.
- 4) Convert the combined NFA to a DFA.
- 5) Each accepting state of the DFA consists of a set of NFA states, some of which are accepting states which we marked by token type in step 2. These marks are used to mark the accepting states of the DFA so each of these will indicate the token types it accepts.

If the same accepting state in the DFA can accept several different token types, it is because these overlap. This is not unusual, as keywords usually overlap with variable names and a description of floating point constants may include integer constants as well. In such cases, we can do one of two things:

- Let the lexer generator generate an error and require the user to make sure the tokens are disjoint.
- Let the user of the lexer generator choose which of the tokens is preferred.

It can be quite difficult (though always possible) with regular expressions to define, *e.g.*, the set of names that are not keywords. Hence, it is common to let the lexer choose according to a prioritised list. Normally, the order in which tokens are defined in the input to the lexer generator indicates priority (earlier defined tokens take precedence over later defined tokens). Hence, keywords are usually defined before variable names, which means that, for example, the string “if” is recognised as a keyword and not a variable name. When an accepting state in a DFA contains accepting NFA states with different marks, the mark corresponding to the highest priority (earliest defined) token is used. Hence, we can simply erase all but one mark from each accepting state. This is a very simple and effective solution to the problem.

When we described minimisation of DFAs, we used two initial groups: One for the accepting states and one for the non-accepting states. As there are now several kinds of accepting states (one for each token), we must use one group for each token, so we will have a total of $n + 1$ initial groups when we have n different tokens.

To illustrate the precedence rule, figure 2.12 shows an NFA made by combining NFAs for variable names, the keyword `if`, integers and floats, as described by the regular expressions in section 2.2.2. The individual NFAs are (simplified versions of) what you get from the method described in section 2.4. When a transition is labelled by a set of characters, it is a shorthand for a set of transitions each labelled by a single character. The accepting states are labelled with token names as described above. The corresponding minimised DFA is shown in figure 2.13. Note that state G is a combination of states 9 and 12 from the NFA, so it can accept both `NUM` and `FLOAT`, but since integers take priority over floats, we have marked G with `NUM` only.

Splitting the input stream

As mentioned, the lexer must cut the input into tokens. This may be done in several ways. For example, the string `if17` can be split in many different ways:

- As one token, which is the variable name `if17`.
- As the variable name `if1` followed by the number 7.
- As the keyword `if` followed by the number 17.
- As the keyword `if` followed by the numbers 1 and 7.
- As the variable name `i` followed by the variable name `f17`.
- And several more.

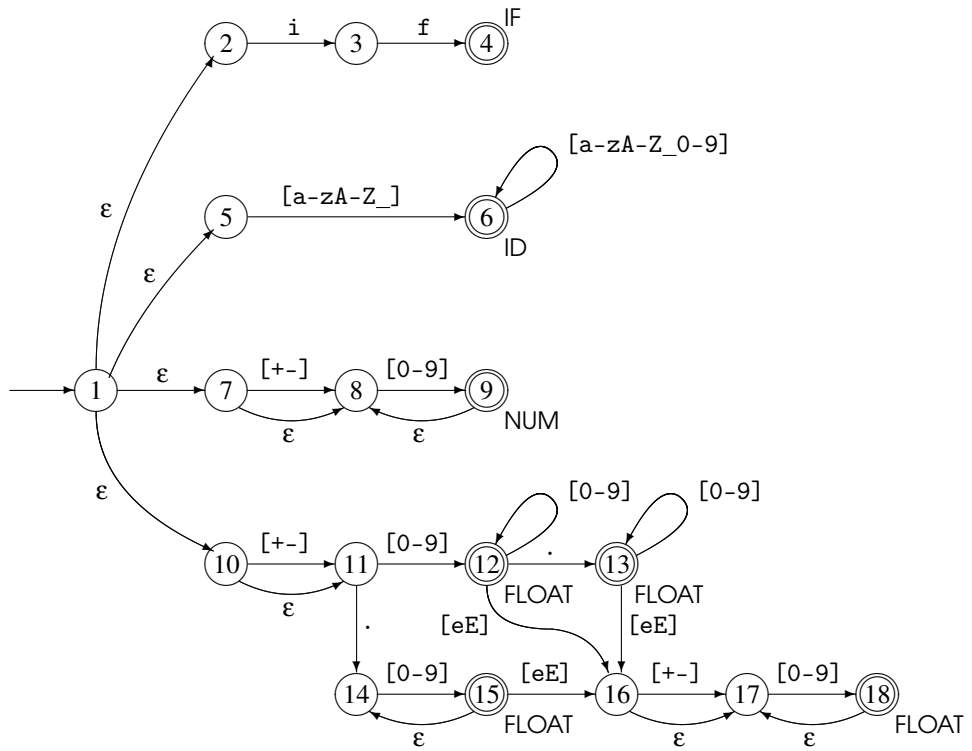


Figure 2.12: Combined NFA for several tokens

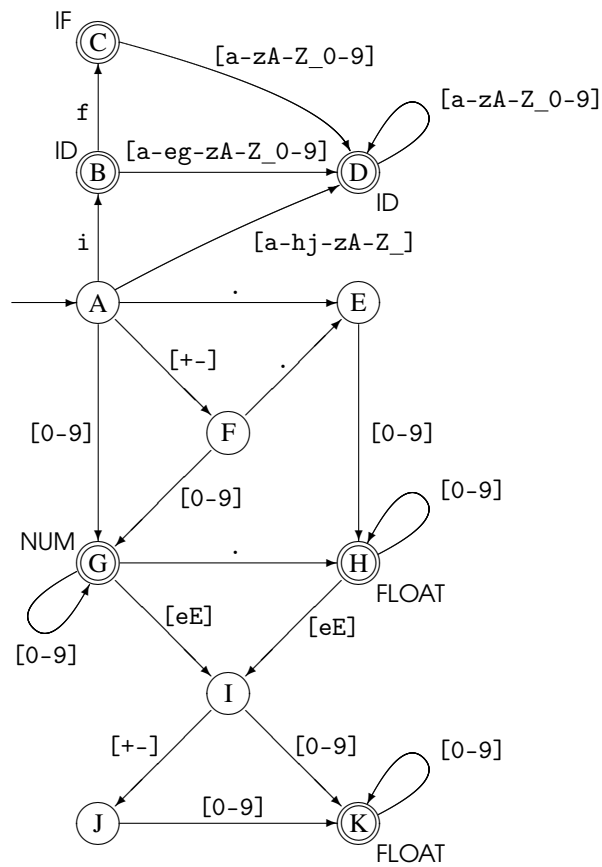


Figure 2.13: Combined DFA for several tokens

A common convention is that it is the longest prefix of the input that matches any token which will be chosen. Hence, the first of the above possible splittings of `if17` will be chosen. Note that the principle of the longest match takes precedence over the order of definition of tokens, so even though the string starts with the keyword `if`, which has higher priority than variable names, the variable name is chosen because it is longer.

Modern languages like C, Java or SML follow this convention, and so do most lexer generators, but some (mostly older) languages like FORTRAN do not. When other conventions are used, lexers must either be written by hand to handle these conventions or the conventions used by the lexer generator must be side-stepped. Some lexer generators allow the user to have some control over the conventions used.

The principle of the longest matching prefix is handled by letting the DFA read as far as it can, until it either reaches the end of the input or no transition is defined on the next input symbol. If the current state at this point is accepting, we are in luck and can simply output the corresponding token. If not, we must go back to the last time we were in an accepting state and output the token indicated by this. The characters read since then are put back in the input stream. The lexer must hence retain the symbols it has read since the last accepting state so it can re-insert these in the input in such situations. If we are not at the end of the input stream, we restart the DFA (in its initial state) on the remaining input to find the next tokens.

As an example, consider lexing of the string `3e-y` with the DFA in figure 2.13. We get to the accepting state G after reading the digit 3. However, we can continue making legal transitions to state I on `e` and then to state J on `-` (as these could be the start of the exponent part of a real number). It is only when we, in state J, find that there is no transition on `y` that we realise that this is not the case. We must now go back to the last accepting state (G) and output the number 3 as the first token and re-insert `-` and `e` in the input stream, so we can continue with `e-y` when we look for the subsequent tokens.

Lexical errors

If no prefix of the input string forms a valid token, a *lexical error* has occurred. When this happens, the lexer will usually report an error. At this point, it may stop reading the input or it may attempt continued lexical analysis by skipping characters until a valid prefix is found. The purpose of the latter approach is to try finding further lexical errors in the same input, so several of these can be corrected by the user before re-running the lexer. Some of these subsequent errors may, however, not be real errors but may be caused by the lexer not skipping enough characters (or skipping too many) after the first error is found. If, for example, the start of a comment is ill-formed, the lexer may try to interpret the contents of the comment

as individual tokens, and if the end of a comment is ill-formed, the lexer will read until the end of the next comment (if any) before continuing, hence skipping too much text.

When the lexer finds an error, the consumer of the tokens that the lexer produces (*e.g.*, the rest of the compiler) can not usually itself produce a valid result. However, the compiler may try to find other errors in the remaining input, again allowing the user to find several errors in one edit-compile cycle. Again, some of the subsequent errors may really be spurious errors caused by lexical error(s), so the user will have to guess at the validity of every error message except the first, as only the first error message is guaranteed to be a real error. Nevertheless, such *error recovery* has, when the input is so large that restarting the lexer from the start of input incurs a considerable time overhead, proven to be an aid in productivity by locating more errors in less time. Less commonly, the lexer may work interactively with a text editor and restart from the point at which an error was spotted after the user has tried to fix the error.

2.9.1 Lexer generators

A lexer generator will typically use a notation for regular expressions similar to the one described in section 2.1, but may require alphabet-characters to be quoted to distinguish them from the symbols used to build regular expressions. For example, an `*` intended to match a multiplication symbol in the input is distinguished from an `*` used to denote repetition by quoting the `*` symbol, *e.g.* as `'*'`. Additionally, some lexer generators extend regular expressions in various ways, *e.g.*, allowing a set of characters to be specified by listing the characters that are *not* in the set. This is useful, for example, to specify the symbols inside a comment up to the terminating character(s).

The input to the lexer generator will normally contain a list of regular expressions that each denote a token. Each of these regular expressions has an associated *action*. The action describes what is passed on to the consumer (*e.g.*, the parser), typically an element from a token data type, which describes the type of token (NUM, ID, *etc.*) and sometimes additional information such as the value of a number token, the name of an identifier token and, perhaps, the position of the token in the input file. The information needed to construct such values is typically provided by the lexer generator through library functions or variables that can be used in the actions.

Normally, the lexer generator requires white-space and comments to be defined by regular expressions. The actions for these regular expressions are typically empty, meaning that white-space and comments are just ignored.

An action can be more than just returning a token. If, for example, a language has a large number of keywords, then a DFA that recognises all of these individu-

ally can be fairly large. In such cases, the keywords are not described as separate regular expressions in the lexer definition but instead treated as special cases of the identifier token. The action for identifiers will then look the name up in a table of keywords and return the appropriate token type (or an identifier token if the name is not a keyword). A similar strategy can be used if the language allows identifiers to shadow keywords.

Another use of non-trivial lexer actions is for nested comments. In principle, a regular expression (or finite automaton) cannot recognise arbitrarily nested comments (see section 2.10), but by using a global counter, the actions for comment tokens can keep track of the nesting level. If escape sequences (for defining, *e.g.*, control characters) are allowed in string constants, the actions for string tokens will, typically, translate the string containing these sequences into a string where they have been substituted by the characters they represent.

Sometimes lexer generators allow several different starting points. In the example in figures 2.12 and 2.13, all regular expressions share the same starting state. However, a single lexer may be used, *e.g.*, for both tokens in the programming language and for tokens in the input to that language. Often, there will be a good deal of sharing between these token sets (the tokens allowed in the input may, for example, be a subset of the tokens allowed in programs). Hence, it is useful to allow these to share a NFA, as this will save space. The resulting DFA will have several starting states. An accepting state may now have more than one token name attached, as long as these come from different token sets (corresponding to different starting points).

In addition to using this feature for several sources of text (program and input), it can be used locally within a single text to read very complex tokens. For example, nested comments and complex-format strings (with nontrivial escape sequences) can be easier to handle if this feature is used.

2.10 Properties of regular languages

We have talked about *regular languages* as the class of languages that can be described by regular expressions or finite automata, but this in itself may not give a clear understanding of what is possible and what is not possible to describe by a regular language. Hence, we will now state a few properties of regular languages and give some examples of some regular and non-regular languages and give informal rules of thumb that can (sometimes) be used to decide if a language is regular.

2.10.1 Relative expressive power

First, we repeat that regular expressions, NFAs and DFAs have exactly the same expressive power: They all can describe all regular languages and only these. Some

languages may, however, have much shorter descriptions in one of these forms than in others.

We have already argued that we from a regular expression can construct an NFA whose size is linear in the size of the regular expression, and that converting an NFA to a DFA can potentially give an exponential increase in size (see below for a concrete example of this). Since DFAs are also NFAs, NFAs are clearly at least as compact as (and sometimes much more compact than) DFAs. Similarly, we can see that NFAs are at least as compact (up to a small constant factor) as regular expressions. But we have not yet considered if the converse is true: Can an NFA be converted to a regular expression of proportional size. The answer is, unfortunately, no: There exist classes of NFAs (and even DFAs) that need regular expressions that are exponentially larger to describe them. This is, however, mainly of academic interest as we rarely have to make conversions in this direction.

If we are only interested in *if* a language is regular rather than the size of its description, however, it does not matter which of the formalisms we choose, so we can in each case choose the formalism that suits us best. Sometimes it is easier to describe a regular language using a DFA or NFA instead of a regular expression. For example, the set of binary number strings that represent numbers that divide evenly by 5 can be described by a 6-state DFA (see exercise 2.9), but it requires a very complex regular expression to do so. For programming language tokens, regular expression are typically quite suitable.

The subset construction (algorithm 2.3) maps sets of NFA states to DFA states. Since there are $2^n - 1$ non-empty sets of n NFA states, the resulting DFA can potentially have exponentially more states than the NFA. But can this potential ever be realised? To answer this, it is not enough to find one n -state NFA that yields a DFA with $2^n - 1$ states. We need to find a family of ever bigger NFAs, all of which yield exponentially-sized DFAs. We also need to argue that the resulting DFAs are minimal. One construction that has these properties is the following: For each integer $n > 1$, construct an n -state NFA in the following way:

1. State 0 is the starting state and state $n - 1$ is accepting.
2. If $0 \leq i < n - 1$, state i has a transition to state $i + 1$ on the symbol a.
3. All states have transitions to themselves *and* to state 0 on the symbol b.

Figure 2.14 shows such an NFA for $n = 4$.

We can represent a set of these states by an n -bit number: Bit i is 1 in the number if and only if state i is in the set. The set that contains only the initial NFA state is, hence, represented by the number 1. We shall see that the way a transition maps a set of states to a new set of states can be expressed as an operation on the number:

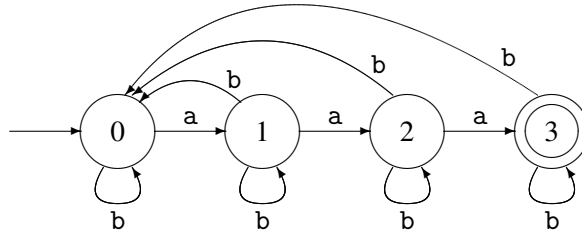


Figure 2.14: A 4-state NFA that gives 15 DFA states

- A transition on a maps the number x to $(2x \bmod (2^n))$.
- A transition on b maps the number x to $(x \text{ or } 1)$, using bitwise or.

This is not hard to verify, so we leave this to the interested reader. It is also easy to see that these two operations can generate any n -bit number from the number 1. Hence, any subset can be reached by a sequence of transitions, which means that the subset-construction will generate a DFA state for every possible non-empty subset of the NFA states.

But is the DFA minimal? If we look at the NFA, we can see that an a leads from state i to $i + 1$ (if $i < n - 1$), so for each NFA state i there is exactly one sequence of as that leads to the accepting state, and that sequence has $n - 1 - i$ as. Hence, a DFA state whose subset contains the NFA state i will lead to acceptance on a string of $n - 1 - i$ as, while a DFA state whose subset does not contain i will not. Hence, for any two different DFA states, we can find an NFA state i that is in one of the sets but not the other and use that to construct a string that will distinguish the DFA states. Hence, all the DFA states are distinct, so the DFA is minimal.

2.10.2 Limits to expressive power

The most basic property of a DFA is that it is *finite*: It has a finite number of states and nowhere else to store information. This means, for example, that any language that requires unbounded counting cannot be regular. An example of this is the language $\{a^n b^n \mid n \geq 0\}$, that is, any sequence of as followed by a sequence of the *same number* of bs. If we must decide membership in this language by a DFA that reads the input from left to right, we must, at the time we have read all the as, know how many there were, so we can compare this to the number of bs. But since a finite automaton cannot count arbitrarily high, the language is not regular. A similar non-regular language is the language of matching parentheses. However, if we limit the nesting depth of parentheses to a constant n , we can recognise this language by a

DFA that has $n + 1$ states (0 to n), where state i corresponds to i unmatched opening parentheses. State 0 is both the starting state and the only accepting state.

Some surprisingly complex languages are regular. As all finite sets of strings are regular languages, the set of all legal Java programs of less than a million lines is a regular language, though it is by no means a simple one. While it can be argued that it would be an acceptable limitation for a language to allow only programs of less than a million lines, it is not practical to describe a programming language as a regular language: The description would be far too large. Even if we ignore such absurdities, we can sometimes be surprised by the expressive power of regular languages. As an example, given any integer constant n , the set of numbers (written in binary or decimal notation) that divide evenly by n is a regular language (see exercise 2.9).

2.10.3 Closure properties

We can also look at closure properties of regular languages. It is clear that regular languages are closed under set union: If we have regular expressions s and t for two languages, the regular expression $s|t$ describes the union of these languages. Similarly, regular languages are closed under concatenation and unbounded repetition, as these correspond to basic operators of regular expressions.

Less obviously, regular languages are also closed under set difference and set intersection. To see this, we first look at set complement: Given a fixed alphabet Σ , the complement of the language L is the set of all strings built from the alphabet Σ , *except* the strings found in L . We write the complement of L as \bar{L} . To get the complement of a regular language L , we first construct a DFA for the language L and make sure that all states have transitions on all characters from the alphabet (as described in section 2.8.2). Now, we simply change every accepting state to non-accepting and *vice versa*, and thus get a DFA for \bar{L} .

We can now (by using the set-theoretic equivalent of De Morgan's law) construct $L_1 \cap L_2$ as $\overline{\bar{L}_1 \cup \bar{L}_2}$. Given this intersection construction, we can now get set difference by $L_1 \setminus L_2 = L_1 \cap \bar{L}_2$.

Regular sets are also closed under a number of common string operations, such as prefix, suffix, subsequence and reversal. The precise meaning of these words in the present context is defined below.

Prefix. A prefix of a string w is any initial part of w , including the empty string and all of w . The prefixes of abc are hence ϵ , a , ab and abc .

Suffix. A suffix of a string is what remains of the string after a prefix has been taken off. The suffixes of abc are hence abc , bc , c and ϵ .

Subsequence. A subsequence of a string is obtained by deleting any number of

symbols from anywhere in the string. The subsequences of *abc* are hence *abc*, *bc*, *ac*, *ab*, *c*, *b*, *a* and ϵ .

Reversal. The reversal of a string is the string read backwards. The reversal of *abc* is hence *cba*.

As with complement, these can be obtained by simple transformations of the DFAs for the language.

Suggested exercises: 2.11.

2.11 Further reading

There are many variants of the method shown in section 2.4. The version presented here has been devised for use in this book in an attempt to make the method easy to understand and manageable to do by hand. Other variants can be found in [5] and [9].

It is possible to convert a regular expression to a DFA directly without going through an NFA. One such method [31] [5] actually at one stage during the calculation computes information equivalent to an NFA (without epsilon-transitions), but more direct methods based on algebraic properties of regular expressions also exist [13, 38]. These, unlike NFA-based methods, generalise fairly easily to handle regular expressions extended with explicit set-intersection and set-difference operators.

A good deal of theoretic information about regular expressions and finite automata can be found in [19]. An efficient DFA minimization algorithm can be found in [24].

Lexer generators can be found for most programming languages. For C, the most common are Lex [28] and Flex [40]. The latter generates the states of the DFA as program code instead of using table-lookup. This makes the generated lexers fast, but can use much more space than a table-driven program.

Finite automata and notation reminiscent of regular expressions are also used to describe behaviour of concurrent systems [33]. In this setting, a state represents the current state of a process and a transition corresponds to an event to which the process reacts by changing state.

Exercises

Exercise 2.1

In the following, a *number-string* is a non-empty sequence of decimal digits, *i.e.*, something in the language defined by the regular expression $[0-9]^+$. The value of

a number-string is the usual interpretation of a number-string as an integer number. Note that leading zeroes are allowed.

Make for each of the following languages a regular expression that describes that language.

- All number-strings that have the value 42.
- All number-strings that *do not* have the value 42.
- All number-strings that have a value that is strictly greater than 42.

Exercise 2.2

Given the regular expression $a^*(a|b)aa$:

- Construct an equivalent NFA using the method in section 2.4.
- convert this NFA to a DFA using algorithm 2.3.

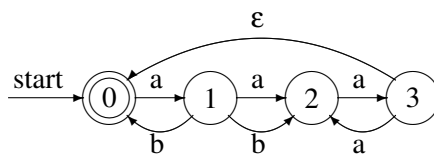
Exercise 2.3

Given the regular expression $((a|b)(a|bb))^*$:

- Construct an equivalent NFA using the method in section 2.4.
- convert this NFA to a DFA using algorithm 2.3.

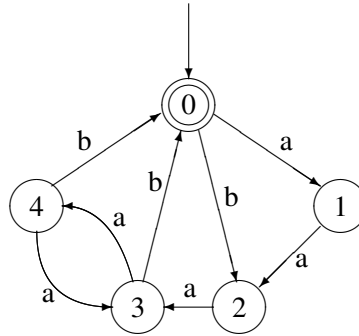
Exercise 2.4

Make a DFA equivalent to the following NFA:

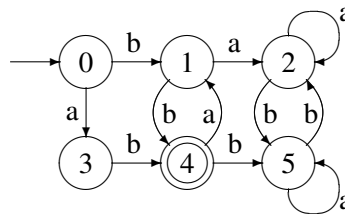


Exercise 2.5

Minimise the following DFA:

**Exercise 2.6**

Minimise the following DFA:

**Exercise 2.7**

Construct DFAs for each of the following regular languages. In all cases the alphabet is $\{a, b\}$.

- The set of strings that has exactly 3 bs (and any number of as).
- The set of strings where the number of bs is a multiple of 3 (and there can be any number of as).
- The set of strings where the difference between the number of as and the number of bs is a multiple of 3.

Exercise 2.8

Construct a DFA that recognises balanced sequences of parenthesis with a maximal nesting depth of 3, *e.g.*, ϵ , $()$, $((()))$ or $((()))()$ but not $((((()))$ or $((()(()))$.

Exercise 2.9

Given that binary number strings are read with the most significant bit first and may have leading zeroes, construct DFAs for each of the following languages:

- a) Binary number strings that represent numbers that are multiples of 4, *e.g.*, 0, 100 and 10100.
- b) Binary number strings that represent numbers that are multiples of 5, *e.g.*, 0, 101, 10100 and 11001.

Hint: Make a state for each possible remainder after division by 5 and then add a state to avoid accepting the empty string.

- c) Given a number n , what is the minimal number of states needed in a DFA that recognises binary numbers that are multiples of n ? Hint: write n as $a * 2^b$, where a is odd.

Exercise 2.10

The empty language, *i.e.*, the language that contains no strings can be recognised by a DFA (any DFA with no accepting states will accept this language), but it can not be defined by any regular expression using the constructions in section 2.2. Hence, the equivalence between DFAs and regular expressions is not complete. To remedy this, a new regular expression ϕ is introduced such that $L(\phi) = \emptyset$.

- a) Argue why each of the following algebraic rules, where s is an arbitrary regular expression, is true:

$$\begin{aligned}\phi|s &= s \\ \phi s &= \phi \\ s\phi &= \phi \\ \phi^* &= \varepsilon\end{aligned}$$

- b) Extend the construction of NFAs from regular expressions to include a case for ϕ .
- c) What consequence will this extension have for converting the NFA to a minimal DFA? Hint: dead states.

Exercise 2.11

Show that regular languages are closed under prefix, suffix, subsequence and reversal, as postulated in section 2.10. Hint: show how an NFA N for a regular language L can be transformed to an NFA N_p for the set of prefixes of strings from L , and similarly for the other operations.

Exercise 2.12

Which of the following statements are true? Argue each answer informally.

- a) Any subset of a regular language is itself a regular language.
- b) Any superset of a regular language is itself a regular language.
- c) The set of anagrams of strings from a regular language forms a regular language. (An anagram of a string is obtained by rearranging the order of characters in the string, but without adding or deleting any. The anagrams of the string `abc` are hence `abc`, `acb`, `bac`, `bca`, `cab` and `cba`).

Exercise 2.13

In figures 2.12 and 2.13 we used character sets on transitions as shorthands for sets of transitions, each with one character. We can, instead, extend the definition of NFAs and DFAs such that such character sets are allowed on a single transition.

For a DFA (to be deterministic), we must require that transitions out of the same state have disjoint character sets.

- a) Sketch how algorithm 2.3 must be modified to handle transitions with sets in such a way that the disjointedness requirement for DFAs are ensured.
- b) Sketch how algorithm 2.4 must be modified to handle character sets. A new requirement for DFA minimality is that the number of transitions as well as the number of states is minimal. How can this be ensured?

Exercise 2.14

As mentioned in section 2.5, DFAs are often implemented by tables where the current state is cross-indexed by the next symbol to find the next state. If the alphabet is large, such a table can take up quite a lot of room. If, for example, 16-bit Unicode is used as the alphabet, there are $2^{16} = 65536$ entries in each row of the table. Even if each entry in the table is only one byte, each row will take up 64KB of memory, which may be a problem.

A possible solution is to split each 16-bit Unicode character c into two 8-bit characters c_1 and c_2 . In the regular expressions, each occurrence of a character c is hence replaced by the regular expression c_1c_2 . This regular expression is then converted to an NFA and then to a DFA in the usual way. The DFA may (and probably will) have more states than the DFA using 16-bit characters, but each state in the new DFA use only 1/256th of the space used by the original DFA.

- a) How much larger is the new NFA compared to the old?

- b) Estimate what the expected size (measured as number of states) of the new DFA is compared to the old. Hint: Some states in the NFA can be reached only after an even number of 8-bit characters are read and the rest only after an odd number of 8-bit characters are read. What does this imply for the sets constructed during the subset construction?
- c) Roughly, how much time does the new DFA require to analyse a string compared to the old?
- d) If space is a problem for a DFA over an 8-bit alphabet, do you expect that a similar trick (splitting each 8-bit character into two 4-bit characters) will help reduce the space requirements? Justify your answer.

Exercise 2.15

If L is a regular language, so is $L \setminus \{\epsilon\}$, i.e., the set of all nonempty strings in L .

So we should be able to transform a regular expression for L into a regular expression for $L \setminus \{\epsilon\}$. We want to do this with a function *nonempty* that is recursive over the structure of the regular expression for L , i.e., of the form:

$$\begin{aligned}
 \text{nonempty}(\epsilon) &= \phi \\
 \text{nonempty}(a) &= \dots && \text{where } a \text{ is an alphabet symbol} \\
 \text{nonempty}(s|t) &= \text{nonempty}(s) | \text{nonempty}(t) \\
 \text{nonempty}(st) &= \dots \\
 \text{nonempty}(s?) &= \dots \\
 \text{nonempty}(s^*) &= \dots \\
 \text{nonempty}(s^+) &= \dots
 \end{aligned}$$

where ϕ is the regular expression for the empty language (see exercise 2.10).

- a) Complete the definition of *nonempty* by replacing the occurrences of “...” in the rules above by expressions similar to those shown in the rules for ϵ and $s|t$.
- b) Use this definition to find *nonempty*(a^*b^*).

Exercise 2.16

If L is a regular language, so is the set of all prefixes of strings in L (see section 2.10.3).

So we should be able to transform a regular expression for L into a regular expression for the set of all prefixes of strings in L . We want to do this with a function *prefixes* that is recursive over the structure of the regular expression for L , i.e., of the form:

$$\begin{aligned} \textit{prefixes}(\epsilon) &= \epsilon \\ \textit{prefixes}(a) &= a? \quad \text{where } a \text{ is an alphabet symbol} \\ \textit{prefixes}(s|t) &= \textit{prefixes}(s) | \textit{prefixes}(t) \\ \textit{prefixes}(st) &= \dots \\ \textit{prefixes}(s^*) &= \dots \\ \textit{prefixes}(s^+) &= \dots \end{aligned}$$

- a) Complete the definition of *prefixes* by replacing the occurrences of “...” in the rules above by expressions similar to those shown in the rules for ϵ , a and $s|t$.
- b) Use this definition to find $\textit{prefixes}(ab^*c)$.

Chapter 3

Syntax Analysis

3.1 Introduction

Where lexical analysis splits the input into tokens, the purpose of syntax analysis (also known as *parsing*) is to recombine these tokens. Not back into a list of characters, but into something that reflects the structure of the text. This “something” is typically a data structure called the *syntax tree* of the text. As the name indicates, this is a tree structure. The leaves of this tree are the tokens found by the lexical analysis, and if the leaves are read from left to right, the sequence is the same as in the input text. Hence, what is important in the syntax tree is how these leaves are combined to form the structure of the tree and how the interior nodes of the tree are labelled.

In addition to finding the structure of the input text, the syntax analysis must also reject invalid texts by reporting *syntax errors*.

As syntax analysis is less local in nature than lexical analysis, more advanced methods are required. We, however, use the same basic strategy: A notation suitable for human understanding is transformed into a machine-like low-level notation suitable for efficient execution. This process is called *parser generation*.

The notation we use for human manipulation is *context-free grammars*¹, which is a recursive notation for describing sets of strings and imposing a structure on each such string. This notation can in some cases be translated almost directly into recursive programs, but it is often more convenient to generate *stack automata*. These are similar to the finite automata used for lexical analysis but they can additionally use a stack, which allows counting and non-local matching of symbols. We shall see two ways of generating such automata. The first of these, LL(1), is relatively simple, but works only for a somewhat restricted class of grammars. The SLR construction, which we present later, is more complex but accepts a wider class of

¹The name refers to the fact that derivation is independent of context.

grammars. Sadly, neither of these work for all context-free grammars. Tools that handle all context-free grammars exist, but they can incur a severe speed penalty, which is why most parser generators restrict the class of input grammars.

3.2 Context-free grammars

Like regular expressions, context-free grammars describe sets of strings, *i.e.*, languages. Additionally, a context-free grammar also defines structure on the strings in the language it defines. A language is defined over some alphabet, for example the set of tokens produced by a lexer or the set of alphanumeric characters. The symbols in the alphabet are called *terminals*.

A context-free grammar recursively defines several sets of strings. Each set is denoted by a name, which is called a *nonterminal*. The set of nonterminals is disjoint from the set of terminals. One of the nonterminals are chosen to denote the language described by the grammar. This is called the *start symbol* of the grammar. The sets are described by a number of *productions*. Each production describes some of the possible strings that are contained in the set denoted by a nonterminal. A production has the form

$$N \rightarrow X_1 \dots X_n$$

where N is a nonterminal and $X_1 \dots X_n$ are zero or more symbols, each of which is either a terminal or a nonterminal. The intended meaning of this notation is to say that the set denoted by N contains strings that are obtained by concatenating strings from the sets denoted by $X_1 \dots X_n$. In this setting, a terminal denotes a singleton set, just like alphabet characters in regular expressions. We will, when no confusion is likely, equate a nonterminal with the set of strings it denotes

Some examples:

$$A \rightarrow a$$

says that the set denoted by the nonterminal A contains the one-character string a .

$$A \rightarrow aA$$

says that the set denoted by A contains all strings formed by putting an a in front of a string taken from the set denoted by A . Together, these two productions indicate that A contains all non-empty sequences of a s and is hence (in the absence of other productions) equivalent to the regular expression a^+ .

We can define a grammar equivalent to the regular expression a^* by the two productions

$$\begin{aligned} B &\rightarrow \\ B &\rightarrow aB \end{aligned}$$

where the first production indicates that the empty string is part of the set B . Compare this grammar with the definition of s^* in figure 2.1.

Productions with empty right-hand sides are called *empty productions*. These are sometimes written with an ϵ on the right hand side instead of leaving it empty.

So far, we have not described any set that could not just as well have been described using regular expressions. Context-free grammars are, however, capable of expressing much more complex languages. In section 2.10, we noted that the language $\{a^n b^n \mid n \geq 0\}$ is not regular. It is, however, easily described by the grammar

$$\begin{aligned} S &\rightarrow \\ S &\rightarrow aSb \end{aligned}$$

The second production ensures that the as and bs are paired symmetrically around the middle of the string, so they occur in equal number.

The examples above have used only one nonterminal per grammar. When several nonterminals are used, we must make it clear which of these is the start symbol. By convention (if nothing else is stated), the nonterminal on the left-hand side of the first production is the start symbol. As an example, the grammar

$$\begin{aligned} T &\rightarrow R \\ T &\rightarrow aTa \\ R &\rightarrow b \\ R &\rightarrow bR \end{aligned}$$

has T as start symbol and denotes the set of strings that start with any number of as followed by a non-zero number of bs and then the same number of as with which it started.

Sometimes, a shorthand notation is used where all the productions of the same nonterminal are combined to a single rule, using the alternative symbol ($|$) from regular expressions to separate the right-hand sides. In this notation, the above grammar would read

$$\begin{aligned} T &\rightarrow R \mid aTa \\ R &\rightarrow b \mid bR \end{aligned}$$

There are still four productions in the grammar, even though the arrow symbol \rightarrow is only used twice.

Form of s_i	Productions for N_i
ϵ	$N_i \rightarrow$
a	$N_i \rightarrow a$
$s_j s_k$	$N_i \rightarrow N_j N_k$
$s_j s_k$	$N_i \rightarrow N_j$ $N_i \rightarrow N_k$
s_j^*	$N_i \rightarrow N_j N_i$ $N_i \rightarrow$
s_j^+	$N_i \rightarrow N_j N_i$ $N_i \rightarrow N_j$
$s_j^?$	$N_i \rightarrow N_j$ $N_i \rightarrow$

Each subexpression of the regular expression is numbered and subexpression s_i is assigned a nonterminal N_i . The productions for N_i depend on the shape of s_i as shown in the table above.

Figure 3.1: From regular expressions to context free grammars

3.2.1 How to write context free grammars

As hinted above, a regular expression can systematically be rewritten as a context free grammar by using a nonterminal for every subexpression in the regular expression and using one or two productions for each nonterminal. The construction is shown in figure 3.1. So, if we can think of a way of expressing a language as a regular expression, it is easy to make a grammar for it. However, we will also want to use grammars to describe non-regular languages. An example is the kind of arithmetic expressions that are part of most programming languages (and also found on electronic calculators). Such expressions can be described by grammar 3.2. Note that, as mentioned in section 2.10, the matching parentheses can not be described by regular expressions, as these can not “count” the number of unmatched opening parentheses at a particular point in the string. However, if we did not have parentheses in the language, it could be described by the regular expression

$$\text{num}((+|-|*|/)\text{num})^*$$

Even so, the regular description is not useful if you want operators to have different precedence, as it treats the expression as a flat string rather than as having structure. We will look at structure in sections 3.3.1 and 3.4.

Most constructions from programming languages are easily expressed by context free grammars. In fact, most modern languages are designed this way.

$$\begin{aligned}
Exp &\rightarrow Exp + Exp \\
Exp &\rightarrow Exp - Exp \\
Exp &\rightarrow Exp * Exp \\
Exp &\rightarrow Exp / Exp \\
Exp &\rightarrow \mathbf{num} \\
Exp &\rightarrow (Exp)
\end{aligned}$$

Grammar 3.2: Simple expression grammar

$$\begin{aligned}
Stat &\rightarrow \mathbf{id} := Exp \\
Stat &\rightarrow Stat ; Stat \\
Stat &\rightarrow \mathbf{if} Exp \mathbf{then} Stat \mathbf{else} Stat \\
Stat &\rightarrow \mathbf{if} Exp \mathbf{then} Stat
\end{aligned}$$

Grammar 3.3: Simple statement grammar

When writing a grammar for a programming language, one normally starts by dividing the constructs of the language into different *syntactic categories*. A syntactic category is a sub-language that embodies a particular concept. Examples of common syntactic categories in programming languages are:

Expressions are used to express calculation of values.

Statements express actions that occur in a particular sequence.

Declarations express properties of names used in other parts of the program.

Each syntactic category is denoted by a main nonterminal, *e.g.*, *Exp* from grammar 3.2. More nonterminals might be needed to describe a syntactic category or provide structure to it, as we shall see, and productions for one syntactic category can refer to nonterminals for other syntactic categories. For example, statements may contain expressions, so some of the productions for statements use the main nonterminal for expressions. A simple grammar for statements might look like grammar 3.3, which refers to the *Exp* nonterminal from grammar 3.2.

Suggested exercises: 3.3 (ignore, for now, the word “unambiguous”), 3.21(a).

3.3 Derivation

So far, we have just appealed to intuitive notions of recursion when we describe the set of strings that a grammar produces. Since the productions are similar to recursive set equations, we might expect to use the techniques from section 2.6.1 to find the set of strings denoted by a grammar. However, though these methods in theory apply to infinite sets by considering limits of chains of sets, they are only practically useful when the sets are finite. Instead, we below introduce the concept of *derivation*. An added advantage of this approach is, as we will later see, that syntax analysis is closely related to derivation.

The basic idea of derivation is to consider productions as rewrite rules: Whenever we have a nonterminal, we can replace this by the right-hand side of any production in which the nonterminal appears on the left-hand side. We can do this anywhere in a sequence of symbols (terminals and nonterminals) and repeat doing so until we have only terminals left. The resulting sequence of terminals is a string in the language defined by the grammar. Formally, we define the derivation relation \Rightarrow by the three rules

1. $\alpha N \beta \Rightarrow \alpha \gamma \beta$ if there is a production $N \rightarrow \gamma$
2. $\alpha \Rightarrow \alpha$
3. $\alpha \Rightarrow \gamma$ if there is a β such that $\alpha \Rightarrow \beta$ and $\beta \Rightarrow \gamma$

where α , β and γ are (possibly empty) sequences of grammar symbols (terminals and nonterminals). The first rule states that using a production as a rewrite rule (anywhere in a sequence of grammar symbols) is a derivation step. The second states that the derivation relation is reflexive, *i.e.*, that a sequence derives itself. The third rule describes transitivity, *i.e.*, that a sequence of derivations is in itself a derivation².

We can use derivation to formally define the language that a context-free grammar generates:

Definition 3.1 *Given a context-free grammar G with start symbol S , terminal symbols T and productions P , the language $L(G)$ that G generates is defined to be the set of strings of terminal symbols that can be obtained by derivation from S using the productions P , *i.e.*, the set $\{w \in T^* \mid S \Rightarrow w\}$.*

As an example, we see that grammar 3.4 generates the string aabbbcc by the derivation shown in figure 3.5. We have, for clarity, in each sequence of symbols underlined the nonterminal that is rewritten in the following step.

²The mathematically inclined will recognise that derivation is a preorder on sequences of grammar symbols.

$$\begin{aligned}
 T &\rightarrow R \\
 T &\rightarrow aTc \\
 R &\rightarrow \\
 R &\rightarrow RbR
 \end{aligned}$$

Grammar 3.4: Example grammar

$$\begin{aligned}
 &\underline{T} \\
 \Rightarrow & a\underline{T}c \\
 \Rightarrow & aa\underline{T}cc \\
 \Rightarrow & aa\underline{R}cc \\
 \Rightarrow & aaRb\underline{R}cc \\
 \Rightarrow & aa\underline{R}bcc \\
 \Rightarrow & aaRb\underline{R}bcc \\
 \Rightarrow & aaRb\underline{R}bRbcc \\
 \Rightarrow & aa\underline{R}bbRbcc \\
 \Rightarrow & aabb\underline{R}bcc \\
 \Rightarrow & aabbbcc
 \end{aligned}$$

Figure 3.5: Derivation of the string aabbbcc using grammar 3.4

$$\begin{aligned}
 &\underline{T} \\
 \Rightarrow & a\underline{T}c \\
 \Rightarrow & aa\underline{T}cc \\
 \Rightarrow & aa\underline{R}cc \\
 \Rightarrow & aa\underline{R}bRcc \\
 \Rightarrow & aa\underline{R}bRbRcc \\
 \Rightarrow & aab\underline{R}bRcc \\
 \Rightarrow & aab\underline{R}bRbRcc \\
 \Rightarrow & aabb\underline{R}bRcc \\
 \Rightarrow & aabbb\underline{R}cc \\
 \Rightarrow & aabbbcc
 \end{aligned}$$

Figure 3.6: Leftmost derivation of the string aabbbcc using grammar 3.4

In this derivation, we have applied derivation steps sometimes to the leftmost nonterminal, sometimes to the rightmost and sometimes to a nonterminal that was neither. However, since derivation steps are local, the order does not matter. So, we might as well decide to always rewrite the leftmost nonterminal, as shown in figure 3.6.

A derivation that always rewrites the leftmost nonterminal is called a *leftmost derivation*. Similarly, a derivation that always rewrites the rightmost nonterminal is called a *rightmost derivation*.

3.3.1 Syntax trees and ambiguity

We can draw a derivation as a tree: The root of the tree is the start symbol of the grammar, and whenever we rewrite a nonterminal we add as its children the symbols on the right-hand side of the production that was used. The leaves of the tree are terminals which, when read from left to right, form the derived string. If a nonterminal is rewritten using an empty production, an ϵ is shown as its child. This is also a leaf node, but is ignored when reading the string from the leaves of the tree.

When we write such a *syntax tree*, the order of derivation is irrelevant: We get the same tree for left derivation, right derivation or any other derivation order. Only the choice of production for rewriting each nonterminal matters.

As an example, the derivations in figures 3.5 and 3.6 yield the same syntax tree, which is shown in figure 3.7.

The syntax tree adds structure to the string that it derives. It is this structure that we exploit in the later phases of the compiler.

For compilation, we do the derivation backwards: We start with a string and want to produce a syntax tree. This process is called *syntax analysis* or *parsing*.

Even though the *order* of derivation does not matter when constructing a syntax tree, the *choice* of production for that nonterminal does. Obviously, different choices can lead to different strings being derived, but it may also happen that several different syntax trees can be built for the same string. As an example, figure 3.8 shows an alternative syntax tree for the same string that was derived in figure 3.7.

When a grammar permits several different syntax trees for some strings we call the grammar *ambiguous*. If our only use of grammar is to describe sets of strings, ambiguity is not a problem. However, when we want to use the grammar to impose structure on strings, the structure had better be the same every time. Hence, it is a desirable feature for a grammar to be unambiguous. In most (but not all) cases, an ambiguous grammar can be rewritten to an unambiguous grammar that generates the same set of strings, or external rules can be applied to decide which of the many possible syntax trees is the “right one”. An unambiguous version of grammar 3.4 is shown in figure 3.9.

$$\begin{aligned}
T &\rightarrow R \\
T &\rightarrow aTc \\
R &\rightarrow \\
R &\rightarrow bR
\end{aligned}$$

Grammar 3.9: Unambiguous version of grammar 3.4

How do we know if a grammar is ambiguous? If we can find a string and show two alternative syntax trees for it, this is a proof of ambiguity. It may, however, be hard to find such a string and, when the grammar is unambiguous, even harder to show that this is the case. In fact, the problem is formally undecidable, *i.e.*, there is no method that for all grammars can answer the question “Is this grammar ambiguous?”.

But in many cases it is not difficult to detect and prove ambiguity. For example, any grammar that has a production of the form

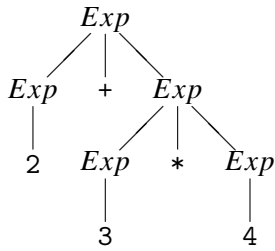
$$N \rightarrow N\alpha N$$

where α is any sequence of grammar symbols, is ambiguous. This is, for example, the case with grammars 3.2 and 3.4.

We will, in sections 3.12 and 3.14, see methods for constructing parsers from grammars. **These methods have the property that they only work on unambiguous grammars, so successful construction of a parser is a proof of unambiguity.** However, the methods may also fail on certain unambiguous grammars, so they can not be used to prove ambiguity.

In the next section, we will see ways of rewriting a grammar to get rid of some sources of ambiguity. These transformations preserve the language that the grammar generates. By using such transformations (and others, which we will see later), we can create a large set of *equivalent* grammars, *i.e.*, grammars that generate the same language (though they may impose different structures on the strings of the language).

Given two grammars, it would be nice to be able to tell if they are equivalent. Unfortunately, no known method is able to decide this in all cases, but, unlike ambiguity, it is not (at the time of writing) known if such a method may or may not theoretically exist. Sometimes, equivalence can be proven *e.g.* by induction over the set of strings that the grammars produce. The converse can be proven by finding an example of a string that one grammar can generate but the other not. But in some cases, we just have to take claims of equivalence on faith or give up on deciding the issue.

Figure 3.10: Preferred syntax tree for $2+3*4$ using grammar 3.2

Suggested exercises: 3.1, 3.2, 3.21(b).

3.4 Operator precedence

As mentioned in section 3.2.1, we can describe traditional arithmetic expressions by grammar 3.2. Note that **num** is a terminal that denotes all integer constants and that, here, the parentheses are terminal symbols (unlike in regular expressions, where they are used to impose structure on the regular expressions).

This grammar is ambiguous, as evidenced by, *e.g.*, the production

$$Exp \rightarrow Exp + Exp$$

which has the form that in section 3.3.1 was claimed to imply ambiguity. This ambiguity is not surprising, as we are used to the fact that an expression like $2+3*4$ can be read in two ways: Either as multiplying the sum of 2 and 3 by 4 or as adding 2 to the product of 3 and 4. Simple electronic calculators will choose the first of these interpretations (as they always calculate from left to right), whereas scientific calculators and most programming languages will choose the second, as they use a hierarchy of *operator precedences* which dictate that the product must be calculated before the sum. The hierarchy can be overridden by explicit parenthesisation, *e.g.*, $(2+3)*4$.

Most programming languages use the same convention as scientific calculators, so we want to make this explicit in the grammar. Ideally, we would like the expression $2+3*4$ to generate the syntax tree shown in figure 3.10, which reflects the operator precedences by grouping of subexpressions: When evaluating an expression, the subexpressions represented by subtrees of the syntax tree are evaluated before the topmost operator is applied.

A possible way of resolving the ambiguity is to use precedence rules during syntax analysis to select among the possible syntax trees. Many parser generators allow this approach, as we shall see in section 3.16. However, some parsing meth-

ods require the grammars to be unambiguous, so we have to express the operator hierarchy in the grammar itself. To clarify this, we first define some concepts:

- An operator \oplus is *left-associative* if the expression $a \oplus b \oplus c$ must be evaluated from left to right, *i.e.*, as $(a \oplus b) \oplus c$.
- An operator \oplus is *right-associative* if the expression $a \oplus b \oplus c$ must be evaluated from right to left, *i.e.*, as $a \oplus (b \oplus c)$.
- An operator \oplus is *non-associative* if expressions of the form $a \oplus b \oplus c$ are illegal.

By the usual convention, $-$ and $/$ are left-associative, as *e.g.*, $2-3-4$ is calculated as $(2-3)-4$. $+$ and $*$ are associative in the mathematical sense, meaning that it does not matter if we calculate from left to right or from right to left. However, to avoid ambiguity we have to choose one of these. By convention (and similarity to $-$ and $/$) we choose to let these be left-associative as well. Also, having a left-associative $-$ and right-associative $+$ would not help resolving the ambiguity of $2-3+4$, as the operators so-to-speak “pull in different directions”.

List construction operators in functional languages, *e.g.*, $::$ and $@$ in SML, are typically right-associative, as are function arrows in types: $a \rightarrow b \rightarrow c$ is read as $a \rightarrow (b \rightarrow c)$. The assignment operator in C is also right-associative: $a=b=c$ is read as $a=(b=c)$.

In some languages (like Pascal), comparison operators (like $<$ or $>$) are non-associative, *i.e.*, you are not allowed to write $2 < 3 < 4$.

3.4.1 Rewriting ambiguous expression grammars

If we have an ambiguous grammar

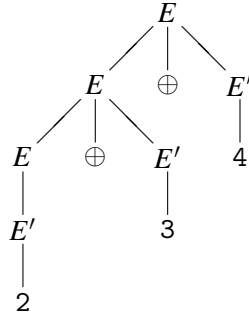
$$\begin{aligned} E &\rightarrow E \oplus E \\ E &\rightarrow \mathbf{num} \end{aligned}$$

we can rewrite this to an unambiguous grammar that generates the correct structure. As this depends on the associativity of \oplus , we use different rewrite rules for different associativities.

If \oplus is left-associative, we make the grammar *left-recursive* by having a recursive reference to the left only of the operator symbol:

$$\begin{aligned} E &\rightarrow E \oplus E' \\ E &\rightarrow E' \\ E' &\rightarrow \mathbf{num} \end{aligned}$$

Now, the expression $2 \oplus 3 \oplus 4$ can only be parsed as



We get a slightly more complex syntax tree than in figure 3.10, but not enormously so.

We handle right-associativity in a similar fashion: We make the offending production *right-recursive*:

$$\begin{aligned}
 E &\rightarrow E' \oplus E \\
 E &\rightarrow E' \\
 E' &\rightarrow \mathbf{num}
 \end{aligned}$$

Non-associative operators are handled by *non-recursive* productions:

$$\begin{aligned}
 E &\rightarrow E' \oplus E' \\
 E &\rightarrow E' \\
 E' &\rightarrow \mathbf{num}
 \end{aligned}$$

Note that the latter transformation actually changes the language that the grammar generates, as it makes expressions of the form $\mathbf{num} \oplus \mathbf{num} \oplus \mathbf{num}$ illegal.

So far, we have handled only cases where an operator interacts with itself. This is easily extended to the case where several operators with the same precedence and associativity interact with each other, as for example + and -:

$$\begin{aligned}
 E &\rightarrow E + E' \\
 E &\rightarrow E - E' \\
 E &\rightarrow E' \\
 E' &\rightarrow \mathbf{num}
 \end{aligned}$$

Operators with the same precedence must have the same associativity for this to work, as mixing left-recursive and right-recursive productions for the same nonterminal makes the grammar ambiguous. As an example, the grammar

$$\begin{aligned}
 E &\rightarrow E + E' \\
 E &\rightarrow E' \oplus E \\
 E &\rightarrow E' \\
 E' &\rightarrow \mathbf{num}
 \end{aligned}$$

$$\begin{aligned}
Exp &\rightarrow Exp + Exp2 \\
Exp &\rightarrow Exp - Exp2 \\
Exp &\rightarrow Exp2 \\
Exp2 &\rightarrow Exp2 * Exp3 \\
Exp2 &\rightarrow Exp2 / Exp3 \\
Exp2 &\rightarrow Exp3 \\
Exp3 &\rightarrow \mathbf{num} \\
Exp3 &\rightarrow (Exp)
\end{aligned}$$

Grammar 3.11: Unambiguous expression grammar

seems like an obvious generalisation of the principles used above, giving $+$ and \oplus the same precedence and different associativity. But not only is the grammar ambiguous, it does not even accept the intended language. For example, the string **num+num \oplus num** is not derivable by this grammar.

In general, there is no obvious way to resolve ambiguity in an expression like $1+2\oplus 3$, where $+$ is left-associative and \oplus is right-associative (or *vice-versa*). Hence, most programming languages (and most parser generators) *require* operators at the same precedence level to have identical associativity.

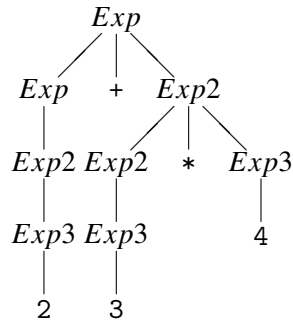
We also need to handle operators with different precedences. This is done by using a nonterminal for each precedence level. The idea is that if an expression uses an operator of a certain precedence level, then its subexpressions cannot use operators of lower precedence (unless these are inside parentheses). Hence, the productions for a nonterminal corresponding to a particular precedence level refers only to nonterminals that correspond to the same or higher precedence levels, unless parentheses or similar bracketing constructs disambiguate the use of these. Grammar 3.11 shows how these rules are used to make an unambiguous version of grammar 3.2. Figure 3.12 show the syntax tree for $2+3*4$ using this grammar.

Suggested exercises: 3.6.

3.5 Other sources of ambiguity

Most of the potential ambiguity in grammars for programming languages comes from expression syntax and can be handled by exploiting precedence rules as shown in section 3.4. Another classical example of ambiguity is the “dangling-else” problem.

Imperative languages like Pascal or C often let the else-part of a conditional be optional, like shown in grammar 3.3. The problem is that it is not clear how to

Figure 3.12: Syntax tree for $2+3*4$ using grammar 3.11

parse, for example,

```
if p then if q then s1 else s2
```

According to the grammar, the `else` can equally well match either `if`. The usual convention is that an `else` matches the closest not previously matched `if`, which, in the example, will make the `else` match the second `if`.

How do we make this clear in the grammar? We can treat `if`, `then` and `else` as a kind of right-associative operators, as this would make them group to the right, making an `if-then` match the closest `else`. However, the grammar transformations shown in section 3.4 can not directly be applied to grammar 3.3, as the productions for conditionals do not have the right form.

Instead we use the following observation: When an `if` and an `else` match, all `ifs` that occur between these must have matching `elses`. This can easily be proven by assuming otherwise and concluding that this leads to a contradiction.

Hence, we make two nonterminals: One for matched (*i.e.* with `else-part`) conditionals and one for unmatched (*i.e.* without `else-part`) conditionals. The result is shown in grammar 3.13. This grammar also resolves the associativity of semicolon (`right`) and the precedence of `if` over semicolon.

An alternative to rewriting grammars to resolve ambiguity is to use an ambiguous grammar and resolve conflicts by using precedence rules during parsing. We shall look into this in section 3.16.

All cases of ambiguity must be treated carefully: It is not enough that we eliminate ambiguity, we must do so in a way that results in the desired structure: The structure of arithmetic expressions is significant, and it makes a difference to which `if` an `else` is matched.

Suggested exercises: 3.3 (focusing now on making the grammar unambiguous).

<i>Stat</i>	→	<i>Stat2 ; Stat</i>
<i>Stat</i>	→	<i>Stat2</i>
<i>Stat2</i>	→	<i>Matched</i>
<i>Stat2</i>	→	<i>Unmatched</i>
<i>Matched</i>	→	if <i>Exp</i> then <i>Matched</i> else <i>Matched</i>
<i>Matched</i>	→	id := <i>Exp</i>
<i>Unmatched</i>	→	if <i>Exp</i> then <i>Matched</i> else <i>Unmatched</i>
<i>Unmatched</i>	→	if <i>Exp</i> then <i>Stat2</i>

Grammar 3.13: Unambiguous grammar for statements

3.6 Syntax analysis

The syntax analysis phase of a compiler will take a string of tokens produced by the lexer, and from this construct a syntax tree for the string by finding a derivation of the string from the start symbol of the grammar.

This can be done by guessing derivations until the right one is found, but random guessing is hardly an effective method. Even so, some parsing techniques are based on “guessing” derivations. However, these make sure, by looking at the string, that they will always guess right. These are called *predictive* parsing methods. Predictive parsers always build the syntax tree from the root down to the leaves and are hence also called (deterministic) top-down parsers.

Other parsers go the other way: They search for parts of the input string that matches right-hand sides of productions and rewrite these to the left-hand nonterminals, at the same time building pieces of the syntax tree. The syntax tree is eventually completed when the string has been rewritten (by inverse derivation) to the start symbol. Also here, we wish to make sure that we always pick the “right” rewrites, so we get deterministic parsing. Such methods are called *bottom-up* parsing methods.

We will in the next sections first look at predictive parsing and later at a bottom-up parsing method called SLR parsing.

3.7 Predictive parsing

If we look at the left-derivation in figure 3.6, we see that, to the left of the rewritten nonterminals, there are only terminals. These terminals correspond to a prefix of the string that is being parsed. In a parsing situation, this prefix will be the part of the input that has already been read. The job of the parser is now to choose the production by which the leftmost unexpanded nonterminal should be rewritten. Our

aim is to be able to make this choice deterministically based on the next unmatched input symbol.

If we look at the third line in figure 3.6, we have already read two *a*s and (if the input string is the one shown in the bottom line) the next symbol is a *b*. Since the right-hand side of the production

$$T \rightarrow aTc$$

starts with an *a*, we obviously can not use this. Hence, we can only rewrite *T* using the production

$$T \rightarrow R$$

We are not quite as lucky in the next step. None of the productions for *R* start with a terminal symbol, so we can not immediately choose a production based on this. As the grammar (grammar 3.4) is ambiguous, it should not be a surprise that we can not always choose uniquely. If we instead use the unambiguous grammar (grammar 3.9) we can immediately choose the second production for *R*. When all the *b*s are read and we are at the following *c*, we choose the empty production for *R* and match the remaining input with the rest of the derived string.

If we can always choose a unique production based on the next input symbol, we are able to do predictive parsing without backtracking.

3.8 Nullable and FIRST

In simple cases, like the above, all productions for a nonterminal start with distinct terminals except at most one production that does not start with a terminal. We chose the latter whenever the input symbol did not match any of the terminal symbols starting the other productions. We can extend the method to work also for grammars where several productions start with nonterminals. We just need to be able to select between them based on the input symbol. In other words, if the strings these productions can derive begin with symbols from disjoint sets, we check if the input symbol is in one of these sets and choose the corresponding production if it is. If not, and there is an empty production, we choose this. Otherwise, we report a syntax error message.

Hence, we define the function *FIRST*, which given a sequence of grammar symbols (*e.g.* the right-hand side of a production) returns the set of symbols with which strings derived from that sequence can begin:

Definition 3.2 A symbol *c* is in *FIRST*(α) if and only if $\alpha \Rightarrow c\beta$ for some (possibly empty) sequence β of grammar symbols.

To calculate *FIRST*, we need an auxiliary function *Nullable*, which for a sequence α of grammar symbols indicates whether or not that sequence can derive the empty string:

Definition 3.3 A sequence α of grammar symbols is *Nullable* (we write this as *Nullable*(α)) if and only if $\alpha \Rightarrow \epsilon$.

A production $N \rightarrow \alpha$ is called nullable if *Nullable*(α). We describe calculation of *Nullable* by case analysis over the possible forms of sequences of grammar symbols:

Algorithm 3.4

$$\begin{aligned} \text{Nullable}(\epsilon) &= \text{true} \\ \text{Nullable}(\mathbf{a}) &= \text{false} \\ \text{Nullable}(\alpha\beta) &= \text{Nullable}(\alpha) \wedge \text{Nullable}(\beta) \\ \text{Nullable}(N) &= \text{Nullable}(\alpha_1) \vee \dots \vee \text{Nullable}(\alpha_n), \\ &\quad \text{where the productions for } N \text{ are} \\ &\quad N \rightarrow \alpha_1, \quad \dots, \quad N \rightarrow \alpha_n \end{aligned}$$

where \mathbf{a} is a terminal, N is a nonterminal, α and β are sequences of grammar symbols and ϵ represents the empty sequence of grammar symbols.

The equations are quite natural: Any occurrence of a terminal on a right-hand side makes *Nullable* false for that right-hand side, but a nonterminal is nullable if any production has a nullable right-hand side.

Note that this is a recursive definition since *Nullable* for a nonterminal is defined in terms of *Nullable* for its right-hand sides, which may contain that same nonterminal. We can solve this in much the same way that we solved set equations in section 2.6.1. We have, however, now booleans instead of sets and several equations instead of one. Still, the method is essentially the same: We have a set of boolean equations:

$$\begin{aligned} X_1 &= F_1(X_1, \dots, X_n) \\ &\vdots \\ X_n &= F_n(X_1, \dots, X_n) \end{aligned}$$

We initially assume X_1, \dots, X_n to be all *false*. We then, in any order, calculate the right-hand sides of the equations and update the variable on the left-hand side by the calculated value. We continue until all equations are satisfied. In appendix A and section 2.6.1, we required the functions to be monotonic with respect to subset. Correspondingly, we now require the boolean functions to be monotonic with respect to truth: If we make more arguments true, the result will also be more true

Right-hand side	Initialisation	Iteration 1	Iteration 2	Iteration 3
R	<i>false</i>	<i>false</i>	<i>true</i>	<i>true</i>
aTc	<i>false</i>	<i>false</i>	<i>false</i>	<i>false</i>
ϵ	<i>false</i>	<i>true</i>	<i>true</i>	<i>true</i>
bR	<i>false</i>	<i>false</i>	<i>false</i>	<i>false</i>
Nonterminal				
T	<i>false</i>	<i>false</i>	<i>true</i>	<i>true</i>
R	<i>false</i>	<i>true</i>	<i>true</i>	<i>true</i>

Figure 3.14: Fixed-point iteration for calculation of *Nullable*

(i.e., it may stay unchanged, change from *false* to *true*, but never change from *true* to *false*).

If we look at grammar 3.9, we get these equations for nonterminals and right-hand sides:

$$\begin{aligned}
\text{Nullable}(T) &= \text{Nullable}(R) \vee \text{Nullable}(aTc) \\
\text{Nullable}(R) &= \text{Nullable}(\epsilon) \vee \text{Nullable}(bR) \\
\\
\text{Nullable}(R) &= \text{Nullable}(R) \\
\text{Nullable}(aTc) &= \text{Nullable}(a) \wedge \text{Nullable}(T) \wedge \text{Nullable}(c) \\
\text{Nullable}(\epsilon) &= \text{true} \\
\text{Nullable}(bR) &= \text{Nullable}(b) \wedge \text{Nullable}(R)
\end{aligned}$$

In a fixed-point calculation, we initially assume that *Nullable* is false for all nonterminals and use this as a basis for calculating *Nullable* for first the right-hand sides and then the nonterminals. We repeat recalculating these until there is no change between two iterations. Figure 3.14 shows the fixed-point iteration for the above equations. In each iteration, we first evaluate the formulae for the right-hand sides and then use the results of this to evaluate the nonterminals. The right-most column shows the final result.

We can calculate *FIRST* in a similar fashion to *Nullable*:

Algorithm 3.5

$$\begin{aligned}
\text{FIRST}(\epsilon) &= \emptyset \\
\text{FIRST}(a) &= \{a\} \\
\text{FIRST}(\alpha\beta) &= \begin{cases} \text{FIRST}(\alpha) \cup \text{FIRST}(\beta) & \text{if } \text{Nullable}(\alpha) \\ \text{FIRST}(\alpha) & \text{if not } \text{Nullable}(\alpha) \end{cases} \\
\text{FIRST}(N) &= \text{FIRST}(\alpha_1) \cup \dots \cup \text{FIRST}(\alpha_n) \\
&\quad \text{where the productions for } N \text{ are} \\
&\quad N \rightarrow \alpha_1, \quad \dots, N \rightarrow \alpha_n
\end{aligned}$$

Right-hand side	Initialisation	Iteration 1	Iteration 2	Iteration 3
R	\emptyset	\emptyset	$\{b\}$	$\{b\}$
aTc	\emptyset	$\{a\}$	$\{a\}$	$\{a\}$
ϵ	\emptyset	\emptyset	\emptyset	\emptyset
bR	\emptyset	$\{b\}$	$\{b\}$	$\{b\}$
Nonterminal				
T	\emptyset	$\{a\}$	$\{a, b\}$	$\{a, b\}$
R	\emptyset	$\{b\}$	$\{b\}$	$\{b\}$

Figure 3.15: Fixed-point iteration for calculation of *FIRST*

where a is a terminal, N is a nonterminal, α and β are sequences of grammar symbols and ϵ represents the empty sequence of grammar symbols.

The only nontrivial equation is that for $\alpha\beta$. Obviously, anything that can start a string derivable from α can also start a string derivable from $\alpha\beta$. However, if α is nullable, a derivation may proceed as $\alpha\beta \Rightarrow \beta \Rightarrow \dots$, so anything in $FIRST(\beta)$ is also in $FIRST(\alpha\beta)$.

The set-equations are solved in the same general way as the boolean equations for *Nullable*, but since we work with sets, we initially assume every set to be empty. For grammar 3.9, we get the following equations:

$$\begin{aligned}
 FIRST(T) &= FIRST(R) \cup FIRST(aTc) \\
 FIRST(R) &= FIRST(\epsilon) \cup FIRST(bR) \\
 \\
 FIRST(R) &= FIRST(R) \\
 FIRST(aTc) &= FIRST(a) \\
 FIRST(\epsilon) &= \emptyset \\
 FIRST(bR) &= FIRST(b)
 \end{aligned}$$

The fixed-point iteration is shown in figure 3.15.

When working with grammars by hand, it is usually quite easy to see for most productions if they are nullable and what their *FIRST* sets are. For example, a production is not nullable if its right-hand side has a terminal anywhere, and if the right-hand side starts with a terminal, the *FIRST* set consists of only that symbol. Sometimes, however, it is necessary to go through the motions of solving the equations. When working by hand, it is often useful to simplify the equations before the fixed-point iteration, *e.g.*, reduce $FIRST(aTc)$ to $\{a\}$.

Suggested exercises: 3.8 (*Nullable* and *FIRST* only).

3.9 Predictive parsing revisited

We are now ready to construct predictive parsers for a wider class of grammars: If the right-hand sides of the productions for a nonterminal have disjoint *FIRST* sets, we can use the next input symbol to choose among the productions.

In section 3.7, we picked the empty production (if any) on any symbol that was not in the *FIRST* sets of the non-empty productions for the same nonterminal. We can extend this, so we in case of no matching *FIRST* sets can select a production if it is *Nullable*. The idea is that a *Nullable* production can derive the empty string, so the input symbol need not be read by the production itself.

But if there are several *Nullable* productions, we have no way of choosing between them. Hence, we do not allow more than one production for a nonterminal to be *Nullable*.

We said in section 3.3.1 that our syntax analysis methods will detect ambiguous grammars. However, this is not true with the method as stated above: We can get unique choice of production even for some ambiguous grammars, including grammar 3.4. The syntax analysis will in this case just choose one of several possible syntax trees for a given input string. In many cases, we do not consider such behaviour acceptable. In fact, we would very much like our parser construction method to tell us if we by mistake write an ambiguous grammar.

Even worse, the rules for predictive parsing as presented here might even for some unambiguous grammars give deterministic choice of production, but reject strings that actually belong to the language described by the grammar. If we, for example, change the second production in grammar 3.9 to

$$T \rightarrow aTb$$

this will not change the choices made by the predictive parser for nonterminal *R*. However, always choosing the last production for *R* on a *b* will lead to erroneous rejection of many strings, including *ab*.

This kind of behaviour is clearly unacceptable. We should, at least, get a warning that this might occur, so we can rewrite the grammar or choose another syntax analysis method.

Hence, we add to our construction of predictive parsers a test that will reject all ambiguous grammars and those unambiguous grammars that can cause the parser to fail erroneously.

We have so far simply chosen a nullable production if and only if no other choice is possible. This is, however, not always the right thing to do, so we must change the rule to say that we choose a production $N \rightarrow \alpha$ on symbol *c* if one of the two conditions below are satisfied:

- 1) $c \in FIRST(\alpha)$.

- 2) α is nullable and the sequence Nc can occur somewhere in a derivation starting from the start symbol of the grammar.

The first rule is obvious, but the second requires a bit of explanation: If α is nullable, we can construct a syntax tree for it without reading any input, so it seems like a nullable production could be a valid choice regardless of the next input symbol. Not only would this give multiple valid choices of production whenever there are both nullable and non-nullable productions for the same nonterminal, but it is also not always correct to choose the nullable production: Predictive parsing makes a leftmost derivation so we always rewrite the leftmost nonterminal N in the current sequence of grammar symbols. So whatever input is not matched by N must be matched by the sequence of grammar symbols that occurs after N in the current sequence. If this is not possible, we have made a bad choice when deriving N . In particular, if N derives to the empty sequence, the next input symbol c should begin the derivation of the sequence that follows N . So at the very least, the sequence of symbols that follow N should have a derivation that begins with c . If we (using a different derivation order) derive the symbols after N before deriving N , we should, hence, see the sequence Nc during the derivation. If we don't, we can not rewrite N to the empty sequence without later getting stuck when rewriting the remaining sequence. Since the derivation order doesn't change the syntax tree, we can see that if the alternative derivation order gets stuck, so will the leftmost derivation order. Hence, we can only rewrite N to the empty sequence if the next input symbol c can occur in combination with N in a legal derivation. We will in the next section see how we can determine this.

Note that a nullable production $N \rightarrow \alpha$ can validly be selected if $c \in FIRST(\alpha)$.

Even with the restriction on choosing nullable productions, we can still have situations where both nullable and non-nullable productions are valid choices. This includes the example above with the modified grammar 3.9 (since Rb can occur in a derivation) and all ambiguous grammars that are not detected as being ambiguous by the original method where we only choose nullable productions if no other choice is possible.

3.10 FOLLOW

To determine when we can select nullable productions during predictive parsing, we introduce *FOLLOW* sets for nonterminals.

Definition 3.6 A terminal symbol a is in $FOLLOW(N)$ if and only if there is a derivation from the start symbol S of the grammar such that $S \Rightarrow \alpha Na\beta$, where α and β are (possibly empty) sequences of grammar symbols.

In other words, a terminal c is in $FOLLOW(N)$ if c may follow N at some point in a derivation. Unlike $FIRST(N)$, this is not a property of the productions for N , but of the productions that (directly or indirectly) use N on their right-hand side.

To correctly handle end-of-string conditions, we want to detect if $S \Rightarrow \alpha N$, *i.e.*, if there are derivations where N can be followed by the end of input. It turns out to be easy to do this by adding an extra production to the grammar:

$$S' \rightarrow S\$$$

where S' is a new nonterminal that replaces S as start symbol and $\$$ is a new terminal symbol representing the end of input. Hence, in the new grammar, $\$$ will be in $FOLLOW(N)$ exactly if $S' \Rightarrow \alpha N\$$ which is the case exactly when $S \Rightarrow \alpha N$.

The easiest way to calculate $FOLLOW$ is to generate a collection of *set constraints*, which are subsequently solved to find the least sets that obey the constraints.

A production

$$M \rightarrow \alpha N \beta$$

generates the constraint $FIRST(\beta) \subseteq FOLLOW(N)$, since β , obviously, can follow N . Furthermore, if $Nullable(\beta)$ the production also generates the constraint $FOLLOW(M) \subseteq FOLLOW(N)$ (note the direction of the inclusion). The reason is that, if a symbol c is in $FOLLOW(M)$, then there (by definition) is a derivation $S' \Rightarrow \gamma M c \delta$. But since $M \rightarrow \alpha N \beta$ and β is nullable, we can continue this by $\gamma M c \delta \Rightarrow \gamma \alpha N c \delta$, so c is also in $FOLLOW(N)$.

If a right-hand side contains several occurrences of nonterminals, we add constraints for all occurrences, *i.e.*, splitting the right-hand side into different α s, N s and β s. For example, the production $A \rightarrow BcB$ generates the constraint $\{c\} \subseteq FOLLOW(B)$ by splitting after the first B and, by splitting after the last B , we also get the constraint $FOLLOW(A) \subseteq FOLLOW(B)$.

We solve the constraints in the following fashion:

We start by assuming empty $FOLLOW$ sets for all nonterminals. We then handle the constraints of the form $FIRST(\beta) \subseteq FOLLOW(N)$: We compute $FIRST(\beta)$ and add this to $FOLLOW(N)$. Then, we handle the second type of constraints: For each constraint $FOLLOW(M) \subseteq FOLLOW(N)$, we add all elements of $FOLLOW(M)$ to $FOLLOW(N)$. We iterate these last steps until no further changes happen.

The steps taken to calculate the follow sets of a grammar are, hence:

1. Add a new nonterminal $S' \rightarrow S\$$, where S is the start symbol for the original grammar. S' is the start symbol for the extended grammar.
2. For each nonterminal N , locate all occurrences of N on the right-hand sides of productions. For each occurrence do the following:

- 2.1 Let β be the rest of the right-hand side after the occurrence of N . Note that β may be empty. In other words, the production is of the form $M \rightarrow \alpha N \beta$, where M is a nonterminal (possibly equal to N) and α and β are (possibly empty) sequences of grammar symbols. Note that if a right-hand-side contains several occurrences of N , we make a split for each occurrence.
 - 2.2 Let $m = FIRST(\beta)$. Add the constraint $m \subseteq FOLLOW(N)$ to the set of constraints. If β is empty, you can omit this constraint, as it does not add anything.
 - 2.3 If $Nullable(\beta)$, find the nonterminal M at the left-hand side of the production and add the constraint $FOLLOW(M) \subseteq FOLLOW(N)$. If $M = N$, you can omit the constraint, as it does not add anything. Note that if β is empty, $Nullable(\beta)$ is true.
3. Solve the constraints using the following steps:
- 3.1 Start with empty sets for $FOLLOW(N)$ for all nonterminals N (not including S').
 - 3.2 For each constraint of the form $m \subseteq FOLLOW(N)$ constructed in step 2.1, add the contents of m to $FOLLOW(N)$.
 - 3.3 Iterating until a fixed-point is reached, for each constraint of the form $FOLLOW(M) \subseteq FOLLOW(N)$, add the contents of $FOLLOW(M)$ to $FOLLOW(N)$.

We can take grammar 3.4 as an example of this. We first add the production

$$T' \rightarrow T\$$$

to the grammar to handle end-of-text conditions. The table below shows the constraints generated by each production

Production	Constraints
$T' \rightarrow T\$$	$\{\$ \} \subseteq FOLLOW(T)$
$T \rightarrow R$	$FOLLOW(T) \subseteq FOLLOW(R)$
$T \rightarrow aTc$	$\{c\} \subseteq FOLLOW(T)$
$R \rightarrow$	
$R \rightarrow RbR$	$\{b\} \subseteq FOLLOW(R), FOLLOW(R) \subseteq FOLLOW(R)$

In the above table, we have already calculated the required *FIRST* sets, so they are shown as explicit lists of terminals. To initialise the *FOLLOW* sets, we first use the constraints that involve these *FIRST* sets:

$$\begin{aligned} FOLLOW(T) &\supseteq \{\$, c\} \\ FOLLOW(R) &\supseteq \{b\} \end{aligned}$$

and then iterate calculation of the subset constraints. The only nontrivial constraint is $FOLLOW(T) \subseteq FOLLOW(R)$, so we get

$$\begin{aligned} FOLLOW(T) &\supseteq \{\$, c\} \\ FOLLOW(R) &\supseteq \{\$, c, b\} \end{aligned}$$

Now all constraints are satisfied, so we can replace subset with equality:

$$\begin{aligned} FOLLOW(T) &= \{\$, c\} \\ FOLLOW(R) &= \{\$, c, b\} \end{aligned}$$

If we return to the question of predictive parsing of grammar 3.4, we see that for the nonterminal R we should choose the empty production on any symbol in $FOLLOW(R)$, i.e., $\{\$, c, b\}$ and choose the non-empty production on the symbols in $FIRST(RbR)$, i.e., $\{b\}$. Since these sets overlap (on the symbol b), we can not uniquely choose a production for R based on the next input symbol. Hence, the revised construction of predictive parsers (see below) will reject this grammar as possibly ambiguous.

3.11 A larger example

The above examples of calculating $FIRST$ and $FOLLOW$ are rather small, so we show a somewhat more substantial example. The following grammar describes even-length strings of as and bs that are not of the form ww where w is any string of as and bs. In other words, the strings can not consist of two identical halves.

$$\begin{aligned} N &\rightarrow AB \\ N &\rightarrow BA \\ A &\rightarrow a \\ A &\rightarrow CAC \\ B &\rightarrow b \\ B &\rightarrow CBC \\ C &\rightarrow a \\ C &\rightarrow b \end{aligned}$$

The idea is that if the string does not consist of two identical halves, there must be a point in the first string that has an a where the equivalent point in the second string has a b or vice-versa. The grammar states that one of these is the case.

We first note that there are is empty production in the grammar, so no production can be *Nullable*. So we immediately set up the equations for *FIRST* for each nonterminal and right-hand side:

$$\begin{aligned}
 FIRST(N) &= FIRST(A B) \cup FIRST(B A) \\
 FIRST(A) &= FIRST(a) \cup FIRST(C A C) \\
 FIRST(B) &= FIRST(b) \cup FIRST(C B C) \\
 FIRST(C) &= FIRST(a) \cup FIRST(b) \\
 \\
 FIRST(A B) &= FIRST(A) \\
 FIRST(B A) &= FIRST(B) \\
 FIRST(a) &= \{a\} \\
 FIRST(C A C) &= FIRST(C) \\
 FIRST(b) &= \{b\} \\
 FIRST(C B C) &= FIRST(C)
 \end{aligned}$$

which we solve by fixed-point iteration. We initially set the *FIRST* sets for the nonterminals to the empty sets, calculate the *FIRST* sets for right-hand sides and then nonterminals, repeating the last two steps until no changes occur:

RHS	Iteration 1	Iteration 2	Iteration 3
$A B$	\emptyset	$\{a\}$	$\{a, b\}$
$B A$	\emptyset	$\{b\}$	$\{a, b\}$
a	$\{a\}$	$\{a\}$	$\{a\}$
$C A C$	\emptyset	$\{a, b\}$	$\{a, b\}$
b	$\{b\}$	$\{b\}$	$\{b\}$
$C B C$	\emptyset	$\{a, b\}$	$\{a, b\}$
Nonterminal			
N	\emptyset	$\{a, b\}$	$\{a, b\}$
A	$\{a\}$	$\{a, b\}$	$\{a, b\}$
B	$\{b\}$	$\{a, b\}$	$\{a, b\}$
C	$\{a, b\}$	$\{a, b\}$	$\{a, b\}$

The next iteration does not add anything, so the fixed-point is reached. We now add the production $N' \rightarrow N\$$ and set up the constraints for calculating *FOLLOW* sets:

Production	Constraints
$N' \rightarrow N\$$	$\{\$ \} \subseteq FOLLOW(N)$
$N \rightarrow A B$	$FIRST(B) \subseteq FOLLOW(A), FOLLOW(N) \subseteq FOLLOW(B)$
$N \rightarrow B A$	$FIRST(A) \subseteq FOLLOW(B), FOLLOW(N) \subseteq FOLLOW(A)$
$A \rightarrow a$	
$A \rightarrow C A C$	$FIRST(A) \subseteq FOLLOW(C), FIRST(C) \subseteq FOLLOW(A), FOLLOW(A) \subseteq FOLLOW(C)$
$B \rightarrow b$	
$B \rightarrow C B C$	$FIRST(B) \subseteq FOLLOW(C), FIRST(C) \subseteq FOLLOW(B), FOLLOW(B) \subseteq FOLLOW(C)$
$C \rightarrow a$	
$C \rightarrow b$	

We first use the constraint $\{\$ \} \subseteq FOLLOW(N)$ and constraints of the form $FIRST(\dots) \subseteq FOLLOW(\dots)$ to get the initial sets:

$$\begin{aligned}
FOLLOW(N) &\subseteq \{\$ \} \\
FOLLOW(A) &\subseteq \{a, b\} \\
FOLLOW(B) &\subseteq \{a, b\} \\
FOLLOW(C) &\subseteq \{a, b\}
\end{aligned}$$

and then use the constraints of the form $FOLLOW(\dots) \subseteq FOLLOW(\dots)$. If we do this in top-down order, we get after one iteration:

$$\begin{aligned}
FOLLOW(N) &\subseteq \{\$ \} \\
FOLLOW(A) &\subseteq \{a, b, \$ \} \\
FOLLOW(B) &\subseteq \{a, b, \$ \} \\
FOLLOW(C) &\subseteq \{a, b, \$ \}
\end{aligned}$$

Another iteration does not add anything, so the final result is

$$\begin{aligned}
FOLLOW(N) &= \{\$ \} \\
FOLLOW(A) &= \{a, b, \$ \} \\
FOLLOW(B) &= \{a, b, \$ \} \\
FOLLOW(C) &= \{a, b, \$ \}
\end{aligned}$$

Suggested exercises: 3.8 (*FOLLOW* only).

3.12 LL(1) parsing

We have, in the previous sections, looked at how we can choose productions based on *FIRST* and *FOLLOW* sets, *i.e.* using the rule that we choose a production $N \rightarrow \alpha$ on input symbol c if

- $c \in FIRST(\alpha)$, or
- $Nullable(\alpha)$ and $c \in FOLLOW(N)$.

If we can always choose a production uniquely by using these rules, this is called LL(1) parsing – the first L indicates the reading direction (left-to-right), the second L indicates the derivation order (left) and the 1 indicates that there is a one-symbol lookahead. A grammar that can be parsed using LL(1) parsing is called an LL(1) grammar.

In the rest of this section, we shall see how we can implement LL(1) parsers as programs. We look at two implementation methods: Recursive descent, where grammar structure is directly translated into the structure of a program, and a table-based approach that encodes the decision process in a table.

3.12.1 Recursive descent

As the name indicates, *recursive descent* uses recursive functions to implement predictive parsing. The central idea is that each nonterminal in the grammar is implemented by a function in the program.

Each such function looks at the next input symbol in order to choose one of the productions for the nonterminal, using the criteria shown in the beginning of section 3.12. The right-hand side of the chosen production is then used for parsing in the following way:

A terminal on the right-hand side is matched against the next input symbol. If they match, we move on to the following input symbol and the next symbol on the right hand side, otherwise an error is reported.

A nonterminal on the right-hand side is handled by calling the corresponding function and, after this call returns, continuing with the next symbol on the right-hand side.

When there are no more symbols on the right-hand side, the function returns.

As an example, figure 3.16 shows pseudo-code for a recursive descent parser for grammar 3.9. We have constructed this program by the following process:

We have first added a production $T' \rightarrow T\$$ and calculated *FIRST* and *FOLLOW* for all productions.

T' has only one production, so the choice is trivial. However, we have added a check on the next input symbol anyway, so we can report an error if it is not in $FIRST(T')$. This is shown in the function `parseT'`.

For the `parseT` function, we look at the productions for T . As $FIRST(R) = \{b\}$, the production $T \rightarrow R$ is chosen on the symbol b . Since R is also *Nullable*, we must choose this production also on symbols in $FOLLOW(T)$, i.e., c or $\$$. $FIRST(aTc) = \{a\}$, so we select $T \rightarrow aTc$ on an a . On all other symbols we report an error.

```

function parseT'() =
  if next = 'a' or next = 'b' or next = '$' then
    parseT() ; match('$')
  else reportError()

function parseT() =
  if next = 'b' or next = 'c' or next = '$' then
    parseR()
  else if next = 'a' then
    match('a') ; parseT() ; match('c')
  else reportError()

function parseR() =
  if next = 'c' or next = '$' then
    (* do nothing *)
  else if next = 'b' then
    match('b') ; parseR()
  else reportError()

```

Figure 3.16: Recursive descent parser for grammar 3.9

For `parseR`, we must choose the empty production on symbols in $FOLLOW(R)$ (c or \$). The production $R \rightarrow bR$ is chosen on input b. Again, all other symbols produce an error.

The function `match` takes as argument a symbol, which it tests for equality with the next input symbol. If they are equal, the following symbol is read into the variable `next`. We assume `next` is initialised to the first input symbol before `parseT'` is called.

The program in figure 3.16 only checks if the input is valid. It can easily be extended to construct a syntax tree by letting the parse functions return the sub-trees for the parts of input that they parse.

3.12.2 Table-driven LL(1) parsing

In table-driven LL(1) parsing, we encode the selection of productions into a table instead of in the program text. A simple non-recursive program uses this table and a stack to perform the parsing.

The table is cross-indexed by nonterminal and terminal and contains for each such pair the production (if any) that is chosen for that nonterminal when that terminal is the next input symbol. This decision is made just as for recursive descent

	a	b	c	\$
T'	$T' \rightarrow T\$$	$T' \rightarrow T\$$		$T' \rightarrow T\$$
T	$T \rightarrow aTc$	$T \rightarrow R$	$T \rightarrow R$	$T \rightarrow R$
R		$R \rightarrow bR$	$R \rightarrow$	$R \rightarrow$

Figure 3.17: LL(1) table for grammar 3.9

parsing: The production $N \rightarrow \alpha$ is in the table at (N, a) if a is in $FIRST(\alpha)$ or if both $Nullable(\alpha)$ and a is in $FOLLOW(N)$.

For grammar 3.9 we get the table shown in figure 3.17.

The program that uses this table is shown in figure 3.18. It uses a stack, which at any time (read from top to bottom) contains the part of the current derivation that has not yet been matched to the input. When this eventually becomes empty, the parse is finished. If the stack is non-empty, and the top of the stack contains a terminal, that terminal is matched against the input and popped from the stack. Otherwise, the top of the stack must be a nonterminal, which we cross-index in the table with the next input symbol. If the table-entry is empty, we report an error. If not, we pop the nonterminal from the stack and replace this by the right-hand side of the production in the table entry. The list of symbols on the right-hand side are pushed such that the first of these will be at the top of the stack.

As an example, figure 3.19 shows the input and stack at each step during parsing of the string `aabbbcc$` using the table in figure 3.17. The top of the stack is to the left.

The program in figure 3.18, like the one in figure 3.16, only checks if the input is valid. It, too, can be extended to build a syntax tree. This can be done by letting each nonterminal on the stack point to its node in the partially built syntax tree. When the nonterminal is replaced by one of its right-hand sides, nodes for the symbols on the right-hand side are added as children to the node.

3.12.3 Conflicts

When a symbol a allows several choices of production for nonterminal N we say that there is a *conflict* on that symbol for that nonterminal. Conflicts may be caused by ambiguous grammars (indeed all ambiguous grammars will cause conflicts) but there are also unambiguous grammars that cause conflicts. An example of this is the unambiguous expression grammar (grammar 3.11). We will in the next section see how we can rewrite this grammar to avoid conflicts, but it must be noted that this is not always possible: There are languages for which there exist unambiguous context-free grammars but where no grammar for the language generates a conflict-free LL(1) table. Such languages are said to be non-LL(1). It is, however, important

```

stack := empty ; push(T',stack)
while stack <> empty do
  if top(stack) is a terminal then
    match(top(stack)) ; pop(stack)
  else if table(top(stack),next) = empty then
    reportError
  else
    rhs := rightHandSide(table(top(stack),next)) ;
    pop(stack) ;
    pushList(rhs,stack)

```

Figure 3.18: Program for table-driven LL(1) parsing

input	stack
aabbbcc\$	T'
aabbbcc\$	$T\$$
aabbbcc\$	a $Tc\$$
abbbcc\$	$Tc\$$
abbbcc\$	a $Tcc\$$
bbbcc\$	$Tcc\$$
bbbcc\$	$Rcc\$$
bbbcc\$	b $Rcc\$$
bbcc\$	$Rcc\$$
bbcc\$	b $Rcc\$$
bcc\$	$Rcc\$$
bcc\$	b $Rcc\$$
cc\$	$Rcc\$$
cc\$	cc\$
c\$	c\$
\$	\$

Figure 3.19: Input and stack during table-driven LL(1) parsing

to note the difference between a non-LL(1) language and a non-LL(1) grammar: A language may well be LL(1) even though the grammar used to describe it is not.

3.13 Rewriting a grammar for LL(1) parsing

In this section we will look at methods for rewriting grammars such that they are more palatable for LL(1) parsing. In particular, we will look at *elimination of left-recursion* and at *left factorisation*.

It must, however, be noted that not all grammars can be rewritten to allow LL(1) parsing. In these cases stronger parsing techniques must be used.

3.13.1 Eliminating left-recursion

As mentioned above, the unambiguous expression grammar (grammar 3.11) is not LL(1). The reason is that all productions in *Exp* and *Exp2* have the same *FIRST* sets. Overlap like this will always happen when there are left-recursive productions in the grammar, as the *FIRST* set of a left-recursive production will include the *FIRST* set of the nonterminal itself and hence be a superset of the *FIRST* sets of all the other productions for that nonterminal. To solve this problem, we must avoid left-recursion in the grammar. We start by looking at direct left-recursion.

When we have a nonterminal with some left-recursive and some productions that are not, *i.e.*,

$$\begin{array}{lcl}
 N & \rightarrow & N\alpha_1 \\
 & \vdots & \\
 N & \rightarrow & N\alpha_m \\
 N & \rightarrow & \beta_1 \\
 & \vdots & \\
 N & \rightarrow & \beta_n
 \end{array}$$

where the β_i do not start with N , we observe that this generates all sequences that start with one of the β_i and continues with any number (including 0) of the α_j . In other words, the grammar is equivalent to the regular expression $(\beta_1 | \dots | \beta_n)(\alpha_1 | \dots | \alpha_m)^*$.

We saw in figure 3.1 a method for converting regular expressions into context-free grammars can generate the same set of strings. By following this procedure and simplifying a bit afterwards, we get this equivalent grammar:

$$\begin{aligned}
Exp &\rightarrow Exp2\ Exp_* \\
Exp_* &\rightarrow +\ Exp2\ Exp_* \\
Exp_* &\rightarrow -\ Exp2\ Exp_* \\
Exp_* &\rightarrow \\
Exp2 &\rightarrow Exp3\ Exp2_* \\
Exp2_* &\rightarrow *\ Exp3\ Exp2_* \\
Exp2_* &\rightarrow /\ Exp3\ Exp2_* \\
Exp2_* &\rightarrow \\
Exp3 &\rightarrow \mathbf{num} \\
Exp3 &\rightarrow (\ Exp\)
\end{aligned}$$

Grammar 3.20: Removing left-recursion from grammar 3.11

$$\begin{aligned}
N &\rightarrow \beta_1 N_* \\
&\vdots \\
N &\rightarrow \beta_n N_* \\
N_* &\rightarrow \alpha_1 N_* \\
&\vdots \\
N_* &\rightarrow \alpha_m N_* \\
N_* &\rightarrow
\end{aligned}$$

where N_* is a new nonterminal that generates a sequence of α s.

Note that, since the β_i do not start with N , there is no direct left-recursion in the first n productions. Since N_* is a new nonterminal, the α_j can not start with this, so the last m productions can't have direct left-recursion either.

There may, however, still be *indirect* left-recursion if any of the α_j are nullable or the β_i can derive something starting with N . We will briefly look at indirect left-recursion below.

While we have eliminated direct left-recursion, we have also changed the syntax trees that are built from the strings that are parsed. Hence, after parsing, the syntax tree must be re-structured to obtain the structure that the original grammar describes. We will return to this in section 3.17.

As an example of left-recursion removal, we take the unambiguous expression grammar 3.11. This has left recursion in both Exp and $Exp2$, so we apply the transformation to both of these to obtain grammar 3.20. The resulting grammar 3.20 is now LL(1), which can be verified by generating an LL(1) table for it.

Indirect left-recursion

The transformation shown in section 3.13.1 is only applicable in the simple case where there is no *indirect left-recursion*. Indirect left-recursion can have several faces:

1. There are mutually left-recursive productions

$$\begin{array}{rcl}
 N_1 & \rightarrow & N_2\alpha_1 \\
 N_2 & \rightarrow & N_3\alpha_2 \\
 & \vdots & \\
 N_{k-1} & \rightarrow & N_k\alpha_{k-1} \\
 N_k & \rightarrow & N_1\alpha_k
 \end{array}$$

2. There is a production $N \rightarrow \alpha N \beta$ where α is *Nullable*.

or any combination of the two. More precisely, a grammar is (directly or indirectly) left-recursive if there is a non-empty derivation sequence $N \Rightarrow N\alpha$, *i.e.*, if a nonterminal derives a sequence of grammar symbols that start by that same nonterminal. If there is indirect left-recursion, we must first rewrite the grammar to make the left-recursion direct and then use the transformation above.

Rewriting a grammar to turn indirect left-recursion into direct left-recursion can be done systematically, but the process is a bit complicated. We will not go into this here, as in practice most cases of left-recursion are direct left-recursion. Details can be found in [4].

3.13.2 Left-factorisation

If two productions for the same nonterminal begin with the same sequence of symbols, they obviously have overlapping *FIRST* sets. As an example, in grammar 3.3 the two productions for *if* have overlapping prefixes. We rewrite this in such a way that the overlapping productions are made into a single production that contains the common prefix of the productions and uses a new auxiliary nonterminal for the different suffixes. See grammar 3.21. In this grammar³, we can uniquely choose one of the productions for *Stat* based on one input token.

For most grammars, combining productions with common prefix will solve the problem. However, in this particular example the grammar still is not LL(1): We can not uniquely choose a production for the auxiliary nonterminal *Elsepart*, since *else* is in *FOLLOW(Elsepart)* as well as in the *FIRST* set of the first production for *Elsepart*. This should not be a surprise to us, since, after all, the grammar

³We have omitted the production for semicolon, as that would only muddle the issue by introducing more ambiguity.

$$\begin{array}{ll}
Stat & \rightarrow \text{ id } := Exp \\
Stat & \rightarrow \text{ if } Exp \text{ then } Stat \text{ Elsepart} \\
\\
Elsepart & \rightarrow \text{ else } Stat \\
Elsepart & \rightarrow
\end{array}$$

Grammar 3.21: Left-factorised grammar for conditionals

is ambiguous and ambiguous grammars can not be LL(1). The equivalent unambiguous grammar (grammar 3.13) can not easily be rewritten to a form suitable for LL(1), so in practice grammar 3.21 is used anyway and the conflict is handled by choosing the non-empty production for *Elsepart* whenever the symbol `else` is encountered, as this gives the desired behaviour of letting an `else` match the nearest `if`. We can achieve this by removing the empty production from the table entry for *Elsepart/else*, so only the non-empty production $Elsepart \rightarrow \text{ else } Stat$ remains.

Very few LL(1) conflicts caused by ambiguity can be removed in this way, however, without also changing the language recognized by the grammar. For example, operator precedence ambiguity can not be resolved by deleting conflicting entries in the LL(1) table.

3.13.3 Construction of LL(1) parsers summarized

1. Eliminate ambiguity
2. Eliminate left-recursion
3. Perform left factorisation where required
4. Add an extra start production $S' \rightarrow S\$$ to the grammar.
5. Calculate *FIRST* for every production and *FOLLOW* for every nonterminal.
6. For nonterminal N and input symbol c , choose production $N \rightarrow \alpha$ when:
 - $c \in FIRST(\alpha)$, or
 - $Nullable(\alpha)$ and $c \in FOLLOW(N)$.

This choice is encoded either in a table or a recursive-descent program.

Suggested exercises: 3.14.

3.14 SLR parsing

A problem with LL(1) parsing is that most grammars need extensive rewriting to get them into a form that allows unique choice of production. Even though this rewriting can, to a large extent, be automated, there are still a large number of grammars that can not be automatically transformed into LL(1) grammars.

LR parsers is a class of bottom-up methods for parsing that accept a much larger class of grammars than LL(1) parsing, though still not all grammars. The main advantage of LR parsing is that less rewriting is required to get a grammar in acceptable form for LR parsing than is the case for LL(1) parsing. Furthermore, as we shall see in section 3.16, LR parsers allow external declaration of operator precedences for resolving ambiguity, instead of requiring the grammars themselves to be unambiguous.

We will look at a simple form of LR-parsing called SLR parsing. The letters “SLR” stand for “Simple”, “Left” and “Right”. “Left” indicates that the input is read from left to right and the “Right” indicates that a rightmost derivation is built.

LR parsers are table-driven bottom-up parsers and use two kinds of “actions” involving the input stream and a stack:

shift: A symbol is read from the input and pushed on the stack.

reduce: The top N elements of the stack hold symbols identical to the N symbols on the right-hand side of a specified production. These N symbols are by the reduce action replaced by the nonterminal at the left-hand side of the specified production. Contrary to LL parsers, the stack holds the right-hand-side symbols such that the *last* symbol on the right-hand side is at the top of the stack.

If the input text does not conform to the grammar, there will at some point during the parsing be no applicable actions and the parser will stop with an error message. Otherwise, the parser will read through all the input and leave a single element (the start symbol of the grammar) on the stack.

LR parsers are also called *shift-reduce parsers*. As with LL(1), our aim is to make the choice of action depend only on the next input symbol and the symbol on top of the stack. To achieve this, we construct a DFA. Conceptually, this DFA reads the contents of the stack, starting from the bottom. If the DFA is in an accepting state when it reaches the top of the stack, the correct action is reduction by a production that is determined by the next input symbol and a mark on the accepting DFA state. If the DFA is not in an accepting state when it reaches the stack top, the correct action is a shift on one of the symbols for which there is an outgoing edge from the DFA state. Hence, at every step, the DFA reads the stack from bottom to top and the action is determined by looking at the DFA state and the next input symbol.

Letting the DFA read the entire stack at every action is, however, not very efficient so, instead, we store with each stack element the state of the DFA when it reads this element. This way, we do not need to start from the bottom of the stack, but can start from the current stack top starting the DFA in the stored state instead of in its initial state.

When the DFA has indicated a shift, the course of action is easy: We get the state from the top of the stack and follow the transition marked with the next input symbol to find the next DFA state, and we store both the symbol and the new state on the stack.

If the DFA indicated a reduce, we pop the symbols corresponding to the right-hand side of the production off the stack. We then read the DFA state from the new stack top. This DFA state should have a transition on the nonterminal that is the left-hand side of the production, so we store both the nonterminal and the state at the end of the transition on the stack.

With these optimisations, the DFA only has to inspect a terminal or nonterminal once, at the time it is pushed on the stack. At all other times, it just needs to read the DFA state that is stored on the top of the stack. We, actually, do not need to store the current input symbol or nonterminal once we have made a transition on it, as no future transitions will depend on it – the stored DFA state is enough. So we can let each stack element just contain the DFA state instead of both the symbol and the state. We still use the DFA to determine the next action, but it now only needs to look at the current state (stored at the top of the stack) and the next input symbol (at a shift action) or nonterminal (at a reduce action).

We represent the DFA as a table, where we cross-index a DFA state with a symbol (terminal or nonterminal) and find one of the following actions:

- shift n*: Read next input symbol and push state *n* on the stack.
- go n*: Push state *n* on the stack.
- reduce p*: Reduce with the production numbered *p*.
- accept*: Parsing has completed successfully..
- error*: A syntax error has been detected.

Note that the current state is always found at the top of the stack. *Shift* and *reduce* actions are used when a state is cross-indexed with a terminal symbol. *Go* actions are used when a state is cross-indexed with a nonterminal. A *Go* action can occur only immediately after a reduce, but we can not in the table combine the *go* actions with the *reduce* actions, as the destination state of a *go* action depends on the state at the top of the stack *after* the right-hand side of the reduced production is popped off: A *reduce* in the current state is immediately followed by a *go* in the state that is found when the stack is popped.

An example SLR table is shown in figure 3.22. The table has been produced

	a	b	c	\$	<i>T</i>	<i>R</i>
0	s3	s4	r3	r3	g1	g2
1				a		
2			r1	r1		
3	s3	s4	r3	r3	g5	g2
4		s4	r3	r3		g6
5			s7			
6			r4	r4		
7			r2	r2		

Figure 3.22: SLR table for grammar 3.9

from grammar 3.9 by the method shown below in section 3.15. The actions have been abbreviated to their first letters and *error* is shown as a blank entry.

The algorithm for parsing a string using the table is shown in figure 3.23. The shown algorithm just determines if a string is in the language generated by the grammar. It can, however, easily be extended to build a syntax tree: Each stack element holds (in addition to the state number) a portion of a syntax tree. When performing a *reduce* action, a new (partial) syntax tree is built by using the non-terminal from the reduced production as root and the syntax trees stored at the popped-off stack elements as children. The new tree and the new state are then pushed (as a single stack element).

Figure 3.24 shows an example of parsing the string aabbbcc using the table in figure 3.22. The sequences of numbers in the “stack” column represent the stack contents with the stack bottom shown to the left and the stack top to the right. At each step, we look at the next input symbol (at the left end of the string in the input column) and the state at the top of the stack (at the right end of the sequence in the stack column). We look up the pair of input symbol and state in the table and find an action, which is shown in the action column. When the shown action is a reduce action, we also show the reduction used (in parentheses) and after a semicolon also the go action that is performed after the reduction.

3.15 Constructing SLR parse tables

An SLR parse table has a DFA as its core. Constructing this DFA from the grammar is similar to constructing a DFA from a regular expression, as shown in chapter 2: We first construct an NFA using techniques similar to those in section 2.4 and then convert this into a DFA using the construction shown in section 2.6.

Before we construct the NFA, we extend the grammar with a new starting production. Doing this to grammar 3.9 yields grammar 3.25.

```

stack := empty ; push(0,stack) ; read(next)
loop
  case table[top(stack),next] of
    shift s:  push(s,stack) ;
              read(next)

    reduce p: n := the left-hand side of production p ;
              r := the number of symbols
                  on the right-hand side of p ;
              pop r elements from the stack ;
              push(s,stack)
              where table[top(stack),n] = go s

    accept:   terminate with success

    error:    reportError
  endloop

```

Figure 3.23: Algorithm for SLR parsing

input	stack	action
aabbbcc\$	0	s3
abbbcc\$	03	s3
bbbcc\$	033	s4
bbcc\$	0334	s4
bcc\$	03344	s4
cc\$	033444	r3 ($R \rightarrow$) ; g6
cc\$	0334446	r4 ($R \rightarrow bR$) ; g6
cc\$	033446	r4 ($R \rightarrow bR$) ; g6
cc\$	03346	r4 ($R \rightarrow bR$) ; g2
cc\$	0332	r1 ($T \rightarrow R$) ; g5
cc\$	0335	s7
c\$	03357	r2 ($T \rightarrow aTc$) ; g5
c\$	035	s7
\$	0357	r2 ($T \rightarrow aTc$) ; g1
\$	01	accept

Figure 3.24: Example SLR parsing

- 0: $T' \rightarrow T$
 1: $T \rightarrow R$
 2: $T \rightarrow aTc$
 3: $R \rightarrow$
 4: $R \rightarrow bR$

Grammar 3.25: Example grammar for SLR-table construction

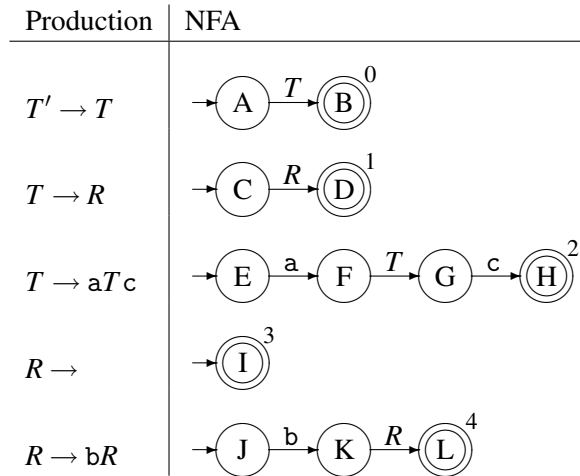


Figure 3.26: NFAs for the productions in grammar 3.25

The next step is to make an NFA for each production. This is done exactly like in section 2.4, treating both terminals and nonterminals as alphabet symbols. The accepting state of each NFA is labelled with the number of the corresponding production. The result is shown in figure 3.26. Note that we have used the optimised construction for ϵ (the empty production) as shown in figure 2.6.

The NFAs in figure 3.26 make transitions both on terminals and nonterminals. Transitions by terminal corresponds to *shift* actions and transitions on nonterminals correspond to *go* actions. A *go* action happens after a reduction, whereby some elements of the stack (corresponding to the right-hand side of a production) are replaced by a nonterminal (corresponding to the left-hand side of that production). However, before we can reduce on a nonterminal, the symbols that form the right-hand side must be on the stack. So we prepare for a later transition on a nonterminal by allowing transitions that, eventually, will leave the symbols forming right-hand side of the production on the stack, so we can reduce these to the nonterminal and then make a transition in this.

state	epsilon-transitions
A	C, E
C	I, J
F	C, E
K	I, J

Figure 3.27: Epsilon-transitions added to figure 3.26

So, whenever a transition by a nonterminal is possible, we also allow transitions on the sequences of symbols on the right-hand sides of the productions for that nonterminal. We achieve this by adding epsilon transitions to the NFAs in figure 3.26: Whenever there is a transition from state s to state t on a nonterminal N , we add epsilon transitions from s to the initial states of all the NFAs for productions with N on the left-hand side. Adding the epsilon transitions as arrows in figure 3.26 would produce a very cluttered picture, so instead we note the transitions in a separate table, shown in figure 3.27. But this is for presentation purposes only: The transitions have the same meaning regardless of whether they are shown in the picture of the DFA or in the table.

Together with these epsilon-transitions, the NFAs in figure 3.26 form a single, combined NFA. This NFA has the starting state A (the starting state of the NFA for the added start production) and an accepting state for each production in the grammar. **We must now convert this NFA into a DFA using** the subset construction shown in section 2.5. Instead of showing the resulting DFA graphically, we construct a table where transitions on terminals are shown as *shift* actions and transitions on nonterminals as *go* actions. This will make the table look similar to figure 3.22, except that no *reduce* or *accept* actions are present yet. Figure 3.28 shows the DFA constructed from the NFA made by adding epsilon-transitions in 3.27 to figure 3.26. The sets of NFA states that form each DFA state is shown in the second column of the table in figure 3.28. We will need these below for adding *reduce* and *accept* actions, but once this is done, we will not need them anymore, so we can remove them from the final table.

To add *reduce* and *accept* actions, we first need to compute the *FOLLOW* sets for each nonterminal, as described in section 3.10. For purpose of calculating *FOLLOW*, we add yet another extra start production: $T'' \rightarrow T'\$,$ to handle end-of-text conditions as described in section 3.10. This gives us the following result:

$$\begin{aligned}
 FOLLOW(T') &= \{\$ \} \\
 FOLLOW(T) &= \{c, \$ \} \\
 FOLLOW(R) &= \{c, \$ \}
 \end{aligned}$$

We then add *reduce* actions by the following rule: If a DFA state s contains the accepting NFA state for a production $p : N \rightarrow \alpha$, we add *reduce p* as action to s on

DFA state	NFA states	Transitions				
		a	b	c	T	R
0	A, C, E, I, J	s3	s4		g1	g2
1	B					
2	D					
3	F, C, E, I, J	s3	s4		g5	g2
4	K, I, J		s4			g6
5	G			s7		
6	L					
7	H					

Figure 3.28: SLR DFA for grammar 3.9

all symbols in $FOLLOW(N)$. Reduction on production 0 (the extra start production that was added before constructing the NFA) is written as *accept*.

In figure 3.28, state 0 contains NFA state I, which accepts production 3. Hence, we add r_3 as actions at the symbols c and $\$$ (as these are in $FOLLOW(R)$). State 1 contains NFA state B, which accepts production 0. We add this at the symbol $\$$ ($FOLLOW(T')$). As noted above, this is written as *accept* (abbreviated to “a”). In the same way, we add reduce actions to state 3, 4, 6 and 7. The result is shown in figure 3.22.

Figure 3.29 summarises the SLR construction.

3.15.1 Conflicts in SLR parse-tables

When *reduce* actions are added to SLR parse-tables, we might add one to a place where there is already a *shift* action, or we may add *reduce* actions for several different productions to the same place. When either of this happens, we no longer have a unique choice of action, *i.e.*, we have a *conflict*. The first situation is called a *shift-reduce conflict* and the other case a *reduce-reduce conflict*. Both may occur in the same place.

Conflicts are often caused by ambiguous grammars, but (as is the case for LL-parsers) even some non-ambiguous grammars may generate conflicts. If a conflict is caused by an ambiguous grammar, it is usually (but not always) possible to find an equivalent unambiguous grammar. Methods for eliminating ambiguity were discussed in sections 3.4 and 3.5. Alternatively, operator precedence declarations may be used to disambiguate an ambiguous grammar, as we shall see in section 3.16.

But even unambiguous grammars may in some cases generate conflicts in SLR-tables. In some cases, it is still possible to rewrite the grammar to get around the problem, but in a few cases the language simply is not SLR. Rewriting an

1. Add the production $S' \rightarrow S$, where S is the start symbol of the grammar.
2. Make an NFA for the right-hand side of each production.
3. If an NFA state s has an outgoing transition on a nonterminal N , add epsilon-transitions from s to the starting states of the NFAs for the right-hand sides of the productions for N .
4. Convert the combined NFA to a DFA. Use the starting state of the NFA for the production added in step 1 as the starting state for the combined NFA.
5. Build a table cross-indexed by the DFA states and grammar symbols (terminals including $\$$ and nonterminals). Add *shift* actions for transitions on terminals and *go* actions for transitions on nonterminals.
6. Calculate *FOLLOW* for each nonterminal. For this purpose, we add one more start production: $S'' \rightarrow S'\$$.
7. When a DFA state contains an accepting NFA state marked with production number p , where the nonterminal for p is N , find the symbols in $FOLLOW(N)$ and add a *reduce* p action in the DFA state at all these symbols. If production p is the production added in step 1, add an *accept* action instead of a *reduce* p action.

Figure 3.29: Summary of SLR parse-table construction

unambiguous grammar to eliminate conflicts is somewhat of an art. Investigation of the NFA states that form the problematic DFA state will often help identifying the exact nature of the problem, which is the first step towards solving it. Sometimes, changing a production from left-recursive to right-recursive may help, even though left-recursion in general is not a problem for SLR-parsers, as it is for LL(1)-parsers.

Suggested exercises: 3.16.

3.16 Using precedence rules in LR parse tables

We saw in section 3.13.2, that the conflict arising from the dangling-else ambiguity could be removed by removing one of the entries in the LL(1) parse table. Resolving ambiguity by deleting conflicting actions can also be done in SLR-tables. In general, there are more cases where this can be done successfully for SLR-parsers than for LL(1)-parsers. In particular, ambiguity in expression grammars like gram-

mar 3.2 can be eliminated this way in an SLR table, but not in an LL(1) table. Most LR-parser generators allow declarations of precedence and associativity for tokens used as infix-operators. These declarations are then used to eliminate conflicts in the parse tables.

There are several advantages to this approach:

- Ambiguous expression grammars are more compact and easier to read than unambiguous grammars in the style of section 3.4.1.
- The parse tables constructed from ambiguous grammars are often smaller than tables produced from equivalent unambiguous grammars.
- Parsing using ambiguous grammars is (slightly) faster, as fewer reductions of the form $Exp2 \rightarrow Exp3$ etc. are required.

Using precedence rules to eliminate conflicts is very simple. Grammar 3.2 will generate several conflicts:

- 1) A conflict between shifting on + and reducing by the production $Exp \rightarrow Exp + Exp$.
- 2) A conflict between shifting on + and reducing by the production $Exp \rightarrow Exp * Exp$.
- 3) A conflict between shifting on * and reducing by the production $Exp \rightarrow Exp + Exp$.
- 4) A conflict between shifting on * and reducing by the production $Exp \rightarrow Exp * Exp$.

And several more of similar nature involving - and /, for a total of 16 conflicts. Let us take each of the four conflicts above in turn and see how precedence rules can be used to eliminate them. We use the rules that + and * are both left-associative and that * binds more strongly than +.

- 1) This conflict arises from expressions like $a+b+c$. After having read $a+b$, the next input symbol is a +. We can now either choose to reduce $a+b$, grouping around the first addition before the second, or shift on the plus, which will later lead to $b+c$ being reduced and hence grouping around the second addition before the first. Since the rules give that + is left-associative, we prefer the first of these options and, hence, eliminate the shift-action from the table and keep the reduce-action.

- 2) The offending expressions here have the form $a*b+c$. Since the rules make multiplication bind stronger than addition, we, again, prefer reduction over shifting.
- 3) In expressions of the form $a+b*c$, the rules, again, make multiplication bind stronger, so we do a shift to avoid grouping around the $+$ operator and, hence, eliminate the reduce-action from the table.
- 4) This case is identical to case 1, where an operator that by the rules is left-associative conflicts with itself. We, like in case 1, handle this by eliminating the shift.

In general, elimination of conflicts by operator precedence declarations can be summarised into the following rules:

- a) If the conflict is between two operators of different priority, eliminate the action with the lowest priority operator in favour of the action with the highest priority. The operator associated with a reduce-action is the operator used in the production that is reduced.
- b) If the conflict is between operators of the same priority, the associativity (which must be the same, as noted in section 3.4.1) of the operators is used: If the operators are left-associative, the shift-action is eliminated and the reduce-action retained. If the operators are right-associative, the reduce-action is eliminated and the shift-action retained. If the operators are non-associative, both actions are eliminated.
- c) If there are several operators with declared precedence in the production that is used in a reduce-action, the last of these is used to determine the precedence of the reduce-action.⁴

Prefix and postfix operators can be handled similarly. Associativity only applies to infix operators, so only the precedence of prefix and postfix operators matters.

Note that only shift-reduce conflicts are eliminated by the above rules. Some parser generators allow also reduce-reduce conflicts to be eliminated by precedence rules (in which case the production with the highest-precedence operator is preferred), but this is not as obviously useful as the above.

The dangling-else ambiguity (section 3.5) can also be eliminated using precedence rules. If we have read *if Exp then Stat* and the next symbol is a *else*, we want to shift on *else*, so the *else* will be associated with the *then*. If we, instead, reduced on the production $Stat \rightarrow \text{if } Exp \text{ then } Stat$, we would lose this association. Giving *else* a higher precedence than *then* or giving them the same

⁴Using several operators with declared priorities in the same production should be done with care.

precedence and making them right-associative will ensure that a shift is made on else when we need it.

Not all conflicts should be eliminated by precedence rules. If you blindly add precedence rules until no conflicts are reported, you risk disallowing also legal strings, so the parser will accept only a subset of the intended language. Normally, you should only use precedence declarations to specify operator hierarchies, unless you have analysed the parser actions carefully and found that there is no undesirable consequences of adding the precedence rules.

Suggested exercises: 3.18.

3.17 Using LR-parser generators

Most LR-parser generators use an extended version of the SLR construction called LALR(1). The “LA” in the abbreviation is short for “lookahead” and the (1) indicates that the lookahead is one symbol, i.e., the next input symbol.

We have chosen to present the SLR construction instead of the LALR(1) construction for several reasons:

- It is simpler.
- In practice, LALR(1) handles only a few more grammars than SLR.
- When a grammar is in the SLR class, the parse-table produced by an SLR parser generator is identical to the table produced by an LALR(1) parser generator.
- Understanding of SLR principles is sufficient to know how to handle a grammar rejected by a LALR(1) parser generator by adding precedence declarations or by rewriting the grammar.

In short, knowledge of SLR parsing is sufficient when using LALR(1) parser generators.

Most LR-parser generators organise their input in several sections:

- Declarations of the terminals and nonterminals used.
- Declaration of the start symbol of the grammar.
- Declarations of operator precedence.
- The productions of the grammar.
- Declaration of various auxiliary functions and data-types used in the actions (see below).

3.17.1 Declarations and actions

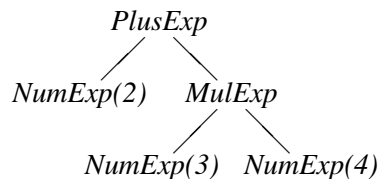
Each nonterminal and terminal is declared and associated with a data-type. For a terminal, the data-type is used to hold the values that are associated with the tokens that come from the lexer, *e.g.*, the values of numbers or names of identifiers. For a nonterminal, the type is used for the values that are built for the nonterminals during parsing (at reduce-actions).

While, conceptually, parsing a string produces a syntax tree for that string, parser generators usually allow more control over what is actually produced. This is done by assigning an *action* to each production. The action is a piece of program text that is used to calculate the value of a production that is being reduced by using the values associated with the symbols on the right-hand side. For example, by putting appropriate actions on each production, the numerical value of an expression may be calculated as the result of parsing the expression. Indeed, compilers can be made such that the value produced during parsing is the compiled code of a program. For all but the simplest compilers it is, however, better to build some kind of syntax representation during parsing and then later operate on this representation.

3.17.2 Abstract syntax

The syntax trees described in section 3.3.1 are not always optimally suitable for compilation. They contain a lot of redundant information: Parentheses, keywords used for grouping purposes only, and so on. They also reflect structures in the grammar that are only introduced to eliminate ambiguity or to get the grammar accepted by a parser generator (such as left-factorisation or elimination of left-recursion). Hence, *abstract syntax* is commonly used.

Abstract syntax keeps the essence of the structure of the text but omits the irrelevant details. An *abstract syntax tree* is a tree structure where each node corresponds to one or more nodes in the (concrete) syntax tree. For example, the concrete syntax tree shown in figure 3.12 may be represented by the following abstract syntax tree:



Here the names *PlusExp*, *MulExp* and *NumExp* may be constructors in a data-type, they may be elements from an enumerated type used as tags in a union-type or they

may be names of subclasses of an *Exp* class. The names indicate which production is chosen, so there is no need to keep the subtrees that are implied by the choice of production, such as the subtree from figure 3.12 that holds the symbol $+$. Likewise, the sequence of nodes *Exp*, *Exp2*, *Exp3*, 2 at the left of figure 3.12 are combined to a single node *NumExp*(2) that includes both the choice of productions for *Exp*, *Exp2* and *Exp3* and the value of the terminal node. In short, each node in the abstract syntax tree corresponds to one or more nodes in the concrete syntax tree.

A designer of a compiler or interpreter has much freedom in the choice of abstract syntax. Some use abstract syntax that retain all of the structure of the concrete syntax trees plus additional positioning information used for error-reporting. Others prefer abstract syntax that contains only the information necessary for compilation or interpretation, skipping parentheses and other (for compilation or interpretation) irrelevant structure, like we did above.

Exactly how the abstract syntax tree is represented and built depends on the parser generator used. Normally, the action assigned to a production can access the values of the terminals and nonterminals on the right-hand side of a production through specially named variables (often called \$1, \$2, *etc.*) and produces the value for the node corresponding to the left-hand-side either by assigning it to a special variable (\$0) or letting it be the value of an action expression.

The data structures used for building abstract syntax trees depend on the language. Most statically typed functional languages support tree-structured datatypes with named constructors. In such languages, it is natural to represent abstract syntax by one datatype per syntactic category (*e.g.*, *Exp* above) and one constructor for each instance of the syntactic category (*e.g.*, *PlusExp*, *NumExp* and *MulExp* above). In Pascal, each syntactic category can be represented by a variant record type and each instance as a variant of that. In C, a syntactic category can be represented by a union of structs, each struct representing an instance of the syntactic category and the union covering all possible instances. In object-oriented languages such as Java, a syntactic category can be represented as an abstract class or interface where each instance in a syntactic category is a concrete class that implements the abstract class or interface.

In most cases, it is fairly simple to build abstract syntax using the actions for the productions in the grammar. It becomes complex only when the abstract syntax tree must have a structure that differs nontrivially from the concrete syntax tree.

One example of this is if left-recursion has been eliminated for the purpose of making an LL(1) parser. The preferred abstract syntax tree will in most cases be similar to the concrete syntax tree of the original left-recursive grammar rather than that of the transformed grammar. As an example, the left-recursive grammar

$$\begin{aligned} E &\rightarrow E + \mathbf{num} \\ E &\rightarrow \mathbf{num} \end{aligned}$$

gets transformed by left-recursion elimination into

$$\begin{aligned} E &\rightarrow \mathbf{num} E' \\ E' &\rightarrow + \mathbf{num} E' \\ E' &\rightarrow \end{aligned}$$

Which yields a completely different syntax tree. We can use the actions assigned to the productions in the transformed grammar to build an abstract syntax tree that reflects the structure in the original grammar.

In the transformed grammar, E' should return an abstract syntax tree with a *hole*. The intention is that this hole will eventually be filled by another abstract syntax tree:

- The second production for E' returns just a hole.
- In the first production for E' , the $+$ and **num** terminals are used to produce a tree for a plus-expression (*i.e.*, a *PlusExp* node) with a hole in place of the first subtree. This tree is used to fill the hole in the tree returned by the recursive use of E' , so the abstract syntax tree is essentially built outside-in. The result is a new tree with a hole.
- In the production for E , the hole in the tree returned by the E' nonterminal is filled by a *NumExp* node with the number that is the value of the **num** terminal.

The best way of building trees with holes depends on the type of language used to implement the actions. Let us first look at the case where a functional language is used.

The actions shown below for the original grammar will build an abstract syntax tree similar to the one shown in the beginning of this section.

$$\begin{aligned} E &\rightarrow E + \mathbf{num} \quad \{ \text{PlusExp}(\$1, \text{NumExp}(\$3)) \} \\ E &\rightarrow \mathbf{num} \quad \{ \text{NumExp}(\$1) \} \end{aligned}$$

We now want to make actions for the transformed grammar that will produce the same abstract syntax trees as this will.

In functional languages, an abstract syntax tree with a hole can be represented by a function. The function takes as argument what should be put into the hole and returns a syntax tree where the hole is filled with this argument. The hole is represented by the argument variable of the function. We can write this as actions to the transformed grammar:

$$\begin{aligned} E &\rightarrow \mathbf{num} E' \quad \{ \$2(\text{NumExp}(\$1)) \} \\ E' &\rightarrow + \mathbf{num} E' \quad \{ \lambda x. \$3(\text{PlusExp}(x, \text{NumExp}(\$2))) \} \\ E' &\rightarrow \quad \{ \lambda x. x \} \end{aligned}$$

where $\lambda x.e$ is a nameless function that takes x as argument and returns the value of the expression e . The empty production returns the identity function, which works like a top-level hole. The non-empty production for E' applies the function $\$3$ returned by the E' on the right-hand side to a subtree, hence filling the hole in $\$3$ by this subtree. The subtree itself has a hole x , which is filled when applying the function returned by the right-hand side. The production for E applies the function $\$2$ returned by E' to a subtree that has no holes and, hence, returns a tree with no holes.

In SML, $\lambda x.e$ is written as `fn x => e`, in Haskell as `\x -> e` and in Scheme as `(lambda (x) e)`.

The imperative version of the actions in the original grammar is

$$\begin{aligned} E &\rightarrow E + \mathbf{num} & \{ \$0 = \text{PlusExp}(\$1, \text{NumExp}(\$3)) \} \\ E &\rightarrow \mathbf{num} & \{ \$0 = \text{NumExp}(\$1) \} \end{aligned}$$

In this setting, `NumExp` and `PlusExp` are not constructors but functions that allocate and build node and return pointers to these. Unnamed functions of the kind used in the above solution for functional languages can not be built in most imperative languages, so holes must be an explicit part of the data-type that is used to represent abstract syntax. These holes will be overwritten when the values are supplied. E' will, hence, return a record holding both an abstract syntax tree (in a field named `tree`) and a pointer to the hole that should be overwritten (in a field named `hole`). As actions (using C-style notation), this becomes

$$\begin{aligned} E &\rightarrow \mathbf{num} E' & \{ \$2 \rightarrow \text{hole} = \text{NumExp}(\$1); \\ & & \$0 = \$2.\text{tree} \} \\ E' &\rightarrow + \mathbf{num} E' & \{ \$0.\text{hole} = \text{makeHole}(); \\ & & \$3 \rightarrow \text{hole} = \text{PlusExp}(\$0.\text{hole}, \text{NumExp}(\$2)); \\ & & \$0.\text{tree} = \$3.\text{tree} \} \\ E' &\rightarrow & \{ \$0.\text{hole} = \text{makeHole}(); \\ & & \$0.\text{tree} = \$0.\text{hole} \} \end{aligned}$$

This may look bad, but left-recursion removal is rarely needed when using LR-parser generators.

An alternative approach is to let the parser build an intermediate (semi-abstract) syntax tree from the transformed grammar, and then let a separate pass restructure the intermediate syntax tree to produce the intended abstract syntax. Some LL(1) parser generators can remove left-recursion automatically and will afterwards restructure the syntax tree so it fits the original grammar.

3.17.3 Conflict handling in parser generators

For all but the simplest grammars, the user of a parser generator should expect conflicts to be reported when the grammar is first presented to the parser generator.

NFA-state	Textual representation
A	T' -> . T
B	T' -> T .
C	T -> . R
D	T -> R .
E	T -> . aTc
F	T -> a . Tc
G	T -> aT . c
H	T -> aTc .
I	R -> .
J	R -> . bR
K	R -> b . R
L	R -> bR .

Figure 3.30: Textual representation of NFA states

These conflicts can be caused by ambiguity or by the limitations of the parsing method. In any case, the conflicts can normally be eliminated by rewriting the grammar or by adding precedence declarations.

Most parser generators can provide information that is useful to locate where in the grammar the problems are. When a parser generator reports conflicts, it will tell in which state in the table these occur. This state can be written out in a (barely) human-readable form as a set of NFA-states. Since most parser generators rely on pure ASCII, they can not actually draw the NFAs as diagrams. Instead, they rely on the fact that each state in the NFA corresponds to a position in a production in the grammar. If we, for example, look at the NFA states in figure 3.26, these would be written as shown in figure 3.30. Note that a ‘.’ is used to indicate the position of the state in the production. State 4 of the table in figure 3.28 will hence be written as

```
R -> b . R
R -> .
R -> . bR
```

The set of NFA states, combined with information about on which symbols a conflict occurs, can be used to find a remedy, *e.g.* by adding precedence declarations.

If all efforts to get a grammar through a parser generator fails, a practical solution may be to change the grammar so it accepts a larger language than the intended language and then post-process the syntax tree to reject “false positives”. This elimination can be done at the same time as type-checking (which, too, may reject programs).

Some languages allow programs to declare precedence and associativity for user-defined operators. This can make it difficult to handle precedence during parsing, as the precedences are not known when the parser is generated. A typical solution is to parse all operators using the same precedence and then restructure the syntax tree afterwards. See exercise 3.20 for other approaches.

3.18 Properties of context-free languages

In section 2.10, we described some properties of regular languages. Context-free languages share some, but not all, of these.

For regular languages, deterministic (finite) automata cover exactly the same class of languages as nondeterministic automata. This is not the case for context-free languages: Nondeterministic stack automata do indeed cover all context-free languages, but deterministic stack automata cover only a strict subset. The subset of context-free languages that can be recognised by deterministic stack automata are called deterministic context-free languages. Deterministic context-free languages can be recognised by LR parsers.

We have noted that the basic limitation of regular languages is finiteness: A finite automaton can not count unboundedly and hence can not keep track of matching parentheses or similar properties. Context-free languages are capable of such counting, essentially using the stack for this purpose. Even so, there are limitations: A context-free language can only keep count of one thing at a time, so while it is possible (even trivial) to describe the language $\{a^n b^n \mid n \geq 0\}$ by a context-free grammar, the language $\{a^n b^n c^n \mid n \geq 0\}$ is not a context-free language. The information kept on the stack follows a strict LIFO order, which further restricts the languages that can be described. It is, for example, trivial to represent the language of palindromes (strings that read the same forwards and backwards) by a context-free grammar, but the language of strings that can be constructed by repeating a string twice is not context-free.

Context-free languages are, as regular languages, closed under union: It is easy to construct a grammar for the union of two languages given grammars for each of these. Context-free languages are also closed under prefix, suffix, subsequence and reversal. Indeed, the language consisting of all subsequences of a context-free language is actually regular. However, context-free languages are *not* closed under intersection or complement. For example, the languages $\{a^n b^n c^m \mid m, n \geq 0\}$ and $\{a^m b^n c^n \mid m, n \geq 0\}$ are both context-free while their intersection $\{a^n b^n c^n \mid n \geq 0\}$ is not, and the complement of the language described by the grammar in section 3.11 is not a context-free language.

3.19 Further reading

Context-free grammars were first proposed as a notation for describing natural languages (*e.g.*, English or French) by the linguist Noam Chomsky [14], who defined this as one of three grammar notations for this purpose. The qualifier “context-free” distinguishes this notation from the other two grammar notations, which were called “context-sensitive” and “unconstrained”. In context-free grammars, derivation of a nonterminal is independent of the context in which the terminal occurs, whereas the context can restrict the set of derivations in a context-sensitive grammar. Unrestricted grammars can use the full power of a universal computer, so these represent all computable languages.

Context-free grammars are actually too weak to describe natural languages, but were adopted for defining the Algol60 programming language [16]. Since then, variants of this notation has been used for defining or describing almost all programming languages.

Some languages have been designed with specific parsing methods in mind: Pascal [20] has been designed for LL(1) parsing while C [25] was originally designed to fit LALR(1) parsing, but this property was lost in subsequent versions of the language.

Most parser generators are based on LALR(1) parsing, but a few use LL(1) parsing. An example of this is ANTLR (<http://www.antlr.org/>).

“The Dragon Book” [4] tells more about parsing methods than the present book.

Several textbooks, *e.g.*, [19] describe properties of context-free languages.

The methods presented here for rewriting grammars based on operator precedence uses only infix operators. If prefix or postfix operators have higher precedence than all infix operators, the method presented here will work (with trivial modifications), but if there are infix operators that have higher precedence than some prefix or postfix operators, it breaks down. A method for handling arbitrary precedences of infix, prefix and postfix operators is presented in [1].

Exercises

Exercise 3.1

Figures 3.7 and 3.8 show two different syntax trees for the string aabbbcc using grammar 3.4. Draw a third, different syntax tree for aabbbcc using the same grammar and show the left-derivation that corresponds to this syntax tree.

Exercise 3.2

Draw the syntax tree for the string aabbbcc using grammar 3.9.

Exercise 3.3

Write an unambiguous grammar for the language of balanced parentheses, *i.e.* the language that contains (among other) the sequences

ε (i.e. the empty string)
 ()
 (())
 (())
 ((()))

but none of the following

(
)
)(
 ((
))

Exercise 3.4

Write grammars for each of the following languages:

- All sequences of as and bs that contain the same number of as and bs (in any order).
- All sequences of as and bs that contain strictly more as than bs.
- All sequences of as and bs that contain a different number of as and bs.
- All sequences of as and bs that contain twice as many as as bs.

Exercise 3.5

We extend the language of balanced parentheses from exercise 3.3 with two symbols: [and]. [corresponds to exactly two normal opening parentheses and] corresponds to exactly two normal closing parentheses. A string of mixed parentheses is legal if and only if the string produced by replacing [by ((and] by)) is a balanced parentheses sequence. Examples of legal strings are

ϵ
 $()()$
 $([]$
 $[]$
 $[]()$
 $[()]$

- Write a grammar that recognises this language.
- Draw the syntax trees for $[]()$ and $[()]$.

Exercise 3.6

Show that the grammar

$$\begin{aligned}
 A &\rightarrow -A \\
 A &\rightarrow A - \mathbf{id} \\
 A &\rightarrow \mathbf{id}
 \end{aligned}$$

is ambiguous by finding a string that has two different syntax trees.

Now make two different unambiguous grammars for the same language:

- One where prefix minus binds stronger than infix minus.
- One where infix minus binds stronger than prefix minus.

Show the syntax trees using the new grammars for the string you used to prove the original grammar ambiguous.

Exercise 3.7

In grammar 3.2, replace the operators $-$ and $/$ by $<$ and $:$. These have the following precedence rules:

$<$ is non-associative and binds less tightly than $+$ but more tightly than $:$.

$:$ is right-associative and binds less tightly than any other operator.

Write an unambiguous grammar for this modified grammar using the method shown in section 3.4.1. Show the syntax tree for $2 : 3 < 4 + 5 : 6 * 7$ using the unambiguous grammar.

Exercise 3.8

Extend grammar 3.13 with the productions

$$\begin{array}{ll} \textit{Exp} & \rightarrow \textbf{id} \\ \textit{Matched} & \rightarrow \end{array}$$

then calculate *Nullable* and *FIRST* for every production in the grammar.

Add an extra start production as described in section 3.10 and calculate *FOLLOW* for every nonterminal in the grammar.

Exercise 3.9

Calculate *Nullable*, *FIRST* and *FOLLOW* for the nonterminals *A* and *B* in the grammar

$$\begin{array}{ll} A & \rightarrow BAa \\ A & \rightarrow \\ B & \rightarrow bBc \\ B & \rightarrow AA \end{array}$$

Remember to extend the grammar with an extra start production when calculating *FOLLOW*.

Exercise 3.10

Eliminate left-recursion from grammar 3.2.

Exercise 3.11

Calculate *Nullable* and *FIRST* for every production in grammar 3.20.

Exercise 3.12

Add a new start production $\textit{Exp}' \rightarrow \textit{Exp}\$$ to the grammar produced in exercise 3.10 and calculate *FOLLOW* for all nonterminals in the resulting grammar.

Exercise 3.13

Make a LL(1) parser-table for the grammar produced in exercise 3.12.

Exercise 3.14

Consider the following grammar for postfix expressions:

$$\begin{aligned} E &\rightarrow EE+ \\ E &\rightarrow EE* \\ E &\rightarrow \text{num} \end{aligned}$$

- a) Eliminate left-recursion in the grammar.
- b) Do left-factorisation of the grammar produced in question a.
- c) Calculate *Nullable*, *FIRST* for every production and *FOLLOW* for every non-terminal in the grammar produced in question b.
- d) Make a LL(1) parse-table for the grammar produced in question b.

Exercise 3.15

Extend grammar 3.11 with a new start production as shown in section 3.15 and calculate *FOLLOW* for every nonterminal. Remember to add an extra start production for the purpose of calculating *FOLLOW* as described in section 3.10.

Exercise 3.16

Make NFAs (as in figure 3.26) for the productions in grammar 3.11 (after extending it as shown in section 3.15) and show the epsilon-transitions as in figure 3.27. Convert the combined NFA into an SLR DFA like the one in figure 3.28. Finally, add reduce and accept actions based on the *FOLLOW* sets calculated in exercise 3.15.

Exercise 3.17

Extend grammar 3.2 with a new start production as shown in section 3.15 and calculate *FOLLOW* for every nonterminal. Remember to add an extra start production for the purpose of calculating *FOLLOW* as described in section 3.10.

Exercise 3.18

Make NFAs (as in figure 3.26) for the productions in grammar 3.2 (after extending it as shown in section 3.15) and show the epsilon-transitions as in figure 3.27. Convert the combined NFA into an SLR DFA like the one in figure 3.28. Add reduce actions based on the *FOLLOW* sets calculated in exercise 3.17. Eliminate the conflicts in the table by using operator precedence rules as described in section 3.16. Compare the size of the table to that from exercise 3.16.

Exercise 3.19

Consider the grammar

$$\begin{aligned} T &\rightarrow T \rightarrow T \\ T &\rightarrow T * T \\ T &\rightarrow \mathbf{int} \end{aligned}$$

where \rightarrow is considered a single terminal symbol.

- a) Add a new start production as shown in section 3.15.
- b) Calculate $FOLLOW(T)$. Remember to add an extra start production.
- c) Construct an SLR parser-table for the grammar.
- d) Eliminate conflicts using the following precedence rules:
 - $*$ binds tighter than \rightarrow .
 - $*$ is left-associative.
 - \rightarrow is right-associative.

Exercise 3.20

In section 3.17.3 it is mentioned that user-defined operator precedences in programming languages can be handled by parsing all operators with a single fixed precedence and associativity and then using a separate pass to restructure the syntax tree to reflect the declared precedences. Below are two other methods that have been used for this purpose:

- a) An ambiguous grammar is used and conflicts exist in the SLR table. Whenever a conflict arises during parsing, the parser consults a table of precedences to resolve this conflict. The precedence table is extended whenever a precedence declaration is read.
- b) A terminal symbol is made for every possible precedence and associativity combination. A conflict-free parse table is made either by writing an unambiguous grammar or by eliminating conflicts in the usual way. The lexical analyser uses a table of precedences to assign the correct terminal symbol to each operator it reads.

Compare all three methods. What are the advantages and disadvantages of each method?.

Exercise 3.21

Consider the grammar

$$\begin{aligned} A &\rightarrow a A a \\ A &\rightarrow b A b \\ A &\rightarrow \end{aligned}$$

- Describe the language that the grammar defines.
- Is the grammar ambiguous? Justify your answer.
- Construct a SLR parse table for the grammar.
- Can the conflicts in the table be eliminated?

Exercise 3.22

The following ambiguous grammar describes boolean expressions:

$$\begin{aligned} B &\rightarrow \mathbf{true} \\ B &\rightarrow \mathbf{false} \\ B &\rightarrow B \vee B \\ B &\rightarrow B \wedge B \\ B &\rightarrow \neg B \\ B &\rightarrow (B) \end{aligned}$$

- Given that negation (\neg) binds tighter than conjunction (\wedge) which binds tighter than disjunction (\vee) and that conjunction and disjunction are both right-associative, rewrite the grammar to be unambiguous.
- Write a grammar that accepts only true boolean expressions. Hint: Use the answer from question a) and add an additional nonterminal F for false boolean expressions.

Chapter 4

Scopes and Symbol Tables

4.1 Introduction

An important concept in programming languages is the ability to *name* objects such as variables, functions and types. Each such named object will have a *declaration*, where the name is defined as a synonym for the object. This is called *binding*. Each name will also have a number of *uses*, where the name is used as a reference to the object to which it is bound.

Often, the declaration of a name has a limited *scope*: a portion of the program where the name will be visible. Such declarations are called *local declarations*, whereas a declaration that makes the declared name visible in the entire program is called *global*. It may happen that the same name is declared in several nested scopes. In this case, it is normal that the declaration closest to a use of the name will be the one that defines that particular use. In this context *closest* is related to the syntax tree of the program: The scope of a declaration will be a sub-tree of the syntax tree and nested declarations will give rise to scopes that are nested sub-trees. The closest declaration of a name is hence the declaration corresponding to the smallest sub-tree that encloses the use of the name. As an example, look at this C statement block:

```
{
    int x = 1;
    int y = 2;
    {
        double x = 3.14159265358979;
        y += (int)x;
    }
    y += x;
}
```

The two lines immediately after the first opening brace declare integer variables x and y with scope until the closing brace in the last line. A new scope is started by the second opening brace and a floating-point variable x with an initial value close to π is declared. This will have scope until the first closing brace, so the original x variable is not visible until the inner scope ends. The assignment $y += (\text{int})x;$ will add 3 to y , so its new value is 5. In the next assignment $y += x;$, we have exited the inner scope, so the original x is restored. The assignment will, hence, add 1 to y , which will have the final value 6.

Scoping based on the structure of the syntax tree, as shown in the example, is called *static* or *lexical* binding and is the most common scoping rule in modern programming languages. We will in the rest of this chapter (indeed, the rest of this book) assume that static binding is used. A few languages have *dynamic* binding, where the declaration that was most recently encountered during execution of the program defines the current use of the name. By its nature, dynamic binding can not be resolved at compile-time, so the techniques that in the rest of this chapter are described as being used in a compiler will have to be used at run-time if the language uses dynamic binding.

A compiler will need to keep track of names and the objects these are bound to, so that any use of a name will be attributed correctly to its declaration. This is typically done using a *symbol table* (or *environment*, as it is sometimes called).

4.2 Symbol tables

A symbol table is a table that binds names to information. We need a number of operations on symbol tables to accomplish this:

- We need an *empty* symbol table, in which no name is defined.
- We need to be able to *bind* a name to a piece of information. In case the name is already defined in the symbol table, the new binding takes precedence over the old.
- We need to be able to *look up* a name in a symbol table to find the information the name is bound to. If the name is not defined in the symbol table, we need to be told that.
- We need to be able to *enter* a new scope.
- We need to be able to *exit* a scope, reestablishing the symbol table to what it was before the scope was entered.

4.2.1 Implementation of symbol tables

There are many ways to implement symbol tables, but the most important distinction between these is how scopes are handled. This may be done using a *persistent* (or *functional*) data structure, or it may be done using an *imperative* (or destructively-updated) data structure.

A persistent data structure has the property that no operation on the structure will destroy it. Conceptually, a new modified copy is made of the data structure whenever an operation updates it, hence preserving the old structure unchanged. This means that it is trivial to reestablish the old symbol table when exiting a scope, as it has been preserved by the persistent nature of the data structure. In practice, only a small portion of the data structure is copied when a symbol table is updated, most is shared with the previous version.

In the imperative approach, only one copy of the symbol table exists, so explicit actions are required to store the information needed to restore the symbol table to a previous state. This can be done by using an auxiliary stack. When an update is made, the old binding of a name that is overwritten is recorded (pushed) on the auxiliary stack. When a new scope is entered, a marker is pushed on the auxiliary stack. When the scope is exited, the bindings on the auxiliary stack (down to the marker) are used to reestablish the old symbol table. The bindings and the marker are popped off the auxiliary stack in the process, returning the auxiliary stack to the state it was in before the scope was entered.

Below, we will look at simple implementations of both approaches and discuss how more advanced approaches can overcome some of the efficiency problems with the simple approaches.

4.2.2 Simple persistent symbol tables

In functional languages like SML, Scheme or Haskell, persistent data structures are the norm rather than the exception (which is why persistent data structures are sometimes called *functional* data structures). For example, when a new element is added to the front of a list or an element is taken off the front of the list, the old list still exists and can be used elsewhere. A list is a natural way to implement a symbol table in a functional language: A binding is a pair of a name and its associated information, and a symbol table is a list of such pairs. The operations are implemented in the following way:

empty: An empty symbol table is an empty list.

binding: A new binding (name/information pair) is added (consed) to the front of the list.

lookup: The list is searched until a pair with a matching name is found. The information paired with the name is then returned. If the end of the list is reached, an indication that this happened is returned instead. This indication can be made by raising an exception or by letting the lookup function return a special value representing “not found”. This requires a type that can hold both normal information and and this special value, *i.e.*, a sum-type.

enter: The old list is remembered, *i.e.*, a reference is made to it.

exit: The old list is recalled, *i.e.*, the above reference is used.

The latter two operations are not really explicit operations, as the variable used to hold the symbol table before entering a new scope will still hold the same symbol table after the scope is exited. So all that is needed is a variable to hold (a reference to) the symbol table.

As new bindings are added to the front of the list and the list is searched from the front to the back, bindings in inner scopes will automatically take precedence over bindings in outer scopes.

Another functional approach to symbol tables is using functions: A symbol table is quite naturally seen as a function from names to information. The operations are:

empty: An empty symbol table is a function that returns an error indication (or raises an exception) no matter what its argument is.

binding: Adding a binding of the name n to the information i in a symbol table t is done by defining a new symbol-table function t' in terms t and the new binding. When t' is called with a name $n1$ as argument, it compares $n1$ to n . If they are equal, t' returns the information i . Otherwise, t' calls t with $n1$ as argument and returns the result that this call yields. In Standard ML, we can define a binding function this way:

```
bind (n,i,t) = fn n1 => if n1=n then i else t n1
```

lookup: The symbol-table function is called with the name as argument.

enter: The old function is remembered (referenced).

exit: The old function is recalled (by using a reference).

Again, the latter two operations are mostly implicit.

4.2.3 A simple imperative symbol table

Imperative symbol tables are natural to use if the compiler is written in an imperative language. A simple imperative symbol table can be implemented as a stack, which works in a way similar to the list-based functional implementation:

empty: An empty symbol table is an empty stack.

binding: A new binding (name/information pair) is pushed on top of the stack.

lookup: The stack is searched top-to-bottom until a matching name is found. The information paired with the name is then returned. If the bottom of the stack is reached, we instead return an error-indication.

enter: We push a marker on the top of the stack.

exit: We pop bindings from the stack until a marker is found. This is also popped from the stack.

Note that since the symbol table is itself a stack, we don't need the auxiliary stack mentioned in section 4.2.1.

This is not quite a persistent data structure, as leaving a scope will destroy its symbol table. For simple languages, this won't matter, as a scope isn't needed again after it is exited. But language features such as classes, modules and lexical closures can require symbol tables to persist after their scope is exited. In these cases, a real persistent symbol table must be used, or the needed parts of the symbol table must be copied and stored for later retrieval before exiting a scope.

4.2.4 Efficiency issues

While all of the above implementations are simple, they all share the same efficiency problem: Lookup is done by linear search, so the worst-case time for lookup is proportional to the size of the symbol table. This is mostly a problem in relation to libraries: It is quite common for a program to use libraries that define literally hundreds of names.

A common solution to this problem is *hashing*: Names are hashed (processed) into integers, which are used to index an array. Each array element is then a linear list of the bindings of names that share the same hash code. Given a large enough hash table, these lists will typically be very short, so lookup time is basically constant.

Using hash tables complicates entering and exiting scopes somewhat. While each element of the hash table is a list that can be handled like in the simple cases, doing this for *all* the array-elements at every entry and exit imposes a major overhead. Instead, it is typical for imperative implementations to use a single auxiliary

stack (as described in section 4.2.1) to record all updates to the table so they can be undone in time proportional to the number of updates that were done in the local scope. Functional implementations typically use persistent hash-tables, which eliminates the problem.

4.2.5 Shared or separate name spaces

In some languages (like Pascal) a variable and a function in the same scope may have the same name, as the context of use will make it clear whether a variable or a function is used. We say that functions and variables have *separate name spaces*, which means that defining a name in one space doesn't affect the same name in the other space. In other languages (*e.g.* C or SML) the context can not (easily) distinguish variables from functions. Hence, declaring a local variable might hide a function declared in an outer scope or vice versa. These languages have a *shared name space* for variables and functions.

Name spaces may be shared or separate for all the kinds of names that can appear in a program, *e.g.*, variables, functions, types, exceptions, constructors, classes, field selectors *etc.* Which name spaces are shared is language-dependent.

Separate name spaces are easily implemented using one symbol table per name space, whereas shared name spaces naturally share a single symbol table. However, it is sometimes convenient to use a single symbol table even if there are separate name spaces. This can be done fairly easily by adding name-space indicators to the names. A name-space indicator can be a textual prefix to the name or it may be a tag that is paired with the name. In either case, a lookup in the symbol table must match both the name and the name-space indicator of the symbol that is looked up with the name and the name-space indicator of the entry in the table.

Suggested exercises: 4.1.

4.3 Further reading

Most algorithms-and-data-structures textbooks include descriptions of methods for hashing strings and implementing hash tables. A description of efficient persistent data structures for functional languages can be found in [37].

Exercises

Exercise 4.1

Pick some programming language that you know well and determine which of the following objects share name spaces: Variables, functions, procedures and types.

If there are more kinds of named objects (labels, data constructors, modules, *etc.*) in the language, include these in the investigation.

Exercise 4.2

Implement, in a programming language of your choice, data structures and operations (empty, binding, lookup, enter and exit) for simple symbol tables.

Exercise 4.3

In some programming languages, identifiers are case-insensitive so, *e.g.*, `size` and `SiZe` refer to the same identifier. Describe how symbol tables can be made case-insensitive.

Chapter 5

Interpretation

5.1 Introduction

After lexing and parsing, we have the abstract syntax tree of a program as a data structure in memory. But a program needs to be executed, and we have not yet dealt with that issue.

The simplest way to execute a program is *interpretation*. Interpretation is done by a program called an *interpreter*, which takes the abstract syntax tree of a program and executes it by inspecting the syntax tree to see what needs to be done. This is similar to how a human evaluates a mathematical expression: We insert the values of variables in the expression and evaluate it bit by bit, starting with the innermost parentheses and moving out until we have the result of the expression. We can then repeat the process with other values for the variables.

There are some differences, however. Where a human being will copy the text of the formula with variables replaced by values and then write a sequence of more and more reduced copies of a formula until it is reduced to a single value, an interpreter will keep the formula (or, rather, the abstract syntax tree of an expression) unchanged and use a symbol table to keep track of the values of variables. Instead of reducing a formula, the interpreter is a function that takes an abstract syntax tree and a symbol table as arguments and returns the value of the expression represented by the abstract syntax tree. The function can call itself recursively on parts of the abstract syntax tree to find the values of subexpressions, and when it evaluates a variable, it can look its value up in the symbol table.

This process can be extended to also handle statements and declarations, but the basic idea is the same: A function takes the abstract syntax tree of the program and, possibly, some extra information about the context (such as a symbol table or the input to the program) and returns the output of the program. Some input and output may be done as side effects by the interpreter.

We will in this chapter assume that the symbol tables are persistent, so no ex-

explicit action is required to restore the symbol table for the outer scope when exiting an inner scope. In the main text of the chapter, we don't need to preserve symbol tables for inner scopes once these are exited (so a stack-like behaviour is fine), but in one of the exercises we will need this.

5.2 The structure of an interpreter

An interpreter will typically consist of one function per syntactic category. Each function will take as arguments an abstract syntax tree from the syntactic category and, possibly, extra arguments such as symbol tables. Each function will return one or more results, which may be the value of an expression or an updated symbol table.

The functions can be implemented in any language that we already have an implementation of. This implementation can also be an interpreter, or it can be a compiler that compiles to some other language. Eventually, we will need to either have an interpreter written in machine language or a compiler that compiles to machine language. Chapter 13 will discuss how this can be done. For the moment, we just write the interpreter functions in a notation reminiscent of a high-level programming language and assume an implementation exists. Additionally, we want to avoid being specific about how abstract syntax is represented, so we will use a notation that looks like concrete syntax to represent abstract syntax.

5.3 A small example language

We will use a small (somewhat contrived) language to show the principles of interpretation. The language is a first-order functional language with recursive definitions. The syntax is given in grammar 5.1. The shown grammar is clearly ambiguous, but we assume that any ambiguities have been resolved during parsing, so we have an unambiguous abstract syntax tree.

In the example language, a program is a list of function declarations. The functions are all mutually recursive, and no function may be declared more than once. Each function declares its result type and the types and names of its arguments. There may not be repetitions in the list of parameters for a function. Functions and variables have separate name spaces. The body of a function is an expression, which may be an integer constant, a variable name, a sum-expression, a comparison, a conditional, a function call or an expression with a local declaration. Comparison is defined both on booleans and integers, but addition only on integers.

A program must contain a function called `main`, which has one integer argument and which returns an integer. Execution of the program is by calling this function with the input (which must be an integer). The result of this function call

$$\begin{aligned}
\textit{Program} &\rightarrow \textit{Funs} \\
\\
\textit{Funs} &\rightarrow \textit{Fun} \\
\textit{Funs} &\rightarrow \textit{Fun Funs} \\
\\
\textit{Fun} &\rightarrow \textit{TypeId} (\textit{TypeIds}) = \textit{Exp} \\
\\
\textit{TypeId} &\rightarrow \text{int } \mathbf{id} \\
\textit{TypeId} &\rightarrow \text{bool } \mathbf{id} \\
\\
\textit{TypeIds} &\rightarrow \textit{TypeId} \\
\textit{TypeIds} &\rightarrow \textit{TypeId} , \textit{TypeIds} \\
\\
\textit{Exp} &\rightarrow \mathbf{num} \\
\textit{Exp} &\rightarrow \mathbf{id} \\
\textit{Exp} &\rightarrow \textit{Exp} + \textit{Exp} \\
\textit{Exp} &\rightarrow \textit{Exp} = \textit{Exp} \\
\textit{Exp} &\rightarrow \text{if } \textit{Exp} \text{ then } \textit{Exp} \text{ else } \textit{Exp} \\
\textit{Exp} &\rightarrow \mathbf{id} (\textit{Exps}) \\
\textit{Exp} &\rightarrow \text{let } \mathbf{id} = \textit{Exp} \text{ in } \textit{Exp} \\
\\
\textit{Exps} &\rightarrow \textit{Exp} \\
\textit{Exps} &\rightarrow \textit{Exp} , \textit{Exps}
\end{aligned}$$

Grammar 5.1: Example language for interpretation

is the output of the program.

5.4 An interpreter for the example language

An interpreter for this language must take the abstract syntax tree of the program and an integer (the input to the program) and return another integer (the output of the program). Since values can be both integers or booleans, the interpreter uses a value type that contains both integers and booleans (and enough information to tell them apart). We will not go into detail about how such a type can be defined but simply assume that there are operations for testing if a value is a boolean or an integer and operating on values known to be integers or booleans. If we during interpretation find that we have to, say, add a boolean to an integer, we stop the interpretation with an error message. We do this by letting the interpreter call a function called **error()**.

We will start by showing how we can evaluate (interpret) expressions, and then extend this to handle the whole program.

5.4.1 Evaluating expressions

When we evaluate expressions, we need, in addition to the abstract syntax tree of the expression, also a symbol table *vtable* that binds variables to their values. Additionally, we need to be able to handle function calls, so we also need a symbol table *ftable* that binds function names to the abstract syntax trees of their declarations. The result of evaluating an expression is the value of the expression.

For terminals (variable names and numeric constants) with attributes, we assume that there are predefined functions for extracting these. Hence, **id** has an associated function *getname*, that extracts the name of the identifier. Similarly, **num** has a function *getvalue*, that returns the value of the number.

Figure 5.2 shows a function $Eval_{Exp}$, that takes an expression *Exp* and symbol tables *vtable* and *ftable* and returns a value, which may be either an integer or a boolean. Also shown is a function $Eval_{Exps}$, that evaluates a list of expressions to a list of values. We also use a function $Call_{Fun}$ that handles function calls. We will return to this later.

The main part of $Eval_{Exp}$ is a case-expression that identifies which kind of expression the top node of the abstract syntax tree represents. The patterns are shown as concrete syntax, but you should think of it as pattern matching on the abstract syntax. The box to the right of a pattern shows the actions needed to evaluate the expression. These actions can refer to parts of the pattern on the left. An action is a sequence of definitions of local variables followed by an expression (in the interpreter) that evaluates to the result of the expression that was given (in abstract syntax) as argument to $Eval_{Exp}$.

$Eval_{Exp}(Exp, vtable, ftable) = \text{case } Exp \text{ of}$	
num	$getvalue(\mathbf{num})$
id	$v = lookup(vtable, getname(\mathbf{id}))$ <i>if</i> $v = \text{unbound}$ <i>then</i> error() <i>else</i> v
$Exp_1 + Exp_2$	$v_1 = Eval_{Exp}(Exp_1, vtable, ftable)$ $v_2 = Eval_{Exp}(Exp_2, vtable, ftable)$ <i>if</i> v_1 and v_2 are integers <i>then</i> $v_1 + v_2$ <i>else</i> error()
$Exp_1 = Exp_2$	$v_1 = Eval_{Exp}(Exp_1, vtable, ftable)$ $v_2 = Eval_{Exp}(Exp_2, vtable, ftable)$ <i>if</i> v_1 and v_2 are both integers or both booleans <i>then</i> <i>if</i> $v_1 = v_2$ <i>then</i> true <i>else</i> false <i>else</i> error()
<i>if</i> Exp_1 <i>then</i> Exp_2 <i>else</i> Exp_3	$v_1 = Eval_{Exp}(Exp_1, vtable, ftable)$ <i>if</i> v_1 is a boolean <i>then</i> <i>if</i> $v_1 = \mathbf{true}$ <i>then</i> $Eval_{Exp}(Exp_2, vtable, ftable)$ <i>else</i> $Eval_{Exp}(Exp_3, vtable, ftable)$ <i>else</i> error()
id ($Exps$)	$def = lookup(ftable, getname(\mathbf{id}))$ <i>if</i> $def = \text{unbound}$ <i>then</i> error() <i>else</i> $args = Eval_{Exps}(Exps, vtable, ftable)$ $Call_{Fun}(def, args, ftable)$
let $\mathbf{id} = Exp_1$ <i>in</i> Exp_2	$v_1 = Eval_{Exp}(Exp_1, vtable, ftable)$ $vtable' = bind(vtable, getname(\mathbf{id}), v_1)$ $Eval_{Exp}(Exp_2, vtable', ftable)$

$Eval_{Exps}(Exps, vtable, ftable) = \text{case } Exps \text{ of}$	
Exp	$[Eval_{Exp}(Exp, vtable, ftable)]$
$Exp, Exps$	$Eval_{Exp}(Exp, vtable, ftable)$ $:: Eval_{Exps}(Exps, vtable, ftable)$

Figure 5.2: Evaluating expressions

We will briefly explain each of the cases handled by $Eval_{Exp}$.

- The value of a number is found as the *value* attribute to the node in the abstract syntax tree.
- The value of a variable is found by looking its name up in the symbol table for variables. If the variable is not found in the symbol table, the lookup-function returns the special value *unbound*. When this happens, an error is reported and the interpretation stops. Otherwise, it returns the value returned by *lookup*.
- At a plus-expression, both arguments are evaluated, then it is checked that they are both integers. If they are, we return the sum of the two values. Otherwise, we report an error (and stop).
- Comparison requires that the arguments have the same type. If that is the case, we compare the values, otherwise we report an error.
- In a conditional expression, the condition must be a boolean. If it is, we check if it is **true**. If so, we evaluate the then-branch, otherwise, we evaluate the else-branch. If the condition is not a boolean, we report an error.
- At a function call, the function name is looked up in the function environment to find its definition. If the function is not found in the environment, we report an error. Otherwise, we evaluate the arguments by calling $Eval_{Exps}$ and call $Call_{Fun}$ to find the result of the call.
- A let-expression declares a new variable with an initial value defined by an expression. The expression is evaluated and the symbol table for variables is extended using the function *bind* to bind the variable to the value. The extended table is used when evaluating the body-expression, which defines the value of the whole expression.

$Eval_{Exps}$ builds a list of the values of the expressions in the expression list. The notation is taken from SML: A list is written in square brackets with commas between the elements. The operator $::$ adds an element to the front of a list.

Suggested exercises: 5.1.

5.4.2 Interpreting function calls

A function declaration explicitly declares the types of the arguments. When a function is called, it must be checked that the number of arguments is the same as the declared number, and that the values of the arguments match the declared types.

$Call_{Fun}(Fun, args, ftable) = \text{case } Fun \text{ of}$	
$TypeId (TypeIds) = Exp$	$(f, t_0) = Get_{TypeId}(TypeId)$ $vtable = Bind_{TypeIds}(TypeIds, args)$ $v_1 = Eval_{Exp}(Exp, vtable, ftable)$ if v_1 is of type t_0 then v_1 else error()

$Get_{TypeId}(TypeId) = \text{case } TypeId \text{ of}$	
int id	$(getname(id), \text{int})$
bool id	$(getname(id), \text{bool})$

$Bind_{TypeIds}(TypeIds, args) = \text{case } (TypeIds, args) \text{ of}$	
$(TypeId, [v])$	$(x, t) = Get_{TypeId}(TypeId)$ if v is of type t then $bind(emptytable, x, v)$ else error()
$(TypeId, TypeIds, (v :: vs))$	$(x, t) = Get_{TypeId}(TypeId)$ $vtable = Bind_{TypeIds}(TypeIds, vs)$ if $lookup(vtable, x) = unbound$ and v is of type t then $bind(vtable, x, v)$ else error()
—	error()

Figure 5.3: Evaluating a function call

If this is the case, we build a symbol table that binds the parameter variables to the values of the arguments and use this in evaluating the body of the function. The value of the body must match the declared result type of the function.

$Call_{Fun}$ is also given a symbol table for functions, which is passed to the $Eval_{Exp}$ when evaluating the body.

$Call_{Fun}$ is shown in figure 5.3, along with the functions for $TypeId$ and $TypeIds$, which it uses. The function Get_{TypeId} just returns a pair of the declared name and type, and $Bind_{TypeIds}$ checks the declared type against a value and builds a symbol table that binds the name to the value if they match (and report an error if they do not). $Binds_{TypeIds}$ also checks if all parameters have different names by seeing if the current name is already bound. $emptytable$ is an empty symbol table. Looking any name up in the empty symbol table returns *unbound*. The underscore used in the last rule for $Bind_{TypeIds}$ is a wildcard pattern that matches anything, so this rule is used when the number of arguments do not match the number of declared parameters.

$Run_{Program}(Program, input) = \text{case } Program \text{ of}$	
$Funs$	$f_{table} = Build_{f_{table}}(Funs)$ $def = lookup(f_{table}, main)$ $\text{if } def = unbound$ $\text{then } \mathbf{error}()$ else $Call_{Fun}(def, [input], f_{table})$
$Build_{f_{table}}(Funs) = \text{case } Funs \text{ of}$	
Fun	$f = Get_{fname}(Fun)$ $bind(emptytable, f, Fun)$
$Fun Funs$	$f = Get_{fname}(Fun)$ $f_{table} = Build_{f_{table}}(Funs)$ $\text{if } lookup(f_{table}, f) = unbound$ $\text{then } bind(f_{table}, f, Fun)$ $\text{else } \mathbf{error}()$
$Get_{fname}(Fun) = \text{case } Fun \text{ of}$	
$TypeId (TypeIds) = Exp$	$(f, t_0) = Get_{TypeId}(TypeId)$ f

Figure 5.4: Interpreting a program

5.4.3 Interpreting a program

Running a program is done by calling the main function with a single argument that is the input to the program. So we build the symbol table for functions, look up `main` in this and call $Call_{Fun}$ with the resulting definition and an argument list containing just the input.

Hence, $Run_{Program}$, which runs the whole program, calls a function $Build_{f_{table}}$ that builds the symbol table for functions. This, in turn, uses a function Get_{fname} that finds the name of a function. All these functions are shown in figure 5.4.

This completes the interpreter for our small example language.

Suggested exercises: 5.5.

5.5 Advantages and disadvantages of interpretation

Interpretation is often the simplest way of executing a program once you have its abstract syntax tree. However, it is also a relatively slow way to do so. When we perform an operation in the interpreted program, the interpreter must first inspect

the abstract syntax tree to see what operation it needs to perform, then it must check that the types of the arguments to the operation match the requirements of the operation, and only then can it perform the operation. Additionally, each value must include sufficient information to identify its type, so after doing, say, an addition, we must add type information to the resulting number.

It should be clear from this that we spend much more time on figuring out what to do and if it is O.K. to do it than on actually doing it.

To get faster execution, we use the observation that a program that only executes each part of the program once will finish quite quickly. In other words, any time-consuming program will contain parts that are executed many times. The idea is that we can do the inspection of the abstract syntax tree and the type checking once for each part and only repeat the actual operations that are performed in this part. Since performing the operations is a small fraction of the total time, we can get a substantial speedup by doing this.

This is the basic idea behind *static type checking* and *compilation*.

Static type checking checks the program for potential mismatches between the types of values and the requirements of operations. It does so for the whole program regardless of whether all parts will actually be executed, so it may report errors even if an interpretation of the program would finish without errors. So static type checking puts extra limitations on programs but reduces the time needed to check errors and, as a bonus, reports potential problems before a program is executed, which can help when debugging a program. We look at static type checking in chapter 6. It does, however, need some time to do the checking before we can start executing the program, so the time from doing an edit in a program to executing it will increase.

Compilation gets rid of the abstract syntax tree of the source program by translating it into a target program (in a language that we already have an implementation of) that only performs the operations in the source program. Usually, the target language is a low-level language such as machine code. Like static checking, compilation must complete before execution can begin, so it adds delay between editing a program and running it.

Usually, static checking and compilation go hand in hand, but there are compilers for languages with dynamic (run-time) type checking and interpreters for statically typed languages.

Some implementations combine interpretation and compilation: The first few times a function is called, it is interpreted, but if it is called sufficiently often, it is compiled and all subsequent calls to the function will execute the compiled code. This is often called *just-in-time compilation*, though this term was originally used for just postponing compilation of a function until the first time it is called, hence reducing the delay from editing a program to its execution but at the cost of adding smaller delays during execution.

5.6 Further reading

A step-by-step construction of an interpreter for a LISP-like language is shown in [41]. A survey of programming language constructs (also for LISP-like languages) and their interpretation is shown in [2].

Exercises

Exercise 5.1

We extend the language from section 5.3 with boolean operators. We add the following productions to grammar 5.1:

$$\begin{aligned} Exp &\rightarrow \text{not } Exp \\ Exp &\rightarrow Exp \text{ and } Exp \end{aligned}$$

When evaluating `not e`, we first evaluate `e` to a value `v` that is checked to be a boolean. If it is, we return $\neg v$, where \neg is logical negation.

When evaluating `e1 and e2`, we first evaluate `e1` and `e2` to values `v1` and `v2` that are both checked to be booleans. If they are, we return `v1 ∧ v2`, where \wedge is logical conjunction.

Extend the interpretation function in figure 5.2 to handle these new constructions as described above.

Exercise 5.2

Add the productions

$$\begin{aligned} Exp &\rightarrow \text{floatconst} \\ TypeId &\rightarrow \text{float id} \end{aligned}$$

to grammar 5.1. This introduces floating-point numbers to the language. The operator `+` is overloaded so it can do integer addition or floating-point addition, and `=` is extended so it can also compare floating point numbers for equality.

- a) Extend the interpretation functions in figures 5.2-5.4 to handle these extensions.
- b) We now add implicit conversion of integers to floats to the language, using the rules: Whenever an operator has one integer argument and one floating-point argument, the integer is converted to a float. Extend the interpretation functions from question a) above to handle this.

Exercise 5.3

The language defined in section 5.3 declares types of parameters and results of functions. The interpreter in section 5.4 adds explicit type information to every value, and checks this before doing any operations on values. So, we could omit type declarations and rely solely on the type information in values.

Replace in grammar 5.1 *TypeId* by **id** and rewrite the interpretation functions in figure 5.3 so they omit checking types of parameters and results, but still check that the number of arguments match the declaration and that no parameter name is repeated.

Exercise 5.4

In the language defined in section 5.3, variables bound in **let**-expressions have no declared type, so it is possible to write a program where the same **let**-bound variable sometimes is bound to an integer and at other times to a boolean value.

Write an example of such a program.

Exercise 5.5

We extend the language from section 5.3 with functional values. These require lexical closures, so we assume symbol tables are fully persistent. We add the following productions to grammar 5.1:

$$\begin{aligned}
 TypeId &\rightarrow \text{fun } \mathbf{id} \\
 Exp &\rightarrow Exp\ Exp \\
 Exp &\rightarrow \text{fn } \mathbf{id} \Rightarrow Exp
 \end{aligned}$$

The notation is taken from Standard ML.

Evaluating **fn** *x* \Rightarrow *e* in an environment *vtable* produces a functional value *f*. When *f* is applied to an argument *v*, it is checked that *v* is an integer. If this is the case, *e* is evaluated in *vtable* extended with a binding that binds *x* to *v*. We then check if the result *w* of this evaluation is an integer, and if so use it as the result of the function application.

When evaluating *e*₁ *e*₂, we evaluate *e*₁ to a functional value *f* and *e*₂ to an integer *v* and then apply *f* to *v* as described above.

Extend the interpreter from figure 5.3 to handle these new constructions as described above. Represent a lexical closures as a pair of (the abstract syntax of) an expression and an environment.

Chapter 6

Type Checking

6.1 Introduction

Lexing and parsing will reject many texts as not being correct programs. However, many languages have well-formedness requirements that can not be handled exclusively by the techniques seen so far. These requirements can, for example, be static type correctness or a requirement that pattern-matching or case-statements are exhaustive.

These properties are most often not context-free, *i.e.*, they can not be checked by membership of a context-free language. Consequently, they are checked by a phase that (conceptually) comes after syntax analysis (though it may be interleaved with it). These checks may happen in a phase that does nothing else, or they may be combined with the actual execution or translation to another language. Often, the translator may exploit or depend on type information, which makes it natural to combine calculation of types with the actual translation. In the chapter 5, we covered type-checking during execution, which is normally called *dynamic typing*. We will in this chapter assume that type checking and related checks are done in a phase previous to execution or translation (*i.e.*, *static typing*), and similarly assume that any information gathered by this phase is available in subsequent phases.

6.2 The design space of types

We have already discussed the difference between static and dynamic typings, *i.e.*, if type checks are made before or during execution of a program. Additionally, we can distinguish *weakly* and *strongly* typed languages.

Strong typing means that the language implementation ensures that whenever an operation is performed, the arguments to the operation are of a type that the operation is defined for, so you, for example, do not try to concatenate a string and

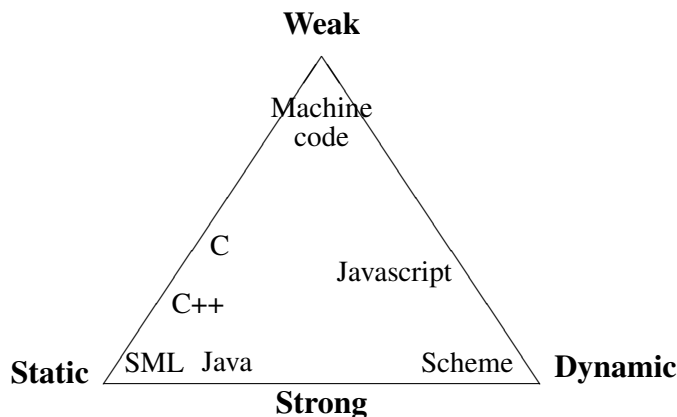


Figure 6.1: The design space of types

a floating-point number. This is independent of whether this is ensured statically (prior to execution) or dynamically (during execution).

In contrast, a weakly typed language gives no guarantee that operations are performed on arguments that make sense for the operation. The archetypical weakly typed language is machine code: Operations are just performed with no checks, and if there is any concept of type at the machine level, it is fairly limited: Registers may be divided into integer, floating point and (possibly) address registers, and memory is (if at all) divided into only code and data. Weakly typed languages are mostly used for system programming, where you need to manipulate move, copy, encrypt or compress data without regard to what the data represents.

Many languages combine both strong and weak typing or both static and dynamic typing: Some types are checked before execution and other during execution, and some types are not checked at all. For example, C is a statically typed language (since no checks are performed during execution), but not all types are checked. For example, you can store an integer in a union-typed variable and read it back as a pointer or floating-point number. Another example is Javascript: If you try to multiply two strings, the interpreter will see if the strings contain sequences of digits and, if so, “read” the strings as numbers and multiply these. This is a kind of weak typing, as the multiplication operation is applied to arguments (strings) where multiplication does not make sense. But instead of, like machine code, blindly trying to multiply the machine representations of the strings as if they were numbers, Javascript performs a dynamic check and conversion to make the values conform to the operation. I will still call this behaviour weak typing, as there is nothing that indicates that converting strings to numbers before multiplication makes any more sense than just multiplying the machine representations of the strings. The main point is that the language, instead of reporting a possible problem, silently does

something that probably makes no sense.

Figure 6.1 shows a diagram of the design space of static vs. dynamic and weak vs. strong typing, placing some well-known programming languages in this design space. Note that the design space is shown as a triangle: If you never check types, you do so neither statically nor dynamically, so at the weak end of the weak vs. strong spectrum, the distinction between static and dynamic is meaningless.

6.3 Attributes

The checking phase operates on the abstract syntax tree of the program and may make several passes over this. Typically, each pass is a recursive walk over the syntax tree, gathering information or using information gathered in earlier passes. Such information is often called *attributes* of the syntax tree. Typically, we distinguish between two types of attributes: *Synthesised attributes* are passed upwards in the syntax tree, from the leaves up to the root. *Inherited attributes* are, conversely, passed downwards in the syntax tree. Note, however, that information that is synthesised in one subtree may be inherited by another subtree or, in a later pass, by the same subtree. An example of this is a symbol table: This is synthesised by a declaration and inherited by the scope of the declaration. When declarations are recursive, the scope may be the same syntax tree as the declaration itself, in which case one pass over this tree will build the symbol table as a synthesised attribute while a second pass will use it as an inherited attribute.

Typically, each *syntactic category* (represented by a type in the data structure for the abstract syntax tree or by a group of related nonterminals in the grammar) will have its own set of attributes. When we write a checker as a set of mutually recursive functions, there will be one or more such functions for each syntactical category. Each of these functions will take inherited attributes (including the syntax tree itself) as arguments and return synthesised attributes as results.

We will, in this chapter, focus on type checking, and only briefly mention other properties that can be checked. The methods used for type checking can in most cases easily be modified to handle such other checks.

We will use the language in section 5.3 as an example for static type checking.

6.4 Environments for type checking

In order to type check the program, we need symbol tables that bind variables and functions to their types. Since there are separate name spaces for variables and functions, we will use two symbol tables, one for variables and one for functions. A variable is bound to one of the two types `int` or `bool`. A function is bound to its type, which consists of the types of its arguments and the type of its result.

Function types are written as a parenthesised list of the argument types, an arrow and the result type, *e.g.*, $(\text{int}, \text{bool}) \rightarrow \text{int}$ for a function taking two parameters (of type `int` and `bool`, respectively) and returning an integer.

We will assume that symbol tables are persistent, so no explicit action is required to restore the symbol table for the outer scope when exiting an inner scope. We don't need to preserve symbol tables for inner scopes once these are exited (so a stack-like behaviour is fine).

6.5 Type checking expressions

When we type check expressions, the symbol tables for variables and functions are inherited attributes. The type (`int` or `bool`) of the expression is returned as a synthesised attribute. To make the presentation independent of any specific data structure for abstract syntax, we will (like in chapter 5) let the type checker function use a notation similar to the concrete syntax for pattern-matching purposes. But you should still think of it as abstract syntax, so all issues of ambiguity, *etc.*, have been resolved.

For terminals (variable names and numeric constants) with attributes, we assume that there are predefined functions for extracting these. Hence, **id** has an associated function *getname*, that extracts the name of the identifier. Similarly, **num** has a function *getvalue*, that returns the value of the number. The latter is not required for static type checking, but we used it in chapter 5 and we will use it again in chapter 7.

For each nonterminal, we define one or more functions that take an abstract syntax subtree and inherited attributes as arguments and return the synthesised attributes.

In figure 6.2, we show the type-checking function for expressions. The function for type checking expressions is called *Check_{Exp}*. The symbol table for variables is given by the parameter *vtable*, and the symbol table for functions by the parameter *ftable*. The function **error** reports a type error. To allow the type checker to continue and report more than one error, we let the error-reporting function return.¹ After reporting a type error, the type checker can make a guess at what the type should have been and return this guess, allowing type checking to continue for the rest of the program. This guess might, however, be wrong, which can cause spurious type errors to be reported later on. Hence, all but the first type error message should be taken with a grain of salt.

We will briefly explain each of the cases handled by *Check_{Exp}*.

- A number has type `int`.

¹Unlike in chapter 5, where the **error** function stops execution.

$Check_{Exp}(Exp, vtable, ftable) = \text{case } Exp \text{ of}$	
num	int
id	$t = \text{lookup}(vtable, \text{getname}(\mathbf{id}))$ if $t = \text{unbound}$ then error() ; int else t
$Exp_1 + Exp_2$	$t_1 = Check_{Exp}(Exp_1, vtable, ftable)$ $t_2 = Check_{Exp}(Exp_2, vtable, ftable)$ if $t_1 = \text{int}$ and $t_2 = \text{int}$ then int else error() ; int
$Exp_1 = Exp_2$	$t_1 = Check_{Exp}(Exp_1, vtable, ftable)$ $t_2 = Check_{Exp}(Exp_2, vtable, ftable)$ if $t_1 = t_2$ then bool else error() ; bool
if Exp_1 then Exp_2 else Exp_3	$t_1 = Check_{Exp}(Exp_1, vtable, ftable)$ $t_2 = Check_{Exp}(Exp_2, vtable, ftable)$ $t_3 = Check_{Exp}(Exp_3, vtable, ftable)$ if $t_1 = \text{bool}$ and $t_2 = t_3$ then t_2 else error() ; t_2
id ($Exps$)	$t = \text{lookup}(ftable, \text{getname}(\mathbf{id}))$ if $t = \text{unbound}$ then error() ; int else $((t_1, \dots, t_n) \rightarrow t_0) = t$ $[t'_1, \dots, t'_m] = Check_{Exps}(Exps, vtable, ftable)$ if $m = n$ and $t_1 = t'_1, \dots, t_n = t'_n$ then t_0 else error() ; t_0
let $\mathbf{id} = Exp_1$ in Exp_2	$t_1 = Check_{Exp}(Exp_1, vtable, ftable)$ $vtable' = \text{bind}(vtable, \text{getname}(\mathbf{id}), t_1)$ $Check_{Exp}(Exp_2, vtable', ftable)$

$Check_{Exps}(Exps, vtable, ftable) = \text{case } Exps \text{ of}$	
Exp	$[Check_{Exp}(Exp, vtable, ftable)]$
$Exp , Exps$	$Check_{Exp}(Exp, vtable, ftable)$ $:: Check_{Exps}(Exps, vtable, ftable)$

Figure 6.2: Type checking of expressions

- The type of a variable is found by looking its name up in the symbol table for variables. If the variable is not found in the symbol table, the lookup-function returns the special value *unbound*. When this happens, an error is reported and the type checker arbitrarily guesses that the type is *int*. Otherwise, it returns the type returned by *lookup*.
- A plus-expression requires both arguments to be integers and has an integer result.
- Comparison requires that the arguments have the same type. In either case, the result is a boolean.
- In a conditional expression, the condition must be of type *bool* and the two branches must have identical types. The result of a condition is the value of one of the branches, so it has the same type as these. If the branches have different types, the type checker reports an error and arbitrarily chooses the type of the *then*-branch as its guess for the type of the whole expression.
- At a function call, the function name is looked up in the function environment to find the number and types of the arguments as well as the return type. The number of arguments to the call must coincide with the expected number and their types must match the declared types. The resulting type is the return-type of the function. If the function name is not found in *fiable*, an error is reported and the type checker arbitrarily guesses the result type to be *int*.
- A *let*-expression declares a new variable, the type of which is that of the expression that defines the value of the variable. The symbol table for variables is extended using the function *bind*, and the extended table is used for checking the body-expression and finding its type, which in turn is the type of the whole expression. A *let*-expression can not in itself be the cause of a type error (though its parts may), so no testing is done.

Since $Check_{Exp}$ mentions the nonterminal $Exps$ and its related type-checking function $Check_{Exps}$, we have included $Check_{Exps}$ in figure 6.2.

$Check_{Exps}$ builds a list of the types of the expressions in the expression list. The notation is taken from SML: A list is written in square brackets with commas between the elements. The operator $::$ adds an element to the front of a list.

Suggested exercises: 6.1.

6.6 Type checking of function declarations

A function declaration explicitly declares the types of the arguments. This information is used to build a symbol table for variables, which is used when type checking

$Check_{Fun}(Fun, ftable) = \text{case } Fun \text{ of}$	
$TypeId (TypeIds) = Exp$	$(f, t_0) = Get_{TypeId}(TypeId)$ $vtable = Check_{TypeIds}(TypeIds)$ $t_1 = Check_{Exp}(Exp, vtable, ftable)$ if $t_0 \neq t_1$ then error ()

$Get_{TypeId}(TypeId) = \text{case } TypeId \text{ of}$	
int id	$(getname(id), \text{int})$
bool id	$(getname(id), \text{bool})$

$Check_{TypeIds}(TypeIds) = \text{case } TypeIds \text{ of}$	
$TypeId$	$(x, t) = Get_{TypeId}(TypeId)$ $bind(emptytable, x, t)$
$TypeId, TypeIds$	$(x, t) = Get_{TypeId}(TypeId)$ $vtable = Check_{TypeIds}(TypeIds)$ if $lookup(vtable, x) = unbound$ then $bind(vtable, x, t)$ else error (); $vtable$

Figure 6.3: Type checking a function declaration

the body of the function. The type of the body must match the declared result type of the function. The type check function for functions, $Check_{Fun}$, has as inherited attribute the symbol table for functions, which is passed down to the type check function for expressions. $Check_{Fun}$ returns no information, it just checks for internal errors. $Check_{Fun}$ is shown in figure 6.3, along with the functions for $TypeId$ and $TypeIds$, which it uses. The function Get_{TypeId} just returns a pair of the declared name and type, and $Check_{TypeIds}$ builds a symbol table from such pairs. $Check_{TypeIds}$ also checks if all parameters have different names. $emptytable$ is an empty symbol table. Looking any name up in the empty symbol table returns *unbound*.

6.7 Type checking a program

A program is a list of functions and is deemed type correct if all the functions are type correct, and there are no two function definitions defining the same function name. Additionally, there must be a function called `main` with one integer argument and integer result.

Since all functions are mutually recursive, each of these must be type checked using a symbol table where all functions are bound to their type. This requires two

passes over the list of functions: One to build the symbol table and one to check the function definitions using this table. Hence, we need two functions operating over *Funs* and two functions operating over *Fun*. We have already seen one of the latter, *Check_{Fun}*. The other, *Get_{Fun}*, returns the pair of the function's declared name and type, which consists of the types of the arguments and the type of the result. It uses an auxiliary function *Get_{Types}* to find the types of the arguments. The two functions for the syntactic category *Funs* are *Get_{Funs}*, which builds the symbol table and checks for duplicate definitions, and *Check_{Funs}*, which calls *Check_{Fun}* for all functions. These functions and the main function *Check_{Program}*, which ties the loose ends, are shown in figure 6.4.

This completes type checking of our small example language.

Suggested exercises: 6.5.

6.8 Advanced type checking

Our example language is very simple and obviously does not cover all aspects of type checking. A few examples of other features and brief explanations of how they can be handled are listed below.

Assignments. When a variable is given a value by an assignment, it must be verified that the type of the value is the same as the declared type of the variable. Some compilers may check if a variable is potentially used before it is given a value, or if a variable is not used after its assignment. While not exactly type errors, such behaviour is likely to be undesirable. Testing for such behaviour does, however, require somewhat more complicated analysis than the simple type checking presented in this chapter, as it relies on non-structural information.

Data structures. A data structure may define a value with several components (*e.g.*, a *struct*, *tuple* or *record*), or a value that may be of different types at different times (*e.g.*, a *union*, *variant* or *sum*). To type check such structures, the type checker must be able to represent their types. Hence, the type checker may need a data structure that describes complex types. This may be similar to the data structure used for the abstract syntax trees of declarations. Operations that build or take apart structured data need to be tested for correctness. If each operation on structured data has well-defined types for its arguments and a type for its result, this can be done in a way similar to how function calls are tested.

Overloading. Overloading means that the same name is used for several different operations over several different types. We saw a simple example of this in the

$Check_{Program}(Program) = \text{case } Program \text{ of}$	
$Funs$	$f_{table} = Get_{Funs}(Funs)$ $Check_{Funs}(Funs, f_{table})$ <i>if</i> $lookup(f_{table}, main) \neq (int) \rightarrow int$ <i>then</i> error()
$Get_{Funs}(Funs) = \text{case } Funs \text{ of}$	
Fun	$(f, t) = Get_{Fun}(Fun)$ $bind(emptytable, f, t)$
$Fun Funs$	$(f, t) = Get_{Fun}(Fun)$ $f_{table} = Get_{Funs}(Funs)$ <i>if</i> $lookup(f_{table}, f) = unbound$ <i>then</i> $bind(f_{table}, f, t)$ <i>else</i> error(); f_{table}
$Get_{Fun}(Fun) = \text{case } Fun \text{ of}$	
$TypeId (TypeIds) = Exp$	$(f, t_0) = Get_{TypeId}(TypeId)$ $[t_1, \dots, t_n] = Get_{Types}(TypeIds)$ $(f, (t_1, \dots, t_n) \rightarrow t_0)$
$Get_{Types}(TypeIds) = \text{case } TypeIds \text{ of}$	
$TypeId$	$(x, t) = Get_{TypeId}(TypeId)$ $[t]$
$TypeId TypeIds$	$(x_1, t_1) = Get_{TypeId}(TypeId)$ $[t_2, \dots, t_n] = Get_{Types}(TypeIds)$ $[t_1, t_2, \dots, t_n]$
$Check_{Funs}(Funs, f_{table}) = \text{case } Funs \text{ of}$	
Fun	$Check_{Fun}(Fun, f_{table})$
$Fun Funs$	$Check_{Fun}(Fun, f_{table})$ $Check_{Funs}(Funs, f_{table})$

Figure 6.4: Type checking a program

example language, where $=$ was used both for comparing integers and booleans. In many languages, arithmetic operators like $+$ and $-$ are defined both over integers and floating point numbers, and possibly other types as well. If these operators are predefined, and there is only a finite number of cases they cover, all the possible cases may be tried in turn, just like in our example.

This, however, requires that the different instances of the operator have disjoint argument types. If, for example, there is a function *read* that reads a value from a text stream and this is defined to read either integers or floating point numbers, the argument (the text stream) alone can not be used to select the right operator. Hence, the type checker must pass the expected type of each expression down as an inherited attribute, so this (possibly in combination with the types of the arguments) can be used to pick the correct instance of the overloaded operator.

It may not always be possible to send down an expected type due to lack of information. In our example language, this is the case for the arguments to $=$ (as these may be either `int` or `bool`) and the first expression in a `let`-expression (since the variable bound in the `let`-expression is not declared to be a specific type). If the type checker for this or some other reason is unable to pick a unique operator, it may report “unresolved overloading” as a type error, or it may pick a default instance.

Type conversion. A language may have operators for converting a value of one type to a value of another type, *e.g.* an integer to a floating point number. Sometimes these operators are explicit in the program and hence easy to check. However, many languages allow implicit conversion of integers to floats, such that, for example, $3 + 3.12$ is well-typed with the implicit assumption that the integer 3 is converted to a float before the addition. This can be handled as follows: If the type checker discovers that the arguments to an operator do not have the correct type, it can try to convert one or both arguments to see if this helps. If there is a small number of predefined legal conversions, this is no major problem. However, a combination of user-defined overloaded operators and user-defined types with conversions can make the type-checking process quite difficult, as the information needed to choose correctly may not be available at compile-time. This is typically the case in object-oriented languages, where method selection is often done at run-time. We will not go into details of how this can be done.

Polymorphism / Generic types. Some languages allow a function to be *polymorphic* or *generic*, that is, to be defined over a large class of similar types, *e.g.* over all arrays no matter what the types of the elements are. A function can explicitly declare which parts of the type is generic/polymorphic or this can be implicit (see below). The type checker can insert the actual types at every use of the generic/polymorphic function to create *instances* of the generic/polymorphic type.

This mechanism is different from overloading as the instances will be related by a common generic type and because a polymorphic/generic function can be instantiated by any type, not just by a limited list of declared alternatives as is the case with overloading.

Implicit types. Some languages (like Standard ML and Haskell) require programs to be well-typed, but do not require explicit type declarations for variables or functions. For such to work, a *type inference* algorithm is used. A type inference algorithm gathers information about uses of functions and variables and uses this information to infer the types of these. If there are inconsistent uses of a variable, a type error is reported.

Suggested exercises: 6.2.

6.9 Further reading

Overloading of operators and functions is described in section 6.5 of [5]. Section 6.7 of same describes how polymorphism can be handled.

Some theory and a more detailed algorithm for inferring types in a language with implicit types and polymorphism can be found in [32].

Exercises

Exercise 6.1

We extend the language from section 5.3 with boolean operators as described in exercise 5.1.

Extend the type-check function in figure 6.2 to handle these new constructions as described above.

Exercise 6.2

We extend the language from section 5.3 with floating-point numbers as described in exercise 5.2.

- a) Extend the type checking functions in figures 6.2-6.4 to handle these extensions.
- b) We now add implicit conversion of integers to floats to the language, using the rules: Whenever an operator has one integer argument and one floating-point argument, the integer is converted to a float. Similarly, if a condition

expression (if-then-else) has one integer branch and one floating-point branch, the integer branch is converted to floating-point. Extend the type checking functions from question a) above to handle this.

Exercise 6.3

The type check function in figure 6.2 tries to guess the correct type when there is a type error. In some cases, the guess is arbitrarily chosen to be `int`, which may lead to spurious type errors later on. A way around this is to have an extra type: `unknown`, which is only used during type checking. If there is a type error and there is no basis for guessing a correct type, `unknown` is returned (the error is still reported, though). If an argument to an operator is of type `unknown`, the type checker should not report this as a type error but continue as if the type is correct. The use of an `unknown` argument to an operator may make the result `unknown` as well, so these can be propagated arbitrarily far.

Change figure 6.2 to use the `unknown` type as described above.

Exercise 6.4

We look at a simple language with an exception mechanism:

$$\begin{aligned} S &\rightarrow \text{throw } \mathbf{id} \\ S &\rightarrow S \text{ catch } \mathbf{id} \Rightarrow S \\ S &\rightarrow S \text{ or } S \\ S &\rightarrow \text{other} \end{aligned}$$

A `throw` statement throws a named exception. This is caught by the nearest enclosing `catch` statement (*i.e.*, where the `throw` statement is in the left sub-statement of the `catch` statement) using the same name, whereby the statement after the arrow in the `catch` statement is executed. An `or` statement is a non-deterministic choice between the two statements, so either one can be executed. `other` is a statement that do not throw any exceptions.

We want the type checker to ensure that all possible exceptions are caught and that no `catch` statement is superfluous, *i.e.*, that the exception it catches can, in fact, be thrown by its left sub-statement.

Write type-check functions that implement these checks. Hint: Let the type of a statement be the set of possible exceptions it can throw.

Exercise 6.5

In exercise 5.5, we extended the example language with closures and implemented these in the interpreter.

Extend the type-checking functions in figures 6.2-6.4 to statically type check the same extensions.

Hint: Check a function definition when it is declared.

Chapter 7

Intermediate-Code Generation

7.1 Introduction

The final goal of a compiler is to get programs written in a high-level language to run on a computer. This means that, eventually, the program will have to be expressed as machine code which can run on the computer. This does not mean that we need to translate directly from the high-level abstract syntax to machine code. Many compilers use a medium-level language as a stepping-stone between the high-level language and the very low-level machine code. Such stepping-stone languages are called *intermediate code*.

Apart from structuring the compiler into smaller jobs, using an intermediate language has other advantages:

- If the compiler needs to generate code for several different machine-architectures, only one translation to intermediate code is needed. Only the translation from intermediate code to machine language (*i.e.*, the *back-end*) needs to be written in several versions.
- If several high-level languages need to be compiled, only the translation to intermediate code need to be written for each language. They can all share the back-end, *i.e.*, the translation from intermediate code to machine code.
- Instead of translating the intermediate language to machine code, it can be *interpreted* by a small program written in machine code or a language for which a compiler or interpreter already exists.

The advantage of using an intermediate language is most obvious if many languages are to be compiled to many machines. If translation is done directly, the number of compilers is equal to the product of the number of languages and the number of machines. If a common intermediate language is used, one front-end (*i.e.*, compiler

to intermediate code) is needed for every language and one back-end (interpreter or code generator) is needed for each machine, making the total number of front-ends and back-ends equal to the sum of the number of languages and the number of machines.

If an interpreter for an intermediate language is written in a language for which there already exist implementations for the target machines, the same interpreter can be interpreted or compiled for each machine. This way, there is no need to write a separate back-end for each machine. The advantages of this approach are:

- No actual back-end needs to be written for each new machine, as long as the machine is equipped with an interpreter or compiler for the implementation language of the interpreter for the intermediate language.
- A compiled program can be distributed in a single intermediate form for all machines, as opposed to shipping separate binaries for each machine.
- The intermediate form may be more compact than machine code. This saves space both in distribution and on the machine that executes the programs (though the latter is somewhat offset by requiring the interpreter to be kept in memory during execution).

The disadvantage is speed: Interpreting the intermediate form will in most cases be a lot slower than executing translated code directly. Nevertheless, the approach has seen some success, *e.g.*, with Java.

Some of the speed penalty can be eliminated by translating the intermediate code to machine code immediately before or during execution of the program. This hybrid form is called *just-in-time compilation* and is often used for executing the intermediate code for Java.

We will in this book, however, focus mainly on using the intermediate code for traditional compilation, where the intermediate form will be translated to machine code by the back-end of the compiler.

7.2 Choosing an intermediate language

An intermediate language should, ideally, have the following properties:

- It should be easy to translate from a high-level language to the intermediate language. This should be the case for a wide range of different source languages.
- It should be easy to translate from the intermediate language to machine code. This should be true for a wide range of different target architectures.

- The intermediate format should be suitable for optimisations.

The first two of these properties can be somewhat hard to reconcile. A language that is intended as target for translation from a high-level language should be fairly close to this. However, this may be hard to achieve for more than a small number of similar languages. Furthermore, a high-level intermediate language puts more burden on the back-ends. A low-level intermediate language may make it easy to write back-ends, but puts more burden on the front-ends. A low-level intermediate language, also, may not fit all machines equally well, though this is usually less of a problem than the similar problem for front-ends, as machines typically are more similar than high-level languages.

A solution that may reduce the translation burden, though it does not address the other problems, is to have two intermediate levels: One, which is fairly high-level, is used for the front-ends and the other, which is fairly low-level, is used for the back-ends. A single shared translator is then used to translate between these two intermediate formats.

When the intermediate format is shared between many compilers, it makes sense to do as many optimisations as possible on the intermediate format. This way, the (often substantial) effort of writing good optimisations is done only once instead of in every compiler.

Another thing to consider when choosing an intermediate language is the “granularity”: Should an operation in the intermediate language correspond to a large amount of work or to a small amount of work?

The first of these approaches is often used when the intermediate language is interpreted, as the overhead of decoding instructions is amortised over more actual work, but it can also be used for compiling. In this case, each intermediate-code operation is typically translated into a sequence of machine-code instructions. When coarse-grained intermediate code is used, there is typically a fairly large number of different intermediate-code operations.

The opposite approach is to let each intermediate-code operation be as small as possible. This means that each intermediate-code operation is typically translated into a single machine-code instruction or that several intermediate-code operations can be combined into one machine-code operation. The latter can, to some degree, be automated as each machine-code instruction can be described as a sequence of intermediate-code instructions. When intermediate-code is translated to machine-code, the code generator can look for sequences that match machine-code operations. By assigning cost to each machine-code operation, this can be turned into a combinatorial optimisation problem, where the least-cost solution is found. We will return to this in chapter 8.

<i>Program</i>	\rightarrow	[<i>Instructions</i>]
<i>Instructions</i>	\rightarrow	<i>Instruction</i>
<i>Instructions</i>	\rightarrow	<i>Instruction</i> , <i>Instructions</i>
<i>Instruction</i>	\rightarrow	LABEL labelid
<i>Instruction</i>	\rightarrow	id := <i>Atom</i>
<i>Instruction</i>	\rightarrow	id := unop <i>Atom</i>
<i>Instruction</i>	\rightarrow	id := id binop <i>Atom</i>
<i>Instruction</i>	\rightarrow	id := <i>M</i> [<i>Atom</i>]
<i>Instruction</i>	\rightarrow	<i>M</i> [<i>Atom</i>] := id
<i>Instruction</i>	\rightarrow	GOTO labelid
<i>Instruction</i>	\rightarrow	IF id relop <i>Atom</i> THEN labelid ELSE labelid
<i>Instruction</i>	\rightarrow	id := CALL functionid (<i>Args</i>)
<i>Atom</i>	\rightarrow	id
<i>Atom</i>	\rightarrow	num
<i>Args</i>	\rightarrow	id
<i>Args</i>	\rightarrow	id , <i>Args</i>

Grammar 7.1: The intermediate language

7.3 The intermediate language

In this chapter we have chosen a fairly low-level fine-grained intermediate language, as it is best suited to convey the techniques we want to cover.

We will not treat translation of function calls until chapter 10, so a “program” in our intermediate language will, for the time being, correspond to the body of a function or procedure in a real program. For the same reason, function calls are initially treated as primitive operations in the intermediate language.

The grammar for the intermediate language is shown in grammar 7.1. A program is a sequence of instructions. The instructions are:

- A label. This has no effect but serves only to mark the position in the program as a target for jumps.
- An assignment of an atomic expression (constant or variable) to a variable.
- A unary operator applied to an atomic expression, with the result stored in a variable.

- A binary operator applied to a variable and an atomic expression, with the result stored in a variable.
- A transfer from memory to a variable. The memory location is an atomic expression.
- A transfer from a variable to memory. The memory location is an atomic expression.
- A jump to a label.
- A conditional selection between jumps to two labels. The condition is found by comparing a variable with an atomic expression by using a relational operator ($=$, \neq , $<$, $>$, \leq or \geq).
- A function call. The arguments to the function call are variables and the result is assigned to a variable. This instruction is used even if there is no actual result (*i.e.*, if a procedure is called instead of a function), in which case the result variable is a dummy variable.

An atomic expression is either a variable or a constant.

We have not specified the set of unary and binary operations, but we expect these to include normal integer arithmetic and bitwise logical operations.

We assume that all values are integers. Adding floating-point numbers and other primitive types is not difficult, though.

7.4 Syntax-directed translation

We will generate code using translation functions for each syntactic category, similarly to the functions we used for interpretation and type checking. We generate code for a syntactic construct independently of the constructs around it, except that the parameters of a translation function may hold information about the context (such as symbol tables) and the result of a translation function may (in addition to the generated code) hold information about how the generated code interfaces with its context (such as which variables it uses). Since the translation closely follows the syntactic structure of the program, it is called *syntax-directed translation*.

Given that translation of a syntactic construct is mostly independent of the surrounding and enclosed syntactic constructs, we might miss opportunities to exploit synergies between these and, hence, generate less than optimal code. We will try to remedy this in later chapters by using various optimisation techniques.

$$\begin{aligned}
Exp &\rightarrow \mathbf{num} \\
Exp &\rightarrow \mathbf{id} \\
Exp &\rightarrow \mathbf{unop} \, Exp \\
Exp &\rightarrow Exp \, \mathbf{binop} \, Exp \\
Exp &\rightarrow \mathbf{id}(Exps) \\
\\
Exps &\rightarrow Exp \\
Exps &\rightarrow Exp \, , \, Exps
\end{aligned}$$

Grammar 7.2: A simple expression language

7.5 Generating code from expressions

Grammar 7.2 shows a simple language of expressions, which we will use as our initial example for translation. Again, we have let the set of unary and binary operators be unspecified but assume that the intermediate language includes all those used by the expression language. We assume that there is a function *transop* that translates the name of an operator in the expression language into the name of the corresponding operator in the intermediate language. The tokens **unop** and **binop** have the names of the actual operators as attributes, accessed by the function *getopname*.

When writing a compiler, we must decide what needs to be done at compile-time and what needs to be done at run-time. Ideally, as much as possible should be done at compile-time, but some things need to be postponed until run-time, as they need the actual values of variables, *etc.*, which are not known at compile-time. When we, below, explain the workings of the translation functions, we might use phrasing like “the expression is evaluated and the result stored in the variable”. This describes actions that are performed at run-time by the code that is generated at compile-time. At times, the textual description may not be 100% clear as to what happens at which time, but the notation used in the translation functions make this clear: Intermediate-language code is executed at run-time, the rest is done at compile time. Intermediate-language instructions may refer to values (constants and register names) that are generated at compile time. When instructions have operands that are written in *italics*, these operands are variables in the compiler that contain compile-time values that are inserted into the generated code. For example, if *place* holds the variable name `t14` and *v* holds the value 42, then the code template `[place := v]` will generate the code `[t14 := 42]`.

When we want to translate the expression language to the intermediate language, the main complication is that the expression language is tree-structured

while the intermediate language is flat, requiring the result of every operation to be stored in a variable and every (non-constant) argument to be in one. We use a function *newvar* to generate new variable names in the intermediate language. Whenever *newvar* is called, it returns a previously unused variable name.

We will describe translation of expressions by a translation function using a notation similar to the notation we used for type-checking functions in chapter 6.

Some attributes for the translation function are obvious: It must return the code as a synthesised attribute. Furthermore, it must translate variables and functions used in the expression language to the names these correspond to in the intermediate language. This can be done by symbol tables *vtable* and *f table* that bind variable and function names in the expression language into the corresponding names in the intermediate language. The symbol tables are passed as inherited attributes to the translation function. In addition to these attributes, the translation function must use attributes to decide where to put the values of sub-expressions. This can be done in two ways:

- 1) The location of the values of a sub-expression can be passed up as a synthesised attribute to the parent expression, which decides on a position for its own value.
- 2) The parent expression can decide where it wants to find the values of its sub-expressions and pass this information down to these as inherited attributes.

Neither of these is obviously superior to the other. Method 1 has a slight advantage when generating code for a variable access, as it does not have to generate any code, but can simply return the name of the variable that holds the value. This, however, only works under the assumption that the variable is not updated before the value is used by the parent expression. If expressions can have side effects, this is not always the case, as the C expression “*x+(x=3)*” shows. Our expression language does not have assignment, but it does have function calls, which may have side effects.

Method 2 does not have this problem: Since the value of the expression is created immediately before the assignment is executed, there is no risk of other side effects between these two points in time. Method 2 also has a slight advantage when we later extend the language to have assignment statements, as we can then generate code that calculates the expression result directly into the desired variable instead of having to copy it from a temporary variable.

Hence, we will choose method 2 for our translation function *Trans_{Exp}*, which is shown in figure 7.3.

The inherited attribute *place* is the intermediate-language variable that the result of the expression must be stored in.

$Trans_{Exp}(Exp, vtable, ftable, place) = \text{case } Exp \text{ of}$	
num	$v = \text{getvalue}(\mathbf{num})$ $[place := v]$
id	$x = \text{lookup}(vtable, \text{getname}(\mathbf{id}))$ $[place := x]$
unop Exp_1	$place_1 = \text{newvar}()$ $code_1 = Trans_{Exp}(Exp_1, vtable, ftable, place_1)$ $op = \text{transop}(\text{getopname}(\mathbf{unop}))$ $code_1 ++ [place := op \ place_1]$
Exp_1 binop Exp_2	$place_1 = \text{newvar}()$ $place_2 = \text{newvar}()$ $code_1 = Trans_{Exp}(Exp_1, vtable, ftable, place_1)$ $code_2 = Trans_{Exp}(Exp_2, vtable, ftable, place_2)$ $op = \text{transop}(\text{getopname}(\mathbf{binop}))$ $code_1 ++ code_2 ++ [place := place_1 \ op \ place_2]$
id ($Exps$)	$(code_1, [a_1, \dots, a_n])$ $\quad = Trans_{Exps}(Exps, vtable, ftable)$ $fname = \text{lookup}(ftable, \text{getname}(\mathbf{id}))$ $code_1 ++ [place := \text{CALL } fname(a_1, \dots, a_n)]$

$Trans_{Exps}(Exps, vtable, ftable) = \text{case } Exps \text{ of}$	
Exp	$place = \text{newvar}()$ $code_1 = Trans_{Exp}(Exp, vtable, ftable, place)$ $(code_1, [place])$
$Exp, Exps$	$place = \text{newvar}()$ $code_1 = Trans_{Exp}(Exp, vtable, ftable, place)$ $(code_2, args) = Trans_{Exps}(Exps, vtable, ftable)$ $code_3 = code_1 ++ code_2$ $args_1 = place :: args$ $(code_3, args_1)$

Figure 7.3: Translating an expression

If the expression is just a number, the value of that number is stored in the *place*.

If the expression is a variable, the intermediate-language equivalent of this variable is found in *vtable* and an assignment copies it into the intended *place*.

A unary operation is translated by first generating a new intermediate-language variable to hold the value of the argument of the operation. Then the argument is translated using the newly generated variable for the *place* attribute. We then use an **unop** operation in the intermediate language to assign the result to the inherited *place*. The operator ++ concatenates two lists of instructions.

A binary operation is translated in a similar way. Two new intermediate-language variables are generated to hold the values of the arguments, then the arguments are translated and finally a binary operation in the intermediate language assigns the final result to the inherited *place*.

A function call is translated by first translating the arguments, using the auxiliary function *Trans_{Exps}*. Then a function call is generated using the argument variables returned by *Trans_{Exps}*, with the result assigned to the inherited *place*. The name of the function is looked-up in *f_{table}* to find the corresponding intermediate-language name.

Trans_{Exps} generates code for each argument expression, storing the results into new variables. These variables are returned along with the code, so they can be put into the argument list of the call instruction.

7.5.1 Examples of translation

Translation of expressions is always relative to symbol tables and a place for storing the result. In the examples below, we assume a variable symbol table that binds *x*, *y* and *z* to *v0*, *v1* and *v2*, respectively and a function table that binds *f* to *_f*. The place for the result is *t0* and we assume that calls to *newvar()* return, in sequence, the variables *t1*, *t2*, *t3*, ...

We start by the simple expression *x-3*. This is a binop-expression, so the first we do is to call *newvar()* twice, giving *place₁* the value *t1* and *place₂* the value *t2*. We then call *Trans_{Exp}* recursively with the expression *x*. When translating this, we first look up *x* in the variable symbol table, yielding *v0*, and then return the code [*t1* := *v0*]. Back in the translation of the subtraction expression, we assign this code to *code₁* and once more call *Trans_{Exp}* recursively, this time with the expression *3*. This is translated to the code [*t2* := *3*], which we assign to *code₂*. The final result is produced by *code₁*++*code₂*++[*t0* := *t1* - *t2*] which yields [*t1* := *v0*, *t2* := *3*, *t0* := *t1* - *t2*]. We have translated the source-language operator - to the intermediate-language operator -.

The resulting code looks quite suboptimal, and could, indeed, be shortened to [*t0* := *v0* - *3*]. When we generate intermediate code, we want, for simplicity, to

$$\begin{aligned}
Stat &\rightarrow Stat ; Stat \\
Stat &\rightarrow \mathbf{id} := Exp \\
Stat &\rightarrow \mathbf{if} Cond \mathbf{then} Stat \\
Stat &\rightarrow \mathbf{if} Cond \mathbf{then} Stat \mathbf{else} Stat \\
Stat &\rightarrow \mathbf{while} Cond \mathbf{do} Stat \\
Stat &\rightarrow \mathbf{repeat} Stat \mathbf{until} Cond \\
\\
Cond &\rightarrow Exp \mathbf{relop} Exp
\end{aligned}$$

Grammar 7.4: Statement language

treat each subexpression independently of its context. This may lead to superfluous assignments. We will look at ways of getting rid of these when we treat machine code generation and register allocation in chapters 8 and 9.

A more complex expression is $3+f(x-y,z)$. Using the same assumptions as above, this yields the code

```

t1 := 3
  t4 := v0
  t5 := v1
  t3 := t4 - t5
  t6 := v2
  t2 := CALL_f(t3,t6)
t0 := t1 + t2

```

We have, for readability, laid the code out on separate lines rather than using a comma-separated list. The indentation indicates the depth of calls to $Trans_{Exp}$ that produced the code in each line.

Suggested exercises: 7.1.

7.6 Translating statements

We now extend the expression language in figure 7.2 with statements. The extensions are shown in grammar 7.4.

When translating statements, we will need the symbol table for variables (for translating assignment), and since statements contain expressions, we also need f_{table} so we can pass it on to $Trans_{Exp}$.

Just like we use *newvar* to generate new unused variables, we use a similar function *newlabel* to generate new unused labels. The translation function for statements is shown in figure 7.5. It uses an auxiliary translation function for conditions shown in figure 7.6.

A sequence of two statements are translated by putting the code for these in sequence.

An assignment is translated by translating the right-hand-side expression using the left-hand-side variable as target location (*place*).

When translating statements that use conditions, we use an auxiliary function *TransCond*. *TransCond* translates the arguments to the condition and generates an IF-THEN-ELSE instruction using the same relational operator as the condition. The target labels of this instruction are inherited attributes to *TransCond*.

An if-then statement is translated by first generating two labels: One for the then-branch and one for the code following the if-then statement. The condition is translated by *TransCond*, which is given the two labels as attributes. When (at run-time) the condition is true, the first of these are selected, and when false, the second is chosen. Hence, when the condition is true, the then-branch is executed followed by the code after the if-then statement. When the condition is false, we jump directly to the code following the if-then statement, hence bypassing the then-branch.

An if-then-else statement is treated similarly, but now the condition must choose between jumping to the then-branch or the else-branch. At the end of the then-branch, a jump bypasses the code for the else-branch by jumping to the label at the end. Hence, there is need for three labels: One for the then-branch, one for the else-branch and one for the code following the if-then-else statement.

If the condition in a while-do loop is true, the body must be executed, otherwise the body is by-passed and the code after the loop is executed. Hence, the condition is translated with attributes that provide the label for the start of the body and the label for the code after the loop. When the body of the loop has been executed, the condition must be re-tested for further passes through the loop. Hence, a jump is made to the start of the code for the condition. A total of three labels are thus required: One for the start of the loop, one for the loop body and one for the end of the loop.

A repeat-until loop is slightly simpler. The body precedes the condition, so there is always at least one pass through the loop. If the condition is true, the loop is terminated and we continue with the code after the loop. If the condition is false, we jump to the start of the loop. Hence, only two labels are needed: One for the start of the loop and one for the code after the loop.

Suggested exercises: 7.2.

<i>Trans_{Stat}(Stat, vtable, ftable) = case Stat of</i>	
<i>Stat₁ ; Stat₂</i>	<i>code₁ = Trans_{Stat}(Stat₁, vtable, ftable)</i> <i>code₂ = Trans_{Stat}(Stat₂, vtable, ftable)</i> <i>code₁ ++ code₂</i>
<i>id := Exp</i>	<i>place = lookup(vtable, getname(id))</i> <i>Trans_{Exp}(Exp, vtable, ftable, place)</i>
<i>if Cond</i> <i>then Stat₁</i>	<i>label₁ = newlabel()</i> <i>label₂ = newlabel()</i> <i>code₁ = Trans_{Cond}(Cond, label₁, label₂, vtable, ftable)</i> <i>code₂ = Trans_{Stat}(Stat₁, vtable, ftable)</i> <i>code₁ ++ [LABEL label₁] ++ code₂</i> <i>++ [LABEL label₂]</i>
<i>if Cond</i> <i>then Stat₁</i> <i>else Stat₂</i>	<i>label₁ = newlabel()</i> <i>label₂ = newlabel()</i> <i>label₃ = newlabel()</i> <i>code₁ = Trans_{Cond}(Cond, label₁, label₂, vtable, ftable)</i> <i>code₂ = Trans_{Stat}(Stat₁, vtable, ftable)</i> <i>code₃ = Trans_{Stat}(Stat₂, vtable, ftable)</i> <i>code₁ ++ [LABEL label₁] ++ code₂</i> <i>++ [GOTO label₃, LABEL label₂]</i> <i>++ code₃ ++ [LABEL label₃]</i>
<i>while Cond</i> <i>do Stat₁</i>	<i>label₁ = newlabel()</i> <i>label₂ = newlabel()</i> <i>label₃ = newlabel()</i> <i>code₁ = Trans_{Cond}(Cond, label₂, label₃, vtable, ftable)</i> <i>code₂ = Trans_{Stat}(Stat₁, vtable, ftable)</i> <i>[LABEL label₁] ++ code₁</i> <i>++ [LABEL label₂] ++ code₂</i> <i>++ [GOTO label₁, LABEL label₃]</i>
<i>repeat Stat₁</i> <i>until Cond</i>	<i>label₁ = newlabel()</i> <i>label₂ = newlabel()</i> <i>code₁ = Trans_{Stat}(Stat₁, vtable, ftable)</i> <i>code₂ = Trans_{Cond}(Cond, label₂, label₁, vtable, ftable)</i> <i>[LABEL label₁] ++ code₁</i> <i>++ code₂ ++ [LABEL label₂]</i>

Figure 7.5: Translation of statements

$Trans_{Cond}(Cond, label_t, label_f, vtable, ftable) = \text{case } Cond \text{ of}$	
$Exp_1 \text{ relop } Exp_2$	$t_1 = \text{newvar}()$ $t_2 = \text{newvar}()$ $code_1 = Trans_{Exp}(Exp_1, vtable, ftable, t_1)$ $code_2 = Trans_{Exp}(Exp_2, vtable, ftable, t_2)$ $op = \text{transop}(\text{getopname}(\text{relop}))$ $code_1 ++ code_2 ++ [\text{IF } t_1 \text{ opt } t_2 \text{ THEN } label_t \text{ ELSE } label_f]$

Figure 7.6: Translation of simple conditions

7.7 Logical operators

Logical conjunction, disjunction and negation are often available for conditions, so we can write, *e.g.*, $x = y$ **or** $y = z$, where **or** is a logical disjunction operator. There are typically two ways to treat logical operators in programming languages:

- 1) Logical operators are similar to arithmetic operators: The arguments are evaluated and the operator is applied to find the result.
- 2) The second operand of a logical operator is not evaluated if the first operand is sufficient to determine the result. This means that a logical **and** will not evaluate its second operand if the first evaluates to **false**, and a logical **or** will not evaluate the second operand if the first is **true**.

The first variant is typically implemented by using bitwise logical operators and uses 0 to represent **false** and a nonzero value (typically 1 or -1) to represent **true**. In C, there is no separate boolean type. The integer 1 is used for logical truth¹ and 0 for falsehood. Bitwise logical operators **&** (bitwise **and**) and **|** (bitwise **or**) are used to implement the corresponding logical operations. Logical negation is *not* handled by bitwise negation, as the bitwise negation of 1 is not 0. Instead, a special logical negation operator **!** is used that maps any non-zero value to 0 and 0 to 1. We assume an equivalent operator is available in the intermediate language.

The second variant is called *sequential logical operators*. In C, these are called **&&** (logical **and**) and **||** (logical **or**).

Adding non-sequential logical operators to our language is not too difficult. Since we have not said exactly which binary and unary operators exist in the intermediate language, we can simply assume these include relational operators, bitwise logical operations and logical negation. We can now simply allow any expression² as a condition by adding the production

¹Actually, any non-zero value is treated as logical truth.

²If it is of boolean type, which we assume has been verified by the type checker.

$$Cond \rightarrow Exp$$

to grammar 7.4. We then extend the translation function for conditions as follows:

$Trans_{Cond}(Cond, label_t, label_f, vtable, ftable) = \text{case } Cond \text{ of}$	
$Exp_1 \text{ relop } Exp_2$	$t_1 = \text{newvar}()$ $t_2 = \text{newvar}()$ $code_1 = Trans_{Exp}(Exp_1, vtable, ftable, t_1)$ $code_2 = Trans_{Exp}(Exp_2, vtable, ftable, t_2)$ $op = \text{transop}(\text{getopname}(\text{relop}))$ $code_1 ++ code_2 ++ [\text{IF } t_1 \text{ op } t_2 \text{ THEN } label_t \text{ ELSE } label_f]$
Exp	$t = \text{newvar}()$ $code_1 = Trans_{Exp}(Exp, vtable, ftable, t)$ $code_1 ++ [\text{IF } t \neq 0 \text{ THEN } label_t \text{ ELSE } label_f]$

We need to convert the numerical value returned by $Trans_{Exp}$ into a choice between two labels, so we generate an IF instruction that does just that.

The rule for relational operators is now actually superfluous, as the case it handles is covered by the second rule (since relational operators are assumed to be included in the set of binary arithmetic operators). However, we can consider it an optimisation, as the code it generates is shorter than the equivalent code generated by the second rule. It will also be natural to keep it separate when we add sequential logical operators.

7.7.1 Sequential logical operators

We will use the same names for sequential logical operators as C, *i.e.*, **&&** for logical **and**, **||** for logical **or** and **!** for logical negation. The extended language is shown in figure 7.7. Note that we allow an expression to be a condition as well as a condition to be an expression. This grammar is highly ambiguous (not least because **binop** overlaps **relop**). As before, we assume such ambiguity to be resolved by the parser before code generation. We also assume that the last productions of Exp and $Cond$ are used as little as possible, as this will yield the best code.

The revised translation functions for Exp and $Cond$ are shown in figure 7.8. Only the new cases for Exp are shown.

As expressions, **true** and **false** are the numbers 1 and 0.

A condition $Cond$ is translated into code that chooses between two labels. When we want to use a condition as an expression, we must convert this choice into a number. We do this by first assuming that the condition is false and hence assign 0 to the target location. We then, if the condition is true, jump to code that assigns 1 to the target location. If the condition is false, we jump around this code, so

$$\begin{aligned} \textit{Exp} &\rightarrow \mathbf{num} \\ \textit{Exp} &\rightarrow \mathbf{id} \\ \textit{Exp} &\rightarrow \mathbf{unop} \textit{Exp} \\ \textit{Exp} &\rightarrow \textit{Exp} \mathbf{binop} \textit{Exp} \\ \textit{Exp} &\rightarrow \mathbf{id}(\textit{Exps}) \\ \textit{Exp} &\rightarrow \mathbf{true} \\ \textit{Exp} &\rightarrow \mathbf{false} \\ \textit{Exp} &\rightarrow \textit{Cond} \\ \\ \textit{Exps} &\rightarrow \textit{Exp} \\ \textit{Exps} &\rightarrow \textit{Exp} , \textit{Exps} \\ \\ \textit{Cond} &\rightarrow \textit{Exp} \mathbf{relop} \textit{Exp} \\ \textit{Cond} &\rightarrow \mathbf{true} \\ \textit{Cond} &\rightarrow \mathbf{false} \\ \textit{Cond} &\rightarrow \mathbf{!} \textit{Cond} \\ \textit{Cond} &\rightarrow \textit{Cond} \mathbf{\&\&} \textit{Cond} \\ \textit{Cond} &\rightarrow \textit{Cond} \mathbf{||} \textit{Cond} \\ \textit{Cond} &\rightarrow \textit{Exp} \end{aligned}$$

Grammar 7.7: Example language with logical operators

$Trans_{Exp}(Exp, vtable, ftable, place) = \text{case } Exp \text{ of}$	
	\vdots
true	$[place := 1]$
false	$[place := 0]$
Cond	$label_1 = \text{newlabel}()$ $label_2 = \text{newlabel}()$ $code_1 = Trans_{Cond}(Cond, label_1, label_2, vtable, ftable)$ $[place := 0] ++ code_1$ $++ [\text{LABEL } label_1, place := 1]$ $++ [\text{LABEL } label_2]$

$Trans_{Cond}(Cond, label_t, label_f, vtable, ftable) = \text{case } Cond \text{ of}$	
$Exp_1 \text{ relop } Exp_2$	$t_1 = \text{newvar}()$ $t_2 = \text{newvar}()$ $code_1 = Trans_{Exp}(Exp_1, vtable, ftable, t_1)$ $code_2 = Trans_{Exp}(Exp_2, vtable, ftable, t_2)$ $op = \text{transop}(\text{getopname}(\text{relop}))$ $code_1 ++ code_2 ++ [\text{IF } t_1 \text{ op } t_2 \text{ THEN } label_t \text{ ELSE } label_f]$
true	$[\text{GOTO } label_t]$
false	$[\text{GOTO } label_f]$
$! Cond_1$	$Trans_{Cond}(Cond_1, label_f, label_t, vtable, ftable)$
$Cond_1 \ \&\& \ Cond_2$	$arg_2 = \text{newlabel}()$ $code_1 = Trans_{Cond}(Cond_1, arg_2, label_f, vtable, ftable)$ $code_2 = Trans_{Cond}(Cond_2, label_t, label_f, vtable, ftable)$ $code_1 ++ [\text{LABEL } arg_2] ++ code_2$
$Cond_1 \ \ Cond_2$	$arg_2 = \text{newlabel}()$ $code_1 = Trans_{Cond}(Cond_1, label_t, arg_2, vtable, ftable)$ $code_2 = Trans_{Cond}(Cond_2, label_t, label_f, vtable, ftable)$ $code_1 ++ [\text{LABEL } arg_2] ++ code_2$
Exp	$t = \text{newvar}()$ $code_1 = Trans_{Exp}(Exp, vtable, ftable, t)$ $code_1 ++ [\text{IF } t \neq 0 \text{ THEN } label_t \text{ ELSE } label_f]$

Figure 7.8: Translation of sequential logical operators

the value remains 0. We could equally well have done things the other way around, *i.e.*, first assign 1 to the target location and modify this to 0 when the condition is false.

It gets a bit more interesting in *Trans_{Cond}*, where we translate conditions. We have already seen how comparisons and expressions are translated, so we move directly to the new cases.

The constant `true` condition just generates a jump to the label for true conditions, and, similarly, `false` generates a jump to the label for false conditions.

Logical negation generates no code by itself, it just swaps the attribute-labels for true and false when translating its argument. This negates the effect of the argument condition.

Sequential logical **and** is translated as follows: The code for the first operand is translated such that if it is false, the second condition is not tested. This is done by jumping straight to the label for false conditions when the first operand is false. If the first operand is true, a jump to the code for the second operand is made. This is handled by using the appropriate labels as arguments to the call to *Trans_{Cond}*. The call to *Trans_{Cond}* for the second operand uses the original labels for true and false. Hence, both conditions have to be true for the combined condition to be true.

Sequential **or** is similar: If the first operand is true, we jump directly to the label for true conditions without testing the second operand, but if it is false, we jump to the code for the second operand. Again, the second operand uses the original labels for true and false.

Note that the translation functions now work even if **binop** and **unop** do not contain relational operators or logical negation, as we can just choose the last rule for expressions whenever the **binop** rules do not match. However, we can not in the same way omit non-sequential (*e.g.*, bitwise) **and** and **or**, as these have a different effect (*i.e.*, they always evaluate both arguments).

We have, in the above, used two different nonterminals for conditions and expressions, with some overlap between these and consequently ambiguity in the grammar. It is possible to resolve this ambiguity by rewriting the grammar and get two non-overlapping syntactic categories in the abstract syntax. Another solution is to join the two nonterminals into one, *e.g.*, *Exp* and use two different translation functions for this: Whenever an expression is translated, the translation function most appropriate for the context is chosen. For example, *if-then-else* will choose a translation function similar to *Trans_{Cond}* while assignment will choose a one similar to the current *Trans_{Exp}*.

Suggested exercises: 7.3.

7.8 Advanced control statements

We have, so far, shown translation of simple conditional statements and loops, but some languages have more advanced control features. We will briefly discuss how such can be implemented.

Goto and labels. Labels are stored in a symbol table that binds each to a corresponding label in the intermediate language. A jump to a label will generate a GOTO statement to the corresponding intermediate-language label. Unless labels are declared before use, an extra pass may be needed to build the symbol table before the actual translation. Alternatively, an intermediate-language label can be chosen and an entry in the symbol table be created at the first occurrence of the label even if it is in a jump rather than a declaration. Subsequent jumps or declarations of that label will use the intermediate-language label that was chosen at the first occurrence. By setting a mark in the symbol-table entry when the label is declared, it can be checked that all labels are declared exactly once.

The scope of labels can be controlled by the symbol table, so labels can be local to a procedure or block.

Break/exit. Some languages allow exiting loops from the middle of the loop-body by a `break` or `exit` statement. To handle these, the translation function for statements must have an extra inherited parameter which is the label that a `break` or `exit` statement must jump to. This attribute is changed whenever a new loop is entered. Before the first loop is entered, this attribute is undefined. The translation function should check for this, so it can report an error if a `break` or `exit` occurs outside loops. This should, rightly, be done during type-checking (see chapter 6), though.

C's `continue` statement, which jumps to the start of the current loop, can be handled similarly.

Case-statements. A case-statement evaluates an expression and selects one of several branches (statements) based on the value of the expression. In most languages, the case-statement will be exited at the end of each of these statements. In this case, the case-statement can be translated as an assignment that stores the value of the expression followed by a nested `if-then-else` statement, where each branch of the case-statement becomes a `then-branch` of one of the `if-then-else` statements (or, in case of the default branch, the final `else-branch`).

In C, the default is that *all* case-branches following the selected branch are executed unless the case-expression (called `switch` in C) is explicitly terminated with a `break` statement (see above) at the end of the branch. In this case, the case-statement can still be translated to a nested `if-then-else`, but the branches of

these are now GOTO's to the code for each case-branch. The code for the branches is placed in sequence after the nested if-then-else, with break handled by GOTO's as described above. Hence, if no explicit jump is made, one branch will fall through to the next.

7.9 Translating structured data

So far, the only values we have used are integers and booleans. However, most programming languages provide floating-point numbers and structured values like arrays, records (structs), unions, lists or tree-structures. We will now look at how these can be translated. We will first look at floats, then at one-dimensional arrays, multi-dimensional arrays and finally other data structures.

7.9.1 Floating-point values

Floating-point values are, in a computer, typically stored in a different set of registers than integers. Apart from this, they are treated the same way we treat integer values: We use temporary variables to store intermediate expression results and assume the intermediate language has binary operators for floating-point numbers. The register allocator will have to make sure that the temporary variables used for floating-point values are mapped to floating-point registers. For this reason, it may be a good idea to let the intermediate code indicate which temporary variables hold floats. This can be done by giving them special names or by using a symbol table to hold type information.

7.9.2 Arrays

We extend our example language with one-dimensional arrays by adding the following productions:

$$\begin{aligned} Exp &\rightarrow Index \\ Stat &\rightarrow Index := Exp \\ Index &\rightarrow \mathbf{id}[Exp] \end{aligned}$$

Index is an array element, which can be used the same way as a variable, either as an expression or as the left part of an assignment statement.

We will initially assume that arrays are zero-based (*i.e.* the lowest index is 0).

Arrays can be allocated statically, *i.e.*, at compile-time, or *dynamically*, *i.e.*, at run-time. In the first case, the *base address* of the array (the address at which index 0 is stored) is a compile-time constant. In the latter case, a variable will contain the base address of the array. In either case, we assume that the symbol table for variables binds an array name to the constant or variable that holds its base address.

$Trans_{Exp}(Exp, vtable, ftable, place) = \text{case } Exp \text{ of}$	
$Index$	$(code_1, address) = Trans_{Index}(Index, vtable, ftable)$ $code_1 ++ [place := M[address]]$

$Trans_{Stat}(Stat, vtable, ftable) = \text{case } Stat \text{ of}$	
$Index := Exp$	$(code_1, address) = Trans_{Index}(Index, vtable, ftable)$ $t = newvar()$ $code_2 = Trans_{Exp}(Exp, vtable, ftable, t)$ $code_1 ++ code_2 ++ [M[address] := t]$

$Trans_{Index}(Index, vtable, ftable) = \text{case } Index \text{ of}$	
$id[Exp]$	$base = lookup(vtable, getname(id))$ $t = newvar()$ $code_1 = Trans_{Exp}(Exp, vtable, ftable, t)$ $code_2 = code_1 ++ [t := t * 4, t := t + base]$ $(code_2, t)$

Figure 7.9: Translation for one-dimensional arrays

Most modern computers are byte-addressed, while integers typically are 32 or 64 bits long. This means that the index used to access array elements must be multiplied by the size of the elements (measured in bytes), *e.g.*, 4 or 8, to find the actual offset from the base address. In the translation shown in figure 7.9, we use 4 for the size of integers. We show only the new parts of the translation functions for *Exp* and *Stat*.

We use a translation function $Trans_{Index}$ for array elements. This returns a pair consisting of the code that evaluates the address of the array element and the variable that holds this address. When an array element is used in an expression, the contents of the address is transferred to the target variable using a memory-load instruction. When an array element is used on the left-hand side of an assignment, the right-hand side is evaluated, and the value of this is stored at the address using a memory-store instruction.

The address of an array element is calculated by multiplying the index by the size of the elements and adding this to the base address of the array. Note that *base* can be either a variable or a constant (depending on how the array is allocated, see below), but since both are allowed as the second operator to a **binop** in the intermediate language, this is no problem.

Allocating arrays

So far, we have only hinted at how arrays are allocated. As mentioned, one possibility is static allocation, where the base-address and the size of the array are

known at compile-time. The compiler, typically, has a large address space where it can allocate statically allocated objects. When it does so, the new object is simply allocated after the end of the previously allocated objects.

Dynamic allocation can be done in several ways. One is allocation local to a procedure or function, such that the array is allocated when the function is entered and deallocated when it is exited. This typically means that the array is allocated on a stack and popped from the stack when the procedure is exited. If the sizes of locally allocated arrays are fixed at compile-time, their base addresses are constant offsets from the stack top (or from the *frame pointer*, see chapter 10) and can be calculated from this at every array-lookup. However, this does not work if the sizes of these arrays are given at run-time. In this case, we need to use a variable to hold the base address of each array. The address is calculated when the array is allocated and then stored in the corresponding variable. This can subsequently be used as described in *Trans_{Index}* above. At compile-time, the array-name will in the symbol table be bound to the variable that at runtime will hold the base-address.

Dynamic allocation can also be done globally, so the array will survive until the end of the program or until it is explicitly deallocated. In this case, there must be a global address space available for run-time allocation. Often, this is handled by the operating system which handles memory-allocation requests from all programs that are running at any given time. Such allocation may fail due to lack of memory, in which case the program must terminate with an error or release memory enough elsewhere to make room. The allocation can also be controlled by the program itself, which initially asks the operating system for a large amount of memory and then administrates this itself. This can make allocation of arrays faster than if an operating system call is needed every time an array is allocated. Furthermore, it can allow the program to use *garbage collection* to automatically reclaim arrays that are no longer in use.

These different allocation techniques are described in more detail in chapter 12.

Multi-dimensional arrays

Multi-dimensional arrays can be laid out in memory in two ways: *row-major* and *column-major*. The difference is best illustrated by two-dimensional arrays, as shown in Figure 7.10. A two-dimensional array is addressed by two indices, *e.g.*, (using C-style notation) as $a[i][j]$. The first index, i , indicates the *row* of the element and the second index, j , indicates the *column*. The first row of the array is, hence, the elements $a[0][0]$, $a[0][1]$, $a[0][2]$, ... and the first column is $a[0][0]$, $a[1][0]$, $a[2][0]$, ...³

In row-major form, the array is laid out one row at a time and in column-major

³Note that the coordinate system, following computer-science tradition, is rotated 90° clockwise compared to mathematical tradition.

	1st column	2nd column	3rd column	...
1st row	a[0][0]	a[0][1]	a[0][2]	...
2nd row	a[1][0]	a[1][1]	a[1][2]	...
3rd row	a[2][0]	a[2][1]	a[2][2]	...
⋮	⋮	⋮	⋮	⋮

Figure 7.10: A two-dimensional array

form it is laid out one column at a time. In a 3×2 array, the ordering for row-major is

$$a[0][0], a[0][1], a[1][0], a[1][1], a[2][0], a[2][1]$$

For column-major the ordering is

$$a[0][0], a[1][0], a[2][0], a[0][1], a[1][1], a[2][1]$$

If the size of an element is *size* and the sizes of the dimensions in an *n*-dimensional array are $dim_0, dim_1, \dots, dim_{n-2}, dim_{n-1}$, then in row-major format an element at index $[i_0][i_1] \dots [i_{n-2}][i_{n-1}]$ has the address

$$base + ((\dots (i_0 * dim_1 + i_1) * dim_2 \dots + i_{n-2}) * dim_{n-1} + i_{n-1}) * size$$

In column-major format the address is

$$base + ((\dots (i_{n-1} * dim_{n-2} + i_{n-2}) * dim_{n-3} \dots + i_1) * dim_0 + i_0) * size$$

Note that column-major format corresponds to reversing the order of the indices of a row-major array. *i.e.*, replacing i_0 and dim_0 by i_{n-1} and dim_{n-1} , i_1 and dim_1 by i_{n-2} and dim_{n-2} , and so on.

We extend the grammar for array-elements to accommodate multi-dimensional arrays:

$$\begin{aligned} Index &\rightarrow \mathbf{id}[Exp] \\ Index &\rightarrow Index[Exp] \end{aligned}$$

and extend the translation functions as shown in figure 7.11. This translation is for row-major arrays. We leave column-major arrays as an exercise.

With these extensions, the symbol table must return both the base-address of the array and a list of the sizes of the dimensions. Like the base-address, the dimension sizes can either be compile-time constants or variables that at run-time will hold the sizes. We use an auxiliary translation function $Calc_{Index}$ to calculate the position of

$Trans_{Exp}(Exp, vtable, ftable, place) = \text{case } Exp \text{ of}$	
$Index$	$(code_1, address) = Trans_{Index}(Index, vtable, ftable)$ $code_1 ++ [place := M[address]]$

$Trans_{Stat}(Stat, vtable, ftable) = \text{case } Stat \text{ of}$	
$Index := Exp$	$(code_1, address) = Trans_{Index}(Index, vtable, ftable)$ $t = \text{newvar}()$ $code_2 = Trans_{Exp}(Exp, vtable, ftable, t)$ $code_1 ++ code_2 ++ [M[address] := t]$

$Trans_{Index}(Index, vtable, ftable) =$	
$(code_1, t, base, []) = Calc_{Index}(Index, vtable, ftable)$ $code_2 = code_1 ++ [t := t * 4, t := t + base]$ $(code_2, t)$	

$Calc_{Index}(Index, vtable, ftable) = \text{case } Index \text{ of}$	
$id[Exp]$	$(base, dims) = \text{lookup}(vtable, \text{getname}(id))$ $t = \text{newvar}()$ $code = Trans_{Exp}(Exp, vtable, ftable, t)$ $(code, t, base, \text{tail}(dims))$
$Index[Exp]$	$(code_1, t_1, base, dims) = Calc_{Index}(Index, vtable, ftable)$ $dim_1 = \text{head}(dims)$ $t_2 = \text{newvar}()$ $code_2 = Trans_{Exp}(Exp, vtable, ftable, t_2)$ $code_3 = code_1 ++ code_2 ++ [t_1 := t_1 * dim_1, t_1 := t_1 + t_2]$ $(code_3, t_1, base, \text{tail}(dims))$

Figure 7.11: Translation of multi-dimensional arrays

an element. In $Trans_{Index}$ we multiply this position by the element size and add the base address. As before, we assume the size of elements is 4.

In some cases, the sizes of the dimensions of an array are not stored in separate variables, but in memory next to the space allocated for the elements of the array. This uses fewer variables (which may be an issue when these need to be allocated to registers, see chapter 9) and makes it easier to return an array as the result of an expression or function, as only the base-address needs to be returned. The size information is normally stored just before the base-address so, for example, the size of the first dimension can be at address $base - 4$, the size of the second dimension at $base - 8$ and so on. Hence, the base-address will always point to the first element of the array no matter how many dimensions the array has. If this strategy is used, the necessary dimension-sizes must be loaded into variables when an index is calculated. Since this adds several extra (somewhat costly) loads, optimising compilers often try to re-use the values of previous loads, *e.g.*, by doing the loading once outside a loop and referring to variables holding the values inside the loop.

Index checks

The translations shown so far do not test if an index is within the bounds of the array. Index checks are fairly easy to generate: Each index must be compared to the size of (the dimension of) the array and if the index is too big, a jump to some error-producing code is made. If the comparison is made on unsigned numbers, a negative index will look like a very large index. Hence, a single conditional jump is inserted at every index calculation.

This is still fairly expensive, but various methods can be used to eliminate some of these tests. For example, if the array-lookup occurs within a `for`-loop, the bounds of the loop-counter may guarantee that array accesses using this variable will be within bounds. In general, it is possible to make an analysis that finds cases where the index-check condition is subsumed by previous tests, such as the exit test for a loop, the test in an `if-then-else` statement or previous index checks. See section 11.4 for details.

Non-zero-based arrays

We have assumed our arrays to be zero-based, *i.e.*, that the indices start from 0. Some languages allow indices to be arbitrary intervals, *e.g.*, -10 to 10 or 10 to 20 . If such are used, the starting-index must be subtracted from each index when the address is calculated. In a one-dimensional array with known size and base-address, the starting-index can be subtracted (at compile-time) from base-address instead. In a multi-dimensional array with known dimensions, the starting-indices can be multiplied by the sizes of the dimensions and added together to form a

single constant that is subtracted from the base-address instead of subtracting each starting-index from each index.

7.9.3 Strings

Strings are usually implemented in a fashion similar to one-dimensional arrays. In some languages (*e.g.* C or pre-ISO-standard Pascal), strings *are* just arrays of characters.

However, strings often differ from arrays in that the length is not static, but can vary at run-time. This leads to an implementation similar to the kind of arrays where the length is stored in memory, as explained in section 7.9.2. Another difference is that the size of a character is typically one byte (unless 16-bit Unicode characters are used), so the index calculation does not multiply the index by the size (as this is 1).

Operations on strings, *e.g.*, concatenation and substring extraction, are typically implemented by calling library functions.

7.9.4 Records/structs and unions

Records (structs) have many properties in common with arrays. They are typically allocated in a similar way (with a similar choice of possible allocation strategies), and the fields of a record are typically accessed by adding an offset to the base-address of the record. The differences are:

- The types (and hence sizes) of the fields may be different.
- The field-selector is known at compile-time, so the offset from the base address can be calculated at this time.

The offset for a field is simply the sum of the sizes of all fields that occur before it. For a record-variable, the symbol table for variables must hold the base-address and the offsets for each field in the record. The symbol table for types must hold the offsets for every record type, such that these can be inserted into the symbol table for variables when a record of this type is declared.

In a union (sum) type, the fields are not consecutive, but are stored at the same address, *i.e.*, the base-address of the union. The size of an union is the maximum of the sizes of its fields.

In some languages, union types include a *tag*, which identifies which variant of the union is stored in the variable. This tag is stored as a separate field before the union-fields. Some languages (*e.g.* Standard ML) enforce that the tag is tested when the union is accessed, others (*e.g.* Pascal) leave this as an option to the programmer.

Suggested exercises: 7.8.

7.10 Translating declarations

In the translation functions used in this chapter, we have several times required that “The symbol table must contain ...”. It is the job of the compiler to ensure that the symbol tables contain the information necessary for translation. When a name (variable, label, type, *etc.*) is declared, the compiler must keep in the symbol-table entry for that name the information necessary for compiling any use of that name. For scalar variables (*e.g.*, integers), the required information is the intermediate-language variable that holds the value of the variable. For array variables, the information includes the base-address and dimensions of the array. For records, it is the offsets for each field and the total size. If a type is given a name, the symbol table must for that name provide a description of the type, such that variables that are declared to be that type can be given the information they need for their own symbol-table entries.

The exact nature of the information that is put into the symbol tables will depend on the translation functions that use these tables, so it is usually a good idea to write first the translation functions for *uses* of names and then translation functions for their declarations.

Translation of function declarations will be treated in chapter 10.

7.10.1 Example: Simple local declarations

We extend the statement language by the following productions:

$$\begin{aligned} Stat &\rightarrow Decl ; Stat \\ Decl &\rightarrow \text{int } \mathbf{id} \\ Decl &\rightarrow \text{int } \mathbf{id}[\mathbf{num}] \end{aligned}$$

We can, hence, declare integer variables and one-dimensional integer arrays for use in the following statement. An integer variable should be bound to a location in the symbol table, so this declaration should add such a binding to *vtable*. An array should be bound to a variable containing its base address. Furthermore, code must be generated for allocating space for the array. We assume arrays are heap allocated and that the intermediate-code variable *HP* points to the first free element of the (upwards growing) heap. Figure 7.12 shows the translation of these declarations. When allocating arrays, no check for heap overflow is done.

7.11 Further reading

A comprehensive discussion about intermediate languages can be found in [35].

$Trans_{Stat}(Stat, vtable, ftable) = \text{case } Stat \text{ of}$	
$Decl ; Stat_1$	$(code_1, vtable_1) = Trans_{Decl}(Decl, vtable)$ $code_2 = Trans_{Stat}(Stat_1, vtable_1, ftable)$ $code_1 ++ code_2$
$Trans_{Decl}(Decl, vtable) = \text{case } Decl \text{ of}$	
$\text{int } id$	$t_1 = newvar()$ $vtable_1 = bind(vtable, getname(id), t_1)$ $([], vtable_1)$
$\text{int } id[num]$	$t_1 = newvar()$ $vtable_1 = bind(vtable, getname(id), t_1)$ $([t_1 := HP, HP := HP + (4 * getvalue(num))], vtable_1)$

Figure 7.12: Translation of simple declarations

Functional and logic languages often use high-level intermediate languages, which are in many cases translated to lower-level intermediate code before emitting actual machine code. Examples of such intermediate languages can be found in [23], [8] and [6].

Another high-level intermediate language is the Java Virtual Machine [29]. This language has single instructions for such complex things as calling virtual methods and creating new objects. The high-level nature of JVM was chosen for several reasons:

- By letting common complex operations be done by single instructions, the code is smaller, which reduces transmission time when sending the code over the Internet.
- JVM was originally intended for interpretation, and the complex operations also helped reduce the overhead of interpretation.
- A program in JVM is *validated* (essentially type-checked) before interpretation or further translation. This is easier when the code is high-level.

Exercises

Exercise 7.1

Use the translation functions in figure 7.3 to generate code for the expression $2+g(x+y, x*y)$. Use a *vtable* that binds x to $v0$ and y to $v1$ and an *ftable* that binds g to $_g$. The result of the expression should be put in the intermediate-code variable r (so the *place* attribute in the initial call to $Trans_{Exp}$ is r).

Exercise 7.2

Use the translation functions in figures 7.5 and 7.6 to generate code for the statement

```
x:=2+y;
if x<y then x:=x+y;
repeat
  y:=y*2;
  while x>10 do x:=x/2
until x<y
```

use the same *vtable* as in exercise 7.1.

Exercise 7.3

Use the translation functions in figures 7.5 and 7.8 to translate the following statement

```
if x<=y && !(x=y || x=1)
then x:=3
else x:=5
```

use the same *vtable* as in exercise 7.1.

Exercise 7.4

De Morgan's law tells us that $!(p \mid\mid q)$ is equivalent to $(!p) \ \&\& \ (!q)$. Show that these generate identical code when compiled with *Trans_{Cond}* from figure 7.8.

Exercise 7.5

Show that, in any code generated by the functions in figures 7.5 and 7.8, every IF-THEN-ELSE instruction will be followed by one of the target labels.

Exercise 7.6

Extend figure 7.5 to include a break-statement for exiting loops, as described in section 7.8, *i.e.*, extend the statement syntax by

$$Stat \rightarrow \text{break}$$

and add a rule for this to *Trans_{Stat}*. Add whatever extra attributes you may need to do this.

Exercise 7.7

We extend the statement language with the following statements:

$$\begin{aligned} Stat &\rightarrow \text{labelid} : \\ Stat &\rightarrow \text{goto labelid} \end{aligned}$$

for defining and jumping to labels.

Extend figure 7.5 to handle these as described in section 7.8. Labels have scope over the entire program (statement) and need not be defined before use. You can assume that there is exactly one definition for each used label.

Exercise 7.8

Show translation functions for multi-dimensional arrays in column-major format.

Hint: Starting from figure 7.11, it may be a good idea to rewrite the productions for *Index* so they are right-recursive instead of left-recursive, as the address formula for column-major arrays groups to the right. Similarly, it is a good idea to reverse the list of dimension sizes, so the size of the rightmost dimension comes first in the list.

Exercise 7.9

When statements are translated using the functions in figure 7.5, it will often be the case that the statement immediately following a label is a GOTO statement, *i.e.*, we have the following situation:

$$\begin{aligned} \text{LABEL } &label_1 \\ \text{GOTO } &label_2 \end{aligned}$$

It is clear that any jump to *label*₁ can be replaced by a jump to *label*₂, and that this will result in faster code. Hence, it is desirable to do so. This is called jump-to-jump optimisation, and can be done after code-generation by a post-process that looks for these situations. However, it is also possible to avoid most of these situations by modifying the translation function.

This can be done by adding an extra inherited attribute *endlabel*, which holds the name of a label that can be used as the target of a jump to the end of the code that is being translated. If the code is immediately followed by a GOTO statement, *endlabel* will hold the target of this GOTO rather than a label immediately preceding this.

- a) Add the *endlabel* attribute to *Trans*_{Stat} from figure 7.5 and modify the rules so *endlabel* is exploited for jump-to-jump optimisation. Remember to set *endlabel* correctly in recursive calls to *Trans*_{Stat}.

b) Use the modified $Trans_{Stat}$ to translate the following statement:

```
while x>0 do {
  x := x-1;
  if x>10 then x := x/2
}
```

The curly braces are used as disambiguators, though they are not part of grammar 7.4.

Use the same $vtable$ as exercise 7.1 and use `endlab` as the *endlabel* for the whole statement.

Exercise 7.10

In figure 7.5, `while` statements are translated in such a way that every iteration of the loop executes an unconditional jump (`GOTO` in addition to the conditional jumps in the loop condition.

Modify the translation so each iteration only executes the conditional jumps in the loop condition, *i.e.*, so an unconditional jump is saved in every iteration. You may have to add an unconditional jump outside the loop.

Exercise 7.11

Logical conjunction is associative: $p \wedge (q \wedge r) \Leftrightarrow (p \wedge q) \wedge r$.

Show that this also applies to the sequential conjunction operator `&&` when translated as in figure 7.8, *i.e.*, that $p \&\&(q \&\&r)$ generates the same code (up to renaming of labels) as $(p \&\&q) \&\&r$.

Exercise 7.12

Figure 7.11 shows translation of multi-dimensional arrays in row-major layout, where the address of each element is found through multiplication and addition. On machines with fast memory access but slow multiplication, an alternative implementation of multi-dimensional arrays is sometimes used: An array with dimensions $dim_0, dim_1, \dots, dim_n$ is implemented as a one-dimensional array of size dim_0 with pointers to dim_0 different arrays each of dimension dim_1, \dots, dim_n , which again are implemented in the same way (until the last dimension, which is implemented as a normal one-dimensional array of values). This takes up more room, as the pointer arrays need to be stored as well as the elements. But array-lookup can be done using only addition and memory accesses.

- a) Assuming pointers and array elements need four bytes each, what is the total number of bytes required to store an array of dimensions $dim_0, dim_1, \dots, dim_n$?
- b) Write translation functions for array-access in the style of figure 7.11 using this representation of arrays. Use addition to multiply numbers by 4 for scaling indices by the size of pointers and array elements.

Chapter 8

Machine-Code Generation

8.1 Introduction

The intermediate language we have used in chapter 7 is quite low-level and similar to the type of machine code you can find on modern RISC processors, with a few exceptions:

- We have used an unbounded number of variables, where a processor will have a bounded number of registers.
- We have used a complex CALL instruction for function calls.
- In the intermediate language, the IF-THEN-ELSE instruction has two target labels, where, on most processors, the conditional jump instruction has only one target label, and simply falls through to the next instruction when the condition is false.
- We have assumed that any constant can be an operand to an arithmetic instruction. Typically, RISC processors allow only small constants as operands.

The problem of mapping a large set of variables to a small number of registers is handled by *register allocation*, as explained in chapter 9. Function calls are treated in chapter 10. We will look at the remaining two problems below.

The simplest solution for generating machine code from intermediate code is to translate each intermediate-language instruction into one or more machine-code instructions. However, it is often possible to find a machine-code instruction that covers two or more intermediate-language instructions. We will in section 8.4 see how we can exploit complex instructions in this way.

Additionally, we will briefly discuss other optimisations.

8.2 Conditional jumps

Conditional jumps come in many forms on different machines. Some conditional jump instructions embody a relational comparison between two registers (or a register and a constant) and are, hence, similar to the IF-THEN-ELSE instruction in our intermediate language. Other types of conditional jump instructions require the condition to be already resolved and stored in special condition registers or flags. However, it is almost universal that conditional jump instructions specify only one target label (or address), typically used when the condition is true. When the condition is false, execution simply continues with the instructions immediately following the conditional jump instruction.

Converting two-way branches to one-way branches is not terribly difficult: IF c THEN l_t ELSE l_f can be translated to

```
branch_if_c   $l_t$ 
jump         $l_f$ 
```

where `branch_if_c` is a conditional instruction that jumps when the condition c is true and `jump` is an unconditional jump.

Often, an IF-THEN-ELSE instruction is immediately followed by one of its target labels. In fact, this will always be the case if the intermediate code is generated by the translation functions shown in chapter 7 (see exercise 7.5). If this label happens to be l_f (the label taken for false conditions), we can simply omit the unconditional jump from the code shown above. If the following label is l_t , we can negate the condition of the conditional jump and make it jump to l_f , *i.e.*, as

```
branch_if_not_c   $l_f$ 
```

where `branch_if_not_c` is a conditional instruction that jumps when the condition c is false.

Hence, the code generator (the part of the compiler that generates machine code) should test which (if any) of the target labels follow an IF-THEN-ELSE instruction and translate it accordingly. Alternatively, a post-processing pass can be made over the generated machine code to remove superfluous jumps.

If the conditional jump instructions in the target machine language do not allow conditions as complex as those used in the intermediate language, code must be generated to first calculate the condition and put the result somewhere where it can be tested by a subsequent conditional jump instruction. In some machine architectures, *e.g.*, MIPS and Alpha, this “somewhere” can be a general-purpose register. Other machines, *e.g.*, PowerPC or IA-64 (also known as Itanium) use special condition registers, while yet others, *e.g.*, IA-32 (also known as x86), Sparc,

PA-RISC and ARM use a single set of arithmetic flags that can be set by comparison or arithmetic instructions. A conditional jump may test various combinations of the flags, so the same comparison instruction can, depending on the subsequent condition, be used for testing equality, signed or unsigned less-than, overflow and several other properties. Usually, any IF-THEN-ELSE instruction can be translated into at most two instructions: One that does the comparison and one that does the conditional jump.

8.3 Constants

The intermediate language allows arbitrary constants as operands to binary or unary operators. This is not always the case in machine code.

For example, MIPS allows only 16-bit constants in operands even though integers are 32 bits (64 bits in some versions of the MIPS architecture). To build larger constants, MIPS includes instructions to load 16-bit constants into the upper half (the most significant bits) of a register. With help of these, an arbitrary 32-bit integer can be entered into a register using two instructions. On the ARM, a constant can be an 8-bit number positioned at any even bit boundary. It may take up to four instructions to build a 32-bit number using these.

When an intermediate-language instruction uses a constant, the code generator must check if it fits into the constant field (if any) of the equivalent machine-code instruction. If it does, the code generator generates a single machine-code instruction. If not, the code generator generates a sequence of instructions that builds the constant in a register, followed by an instruction that uses this register in place of the constant. If a complex constant is used inside a loop, it may be a good idea to move the code for generating this outside the loop and keep it in a register inside the loop. This can be done as part of a general optimisation to move code out of loops, see section 8.5.

8.4 Exploiting complex instructions

Most instructions in our intermediate language are *atomic*, in the sense that each instruction corresponds to a single operation which can not sensibly be split into smaller steps. The exceptions to this rule are the instructions IF-THEN-ELSE, which we in section 8.2 described how to handle, and CALL, which will be detailed in chapter 10.

CISC (Complex Instruction Set Computer) processors like IA-32 have composite (*i.e.*, non-atomic) instructions in abundance. And while the philosophy behind RISC (Reduced Instruction Set Computer) processors like MIPS and ARM advocates that machine-code instructions should be simple, most RISC processors

include at least a few non-atomic instructions, typically for memory-access instructions.

We will in this chapter use a subset of the MIPS instruction set as an example. A description of the MIPS instruction set can be found Appendix A of [39], which is available online [27]. If you are not already familiar with the MIPS instruction set, it would be a good idea to read the description before continuing.

To exploit composite instructions, several intermediate-language instructions can be grouped together and translated into a single machine-code instruction. For example, the intermediate-language instruction sequence

$$\begin{aligned} t_2 &:= t_1 + 116 \\ t_3 &:= M[t_2] \end{aligned}$$

can be translated into the single MIPS instruction

lw r3, 116(r1)

where r1 and r3 are the registers chosen for t_1 and t_3 , respectively. However, it is only possible to combine the two instructions if the value of the intermediate value t_2 is not required later, as the combined instruction does not store this value anywhere.

We will, hence, need to know if the contents of a variable is required for later use, or if it is *dead* after a particular use. When generating intermediate code, most of the temporary variables introduced by the compiler will be single-use and can be marked as such. Any use of a single-use variable will, by definition, be the last use. Alternatively, last-use information can be obtained by analysing the intermediate code using a *liveness analysis*, which we will describe in chapter 9. For now, we will just assume that the last use of any variable is marked in the intermediate code. We assume this is done, and the last use of any variable in the intermediate code is marked by *last*, such as t^{last} , which indicates that this is the last use of the variable t .

Our next step is to describe each machine-code instruction in terms of one or more intermediate-language instructions. We call the sequence of intermediate-language instructions a *pattern*, and the corresponding machine-code instruction its *replacement*, since the idea is to find sequences in the intermediate code that matches the pattern and replace these sequences by instances of the replacement. When a pattern uses variables such as k , t or r_d , these can match any intermediate-language constants, variables or labels, and when the same variable is used in both pattern and replacement, it means that the corresponding intermediate-language constant or variable/label name is copied to the machine-code instruction, where it will represent a constant, a named register or a machine-code label.

For example, the MIPS lw (load word) instruction can be described by the pattern/replacement pair

$t := r_s + k$ $r_t := M[t^{last}]$	$\text{lw } r_t, k(r_s)$
--	--------------------------

where t^{last} in the pattern indicates that the contents of t must not be used afterwards, *i.e.*, that the intermediate-language variable that is matched against t must have a *last* annotation at this place. A pattern can only match a piece of intermediate code if all *last* annotations in the pattern are matched by *last* annotations in the intermediate code. The converse, however, need not hold: It is not harmful to store a value in a register even if it is not used later, so a *last* annotation in the intermediate code need not be matched by a *last* annotation in the pattern.

The list of patterns that in combination describe the machine-code instruction set must cover the intermediate language in full (excluding function calls, which we handle in chapter 10). In particular, each single intermediate-language instruction (with the exception of CALL, which we handle separately in chapter 10) must be covered by at least one pattern. This means that we must include the MIPS instruction $\text{lw } r_t, 0(r_s)$ to cover the intermediate-code instruction $r_t := M[r_s]$, even though we have already listed a more general form of lw . If there is an intermediate-language instruction for which there are no equivalent single machine-code instruction, a sequence of machine-code instructions must be given for this. Hence, an instruction-set description is a list of pairs, where each pair consists of a *pattern* (a sequence of intermediate-language instructions) and a *replacement* (a sequence of machine-code instructions).

When translating a sequence of intermediate-code instructions, the code generator can look at the patterns and pick the replacement that covers the largest prefix of the intermediate code. A simple way of ensuring that the longest prefix is matched is to list the pairs so longer patterns are listed before shorter patterns. The first pattern in the list that matches a prefix of the intermediate code will now also be the longest matching pattern.

This kind of algorithm is called *greedy*, because it always picks the choice that is best for immediate profit, *i.e.*, the sequence that “eats” most of the intermediate code in one bite. It will, however, not always yield the best possible solution for the total sequence of intermediate-language instructions.

If costs are given for each machine-code instruction sequence in the pattern/replacement pairs, optimal (*i.e.*, least-cost) solutions can be found for straight-line (*i.e.*, jump-free) code sequences. The least-cost sequence that covers the intermediate code can be found, *e.g.*, using a dynamic-programming algorithm. For RISC processors, a greedy algorithm will typically get close to optimal solutions, so the gain from using a better algorithm is small. Hence, we will go into detail only for the greedy algorithm.

As an example, figure 8.1 describes a subset of the instructions for the MIPS

microprocessor architecture in terms of the intermediate language as a set of pattern/replacement pairs. Note that we exploit the fact that register 0 is hardwired to be the value 0 to, *e.g.*, use the `addi` instruction to generate a constant. We assume that we, at this point, have already handled the problem of too-large constants, so any constant that now remains in the intermediate code can be used as an immediate constant in an instruction such as `addi`. Note that we make special cases for IF-THEN-ELSE when one of the labels immediately follows the test. Note, also, that we need (at least) two instructions from our MIPS subset to implement an IF-THEN-ELSE instruction that uses `<` as the relational operator, while we need only one for comparison by `=`. Figure 8.1 does not cover all of the intermediate language, but it can fairly easily be extended to do so. It is also possible to add more special cases to exploit a larger subset of the MIPS instruction set.

The instructions in figure 8.1 are listed so that, when two patterns overlap, the longest of these is listed first. Overlap can happen if the pattern in one pair is a prefix of the pattern for another pair, as is the case with the pairs involving `addi` and `lw/sw` and for the different instances of `beq/bne` and `slt`.

We can try to use figure 8.1 to select MIPS instructions for the following sequence of intermediate-language instructions:

```

a := a + blast
d := c + 8
M[dlast] := a
IF a = c THEN label1 ELSE label2
LABEL label2

```

Only one pattern (for the `add` instruction) in figure 8.1 matches a prefix of this code, so we generate an `add` instruction for the first intermediate instruction. We now have two matches for prefixes of the remaining code: One using `sw` and one using `addi`. Since the pattern using `sw` is listed first in the table, we choose this to replace the next two intermediate-language instructions. Finally, a `beq` instruction matches the last two instructions. Hence, we generate the code

```

add  a, a, b
sw   a, 8(c)
beq  a, c, label1
label2 :

```

Note that we retain `label2` even though the resulting sequence does not refer to it, as some other part of the code might jump to it. We could include single-use annotations for labels like we use for variables, but it is hardly worth the effort, as labels do not generate actual code and hence cost nothing¹.

¹This is, strictly speaking, not entirely true, as superfluous labels might inhibit later optimisations.

$t := r_s + k,$ $r_t := M[t^{last}]$	lw $r_t, k(r_s)$
$r_t := M[r_s]$	lw $r_t, 0(r_s)$
$r_t := M[k]$	lw $r_t, k(R0)$
$t := r_s + k,$ $M[t^{last}] := r_t$	sw $r_t, k(r_s)$
$M[r_s] := r_t$	sw $r_t, 0(r_s)$
$M[k] := r_t$	sw $r_t, k(R0)$
$r_d := r_s + r_t$	add r_d, r_s, r_t
$r_d := r_t$	add $r_d, R0, r_t$
$r_d := r_s + k$	addi r_d, r_s, k
$r_d := k$	addi $r_d, R0, k$
GOTO $label$	j $label$
IF $r_s = r_t$ THEN $label_t$ ELSE $label_f$, LABEL $label_f$	beq $r_s, r_t, label_t$ $label_f:$
IF $r_s = r_t$ THEN $label_t$ ELSE $label_f$, LABEL $label_t$	bne $r_s, r_t, label_f$ $label_t:$
IF $r_s = r_t$ THEN $label_t$ ELSE $label_f$	beq $r_s, r_t, label_t$ j $label_f$
IF $r_s < r_t$ THEN $label_t$ ELSE $label_f$, LABEL $label_f$	slt r_d, r_s, r_t bne $r_d, R0, label_t$ $label_f:$
IF $r_s < r_t$ THEN $label_t$ ELSE $label_f$, LABEL $label_t$	slt r_d, r_s, r_t beq $r_d, R0, label_f$ $label_t:$
IF $r_s < r_t$ THEN $label_t$ ELSE $label_f$	slt r_d, r_s, r_t bne $r_d, R0, label_t$ j $label_f$
LABEL $label$	$label:$

Figure 8.1: Pattern/replacement pairs for a subset of the MIPS instruction set

8.4.1 Two-address instructions

In the above we have assumed that the machine code is three-address code, *i.e.*, that the destination register of an instruction can be distinct from the two operand registers. It is, however, not uncommon that processors use two-address code, where the destination register is the same as the first operand register. To handle this, we use pattern/replacement pairs like these:

$r_t := r_s$	mov	r_t, r_s
$r_t := r_t + r_s$	add	r_t, r_s
$r_d := r_s + r_t$	move	r_d, r_s
	add	r_d, r_t

that add copy instructions in the cases where the destination register is not the same as the first operand. As we will see in chapter 9, the register allocator will often be able to remove the added copy instruction by allocating r_d and r_s in the same register.

Processors that divide registers into data and address registers or integer and floating-point registers can be handled in a similar way: Add instructions that copy to new registers before operations and let register allocation allocate these to the right type of registers (and eliminate as many of the moves as possible).

Suggested exercises: 8.2.

8.5 Optimisations

Optimisations can be done by a compiler in three places: In the source code (*i.e.*, on the abstract syntax), in the intermediate code, and in the machine code. Some optimisations can be specific to the source language or the machine language, but it makes sense to perform optimisations mainly in the intermediate language, as the optimisations hence can be shared among all compilers that use the same intermediate language. Also, the intermediate language is typically simpler than both the source language and the machine language, making the effort of doing optimisations smaller.

Optimising compilers have a wide array of optimisations that they can employ, but we will mention only a few and just hint at how they can be implemented.

Common subexpression elimination. In the statement $a[i] := a[i] + 2$, the address for $a[i]$ is calculated twice. This double calculation can be eliminated by storing the address in a temporary variable when the address is first calculated, and then use this variable instead of calculating the address again. Simple methods for common subexpression elimination work on *basic blocks*, *i.e.*, straight-line

code without jumps or labels, but more advanced methods can eliminate duplicated calculations even across jumps.

Code hoisting. If part of the computation inside a loop is independent of the variables that change inside the loop, it can be moved outside the loop and only calculated once. For example, in the loop

```
while (j < k) {
    sum = sum + a[i][j];
    j++;
}
```

a large part of the address calculation for `a[i][j]` can be done without knowing `j`. This part can be moved outside the loop so it will only be calculated once. Note that this optimisation can not be done on source-code level, as the address calculations are not visible there. For the same reason, the optimised version is not shown here.

If `k` may be less than or equal to `j`, the loop body may never be entered and we may, hence, unnecessarily execute the code that was moved out of the loop. This might even generate a run-time error. Hence, we can unroll the loop once to

```
if (j < k) {
    sum = sum + a[i][j];
    j++;
    while (j < k) {
        sum = sum + a[i][j];
        j++;
    }
}
```

The loop-independent part(s) may now without risk be calculated in the unrolled part and reused in the non-unrolled part. Again, this optimisation is not shown.

Constant propagation. A variable may, at some points in the program, have a value that is always equal to a known constant. When such a variable is used in a calculation, this calculation can often be simplified after replacing the variable by the constant that is guaranteed to be its value. Furthermore, the variable that holds the results of this computation may now also become constant, which may enable even more compile-time reduction.

Constant-propagation algorithms first trace the flow of constant values through the program, and then reduce calculations. More advanced methods also look at conditions, so they can exploit that after a test on, *e.g.*, `x = 0`, `x` is, indeed, the constant 0.

Index-check elimination. As mentioned in chapter 7, some compilers insert run-time checks to catch cases when an index is outside the bounds of the array. Some of these checks can be removed by the compiler. One way of doing this is to see if the tests on the index are subsumed by earlier tests or ensured by assignments. For example, assume that, in the loop shown above, a is declared to be a $k \times k$ array. This means that the entry test for the loop will ensure that j is always less than the upper bound on the array, so this part of the index test can be eliminated. If j is initialised to 0 before entering the loop, we can use this to conclude that we do not need to check the lower bound either.

8.6 Further reading

Code selection by pattern matching normally uses a tree-structured intermediate language instead of the linear instruction sequences we use in this book. This can avoid some problems where the order of unrelated instructions affect the quality of code generation. For example, if the two first instructions in the example at the end of section 8.4 are interchanged, our simple prefix-matching algorithm will not include the address calculation in the `sw` instruction and, hence, needs one more instruction. If the intermediate code is tree-structured, the order of independent instructions is left unspecified, and the code generator can choose whichever ordering gives the best code. See [35] or [9] for more details.

Descriptions of and methods for a large number of different optimisations can be found in [5], [35] and [9].

The instruction set of (one version of) the MIPS microprocessor architecture is described in [39]. This description is also available online [27].

Chapter 11 describes optimisation in more detail.

Exercises

Exercise 8.1

Add extra inherited attributes to $Trans_{Cond}$ in figure 7.8 that, for each of the two target labels, indicates if this label immediately follows the code for the condition, *i.e.*, a boolean-valued attribute for each of the two labels. Use this information to make sure that the false-destination labels of an IF-THEN-ELSE instruction follow immediately after the IF-THEN-ELSE instruction.

You can use the function *negate* to negate relational operators so, *i.e.*, $negate(<) = \geq$.

Make sure the new attributes are maintained in recursive calls and modify $Trans_{Stat}$ in figure 7.5 so it sets these attributes when calling $Trans_{Cond}$.

Exercise 8.2

Use figure 8.1 and the method described in section 8.4 to generate code for the following intermediate code sequence:

$$\begin{aligned} d &:= c + 8 \\ &:= a + b^{last} \\ M[d^{last}] &:= a \\ \text{IF } a < c &\text{ THEN } label_1 \text{ ELSE } label_2 \\ \text{LABEL } &label_1 \end{aligned}$$

Compare this to the example in section 8.4.

Exercise 8.3

In figures 7.3 and 7.5, identify guaranteed last-uses of temporary variables, *i.e.*, places where *last* annotations can be inserted safely.

Exercise 8.4

Choose an instruction set (other than MIPS) and make patterns for the same subset of the intermediate language as covered by figure 8.1. Use this to translate the intermediate-code example from section 8.4.

Exercise 8.5

In some microprocessors, arithmetic instructions use only two registers, as the destination register is the same as one of the argument registers. As an example, copy and addition instructions of such a processor can be described as follows (using notation like in figure 8.1):

$r_d := r_t$	MOV	r_d, r_t
$r_d := r_d + r_t$	ADD	r_d, r_t
$r_d := r_d + k$	ADDI	r_d, k

As in MIPS, register 0 (R0) is hardwired to the value 0.

Add to the above table pattern/replacement pairs sufficient to translate the following intermediate-code instructions to sequences of machine-code instructions using only MOV, ADD and ADDI instructions in the replacement sequences:

$$\begin{aligned} r_d &:= k \\ r_d &:= r_s + r_t \\ r_d &:= r_s + k \end{aligned}$$

Note that neither r_s nor r_t have the *last* annotation, so their values must be preserved. Note, also, that the intermediate-code instructions above are not a sequence, but a list of separate instructions, so you should generate code separately for each instruction.

Chapter 9

Register Allocation

9.1 Introduction

When generating intermediate code in chapter 7, we have freely used as many variables as we found convenient. In chapter 8, we have simply translated variables in the intermediate language one-to-one into registers in the machine language. Processors, however, do not have an unlimited number of registers, so we need *register allocation* to handle this conflict. The purpose of register allocation is to map a large number of variables into a small(ish) number of registers. This can often be done by letting several variables share a single register, but sometimes there are simply not enough registers in the processor. In this case, some of the variables must be temporarily stored in memory. This is called *spilling*.

Register allocation can be done in the intermediate language prior to machine-code generation, or it can be done in the machine language. In the latter case, the machine code initially uses symbolic names for registers, which the register allocation turns into register numbers. Doing register allocation in the intermediate language has the advantage that the same register allocator can easily be used for several target machines (it just needs to be parameterised with the set of available registers).

However, there may be advantages to postponing register allocation to after machine code has been generated. In chapter 8, we saw that several instructions may be combined to a single instruction, and in the process a variable may disappear. There is no need to allocate a register to this variable, but if we do register allocation in the intermediate language we will do so. Furthermore, when an intermediate-language instruction needs to be translated into a sequence of machine-code instructions, the machine code may need an extra register (or two) for storing temporary values (such as the register needed to store the result of the SLT instruction when translating a jump on $<$ to MIPS code). Hence, the register allocator must make sure that there is always at least one spare register for temporary storage.

The techniques used for register allocation are more or less the same regardless of whether register allocation is done on intermediate code or on machine code. So, in this chapter, we will describe register allocation in terms of the intermediate language introduced in chapter 7.

As in chapter 7, we operate on the body of a single procedure or function, so when we below use the word “program”, we mean it to be such a body. In chapter 10, we will look at how to handle programs consisting of several functions that can call each other.

9.2 Liveness

In order to answer the question “When can two variables share a register?”, we must first define the concept of *liveness*:

Definition 9.1 *A variable is live at some point in the program if the value it contains at that point might conceivably be used in future computations. Conversely, it is dead if there is no way its value can be used in the future.*

We have already hinted at this concept in chapter 8, when we talked about last-uses of variables.

Loosely speaking, two variables may share a register if there is no point in the program where they are both live. We will make a more precise definition later.

We can use some rules to determine when a variable is live:

- 1) If an instruction uses the contents of a variable, that variable is *live* at the start of that instruction.
- 2) If a variable is assigned a value in an instruction, and the same variable is not used as an operand in that instruction, then the variable is *dead* at the start of the instruction, as the value it has at this time is not used before it is overwritten.
- 3) If a variable is live at the end of an instruction and that instruction does not assign a value to the variable, then the variable is also live at the start of the instruction.
- 4) A variable is live at the end of an instruction if it is live at the start of any of the immediately succeeding instructions.

Rule 1 tells how liveness is *generated*, rule 2 how liveness is *killed*, and rules 3 and 4 how liveness is *propagated*.

9.3 Liveness analysis

We can formalise the above rules as equations over sets of variables. The process of solving these equations is called *liveness analysis*, and will at any given point in the program determine which variables are live at this point. To better speak of points in a program, we number all instructions as in figure 9.2.

For every instruction in the program, we have a set of *successors*, *i.e.*, instructions that may immediately follow the instruction during execution. We denote the set of successors to the instruction numbered i as $succ[i]$. We use the following rules to find $succ[i]$:

- 1) The instruction numbered j (if any) that is listed just after instruction number i is in $succ[i]$, unless i is a GOTO or IF-THEN-ELSE instruction. If instructions are numbered consecutively, $j = i + 1$.
- 2) If instruction number i is of the form GOTO l , (the number of) the instruction LABEL l is in $succ[i]$. Note that there in a correct program will be exactly one LABEL instruction with the label used by the GOTO instruction.
- 3) If instruction i is IF p THEN l_t ELSE l_f , (the numbers of) the instructions LABEL l_t and LABEL l_f are in $succ[i]$.

Note that we assume that both outcomes of an IF-THEN-ELSE instruction are possible. If this happens not to be the case (*i.e.*, if the condition is always true or always false), our liveness analysis may claim that a variable is live when it is in fact dead. This is no major problem, as the worst that can happen is that we use a register for a variable that is not going to be used after all. The converse (claiming a variable dead when it is, in fact, live) is worse, as we may overwrite a value that could be used later on, and hence get wrong results from the program. Precise liveness information depends on knowing exactly which paths a program may take through the code when executed, and this is not possible to compute exactly (it is a formally undecidable problem), so it is quite reasonable to allow imprecise results from a liveness analysis, as long as we err on the side of safety, *i.e.*, calling a variable live unless we can prove it to be dead.

For every instruction i , we have a set $gen[i]$, which lists the variables that may be read by instruction i and, hence, are live at the start of the instruction. In other words, $gen[i]$ is the set of variables that instruction i *generates* liveness for. We also have a set $kill[i]$ that lists the variables that may be assigned a value by the instruction. Figure 9.1 shows which variables are in $gen[i]$ and $kill[i]$ for the types of instruction found in intermediate code. x , y and z are (possibly identical) variables and k denotes a constant.

Instruction i	$gen[i]$	$kill[i]$
LABEL l	\emptyset	\emptyset
$x := y$	$\{y\}$	$\{x\}$
$x := k$	\emptyset	$\{x\}$
$x := \mathbf{unop} \ y$	$\{y\}$	$\{x\}$
$x := \mathbf{unop} \ k$	\emptyset	$\{x\}$
$x := y \ \mathbf{binop} \ z$	$\{y, z\}$	$\{x\}$
$x := y \ \mathbf{binop} \ k$	$\{y\}$	$\{x\}$
$x := M[y]$	$\{y\}$	$\{x\}$
$x := M[k]$	\emptyset	$\{x\}$
$M[x] := y$	$\{x, y\}$	\emptyset
$M[k] := y$	$\{y\}$	\emptyset
GOTO l	\emptyset	\emptyset
IF $x \ \mathbf{relop} \ y$ THEN l_t ELSE l_f	$\{x, y\}$	\emptyset
$x := \mathbf{CALL} \ f(args)$	$args$	$\{x\}$

Figure 9.1: Gen and kill sets

For each instruction i , we use two sets to hold the actual liveness information: $in[i]$ holds the variables that are live at the start of i , and $out[i]$ holds the variables that are live at the end of i . We define these by the following equations:

$$in[i] = gen[i] \cup (out[i] \setminus kill[i]) \quad (9.1)$$

$$out[i] = \bigcup_{j \in succ[i]} in[j] \quad (9.2)$$

These equations are recursive. We solve these by fixed-point iteration, as shown in appendix A: We initialise all $in[i]$ and $out[i]$ to be empty sets and repeatedly calculate new values for these until no changes occur. This will eventually happen, since we work with sets with finite support (*i.e.*, a finite number of possible values) and because adding elements to the sets $out[i]$ or $in[j]$ on the right-hand sides of the equations can not reduce the number of elements in the sets on the left-hand sides. Hence, each iteration will either add elements to some set (which we can do only a finite number of times) or leave all sets unchanged (in which case we are done). It is also easy to see that the resulting sets form a solution to the equation – the last iteration essentially verifies that all equations hold. This is a simple extension of the reasoning used in section 2.6.1.

The equations work under the assumption that all uses of a variable are visible in the code that is analysed. If a variable contains, *e.g.*, the output of the program,

```

1:   $a := 0$ 
2:   $b := 1$ 
3:   $z := 0$ 
4:  LABEL loop
5:  IF  $n = z$  THEN end ELSE body
6:  LABEL body
7:   $t := a + b$ 
8:   $a := b$ 
9:   $b := t$ 
10:  $n := n - 1$ 
11:  $z := 0$ 
12: GOTO loop
13: LABEL end

```

Figure 9.2: Example program for liveness analysis and register allocation

it will be used after the program finishes, even if this is not visible in the code of the program itself. So we must ensure that the analysis makes this variable live at the end of the program.

Equation 9.2, similarly, is ill-defined if $\text{succ}[i]$ is the empty set (which is, typically, the case for any instruction that ends the program), so we make a special case: $\text{out}[i]$, where i has no successor, is defined to be the set of all variables that are live at the end of the program. This definition replaces (for these instructions only) equation 9.2.

Figure 9.2 shows a small program that we will calculate liveness for. Figure 9.3 shows succ , gen and kill sets for the instructions in the program.

The program in figure 9.2 calculates the Nth Fibonacci number (where N is given as input by initialising n to N prior to execution). When the program ends (by reaching instruction 13), a will hold the Nth fibonacci number, so a is live at the end of the program. Instruction 13 has no successors ($\text{succ}[13] = \emptyset$), so we set $\text{out}[13] = \{a\}$. The other out sets are defined by equation 9.2 and all in sets are defined by equation 9.1. We initialise all in and out sets to the empty set and iterate until we reach a fixed point.

The order in which we treat the instructions does not matter for the final result of the iteration, but it may influence how quickly we reach the fixed-point. Since the information in equations 9.1 and 9.2 flow backwards through the program, it is a good idea to do the evaluation in reverse instruction order and to calculate $\text{out}[i]$ before $\text{in}[i]$. In the example, this means that we will in each iteration calculate the sets in the order

i	$succ[i]$	$gen[i]$	$kill[i]$
1	2		a
2	3		b
3	4		z
4	5		
5	6, 13	n, z	
6	7		
7	8	a, b	t
8	9	b	a
9	10	t	b
10	11	n	n
11	12		z
12	4		
13			

Figure 9.3: $succ$, gen and $kill$ for the program in figure 9.2

$$out[13], in[13], out[12], in[12], \dots, out[1], in[1]$$

Figure 9.4 shows the fixed-point iteration using this backwards evaluation order. Note that the most recent values are used when calculating the right-hand sides of equations 9.1 and 9.2, so, when a value comes from a higher instruction number, the value from the same column in figure 9.4 is used.

We see that the result after iteration 3 is the same as after iteration 2, so we have reached a fixed point. We note that n is live at the start of the program, which is to be expected, as n is expected to hold the input to the program. If a variable that is not expected to hold input is live at the start of a program, it might in some executions of the program be used before it is initialised, which is generally considered an error (since it can lead to unpredictable results and even security holes). Some compilers issue warnings about uninitialised variables and some compilers add instructions to initialise such variables to a default value (usually 0).

Suggested exercises: 9.1(a,b).

9.4 Interference

We can now define precisely the condition needed for two variables to share a register. We first define *interference*:

i	Initial		Iteration 1		Iteration 2		Iteration 3	
	out[i]	in[i]	out[i]	in[i]	out[i]	in[i]	out[i]	in[i]
1			n, a	n	n, a	n	n, a	n
2			n, a, b	n, a	n, a, b	n, a	n, a, b	n, a
3			n, z, a, b	n, a, b	n, z, a, b	n, a, b	n, z, a, b	n, a, b
4			n, z, a, b	n, z, a, b	n, z, a, b	n, z, a, b	n, z, a, b	n, z, a, b
5			a, b, n	n, z, a, b	a, b, n	n, z, a, b	a, b, n	n, z, a, b
6			a, b, n	a, b, n	a, b, n	a, b, n	a, b, n	a, b, n
7			b, t, n	a, b, n	b, t, n	a, b, n	b, t, n	a, b, n
8			t, n	b, t, n	t, n, a	b, t, n	t, n, a	b, t, n
9			n	t, n	n, a, b	t, n, a	n, a, b	t, n, a
10				n	n, a, b	n, a, b	n, a, b	n, a, b
11					n, z, a, b	n, a, b	n, z, a, b	n, a, b
12					n, z, a, b	n, z, a, b	n, z, a, b	n, z, a, b
13			a	a	a	a	a	a

Figure 9.4: Fixed-point iteration for liveness analysis

Definition 9.2 A variable x interferes with a variable y if $x \neq y$ and there is an instruction i such that $x \in \text{kill}[i]$, $y \in \text{out}[i]$ and instruction i is not $x := y$.

Two different variables can share a register precisely if neither interferes with the other. This is almost the same as saying that they should not be live at the same time, but there are small differences:

- After $x := y$, x and y may be live simultaneously, but as they contain the same value, they can still share a register.
- It may happen that x is not in $\text{out}[i]$ even if x is in $\text{kill}[i]$, which means that we have assigned to x a value that is definitely not read from x later on. In this case, x is not technically live after instruction i , but it still interferes with any y in $\text{out}[i]$. This interference prevents an assignment to x overwriting a live variable y .

The first of these differences is essentially an optimisation that allows more sharing than otherwise, but the latter is important for preserving correctness. In some cases, assignments to dead variables can be eliminated, but in other cases the instruction may have another visible effect (*e.g.*, setting condition flags or accessing memory) and hence can not be eliminated without changing program behaviour.

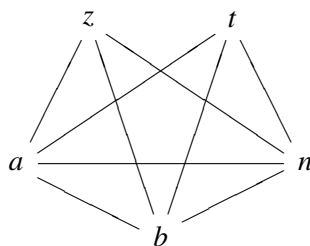


Figure 9.5: Interference graph for the program in figure 9.2

We can use definition 9.2 to generate interference for each assignment statement in the program in figure 9.2:

Instruction	Left-hand side	Interferes with
1	a	n
2	b	n, a
3	z	n, a, b
7	t	b, n
8	a	t, n
9	b	n, a
10	n	a, b
11	z	n, a, b

We will do *global register allocation*, i.e., find for each variable a register that it can stay in at all points in the program (procedure, actually, since a “program” in terms of our intermediate language corresponds to a procedure in a high-level language). This means that, for the purpose of register allocation, two variables interfere if they do so at *any* point in the program. Also, even though interference is defined in an assymetric way in definition 9.2, the conclusion that the two involved variables cannot share a register is symmetric, so interference defines a symmetric relation between variables. A variable can never interfere with itself, so the relation is not reflective.

We can draw interference as an undirected graph, where each node in the graph is a variable, and there is an edge between nodes x and y if x interferes with y (or *vice versa*, as the relation is symmetric). The *interference graph* for the program in figure 9.2 is shown in figure 9.5.

9.5 Register allocation by graph colouring

Two variables can share a register if they are not connected by an edge in the interference graph. Hence, we must assign to each node in the interference graph a register number such that:

- 1) Two nodes that share an edge have different register numbers.
- 2) The total number of different register numbers is no higher than the number of available registers.

This problem is well-known in graph theory, where it is called *graph colouring* (in this context a “colour” is a register number). It is known to be NP-complete, which means that no effective (*i.e.*, polynomial-time) method for doing this optimally is known. In practice, this means that we need to use a heuristic method, which will often find a solution but may give up in some cases even when a solution does exist. This is no great disaster, as we must deal with non-colourable graphs anyway (by moving some variables to memory), so at worst we get slightly slower programs than we would get if we could colour the interference graphs optimally.

The basic idea of the heuristic method we use is simple: If a node in the graph has strictly fewer than N edges, where N is the number of available colours (*i.e.*, registers), we can set this node aside and colour the rest of the graph. When this is done, the (at most $N - 1$) nodes connected by edges to the selected node can not possibly use all N colours, so we can always pick a colour for the selected node from the remaining colours.

We can use this method to four-colour the interference graph from figure 9.5:

- 1) z has three edges, which is strictly less than four. Hence, we remove z from the graph.
- 2) Now, a has less than four edges, so we also remove this.
- 3) Only three nodes are now left (b , t and n), so we can give each of these a number, *e.g.*, 1, 2 and 3 respectively for nodes b , t and n .
- 4) Since three nodes (b , t and n) are connected to a , and these use colours 1, 2 and 3, we must choose a fourth colour for a , *e.g.*, 4.
- 5) z is connected to a , b and n , so we choose a colour that is different from 4, 1 and 3. Giving z colour 2 works.

The problem comes if there are no nodes that have less than N edges. This in itself does not imply that the graph is uncolourable. As an example, a graph with four nodes arranged and connected as the corners of a square can, even though all nodes

have two neighbours, be coloured with two colours by giving opposite corners the same colour. This leads to the following so-called “optimistic” colouring heuristics:

Algorithm 9.3

initialise: *Start with an empty stack.*

simplify: *If there is a node with less than N edges, put this on the stack along with a list of the nodes it is connected to, and remove it and its edges from the graph.*

If there is no node with less than N edges, pick any node and do as above.

*If there are more nodes left in the graph, continue with **simplify**, otherwise go to **select**.*

select: *Take a node and its list of connected nodes from the stack. If possible, give the node a colour that is different from the colours of the connected nodes (which are all coloured at this point). If this is not possible, colouring fails and we mark the node for spilling (see below).*

*If there are more nodes on the stack, continue with **select**.*

The idea in this algorithm is that, even though a node has N or more edges, some of the nodes it is connected to may have been given identical colours, so the total number of colours used for these nodes is less than N . If this is the case, we can use one of the unused colours. If not, we must mark the node for spill.

There are several things left unspecified by algorithm 9.3:

- Which node to choose in **simplify** when none have less than N edges, and
- Which colour to choose in **select** if there are several choices.

If we choose perfectly in both cases, algorithm 9.3 will do optimal colouring. But perfect choices are costly to compute so, in practice, we will sometimes have to guess. We will, in section 9.7, look at some ideas for making qualified guesses. For now, we just make arbitrary choices.

Suggested exercises: 9.1(c,d).

9.6 Spilling

If the **select** phase is unable to find a colour for a node, algorithm 9.3 cannot colour the graph. This means we must give up on keeping all variables in registers throughout the program. We must, hence, select some variables that will reside in memory

(except for brief periods). This process is called *spilling*. Obvious candidates for spilling are variables at nodes that are not given colours by **select**. We simply mark these as *spilled* and continue **select** with the rest of the stack, ignoring spilled nodes when selecting colours for the remaining nodes. When we finish algorithm 9.3, several variables may be marked as spilled.

When we have chosen one or more variables for spilling, we change the program so these are kept in memory. To be precise, for each spilled variable x we:

- 1) Choose a memory address $address_x$ where the value of x is stored.
- 2) In every instruction i that reads or assigns x , we locally in this instruction rename x to x_i .
- 3) Before an instruction i that reads x_i , insert the instruction $x_i := M[address_x]$.
- 4) After an instruction i that assigns x_i , insert the instruction $M[address_x] := x_i$.
- 5) If x is live at the start of the program, add an instruction $M[address_x] := x$ to the start of the program. Note that we use the original name for x here.
- 6) If x is live at the end of the program, add an instruction $x := M[address_x]$ to the end of the program. Note that we use the original name for x here.

After this rewrite of the program, we do register allocation again. This includes re-doing the liveness analysis, since we have added new variables x_i and changed the liveness of x . We may optimise this a bit by repeating the liveness analysis only for the affected variables (x_i and x), as the results will not change for the other variables.

It may happen that the subsequent new register allocation will generate additional spilled variables. There are several reasons why this may be:

- We have ignored spilled variables when selecting colours for a node in the **select** phase. When the spilled variables are replaced by new variables, these may use colours that would otherwise be available, so we may end up with no choices where we originally had one or more colours available.
- The choices of nodes to remove from the graph in the **simplify** phase and the colours to assign in the **select phase** can change, and we might be less lucky in our choices, so we get more spills.

If we have at least as many registers as the number of variables used in a single instruction, all variables can be loaded just before the instruction, and the result can be saved immediately afterwards, so we will eventually be able to find a colouring

Node	Neighbours	Colour
<i>n</i>		1
<i>t</i>	<i>n</i>	2
<i>b</i>	<i>t, n</i>	3
<i>a</i>	<i>b, n, t</i>	<i>spill</i>
<i>z</i>	<i>a, b, n</i>	2

Figure 9.6: Algorithm 9.3 applied to the graph in figure 9.5

by repeated spilling. If we ignore the CALL instruction, no instruction in the intermediate language uses more than two variables, so this is the minimum number of registers that we need. A CALL instruction can use an unbounded number of variables as arguments, possibly even more than the total number of registers available, so it is unrealistic to expect all arguments to function calls to be in registers. We will look at this issue in chapter 10.

If we take our example from figure 9.2, we can attempt to colour its interference graph (figure 9.5) with only three colours. The stack built by the **simplify** phase of algorithm 9.3 and the colours chosen for these nodes in the **select** phase are shown in figure 9.6. The stack grows upwards, so the first node chosen by **simplify** is at the bottom. The colours (numbers) are, conversely, chosen top-down as the stack is popped. We can choose no colour for *a*, as all three available colours are in use by the neighbours *b, n* and *t*. Hence, we mark *a* as spilled. Figure 9.7 shows the program after spill code has been inserted. Note that, since *a* is live at the end of the program, we have inserted a load instruction at the end of the program. Figure 9.8 shows the interference graph for the program in figure 9.7 and figure 9.9 shows the stack used by algorithm 9.3 for colouring this graph, showing that colouring with three colours is now possible.

Suggested exercises: 9.1(e).

9.7 Heuristics

When the **simplify** phase of algorithm 9.3 is unable to find a node with less than N edges, some other node is chosen. So far, we have chosen arbitrarily, but we may apply some heuristics (qualified guessing) to the choice in order to make colouring more likely or reduce the number of spilled variables:

- We may choose a node with close to N neighbours, as this is likely to be colourable in the **select** phase anyway. For example, if a node has exactly N

```

1:  $a_1 := 0$ 
    $M[address_a] := a_1$ 
2:  $b := 1$ 
3:  $z := 0$ 
4: LABEL loop
5: IF  $n = z$  THEN end ELSE body
6: LABEL body
    $a_7 := M[address_a]$ 
7:  $t := a_7 + b$ 
8:  $a_8 := b$ 
    $M[address_a] := a_8$ 
9:  $b := t$ 
10:  $n := n - 1$ 
11:  $z := 0$ 
12: GOTO loop
13: LABEL end
    $a := M[address_a]$ 

```

Figure 9.7: Program from figure 9.2 after spilling variable a

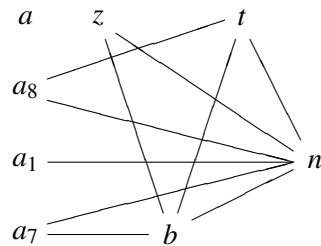


Figure 9.8: Interference graph for the program in figure 9.7

Node	Neighbours	Colour
n		1
t	n	2
a_8	t, n	3
b	t, n	3
a_7	b, n	2
z	b, n	2
a_1	n	2
a		1

Figure 9.9: Colouring of the graph in figure 9.8

neighbours, it will be colourable if just two of its neighbours get the same colour.

- We may choose a node with many neighbours that have close to N neighbours of their own, as spilling this node may allow many of these neighbours to be coloured.
- We may look at the program and select a variable that does not cost so much to spill, *e.g.*, a variable that is not used inside a loop.

These criteria (and maybe others as well) may be combined into a single heuristic by giving numeric values describing how well a variable fits each criterion, multiplying each with a weight and then adding the results to give a weighted sum.

We have also made arbitrary choices when we pick colours for nodes in the **select** phase. We can try to make it more likely that the rest of the graph can be coloured by choosing a colour that is already used elsewhere in the graph instead of picking a colour that is used nowhere else. This will make it less likely that the nodes connected to an as yet uncoloured node will use all the available colours. A simple instance of this idea is to always use the lowest-numbered available colour.

A more advanced variant of this idea is to look at the uncoloured nodes connected to the node we are about to colour. If we have several choices of colour for the current node, we would like to choose a colour that makes it more likely that its uncoloured neighbours can later be coloured. If an uncoloured neighbour has neighbours of its own that are already coloured, we would like to use one of the colours used among these, as this will not increase the number of colours for nodes that neighbour the uncoloured neighbour, so we will not make it any harder to colour this later on. If the current node has several uncoloured neighbours, we can find the set of neighbour-colours for each of these and select a colour that occurs in as many of these sets as possible.

9.7.1 Removing redundant moves

An assignment of the form $x := y$ can be removed from the code if x and y use the same register (as the instruction will have no effect). Most register allocators remove such redundant move instructions, and some even try to increase the number of assignments that can be removed by trying to allocate x and y in the same register whenever possible.

If x has already been given a colour by the time we need to select a colour for y , we can choose the same colour for y , as long as it is not used by any variable that y interferes with (including, possibly, x). Similarly, if x is uncoloured, we can give it the same colour as y if this colour is not used for a variable that interferes with x (including y itself). This is called *biased colouring*.

Another method of achieving the same goal is to combine x and y (if they do not interfere) into a single node before colouring the graph, and only split the combined node if the **simplify** phase can not otherwise find a node with less than N edges. This is called *coalescing*.

The converse of coalescing (called *live-range splitting*) can be used as well: Instead of spilling a variable, we can split its node by giving each occurrence of the variable a different name and inserting assignments between these when necessary. This is not quite as effective at increasing the chance of colouring as spilling, but the cost of the extra assignments is likely to be less than the cost of the loads and stores inserted by spilling.

9.7.2 Using explicit register numbers

Some operations may require their arguments or results to be in specific registers. For example, the integer multiplication instruction in Intel's IA-32 (x86) processors require the first argument to be in the `eax` register and puts the 64-bit result in the `eax` and `edx` registers. Also, as we shall see in chapter 10, function calls can require arguments and results to be in specific registers.

Variables used as arguments results to such operations must, hence, be assigned to these registers *a priori*, before the register allocation begins. We call these precoloured nodes in the interference graph. If two nodes that are precoloured to the same register interfere, we can not make a legal colouring of the graph. One solution would be to spill one or both so they no longer interfere, but that is rather costly.

A better solution is to insert move instructions that move the variables to and from the required registers before and after an instruction that requires specific registers. The specific registers must still be included as precoloured nodes in the interference graph, but are not removed from it in the **simplify** phase. Once only precoloured nodes remain in the graph, the **select** phase starts. When the **select** phase needs to colour a node, it must avoid colours used by all neighbours to the

node – whether they are precoloured or just coloured earlier in the **select** phase. The register allocator can try to remove some of the inserted moves by using the techniques described in section 9.7.1.

9.8 Further reading

Preston Briggs' Ph.D. thesis [12] shows several variants of the register-allocation algorithm shown here, including many optimisations and heuristics as well as considerations about how the various phases can be implemented efficiently. The compiler textbooks [35] and [9] show some other variants and a few newer developments. A completely different approach to register allocation that exploits the structure of a program is suggested in [42].

Exercises

Exercise 9.1

Given the following program:

```

1: LABEL start
2: IF  $a < b$  THEN next ELSE swap
3: LABEL swap
4:  $t := a$ 
5:  $a := b$ 
6:  $b := t$ 
7: LABEL next
8:  $z := 0$ 
9:  $b := b \bmod a$ 
10: IF  $b = z$  THEN end ELSE start
11: LABEL end

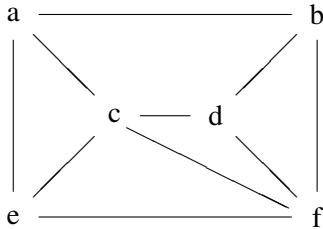
```

- Show *succ*, *gen* and *kill* for every instruction in the program.
- Assuming a is live at the end of the program, *i.e.*, $out[11] = \{a\}$, calculate *in* and *out* for every instruction in the program. Show the iteration as in figure 9.4.
- Draw the interference graph for a, b, t and z .
- Make a three-colouring of the interference graph. Show the stack as in figure 9.6.

- e) Attempt, instead, a two-colouring of the graph. Select variables for spill, do the spill-transformation as shown in section 9.6 and redo the complete register allocation process on the transformed program. If necessary, repeat the process until register allocation is successful.

Exercise 9.2

Three-colour the following graph. Show the stack as in figure 9.6. The graph *is* three-colour-able, so try making different choices if you get spill.



Exercise 9.3

Combine the heuristics suggested in section 9.7 for selecting nodes in the **simplify** phase of algorithm 9.3 into a formula that gives a single numerical score for each node, such that a higher score implies a stronger candidate for spill.

Exercise 9.4

Some processors (such as Motorola 68000) have two types of registers: data registers and address registers. Some instructions (such as load and store) expect their arguments or put their results in address registers while other instructions (such as multiplication and division) expect their arguments or put their results in data registers. Some operations (like addition and subtraction) can use either type of register. There are instructions for moving between address and data registers.

By adding the registers as nodes in the interference graph, a variable can be prevented from being allocated in a specific register by making it interfere with it.

- Describe how instructions that require argument or result variables to be in a specific type of register can ensure this by adding interference for its argument and result variables.
- The answer above is likely to cause spilling of variables that are used as both address and data (as they interfere with all registers). Describe how this can be avoided by taking care of this situation in the spill phase. Hint: Add register-to-register move instructions to the program.

- c) If there are not enough registers of one type, but there are still available registers of the other type, describe how you can spill a variable to a register of the other type instead of to memory.

Exercise 9.5

Some processors have instructions that operate on values that require two registers to hold. Such processors usually require these values to be held in pairs of adjacent registers, so the instructions only need specify one register number per value (as the other part of the value is implicitly stored in the following register).

We will now look at register allocation where some values must be allocated in register pairs. We note that live-ness analysis is unaffected, so only colouring and spill is affected. Hence, we start with an interference graph where some nodes are marked as requiring register pairs.

- a) Modify algorithm 9.3 to take register pairs into account. Focus on correctness, not efficiency. You can assume “colours” are numbers, so you can talk about adjacent colours, the next colour, *etc.*
- b) Describe for the **simplify** phase of algorithm 9.3 heuristics that take into account that some nodes require two registers.
- c) Describe for the **select** phase of algorithm 9.3 heuristics that take into account that some nodes require two registers.

Chapter 10

Function calls

10.1 Introduction

In chapter 7 we have shown how to translate the body of a single function. Function calls were left (mostly) untranslated by using the CALL instruction in the intermediate code. Nor did we in chapter 8 show how the CALL instruction should be translated.

We will, in this chapter, remedy these omissions. We will initially assume that all variables are local to the procedure or function that access them and that parameters are *call-by-value*, meaning that the *value* of an argument expression is passed to the called function. This is the default parameter-passing mechanism in most languages, and in many languages (*e.g.*, C or SML) it is the only one.

10.1.1 The call stack

A single procedure body uses (in most languages) a finite number of variables. We have seen in chapter 9 that we can map these variables into a (possibly smaller) set of registers. A program that uses recursive procedures or functions may, however, use an unbounded number of variables, as each recursive invocation of the function has its own set of variables, and there is no bound on the recursion depth. We can not hope to keep all these variables in registers, so we will use memory for some of these. The basic idea is that only variables that are local to the active (most recently called) function will be kept in registers. All other variables will be kept in memory.

When a function is called, all the live variables of the calling function (which we will refer to as the *caller*) will be stored in memory so the registers will be free for use by the called function (the *callee*). When the callee returns, the stored variables are loaded back into registers. It is convenient to use a stack for this storing and loading, pushing register contents on the stack when they must be saved

and popping them back into registers when they must be restored. Since a stack is (in principle) unbounded, this fits well with the idea of unbounded recursion.

The stack can also be used for other purposes:

- Space can be set aside on the stack for variables that need to be spilled to memory. In chapter 9, we used a constant address ($address_x$) for spilling a variable x . When a stack is used, $address_x$ is actually an offset relative to a stack-pointer. This makes the spill-code slightly more complex, but has the advantage that spilled registers are already saved on the stack when or if a function is called, so they do not need to be stored again.
- Parameters to function calls can be passed on the stack, *i.e.*, written to the top of the stack by the caller and read therefrom by the callee.
- The address of the instruction where execution must be resumed after the call returns (the *return address*) can be stored on the stack.
- Since we decided to keep only local variables in registers, variables that are in scope in a function but not declared locally in that function must reside in memory. It is convenient to access these through the stack.
- Arrays and records that are allocated locally in a function can be allocated on the stack, as hinted in section 7.9.2.

We shall look at each of these in more detail later on.

10.2 Activation records

Each function invocation will allocate a chunk of memory on the stack to cover all of the function's needs for storing values on the stack. This chunk is called the *activation record* or *frame* for the function invocation. We will use these two names interchangeably. Activation records will typically have the same overall structure for all functions in a program, though the sizes of the various fields in the records may differ. Often, the machine architecture (or operating system) will dictate a *calling convention* that standardises the layout of activation records. This allows a program to call functions that are compiled with another compiler or even written in a different language, as long as both compilers follow the same calling convention.

We will start by defining very simple activation records and then extend and refine these later on. Our first model uses the assumption that all information is stored in memory when a function is called. This includes parameters, return address and the contents of registers that need to be preserved. A possible layout for such an activation record is shown in figure 10.1.

	...
	Next activation records
	Space for storing local variables for spill or preservation across function calls
	Remaining incoming parameters
	First incoming parameter / return value
FP →	Return address
	Previous activation records
	...

Figure 10.1: Simple activation record layout

FP is shorthand for “Frame pointer” and points to the first word of the activation record. In this layout, the first word holds the return address. Above this, the incoming parameters are stored. The function will typically move the parameters to registers (except for parameters that have been spilled by the register allocator) before executing its body. The space used for the first incoming parameter is also used for storing the return value of the function call (if any). Above the incoming parameters, the activation record has space for storing other local variables, *e.g.*, for spilling or for preserving across later function calls.

10.3 Prologues, epilogues and call-sequences

In chapter 7, we generated code corresponding to the body of a single function. We assumed parameters to this function were already available in variables and the generated code would just put the function result in a variable.

But, since parameters and results are passed through the activation record, we need in front of the code for the function body to add code that reads parameters from the activation record into variables. This code is called the *prologue* of the function. Likewise, after the code for the function body, we need code to store the calculated return value in the activation record and jump to the return address that was stored in the activation record by the caller. This is called the *epilogue* of the function. For the activation-record layout shown in figure 10.1, a suitable prologue and epilogue is shown in figure 10.2. Note that, though we have used a notation similar to the intermediate language introduced in chapter 7, we have extended this a bit: We have used $M[]$ and GOTO with general expressions as arguments.

We use the names $parameter_1, \dots, parameter_n$ for the intermediate-language variables used in the function body for the n parameters. *result* is the intermediate-language variable that holds the result of the function after the body have been executed.

$$\begin{array}{l}
 \text{Prologue} \quad \left\{ \begin{array}{l} \text{LABEL } \textit{function-name} \\ \textit{parameter}_1 := M[FP + 4] \\ \dots \\ \textit{parameter}_n := M[FP + 4 * n] \end{array} \right. \\
 \\
 \text{code for the function body} \\
 \\
 \text{Epilogue} \quad \left\{ \begin{array}{l} M[FP + 4] := \textit{result} \\ \text{GOTO } M[FP] \end{array} \right.
 \end{array}$$

Figure 10.2: Prologue and epilogue for the frame layout shown in figure 10.1

In chapter 7, we used a single intermediate-language instruction to implement a function call. This function-call instruction must be translated into a *call-sequence* of instructions that will save registers, put parameters in the activation record, *etc.* A call-sequence suitable for the activation-record layout shown in figure 10.1 is shown in figure 10.3. The code is an elaboration of the intermediate-language instruction $x := \text{CALL } f(a_1, \dots, a_n)$.

First, all registers that can be used to hold variables are stored in the frame. In figure 10.3, $R0-Rk$ are assumed to hold variables. These are stored in the activation record just above the calling function's own m incoming parameters. Then, the frame-pointer is advanced to point to the new frame and the parameters and the return address are stored in the prescribed locations in the new frame. Finally, a jump to the function is made. When the function call returns, the result is read from the frame into the variable x , FP is restored to its former value and the saved registers are read back from the old frame.

Keeping all the parameters in register-allocated variables until just before the call, and only then storing them in the frame can require a lot of registers to hold the parameters (as these are all live up to the point where they are stored). An alternative is to store each parameter in the frame as soon as it is evaluated. This way, only one of the variables a_1, \dots, a_n will be live at any one time. However, this can go wrong if a later parameter-expression contains a function call, as the parameters to this call will overwrite the parameters of the outer call. Hence, this optimisation must only be used if no parameter-expressions contain function calls or if nested calls use stack-locations different from those used by the outer call.

In this simple call-sequence, we save on the stack all registers that can potentially hold variables, so these are preserved across the function call. This may save more registers than needed, as not all registers will hold values that are required after the call (*i.e.*, they may be dead). We will return to this issue in section 10.6.


```

 $M[FP + 4 * m + 4] := R0$ 
...
 $M[FP + 4 * m + 4 * (k + 1)] := Rk$ 
 $FP := FP + framesize$ 
 $M[FP + 4] := a_1$ 
...
 $M[FP + 4 * n] := a_n$ 
 $M[FP] := returnaddress$ 
GOTO  $f$ 
LABEL  $returnaddress$ 
 $x := M[FP + 4]$ 
 $FP := FP - framesize$ 
 $R0 := M[FP + 4 * m + 4]$ 
...
 $Rk := M[FP + 4 * m + 4 * (k + 1)]$ 

```

Figure 10.3: Call sequence for $x := \text{CALL } f(a_1, \dots, a_n)$ using the frame layout shown in figure 10.1

Suggested exercises: 10.1.

10.4 Caller-saves versus callee-saves

The convention used by the activation record layout in figure 10.1 is that, before a function is called, the caller saves all registers that must be preserved. Hence, this strategy is called *caller-saves*. An alternative strategy is that the called function saves the contents of the registers that need to be preserved and restores these immediately before the function returns. This strategy is called *callee-saves*.

Stack-layout, prologue/epilogue and call sequence for the callee-saves strategy are shown in figures 10.4, 10.5 and 10.6.

Note that it may not be necessary to store *all* registers that may potentially be used to hold variables, only those that the function actually uses to hold its local variables. We will return to this issue in section 10.6.

So far, the only difference between caller-saves and callee-saves is *when* registers are saved. However, once we refine the strategies to save only a subset of the registers that may potentially hold variables, other differences emerge: Caller-saves need only save the registers that hold *live* variables and callee-saves need only save the registers that the function actually uses. We will in section 10.6 return to how this can be achieved, but at the moment just assume these optimisations are made.

	...
	Next activation records
	Space for storing local variables for spill
	Space for storing registers that need to be preserved
	Remaining incoming parameters
	First incoming parameter / return value
FP →	Return address
	Previous activation records
	...

Figure 10.4: Activation record layout for callee-saves

$$\begin{array}{l}
 \text{Prologue} \left\{ \begin{array}{l} \text{LABEL } \textit{function-name} \\ M[FP + 4 * n + 4] := R0 \\ \dots \\ M[FP + 4 * n + 4 * (k + 1)] := Rk \\ \textit{parameter}_1 := M[FP + 4] \\ \dots \\ \textit{parameter}_n := M[FP + 4 * n] \end{array} \right. \\
 \\
 \textit{code for the function body} \\
 \\
 \text{Epilogue} \left\{ \begin{array}{l} M[FP + 4] := \textit{result} \\ R0 := M[FP + 4 * n + 4] \\ \dots \\ Rk := M[FP + 4 * n + 4 * (k + 1)] \\ \text{GOTO } M[FP] \end{array} \right.
 \end{array}$$

Figure 10.5: Prologue and epilogue for callee-saves

```

     $FP := FP + framesize$ 
     $M[FP + 4] := a_1$ 
    ...
     $M[FP + 4 * n] := a_n$ 
     $M[FP] := returnaddress$ 
    GOTO  $f$ 
    LABEL  $returnaddress$ 
     $x := M[FP + 4]$ 
     $FP := FP - framesize$ 

```

Figure 10.6: Call sequence for $x := \text{CALL } f(a_1, \dots, a_n)$ for callee-saves

Caller-saves and callee-saves each have their advantages (described above) and disadvantages: When caller-saves is used, we might save a live variable in the frame even though the callee does not use the register that holds this variable. On the other hand, with callee-saves we might save some registers that do not actually hold live values. We can not avoid these unnecessary saves, as each function is compiled independently and hence do not know the register usage of their callers/callees. We can, however, try to reduce unnecessary saving of registers by using a mixed caller-saves and callee-saves strategy:

Some registers are designated caller-saves and the rest as callee-saves. If any live variables are held in caller-saves registers, it is the caller that must save these to its own frame (as in figure 10.3, though only registers that are both designated caller-saves *and* hold live variables are saved). If a function uses any callee-saves registers in its body, it must save these before using them, as in figure 10.5. Only callee-saves registers that are actually used in the body need to be saved.

Calling conventions typically specify which registers are caller-saves and which are callee-saves, as well as the layout of the activation records.

10.5 Using registers to pass parameters

In both call sequences shown (in figures 10.3 and 10.6), parameters are stored in the frame, and in both prologues (figures 10.2 and 10.5) most of these are immediately loaded back into registers. It will save a good deal of memory traffic if we pass the parameters in registers instead of memory.

Normally, only a few (4-8) registers are used for parameter passing. These are used for the first parameters of a function, while the remaining parameters are passed on the stack, as we have done above. Since most functions have fairly short parameter lists, most parameters will normally be passed in registers. The registers

Register	Saved by	Used for
0	caller	parameter 1 / result / local variable
1-3	caller	parameters 2 - 4 / local variables
4-12	callee	local variables
13	caller	temporary storage (unused by register allocator)
14	callee	FP
15	callee	return address

Figure 10.7: Possible division of registers for 16-register architecture

	...
	Next activation records
	Space for storing local variables for spill and for storing live variables allocated to caller-saves registers across function calls
	Space for storing callee-saves registers that are used in the body
	Incoming parameters in excess of four
FP →	Return address
	Previous activation records
	...

Figure 10.8: Activation record layout for the register division shown in figure 10.7

used for parameter passing are typically a subset of the caller-saves registers, as parameters are not live after the call and hence do not have to be preserved.

A possible division of registers for a 16-register architecture is shown in figure 10.7. Note that the return address is also passed in a register. Most RISC architectures have jump-and-link (function-call) instructions, which leaves the return address in a register, so this is only natural. However, if a function call is made inside the body, this register is overwritten, so the return address must be saved in the activation record before any calls. The return-address register is marked as callee-saves in figure 10.7. In this manner, the return-address register is just like any other variable that must be preserved in the frame if it is used in the body (which it is if a function call is made). Strictly speaking, we do not need the return address after the call has returned, so we can also argue that R15 is a caller-saves register. If so, the caller must save R15 prior to any call, *e.g.*, by spilling it.

Activation record layout, prologue/epilogue and call sequence for a calling convention using the register division in figure 10.7 are shown in figures 10.8, 10.9 and 10.10.

Prologue	{	LABEL <i>function-name</i> $M[FP + offset_{R4}] := R4$ (if used in body) ... $M[FP + offset_{R12}] := R12$ (if used in body) $M[FP] := R15$ (if used in body) $parameter_1 := R0$ $parameter_2 := R1$ $parameter_3 := R2$ $parameter_4 := R3$ $parameter_5 := M[FP + 4]$... $parameter_n := M[FP + 4 * (n - 4)]$
<i>code for the function body</i>		
Epilogue	{	$R0 := result$ $R4 := M[FP + offset_{R4}]$ (if used in body) ... $R12 := M[FP + offset_{R12}]$ (if used in body) $R15 := M[FP]$ (if used in body) GOTO R15

Figure 10.9: Prologue and epilogue for the register division shown in figure 10.7

```

 $M[FP + offset_{live_1}] := live_1$     (if allocated to a caller-saves register)
...
 $M[FP + offset_{live_k}] := live_k$     (if allocated to a caller-saves register)
 $FP := FP + framesize$ 
 $R0 := a_1$ 
...
 $R3 := a_4$ 
 $M[FP + 4] := a_5$ 
...
 $M[FP + 4 * (n - 4)] := a_n$ 
 $R15 := returnaddress$ 
GOTO  $f$ 
LABEL  $returnaddress$ 
 $x := R0$ 
 $FP := FP - framesize$ 
 $live_1 := M[FP + offset_{live_1}]$     (if allocated to a caller-saves register)
...
 $live_k := M[FP + offset_{live_k}]$     (if allocated to a caller-saves register)

```

Figure 10.10: Call sequence for $x := \text{CALL } f(a_1, \dots, a_n)$ for the register division shown in figure 10.7

Note that the offsets for storing registers are not simple functions of their register numbers, as only a subset of the registers need to be saved. R15 (which holds the return address) is, like any other callee-saves register, saved in the prologue and restores in the epilogue if it is used inside the body (*i.e.*, if the body makes a function call). Its offset is 0, as the return address is stored at offset 0 in the frame.

In a call-sequence, the instructions

```
R15 := returnaddress  
GOTO f  
LABEL returnaddress
```

can on most RISC processors be implemented by a jump-and-link instruction.

10.6 Interaction with the register allocator

As we have hinted above, the register allocator can be used to optimise function calls, as it can provide information about which registers need to be saved.

The register allocator can tell which variables are live after the function call. In a caller-saves strategy (or for caller-saves registers in a mixed strategy), only the (caller-saves) registers that hold such variables need to be saved before the function call.

Likewise, the register allocator can return information about which registers are used by the function body, so only these need to be saved in a callee-saves strategy.

If a mixed strategy is used, variables that are live across a function call should, if possible, be allocated to callee-saves registers. This way, the caller does not have to save these and, with luck, they do not have to be saved by the callee either (if the callee does not use these registers in its body). If all variables that are live across function calls are made to interfere with all caller-saves registers, the register allocator will not allocate these variables in caller-saves registers, which achieves the desired effect. If no callee-saves register is available, the variable will be spilled and hence, effectively, be saved across the function call. This way, the call sequence will not need to worry about saving caller-saves registers, this is all done by the register allocator.

As spilling may be somewhat more costly than local save/restore around a function call, it is a good idea to have plenty of callee-saves registers for holding variables that are live across function calls. Hence, most calling conventions specify more callee-saves registers than caller-saves registers.

Note that, though the prologues shown in figures 10.2, 10.5 and 10.9 load all stack-passed parameters into registers, this should actually only be done for parameters that are not spilled. Likewise, a register-passed parameter that needs to be spilled should be transferred to a stack location instead of to a symbolic register (*parameter_i*).

In figures 10.2, 10.5 and 10.9, we have moved register-passed parameters from the numbered registers or stack locations to named registers, to which the register allocator must assign numbers. Similarly, in the epilogue we move the function result from a named variable to *R0*. This means that these parts of the prologue and epilogue must be included in the body when the register allocator is called (so the named variables will be replaced by numbers). This will also automatically handle the issue about spilled parameters mentioned above, as spill-code is inserted immediately after the parameters are (temporarily) transferred to registers. This may cause some extra memory transfers when a spilled stack-passed parameter is first loaded into a register and then immediately stored back again. This problem is, however, usually handled by later optimisations.

It may seem odd that we move register-passed parameters to named registers instead of just letting them stay in the registers they are passed in. But these registers may be needed for other function calls, which gives problems if a parameter allocated to one of these needs to be preserved across the call (as mentioned above, variables that are live across function calls should not be allocated to caller-saves registers). By moving the parameters to named registers, the register allocator is free to allocate these to callee-saves registers if needed. If this is not needed, the register allocator may allocate the named variable to the same register as the parameter was passed in and eliminate the (superfluous) register-to-register move. As mentioned in section 9.7, modern register allocators will eliminate most such moves anyway, so we might as well exploit this.

In summary, given a good register allocator, the compiler needs to do the following to compile a function:

- 1) Generate code for the body of the function, using symbolic names for variables (except precoloured temporary variables used for parameter-passing in call sequences or for instructions that require specific registers, see section 9.7.2).
- 2) Add code for moving parameters from numbered registers and stack locations into the named variables used for accessing the parameters in the body of the function, and for moving the function-result from a named register to the register used for function results.
- 3) Call the register allocator with this extended function body. The register allocator should be aware of the register division (caller-saves/callee-saves split) and allocate variables that are live across function calls only to callee-saves registers and should return both the set of used callee-saves registers and the set of spilled variables.
- 4) To the register-allocated code, add code for saving and restoring the callee-saves registers that the register allocator says have been used in the extended

function body and for updating the frame pointer with the size of the frame (including space for saved registers and spilled variables).

- 5) Add a function label at the beginning of the code and a return jump at the end.

10.7 Accessing non-local variables

We have up to now assumed that all variables used in a function are local to that function, but most high-level languages also allow functions to access variables that are not declared locally in the functions themselves.

10.7.1 Global variables

In C, variables are either global or local to a function. Local variables are treated exactly as we have described, *i.e.*, typically stored in a register. Global variables will, on the other hand, be stored in memory. The location of each global variable will be known at compile-time or link-time. Hence, a use of a global variable x generates the code

$$x := M[address_x]$$

instruction that uses x

The global variable is loaded into a (register-allocated) temporary variable and this will be used in place of the global variable in the instruction that needs the value of the global variable.

An assignment to a global variable x is implemented as

$$x := \text{the value to be stored in } x$$

$$M[address_x] := x$$

Note that global variables are treated almost like spilled variables: Their value is loaded from memory into a register immediately before any use and stored from a register into memory immediately after an assignment. Like with spill, it is possible to use different register-allocated variables for each use of x .

If a global variable is used often within a function, it can be loaded into a local variable at the beginning of the function and stored back again when the function returns. However, a few extra considerations need to be made:

- The variable must be stored back to memory whenever a function is called, as the called function may read or change the global variable. Likewise, the global variable must be read back from memory after the function call, so any changes will be registered in the local copy. Hence, it is best to allocate local copies of global variables in caller-saves registers.

- If the language allows *call-by-reference* parameters (see below) or pointers to global variables, there may be more than one way to access a global variable: Either through its name or via a call-by-reference parameter or pointer. If we cannot exclude the possibility that a call-by-reference parameter or pointer can access a global variable, it must be stored/retrieved before/after any access to a call-by-reference parameter or any access through a pointer. It is possible to make a global *alias analysis* that determines if global variables, call-by-reference parameters or pointers may point to the same location (*i.e.*, may be *aliased*). However, this is a fairly complex analysis, so many compilers simply assume that a global variable may be aliased with *any* call-by-reference parameter or pointer and that any two of the latter may be aliased.

The above tells us that accessing local variables (including call-by-value parameters) is faster than accessing global variables. Hence, good programmers will use global variables sparingly.

10.7.2 Call-by-reference parameters

Some languages, *e.g.*, Pascal (which uses the term **var**-parameters), allow parameters to be passed by *call-by-reference*. A parameter passed by call-by-reference must be a variable, an array element, a field in a record or, in general, anything that is allowed at the left-hand-side of an assignment statement. Inside the function that has a call-by-reference parameter, values can be assigned to the parameter and these assignments actually update the variable, array element or record-field that was passed as parameter such that the changes are visible to the caller. This differs from assignments to call-by-value parameters in that these update only a local copy.

Call-by-reference is implemented by passing the address of the variable, array element or whatever that is given as parameter. Any access (use or definition) to the call-by-reference parameter must be through this address.

In C, there are no explicit call-by-reference parameters, but it is possible to explicitly pass pointers to variables, array-elements, *etc.* as parameters to a function by using the `&` (address-of) operator. When the value of the variable is used or updated, this pointer must be explicitly followed, using the `*` (de-reference) operator. So, apart from notation and a higher potential for programming errors, this is not significantly different from “real” call-by-reference parameters.

In any case, a variable that is passed as a call-by-reference parameter or has its address passed via a `&` operator, must reside in memory. This means that it must be spilled at the time of the call or allocated to a caller-saves register, so it will be stored before the call and restored afterwards.

It also means that passing a result back to the caller by call-by-reference or pointer parameters can be slower than using the function’s return value, as the

```
procedure f (x : integer);  
  var y : integer;  
  function g(p : integer);  
  var q : integer;  
  begin  
    if p<10 then y := g(p+y)  
    else q := p+y;  
    if (y<20) then f(y);  
    g := q  
  end;  
begin  
  y := x+x;  
  writeln(g(y),y)  
end;
```

Figure 10.11: Example of nested scopes in Pascal

return value can be passed in a register. Hence, like global variables, call-by-reference and pointer parameters should be used sparingly.

Either of these on their own have the same aliasing problems as when combined with global variables.

10.7.3 Nested scopes

Some languages, *e.g.*, Pascal and SML, allow functions to be declared locally within other functions. A local function typically has access to variables declared in the function in which it itself is declared. For example, figure 10.11 shows a fragment of a Pascal program. In this program, *g* can access *x* and *y* (which are declared in *f*) as well as its own local variables *p* and *q*.

Note that, since *f* and *g* are recursive, there can be many instances of their variables in different activation records at any one time.

When *g* is called, its own local variables (*p* and *q*) are held in registers, as we have described above. All other variables (*i.e.*, *x* and *y*) reside in the activation records of the procedures/functions in which they are declared (in this case *f*). It is no problem for *g* to know the offsets for *x* and *y* in the activation record for *f*, as *f* can be compiled before *g*, so full information about *f*'s activation record layout is available for the compiler when it compiles *g*. However, we will not at compile-time know the position of *f*'s activation record on the stack. *f*'s activation record will not always be directly below that of *g*, since there may be several recursive invocations of *g* (each with its own activation record) above the last activation record for *f*. Hence, a pointer to *f*'s activation record will be given as parameter to

```

function g(var fFrame : fRecord, p : integer);
var q : integer;
begin
    if p<10 then fFrame.y := g(fFrame,p+fFrame.y)
    else q := p+fFrame.y;
    if (fFrame.y<20) then f(fFrame.y);
    g := q
end;

procedure f (x : integer);
    var y : integer;
begin
    y := x+x;
    writeln(g(FP,y),y)
end;

```

Figure 10.12: Adding an explicit frame-pointer to the program from figure 10.11

g when it is called. When f calls g, this pointer is just the contents of FP, as this, by definition, points to the activation record of the active function (*i.e.*, f). When g is called recursively from g itself, the incoming parameter that points to f's activation record is passed on as a parameter to the new call, so every instance of g will have its own copy of this pointer.

To illustrate this, we have in figure 10.12 added this extra parameter explicitly to the program from figure 10.11. Now, g accesses all non-local variables through the fFrame parameter, so it no longer needs to be declared locally inside f. Hence, we have moved it out. We have used record-field-selection syntax in g for accessing f's variables through fFrame. Note that fFrame is a call-by-reference parameter (indicated by the var keyword), as g can update f's variables (*i.e.*, y). In f, we have used FP to refer to the current activation record. Normally, a function in a Pascal program will not have access to its own frame, so this is not quite standard Pascal.

It is sometimes possible to make the transformation entirely in the source language (*e.g.*, Pascal), but the extra parameters are usually not added until the intermediate code, where FP is made explicit, has been generated. Hence, figure 10.12 mainly serves to illustrate the idea, not as a suggestion for implementation.

Note that all variables that can be accessed in inner scopes need to be stored in memory when a function is called. This is the same requirement as was made for call-by-reference parameters, and for the same reason. This can, in the same way, be handled by allocating such variables in caller-saves registers.

	...
	Next activation records
	Space for storing local variables for spill and for storing live variables allocated to caller-saves registers across function calls
	Space for storing callee-saves registers that are used in the body
	Incoming parameters in excess of four
	Return address
FP →	Static link (SL)
	Previous activation records
	...

Figure 10.13: Activation record with static link

	y		q
	x		p
	Return address		Return address
FP →	SL (null)	FP →	SL (to f)

Figure 10.14: Activation records for f and g from figure 10.11

Static links

If there are more than two nested scopes, pointers to all outer scopes need to be passed as parameters to locally declared functions. If, for example, g declared a local function h, h would need pointers to both f's and g's activation records. If there are many nested scopes, this list of extra parameters can be quite long. Typically, a single parameter is instead used to hold a linked list of the frame pointers for the outer scopes. This is normally implemented by putting the links in the activation records themselves. Hence, the first field of an activation record (the field that FP points to) will point to the activation record of the next outer scope. This is shown in figure 10.13. The pointer to the next outer scope is called the *static link*, as the scope-nesting is static as opposed to the actual sequence of run-time calls that determine the stacking-order of activation records¹. The layout of the activation records for f and g from figure 10.11 is shown in figure 10.14.

g's static link will point to the most recent activation record for f. To read y, g will use the code

¹Sometimes, the return address is referred to as the *dynamic link*.

$FP_f := M[FP]$	Follow g 's static link
$address := FP_f + 12$	Calculate address of y
$y := M[address]$	Get y 's value

where y afterwards holds the value of y . To write y , g will use the code

$FP_f := M[FP]$	Follow g 's static link
$address := FP_f + 12$	Calculate address of y
$M[address] := y$	Write to y

where y holds the value that is written to y . If a function h was declared locally inside g , it would need to follow two links to find y :

$FP_g := M[FP]$	Follow h 's static link
$FP_f := M[FP_g]$	Follow g 's static link
$address := FP_f + 12$	Calculate address of y
$y := M[address]$	Get y 's value

This example shows why the static link is put in the first element of the activation record: It makes following a chain of links easier, as no offsets have to be added in each step.

Again, we can see that a programmer should keep variables as local as possible, as non-local variables take more time to access.

10.8 Variants

We have so far seen fixed-size activation records on stacks that grow upwards in memory, and where FP points to the first element of the frame. There are, however, reasons why you sometimes may want to change this.

10.8.1 Variable-sized frames

If arrays are allocated on the stack, the size of the activation record depends on the size of the arrays. If these sizes are not known at compile-time, neither will the size of the activation records. Hence, we need a run-time variable to point to the end of the frame. This is typically called the *stack pointer*, because the end of the frame is also the top of the stack. When setting up parameters to a new call, these are put at places relative to SP rather than relative to FP . When a function is called, the new FP takes the value of the old SP , but we now need to store the old value of FP , as we no longer can restore it by subtracting a constant from the current FP . Hence, the old FP is passed as a parameter (in a register or in the frame) to the new function, which restores FP to this value just before returning.

If arrays are allocated on a separate stack, frames can be of fixed size, but a separate stack-pointer is now needed for allocating/deallocating arrays.

If two stacks are used, it is customary to let one grow upwards and the other downwards, such that they grow towards each other. This way, stack-overflow tests on both stacks can be replaced by a single test on whether the stack-tops meet. It also gives a more flexible division of memory between the two stacks than if each stack is allocated its own fixed-size memory segment.

10.8.2 Variable number of parameters

Some languages (*e.g.*, C and LISP) allow a function to have a variable number of parameters. This means that the function can be called with a different number of parameters at each call. In C, the `printf` function is an example of this.

The layouts we have shown in this chapter all assume that there is a fixed number of arguments, so the offsets to, *e.g.*, local variables are known. If the number of parameters can vary, this is no longer true.

One possible solution is to have two frame pointers: One that shows the position of the first parameter and one that points to the part of the frame that comes after the parameters. However, manipulating two FP's is somewhat costly, so normally another trick is used: The FP points to the part of the frame that comes after the parameters. Below this, the parameters are stored at negative offsets from FP, while the other parts of the frame are accessed with (fixed) positive offsets. The parameters are stored such that the first parameter is closest to FP and later parameters further down the stack. This way, parameter number k will be a fixed offset ($-4 * k$) from FP.

When a function call is made, the number of arguments to the call is known to the caller, so the offsets (from the old FP) needed to store the parameters in the new frame will be fixed at this point.

Alternatively, FP can point to the top of the frame and all fields can be accessed by fixed negative offsets. If this is the case, FP is sometimes called SP, as it points to the top of the stack.

10.8.3 Direction of stack-growth and position of FP

There is no particular reason why a stack has to grow upwards in memory. It is, in fact, more common that call stacks grow downwards in memory. Sometimes the choice is arbitrary, but at other times there is an advantage to have the stack growing in a particular direction. Some instruction sets have memory-access instructions that include a constant offset from a register-based address. If this offset is unsigned (as it is on, *e.g.*, IBM System/370), it is an advantage that all fields in the activation record are at non-negative offsets. This means that either FP must point to the

bottom of the frame and the stack grow upwards, or FP must point to the top of the frame and the stack grow downwards.

If, on the other hand, offsets are signed but have a small range (as on Digital's Vax, where the range is $-128 - +127$), it is an advantage to use both positive and negative offsets. This can be done, as suggested in section 10.8.2, by placing FP after the parameters but before the rest of the frame, so parameters are addressed by negative offsets and the rest by positive. Alternatively, FP can be positioned k bytes above the bottom of the frame, where $-k$ is the largest negative offset.

10.8.4 Register stacks

Some processors, *e.g.*, Suns Sparc and Intels IA-64 have on-chip stacks of registers. The intention is that frames are kept in registers rather than on a stack in memory. At call or return of a function, the register stack is adjusted. Since the register stack has a finite size, which is often smaller than the total size of the call stack, it may overflow. This is trapped by the operating system which stores part of the stack in memory and shifts the rest down (or up) to make room for new elements. If the stack underflows (at a pop from an empty register stack), the OS will restore earlier saved parts of the stack.

10.8.5 Functions as values

If you can pass a function as a parameter to another function, you need to pass its static link as well, so the function can access its non-local variables when it is called. If there are no local function definitions (and, hence, no need for static links) it is enough to pass the address of the function. This applies, for example, to the language C.

In Pascal, functions can be passed as parameters, but can not be returned as function results. When you pass a function as parameter, you pass a pair consisting of the address of the function's code and its static link. This pair is called a *thunk* or a *closure*. When the function is called and its frame is put on the stack, the static link is taken from this pair. Since you cannot return a functional value, the static link will always point to a frame that is further down the stack.

However, if you can return a functional value as the result of a function (as is possible in most functional languages), the frame that the static link points to can have been unstacked by the time the functional value is used. To prevent this, the part of the frame that contains variables is allocated on the heap instead of the stack and the static link in the closure points to the heap-allocated part of the frame.

Suggested exercises: 10.2.

10.9 Further reading

Calling conventions for various architectures are usually documented in the manuals provided by the vendors of these architectures. Additionally, the calling convention for the MIPS microprocessor is shown in [39].

In figure 10.12, we showed in source-language terms how an extra parameter can be added for accessing non-local parameters, but stated that this was for illustrative purposes only, and that the extra parameters are not normally added at source-level. However, [8] argues that it *is*, actually, a good idea to do this, and goes on to show how many advanced features regarding nested scopes, higher-order functions and even register allocation can be implemented mostly by source-level transformations.

Section 11.7 describes some optimisations that can be used to make function calls faster or use less stack space.

Exercises

Exercise 10.1

In section 10.3 an optimisation is mentioned whereby parameters are stored in the new frame as soon as they are evaluated instead of just before the call. It is warned that this will go wrong if any of the parameter-expressions themselves contain function calls. Argue that the *first* parameter-expression of a function call can contain other function calls without causing the described problem.

Exercise 10.2

Section 10.8.2 suggests that a variable number of arguments can be handled by storing parameters at negative offsets from FP and the rest of the frame at non-negative offsets from FP. Modify figures 10.8, 10.9 and 10.10 to follow this convention.

Exercise 10.3

Find documentation for the calling convention of a processor of your choice and modify figures 10.7, 10.8, 10.9 and 10.10 to follow this convention.

Exercise 10.4

Many functions have a body consisting of an if-then-else statement (or expression), where one or both branches use only a subset of the variables used in the body as a whole. As an example, assume the body is of the form

```
IF cond THEN label1 ELSE label2  
LABEL label1  
code1  
GOTO label3  
LABEL label2  
code2  
LABEL label3
```

The condition *cond* is a simple comparison between variables (which may or may not be callee-saves).

A normal callee-saves strategy will in the prologue save (and in the epilogue restore) all callee-saves registers used in the body. But since only one branch of the if-then-else is taken, some registers are saved and restored unnecessarily.

We can, as usual, get information about variable use in the different parts of the body (*i.e.*, *cond*, *code*₁ and *code*₂) from the register allocator.

We will now attempt to combine the prologue and epilogue with a function body of the above form in order to reduce the number of *callee-saves* registers saved.

Replace in figure 10.9 the text “*code for function body*” by the above body. Then modify the combined code so parts of saving and restoring registers *R4* – *R12* and *R15* is moved into the branches of the if-then-else structure. Be precise about which registers are saved and restored where. You can use clauses like “if used in *code*₁”.

Chapter 11

Analysis and optimisation

In chapter 8, we briefly mentioned optimisations without going into detail about how they are done. We will remedy this in this chapter.

An optimisation, generally, is about recognising instructions that form a specific pattern that can be replaced by a smaller or faster pattern of new instructions. In the simplest case, the pattern is just a short sequence of instructions that can be replaced by another short sequence of instructions. In chapter 8, we replaced sequences of intermediate-code instructions by sequences of machine-code instructions, but the same idea can be applied to replacing sequences of machine-code instructions by sequences of machine-code instructions or sequences of intermediate-code instructions by other sequences of intermediate-code instructions. This kind of optimisation is called *peephole optimisation*, because we look at the code through a small hole that just allows us to see short sequences of instructions.

Another thing to note about the patterns we used in chapter 8 is that they sometimes required some of the variables involved to have no subsequent uses. This is a non-local property that requires looking at an arbitrarily large context of the instructions, sometimes the entire procedure or function in which the instructions appear. A variable having possible subsequent uses is called *live*. In chapter 9 we looked at how *liveness analysis* could determine which variables are live at each instruction. We used this to determine interference between variables, but the same information can also be used to enable optimisations, such as replacing a sequence of instructions by a simpler sequence.

Many optimisations are like this: We have a pattern of instructions and some requirements about the context in which the pattern appears. These contextual requirements are often found by analyses similar to liveness analysis, collectively called *data-flow analyses*.

We will now look at optimisations that can be enabled by such analyses. We will later present a few other types of optimisations.

11.1 Data-flow analysis

As the name indicates, data-flow analysis attempts to discover how information flows through a program. We already discussed liveness analysis, where the information that a variable is live flows through the program in the opposite order of the flow of values: A value flows from an assignment to a variable to its uses, but liveness information flows from a use of a variable back to its assignments. Liveness analysis is, hence, called a backwards analysis. In other analyses information flows in the same order as values. For example, an analysis might try to approximate the set of possible values that a variable can hold, and here the flow is naturally from assignments to uses.

The liveness analysis presented in chapter 9 consisted of four things:

1. Information about which instructions can follow others, *i.e.*, the successors of each instruction.
2. For each instruction *gen* and *kill* sets that describe how data-flow information is created and destroyed by the instruction.
3. Equations that define *in* and *out* sets by describing how data-flow information flow between instructions.
4. Initialisation of the *in* and *out* sets for a fixed-point iteration that solve the data-flow equations.

We will use the same template for other data-flow analyses, but the details might differ. For example:

1. Forwards analyses require information about the predecessors of an instruction instead of its successors.
2. Where liveness analysis uses sets of variables, other analyses might use sets of instructions or sets of variable/value pairs.
3. The equations for *in* and *out* sets might differ. For example, they may use intersection instead of union to combine information from several successors or predecessors of an instruction.
4. Where liveness analysis initialises all *in* and *out* sets to empty sets (except for the *out* set of the last instruction in the function, which is initialised to the set of variables live at the exit of the function), other analyses might initialise the sets to, for example, the set of all instructions in the function. If we want the minimal solution to the equations, we initialise with empty sets, but if we want the maximal solution to the equations, we initialise with the set of all

relevant values. See appendix A for more details about minimal and maximal solutions to set equations.

We will in the following sections show some examples of optimisations and the data-flow analyses required to find the contextual information required to determine if the optimisation is applicable. As the optimisations we show are not specific to any particular processor, we will show them in terms of intermediate code, but the same principles can be applied to analysis and optimisations of machine-code.

11.2 Common subexpression elimination

After translation to intermediate code, there might be several occurrences of the same calculations, even when this is not the case in the source program. For example, the assignment `a[i] := a[i]+1` can in many languages not be simplified at the source level, but the assignment might be translated to the following intermediate-code sequence:

```
t1 := 4*i
t2 := a+t1
t3 := M[t2]
t4 := t3+1
t5 := 4*i
t6 := a+t5
M[t6] := t4
```

Note that both the multiplication by 4 and the addition of a is repeated, but that while we can see that the expressions `a+t1` and `a+t5` have the same value, they are not textually identical.

Our ultimate goal is to eliminate both of these redundancies, but we will start by a simpler analysis that only catches the second occurrence of `4*i` and then discuss how it can be extended to also eliminate `a+t5`.

11.2.1 Available assignments

If we want to replace an expression by a variable that holds the value of the expression, we need to keep track of which expressions we have available and which variables hold their values. So we want at each program point a set of pairs of variables and expressions. Since each such pair originates in an assignment of the expression to the variable, we can, equivalently, use a set of assignment instructions that occur in the program. This makes things slightly simpler later on. We call this analysis *available assignments*.

It is clear that information flows from assignment forwards, so for each instruction in the program, the *in* set is the set of available assignments before the instruction is executed and the *out* set is the set of assignments available afterwards. The *gen* and *kill* sets should, hence, describe which new assignments become available and which are no longer available. An assignment makes itself available, *unless* the variable on the left-hand side also occurs on the right-hand side (because the assignment would make a new occurrence of the expression have a different value). All other assignments where the variable on the left-hand side of the instruction occurs (on either side) are invalidated: If the variable occurs on the left-hand side on another assignment, the variable no longer holds the value of the expression on the right-hand side, and if the variable occurs in an expression, the expression changes value, so the left-hand-side variable no longer holds the current value of the expression. Figure 11.1 shows the *gen* and *kill* sets for each kind of instruction in the intermediate language.

Note that a copy instruction $x := y$ does not generate any available assignment, as nothing is gained by replacing an occurrence of y by x . Note, also, that any store of a value to memory kills all instructions that load from memory. We need to make this rather conservative assumption because we do not know where in memory loads and stores go.

The next step is to define the equations for *in* and *out* sets:

$$out[i] = gen[i] \cup (in[i] \setminus kill[i]) \quad (11.1)$$

$$in[i] = \bigcap_{j \in pred[i]} out[j] \quad (11.2)$$

As mentioned above, the assignments that are available after an instruction (*i.e.*, in the *out* set) are those that are generated by the instruction and those that were available before, except for those that are killed by the instruction. The available assignments before an instruction (*i.e.*, in the *in* set) are those that are available at all predecessors, so we take the intersection of the sets available at the predecessors.

The predecessors $pred[i]$ of instruction i is $\{i - 1\}$, except when i is a LABEL instruction, in which case we also add all GOTO and IF instructions that can jump to i , or when i is the first instruction in the program (*i.e.*, when $i = 1$), in which case $i - 1$ is not in the predecessors of i .

We need to initialise the *in* and *out* sets for the fixed-point iteration. We want to find the maximal solution to the equations: Consider a loop where an assignment is available before the loop and no variable in the assignment is changed inside the loop. We would want the assignment to be available inside the loop, but if we initialise the *out* set of the jump from the end of the loop to its beginning to the empty set, the intersection of the assignments available before the loop and those available at the end of the loop will be empty, and remain that way throughout the

Instruction i	$gen[i]$	$kill[i]$
LABEL l	\emptyset	\emptyset
$x := y$	\emptyset	$assg(x)$
$x := k$	$\{x := k\}$	$assg(x)$
$x := \mathbf{unop} \ y$ where $x \neq y$	$\{x := \mathbf{unop} \ y\}$	$assg(x)$
$x := \mathbf{unop} \ x$	\emptyset	$assg(x)$
$x := \mathbf{unop} \ k$	$\{x := \mathbf{unop} \ k\}$	$assg(x)$
$x := y \ \mathbf{binop} \ z$ where $x \neq y$ and $x \neq z$	$\{x := y \ \mathbf{binop} \ z\}$	$assg(x)$
$x := y \ \mathbf{binop} \ z$ where $x = y$ or $x = z$	\emptyset	$assg(x)$
$x := y \ \mathbf{binop} \ k$ where $x \neq y$	$\{x := y \ \mathbf{binop} \ k\}$	$assg(x)$
$x := x \ \mathbf{binop} \ k$	\emptyset	$assg(x)$
$x := M[y]$ where $x \neq y$	$\{x := M[y]\}$	$assg(x)$
$x := M[x]$	\emptyset	$assg(x)$
$x := M[k]$	$\{x := M[k]\}$	$assg(x)$
$M[x] := y$	\emptyset	$loads$
$M[k] := y$	\emptyset	$loads$
GOTO l	\emptyset	\emptyset
IF $x \ \mathbf{relop} \ y$ THEN l_t ELSE l_f	\emptyset	\emptyset
$x := \mathbf{CALL} \ f(args)$	\emptyset	$assg(x)$

where $assg(x)$ is the set of all assignments in which x occurs on either left-hand or right-hand side, and $loads$ is the set of all assignments of the form $x := M[\cdot]$.

Figure 11.1: Gen and kill sets for available assignments

```

1:  $i := 0$ 
2:  $a := n * 3$ 
3: IF  $i < a$  THEN loop ELSE end
4: LABEL loop
5:  $b := i * 4$ 
6:  $c := p + b$ 
7:  $d := M[c]$ 
8:  $e := d * 2$ 
9:  $f := i * 4$ 
10:  $g := p + f$ 
11:  $M[g] := e$ 
12:  $i := i + 1$ 
13:  $a := n * 3$ 
14: IF  $i < a$  THEN loop ELSE end
15: LABEL end

```

Figure 11.2: Example program for available-assignments analysis

iteration. So, instead we initialise the *in* and *out* sets for all instructions except the first to the set of all assignments in the program (so we find the maximal solution). The *in* set for the first instruction remains empty, as no assignments is available at the beginning of the program.

When an analysis takes the intersection of the values for the predecessors, we will always initialise sets (except for the first instruction) to the largest possible, so we do not get overly conservative results for loops.

11.2.2 Example of available-assignments analysis

Figure 11.2 shows a program that doubles elements of an array p

Figure 11.3 shows *pred*, *gen* and *kill* sets for each instruction in this program. We represent an assignment by the number of the assignment instruction, so *gen* and *kill* sets are sets of numbers. We will, however, identify identical assignments with the same number, so both the assignment in instruction 2 and the assignment in instruction 13 are identified with the number 2, as can be seen in the *gen* set of instruction 13.

Note that each assignment kills itself, but since it also (in most cases) generates itself, the net effect is to remove all conflicting assignments (including itself) and then adding itself. Assignment 12 ($i := i + 1$) does not generate itself, since i also occurs on the right-hand side. Note, also, that the write to memory in instruction 11 kills instruction 7, as this loads from memory.

i	$pred[i]$	$gen[i]$	$kill[i]$
1		1	1, 5, 9, 12
2	1	2	2
3	2		
4	3, 14		
5	4	5	5, 6
6	5	6	6, 7
7	6	7	7, 8
8	7	8	8
9	8	9	9, 10
10	9	10	10
11	10		7
12	11		1, 5, 9, 12
13	12	2	2
14	13		
15	3, 14		

Figure 11.3: $pred$, gen and $kill$ for the program in figure 11.2

For the fixed-point iteration we initialise the in set of instruction 1 to the empty set and all other in and out sets to the set of all assignments. We need only the assignments that are actually generated by some instructions, *i.e.*, $\{1, 2, 5, 6, 7, 8, 9, 10\}$. We then iterate equations 11.2 and 11.1 as assignments until we reach a fixed-point. Since information flow is forwards, we process the instructions by increasing number and calculate $in[i]$ before $out[i]$. The iteration is shown in figure 11.4. For space reasons, the table is shown sideways and the final iteration (which is identical to iteration 2) is not shown.

11.2.3 Using available assignment analysis for common subexpression elimination

If instruction i is of the form $x := e$ for some expression e and $in[i]$ contains an assignment $y := e$, then we can replace $x := e$ by $x := y$.

If we apply this idea to the program in figure 11.2, we see that at instruction 9 ($f := i * 4$), we have assignment 5 ($b := i * 4$) available, so we can replace instruction 9 by ($f := b$). At instruction 13 ($a := n * 3$), we have assignment 2 ($a := n * 3$) available, so we can replace instruction 13 by $a := a$, which we can omit entirely as it is a no-operation. The optimised program is shown in figure 11.5.

Note that while we could eliminate identical expressions, we could not eliminate

i	Initialisation		Iteration 1		Iteration 2	
	$in[i]$	$out[i]$	$in[i]$	$out[i]$	$in[i]$	$out[i]$
1		1, 2, 5, 6, 7, 8, 9, 10		1		1
2	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1	1, 2	1	1, 2
3	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1, 2	1, 2	1, 2	1, 2
4	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1, 2	1, 2	2	2
5	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1, 2	1, 2, 5	2	2, 5
6	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5	1, 2, 5, 6	2, 5	2, 5, 6
7	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6	1, 2, 5, 6, 7	2, 5, 6	2, 5, 6, 7
8	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7	1, 2, 5, 6, 7, 8	2, 5, 6, 7	2, 5, 6, 7, 8
9	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8	1, 2, 5, 6, 7, 8, 9	2, 5, 6, 7, 8	2, 5, 6, 7, 8, 9
10	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9	1, 2, 5, 6, 7, 8, 9, 10	2, 5, 6, 7, 8, 9	2, 5, 6, 7, 8, 9, 10
11	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	2, 5, 6, 7, 8, 9, 10	2, 5, 6, 7, 8, 9, 10
12	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 8, 9, 10	2, 6, 8, 10	2, 5, 6, 8, 9, 10	2, 6, 8, 10
13	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	2, 6, 8, 10	2, 6, 8, 10	2, 6, 8, 10	2, 6, 8, 10
14	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	2, 6, 8, 10	2, 6, 8, 10	2, 6, 8, 10	2, 6, 8, 10
15	1, 2, 5, 6, 7, 8, 9, 10	1, 2, 5, 6, 7, 8, 9, 10	2	2	2	2

Figure 11.4: Fixed-point iteration for available-assignment analysis

```

1:  $i := 0$ 
2:  $a := n * 3$ 
3: IF  $i < a$  THEN loop ELSE end
4: LABEL loop
5:  $b := i * 4$ 
6:  $c := p + b$ 
7:  $d := M[c]$ 
8:  $e := d * 2$ 
9:  $f := b$ 
10:  $g := p + f$ 
11:  $M[g] := e$ 
12:  $i := i + 1$ 
14: IF  $i < a$  THEN loop ELSE end
15: LABEL end

```

Figure 11.5: The program in figure 11.2 after common subexpression elimination.

the expression $p + f$ in instruction 10 even though it has the same value as the right-hand side of the available assignment 6 ($c := p + b$), since $b = f$. A way of achieving this is to first replace all uses of f by b (which we can do since the only assignment to f is $f := b$) and then repeat common subexpression elimination on the resulting program. If a large expression has two occurrences, we might have to repeat this a large number of times to get the optimal result. An alternative is to keep track of sets of variables that have the same value (a technique called *value numbering*), which allows large common subexpressions to be eliminated in one pass.

Another limitation of the available assignment analysis is when two different predecessors to an instruction have the same expression available but in different variables, *e.g.*, if one predecessor of instruction i has the available assignments $\{x := a + b\}$ and the other predecessor has the available assignments $\{y := a + b\}$. These have an empty intersection, so the analysis would have no available assignments at the entry of instruction i . It is possible to make common subexpression elimination make the expression $a + b$ available in this situation, but this makes analysis and transformation more complex.

Suggested exercises: 11.1.

11.3 Jump-to-jump elimination

When we have an instruction sequence like

```

LABEL  $l_1$ 
GOTO  $l_2$ 

```

we would like to replace all jumps to l_1 by jumps to l_2 . However, there might be chains of such jumps, *e.g.*,

```

LABEL  $l_1$ 
GOTO  $l_2$ 
...
LABEL  $l_2$ 
GOTO  $l_3$ 
...
LABEL  $l_3$ 
GOTO  $l_4$ 

```

We want, in this example, to replace a jump to l_1 with a jump to l_4 directly. To do this, we make a data-flow analysis that for each jump finds its ultimate destination. The analysis is a backwards analysis, as information flows from a label back to the instruction that jumps to it. The *gen* and *kill* sets are quite simple:

instruction	<i>gen</i>	<i>kill</i>
LABEL l	$\{l\}$	\emptyset
GOTO l	\emptyset	\emptyset
IF c THEN l_1 ELSE l_2	\emptyset	\emptyset
any other	\emptyset	the set of all labels

The equations for *in* and *out* are:

$$in[i] = \begin{cases} gen[i] \setminus kill[i] & \text{if } out[i] \text{ is empty} \\ out[i] \setminus kill[i] & \text{if } out[i] \text{ is non-empty} \end{cases} \quad (11.3)$$

$$out[i] = \bigcap_{j \in succ[i]} in[j] \quad (11.4)$$

We initialise all *out* sets to the set of all labels, except the *out* set of the last instruction (which has no successors), which is initialised to empty.

Note that, after the fixed-point iteration, $in[i]$ can only have more than one element if i is part of or jumps to a loop consisting entirely of GOTO instructions and labels. Such infinite loops can occur in valid programs (*e.g.*, as an idle-loop that waits for an interrupt), so we allow this.

When we have reached a fixed-point, we can do the following optimisations:

- At GOTO i , we can replace i by any element in $in[i]$.
- If instruction i is LABEL l and $l \notin in[i]$, there will be no jumps to l after the optimisation, so we can remove instruction i and all following instructions up to just before the next label. This is called *dead code elimination*.
- At IF c THEN i ELSE j , we can replace i by any element in $in[i]$ and j by any element in $in[j]$. If $in[i] = in[j]$, the test is redundant, and we can replace IF c THEN i ELSE j by GOTO l , where l is any element in $in[i]$.

Suggested exercises: 11.2.

11.4 Index-check elimination

When a language requires bounds-checking of array accesses, the compiler must insert tests before each array access to verify that the index is in range. Index-check elimination aims to remove these checks when they are redundant.

To find this, we collect for each program point a set of inequalities that hold at this point. If a bounds check is implied by these inequalities, we can eliminate it.

We will use the conditions in IF-THEN-ELSE instructions as our source for inequalities. Note that, since bounds checks are translated into such instructions, this set of conditions includes all bounds checks.

Conditions all have the form x **relop** p , where x is a variable and p is either a constant or a variable. We only care about inequalities, *i.e.*, conditions of the form $p < q$ or $p \leq q$, where p and q are either variables or constants (but they can not both be constants). Each condition that occurs in the program is translated into a set of inequalities of this form. For example $x = y$ is translated into $x \leq y$ and $y \leq x$. We also generate inequalities for the negation of each condition in the program, so a condition like $x < 10$ generates the inequalities $x < 10$ and $10 \leq x$. A condition $x = y$ generates no inequalities for its negation, as $x \neq y$ can not be expressed as a conjunction of inequalities. This gives us a universe Q of inequalities that we work on. At each point in the program, we want to find which of these inequalities are guaranteed to hold at this point, *i.e.*, a subset of Q .

The idea is that, after executing the instruction IF c THEN l_1 ELSE l_2 , the conditions derived from c will be true at l_1 and the conditions derived from the negation of c will be true at l_2 , *assuming* there are no other jumps to l_1 or l_2 . To ensure this assumption, we insert extra labels and jumps in the program: each instruction of the form IF c THEN l_1 ELSE l_2 is replaced by the sequence

$$in[i] = \begin{cases} \bigcap_{j \in pred[i]} in[j] & \text{if } pred[i] \text{ has more than one element} \\ in[pred[i]] \cup when(c) & \text{if } pred[i] \text{ is IF } c \text{ THEN } i \text{ ELSE } j \\ in[pred[i]] \cup whennot(c) & \text{if } pred[i] \text{ is IF } c \text{ THEN } j \text{ ELSE } i \\ (in[pred[i]] \setminus conds(Q, x)) \cup equal(Q, x, p) & \text{if } pred[i] \text{ is of the form } x := p \\ in[pred[i]] \setminus upper(Q, x) & \text{if } pred[i] \text{ is of the form } x := x + k \text{ where } k \geq 0 \\ in[pred[i]] \setminus lower(Q, x) & \text{if } pred[i] \text{ is of the form } x := x - k \text{ where } k \geq 0 \\ in[pred[i]] \setminus conds(Q, x) & \text{if } pred[i] \text{ is of a form } x := e \text{ not covered above} \\ in[pred[i]] & \text{otherwise} \end{cases}$$

Figure 11.6: Equations for index-check elimination

```

      IF  $c$  THEN  $t$  ELSE  $f$ 
      LABEL  $t$ 
      GOTO  $l_1$ 
      LABEL  $f$ 
      GOTO  $l_2$ 

```

where t and f are new labels that do not occur anywhere else.

This way, if a label has more than one predecessor, these will all be GOTO statements (that do not alter the set of valid inequalities). We can later remove the added code by jump-to-jump elimination as described in section 11.3.

We do not use *out* sets, but write the equations for *in* sets in terms of the *in* sets of the predecessors and the type of instruction of the predecessor, as shown in figure 11.6. Note that information flows forwards: from predecessors to successors. We use the following auxiliary definitions:

- $when(c)$ is the set of inequalities implied by the condition c . These will be elements of Q (the universe of inequalities), since Q was constructed to include all these.
- $whennot(c)$ is the set of inequalities implied by the negation of the condition c . These will, likewise, be elements of Q .

- $conds(Q, x)$ is the set of inequalities from Q that involve x .
- $equal(Q, x, p)$, where p is a variable or a constant, is the set of inequalities from Q that are implied by the equality $x = p$. For example, if $Q = \{x < 10, 10 \leq x, 0 < x, x \leq 0\}$ then $equal(Q, x, 7) = \{x < 10, 0 < x\}$.
- $upper(Q, x)$ is the set of inequalities from Q that have the form $p < x$ or $p \leq x$, where p is a variable or a constant.
- $lower(Q, x)$ is the set of inequalities from Q that have the form $x < p$ or $x \leq p$, where p is a variable or a constant.

Normally, any assignment to a variable invalidates all inequalities involving that variable, but we have made some exceptions: If we assign a constant or variable to a variable, we check all the possible inequalities and add those that are implied by the assignment. Also, if x increases, we invalidate all inequalities that bound x from above but keep those that bound x from below, and if x decreases, we invalidate the inequalities that bound x from below but keep those that bound x from above. We can add more special cases to make the analysis more precise, but the above are sufficient for the most common cases.

We initialise all *in* sets to Q , except the *in* set for the first instruction, which is initialised to the empty set.

After the data-flow analysis reaches a fixed-point, the inequalities in $in[i]$ are guaranteed to hold at instruction i . So, if we have an instruction i of the form IF c THEN l_t ELSE l_f and c is implied by the inequalities in $in[i]$, we can replace the instruction by GOTO l_t . If the negation of c is implied by the inequalities in $in[i]$, we can replace the instruction by GOTO l_f .

We illustrate the analysis by an example. Consider the following for-loop and assume that the array a is declared to go from 0 to 10.

```
for i:=0 to 9 do
  a[i] := 0;
```

This loop can be translated (with index check) into the intermediate code shown in figure 11.7.

The set Q of possible inequalities in the program are derived from the conditions in the three IF-THEN-ELSE instructions and their negations, *i.e.*, $Q = \{i \leq 9, 9 < i, i < 0, 0 \leq i, 10 < i, i \leq 10\}$.

We leave the fixed-point iteration and check elimination as an exercise to the reader, but note that the assignment $i := 0$ in instruction 1 implies the inequalities $\{i \leq 9, 0 \leq i, i \leq 10\}$ and that the assignment $i := i + 1$ in instruction 11 preserves $0 \leq i$ but invalidates $i \leq 9$ and $i \leq 10$.

```

1:  $i := 0$ 
2: LABEL for1
3: IF  $i \leq 9$  THEN for2 ELSE for3
4: LABEL for2
5: IF  $i < 0$  THEN error ELSE ok1
6: LABEL ok1
7: IF  $i > 10$  THEN error ELSE ok2
8: LABEL ok2
9:  $t := i * 4$ 
10:  $t := a + t$ 
11:  $M[t] := 0$ 
12:  $i := i + 1$ 
13: GOTO for1
14: LABEL for3

```

Figure 11.7: Intermediate code for for-loop with index check

Suggested exercises: 11.3.

11.5 Limitations of data-flow analyses

All of the data-flow analyses we have seen above are approximations: They will not always accurately reflect what happens at runtime: The index-check analysis may fail to remove a redundant index check, and the available assignment analysis may say an assignment is unavailable when, in fact, it is available.

In all cases, the approximations err on the safe side: It is better to miss an opportunity for optimisation than to make an incorrect optimisation. For liveness analysis, this means that if you are in doubt about a variable being live, you had better say that it is, as assuming it dead might cause its value to be overwritten. When available assignment analysis is used for common subexpression elimination, saying that an assigning is available when it is not may make the optimisation replace an expression by a variable that does not always hold the same value as the expression, so it is better to leave an assignment out of the set if you are in doubt.

It can be shown that no compile-time analysis that seeks to uncover nontrivial information about the run-time behaviour of programs can ever be completely exact. You can make more and more complex analyses that get closer and closer to the exact result, but there will always be programs where the analysis is not precise. So a compiler writer will have to be satisfied with analyses that find most cases where an optimisation can be applied, but misses some.

11.6 Loop optimisations

Since many programs spend most of their time in loops, it is worthwhile to study optimisations specific for loops.

11.6.1 Code hoisting

One such optimisation is recognising computations that are repeated in every iteration of the loop without changing the values involved, *i.e.*, loop-invariant computations. We want to lift such computations outside the loop, so they are performed only once. This is called *code hoisting*.

We saw an example of this in section 11.2.3, where calculation of $n * 3$ was done once before the loop and subsequent re-computations were replaced by a reference to the variable a that holds the value of $n * 3$ computed before the loop. However, it is only because there was an explicit computation of $n * 3$ before the loop, that we could avoid re-computation inside the loop: Otherwise, the occurrence of $n * 3$ inside the loop would not have any available assignment that can replace the calculation.

So our aim is to move or copy loop-invariant assignments to before the loop, so their result can be reused inside the loop. Moving a computation to before the loop may, however, cause it to be computed even when the loop is not entered. In addition to causing unnecessary computation (which goes against the wish for optimisation), such computations can cause errors when the precondition (the loop condition) is not satisfied. For example, if the invariant computation is a memory access, the address may be valid only if the loop is entered.

A common solution to this problem is to unroll the loop once: A loop of the form (using C-like syntax):

```
while (cond) {  
    body  
}
```

is transformed to

```
if (cond) then {  
    body  
    while (cond) {  
        body  
    }  
}
```

Similarly, a test-at-bottom loop of the form

```

do
    body
while (cond)

```

can be unrolled to

```

body
while (cond) {
    body
}

```

Now, we can safely calculate the invariant parts in the first copy of the body and reuse the results in the loop. If the compiler does common subexpression elimination, this unrolling is all that is required to do code hoisting – assuming the unrolling is done before common-subexpression elimination. Unrolling of loops is most easily done at source-code level (*i.e.*, on the abstract syntax tree), so this is no problem. This unrolling will, of course, increase the size of the compiled program, so it should be done with care if the loop body is large.

11.6.2 Memory prefetching

If a loop goes through a large array, it is likely that parts of the array will not be in the cache of the processor. Since access to non-cached memory is *much* slower than access to cached memory, we would like to avoid this.

Many modern processors have *memory prefetch instructions* that tell the processor to load the contents of an address into cache, but unlike a normal load, a memory prefetch does not cause errors if the address is invalid, and it returns immediately without waiting for the load to complete. So a way to ensure that an array element is in the cache is to issue a prefetch of the array element well in advance of its use, but not so well in advance that it is likely that it will be evicted from the cache between the prefetch and the use. Given modern cache sizes and timings, 25 to 10000 cycles ahead of the use is a reasonable time for prefetching – less than 25 increases the risk that the prefetch is not completed before the use, and more than 10000 increases the chance that the value will be evicted from the cache before use.

A prefetch instruction usually loads an entire cache line, which is typically four or eight words, so we do not have to explicitly prefetch every array element – every fourth element is enough.

So, assume we have a loop that adds the elements of an array:

```

sum = 0;
for (i=0; i<100000; i++)
    sum += a[i];
}

```

we can rewrite this to

```
sum = 0;
for (i=0; i<100000; i++) {
    if (i&3 == 0) prefetch a[i+32];
    sum += a[i];
}
```

where `prefetch a[i+32]` prefetches the element of `a` that is 32 places after the current element. The number 32 is rather arbitrary, but makes the number of cycles between prefetch and use lie in the interval mentioned above. Note that we used the test `i&3==0`, which is equivalent to `i%4==0`, but somewhat faster.

We do not have to worry about prefetching past the end of the array – prefetching will never cause runtime errors, so at worst we prefetch something that we do not need.

While this transformation adds a test (that takes time), the potential savings by having all array elements in cache before use are much larger. The overhead of testing can be reduced by unrolling the loop body:

```
sum = 0;
for (i=0; i<100000; i++) {
    prefetch a[i+32];
    sum += a[i];
    i++;
    sum += a[i];
    i++;
    sum += a[i];
    i++;
    sum += a[i];
}
```

This should, of course, only be done if the loop body is small. We have exploited that the number of iterations is a multiple of 4, so the exit test is not needed at every increment of `i`. If we do not know this, the exit test must be replicated after each increase of `i`, like shown here:

```

sum = 0;
for (i=0; i<n; i++) {
    prefetch a[i+32];
    sum += a[i];
    if (++i < n) {
        sum += a[i];
        if (++i < n) {
            sum += a[i];
            if (++i < n) {
                sum += a[i];
            }
        }
    }
}

```

In a nested loop that accesses a multi-dimensional array, you can prefetch the next row while processing the current. For example, the loop

```

sum = 0;
for (i=0; i<1000; i++)
    for (j=0; j<1000; j++)
        sum += a[i][j];
    }
}

```

can be transformed to

```

sum = 0;
for (i=0; i<1000; i++)
    for (j=0; j<1000; j++)
        if (j&3 == 0) prefetch a[i+1][j];
        sum += a[i][j];
    }
}

```

Again, we can unroll the body of the inner loop to reduce the overhead.

11.7 Optimisations for function calls

Modern coding styles use frequent function (or method) calls, so optimising function calls is as worthwhile as optimising loops.

Basically, optimisation of function calls attempt to reduce the overhead associated with call sequences, prologues and epilogues (see chapter 10). We will see a few ways of doing this below.

11.7.1 Inlining

Inlining means replacing a function call by a copy of the body of the function, with some glue code to replace parameter and result passing.

If we have a call (using C-style syntax)

$$x = f(exp_1, \dots, exp_n);$$

and the function f is defined as

```

type0 f(type1 x1, ..., typen xn)
{
    body
    return(exp);
}

```

where *body* represents the body of the function (apart from the `return` statement) and the shown `return` statement is the only exit point of the function, we can replace the call by the block

```

{
    type1 x1 = exp1;
    ...
    typen xn = expn;
    body
    x = exp;
}

```

Note that if, say, exp_n refers to a variable with the same name as x_1 , it will refer to a different instance of x_1 after the transformation. So we need to avoid variables in the inlined function shadowing variables used in the function call statement. We can achieve this by renaming the variables in the inlined function to new, previously unused, names before it is inlined.

Note that, unless the body of the inlined function is very small, inlining causes the program to grow in size. Hence, it is common to put a limit on the size of functions that are inlined. What the limit is depends on the desired balance between speed and size. But even when optimising for speed, care should be taken not to inline too much, as large programs can run slower than small programs due to worse cache behaviour.

Care must be taken if you inline calls recursively: If the inlined body contains a call that is also inlined, and this again contains a call that is inlined, and so on, we might continue inlining forever. So it is common to limit inlining to only one or two levels deep or treat (mutually) recursive functions as special cases.

11.7.2 Tail-call optimisation

A *tail call* is a call that happens just before a return.

As an example, assume we in a function f have (using C-style notation) a statement $\text{return}(g(x, y));$. Clearly, f returns just after g returns, and the result of f is the result of g . We want to combine the call sequence for the call to g with the epilogue of f . We call this *tail-call optimisation*.

We will exploit the following observations:

- No variables in f are live after the call to g (since there are not any uses of variables after the call), so there will be no need to save and restore caller-saves variables around the call.
- If we can eliminate all of the epilogue of f except for the return-jump to f 's caller, we can instead make g return directly to f 's caller, hence skipping f 's epilogue entirely.
- If f 's frame holds no useful information at the time we call g (or we can be sure that g does not overwrite any useful information in the frame), we can reuse the frame for f as the frame for g .

We will look at this in more detail below.

If we assume a simple stack-based caller-saves strategy like the one shown in figures 10.2 and 10.3, the combined call sequence of the call to g and the epilogue of f becomes:

```

 $M[FP + 4 * m + 4] := R0$ 
...
 $M[FP + 4 * m + 4 * (k + 1)] := Rk$ 
 $FP := FP + framesize$ 
 $M[FP + 4] := a_1$ 
...
 $M[FP + 4 * n] := a_n$ 
 $M[FP] := returnaddress$ 
GOTO  $g$ 
LABEL  $returnaddress$ 
 $result := M[FP + 4]$ 
 $FP := FP - framesize$ 
 $R0 := M[FP + 4 * m + 4]$ 
...
 $Rk := M[FP + 4 * m + 4 * (k + 1)]$ 
 $M[FP + 4] := result$ 
GOTO  $M[FP]$ 

```

Since there are no live variables after the call to *g*, we can eliminate the saving and restoring of R_0, \dots, R_k , yielding the following simplified code:

```

     $FP := FP + framesize$ 
     $M[FP + 4] := a_1$ 
    ...
     $M[FP + 4 * n] := a_n$ 
     $M[FP] := returnaddress$ 
    GOTO g
    LABEL returnaddress
     $result := M[FP + 4]$ 
     $FP := FP - framesize$ 
     $M[FP + 4] := result$ 
    GOTO  $M[FP]$ 

```

We now see that all that happens after we return is an adjustment of *FP*, a copy of the result from *g*'s frame to *f*'s frame and a jump to the return address stored in *f*'s frame.

What we now want is to reuse *f*'s frame for *g*'s frame. We do this by not adding and subtracting *framesize* from *FP*, so we get the following simplified code:

```

     $M[FP + 4] := a_1$ 
    ...
     $M[FP + 4 * n] := a_n$ 
     $M[FP] := returnaddress$ 
    GOTO g
    LABEL returnaddress
     $result := M[FP + 4]$ 
     $M[FP + 4] := result$ 
    GOTO  $M[FP]$ 

```

It is immediately evident that the two instructions that copy the result from and to the frame cancel out, so we can simplify further to

```

     $M[FP + 4] := a_1$ 
    ...
     $M[FP + 4 * n] := a_n$ 
     $M[FP] := returnaddress$ 
    GOTO g
    LABEL returnaddress
    GOTO  $M[FP]$ 

```

We also see an unfortunate problem: Just before the call to *g*, we overwrite *f*'s return address in $M[FP]$ by *g*'s return address, so we will not return correctly to *f*'s

caller (we will, instead, get an infinite loop). However, since all that happens after we return from *g* is a return from *f*, we can make *g* return directly to *f*'s caller. We do this simply by not overwriting *f*'s return address. This makes the instructions after the jump to *g* unreachable, so we can just delete them. This results in the following code:

$$\begin{aligned} M[FP + 4] &:= a_1 \\ \dots \\ M[FP + 4 * n] &:= a_n \\ \text{GOTO } g \end{aligned}$$

With this tail-call optimisation, we have eliminated the potential ill effects of reusing *f*'s frame for *g*. Not only have we shortened the combined call sequence for *g* and epilogue for *f* considerably, we have also saved stack space. Functional programming languages rely on this space saving, as they often use recursive tail calls where imperative languages use loops. By doing tail-call optimisation, an arbitrarily long sequence of recursive tail calls can share the same stack frame, so only a constant amount of stack space is used.

In the above, we relied on a pure caller-saves strategy, since the absence of live variables after the call meant that there would be no saving and restoring of caller-saves registers around the call. If a callee-saves strategy or a mixed caller-saves/callee-saves strategy is used, there will still be no saving and restoring around the call, but *f*'s epilogue would restore the callee-saves registers that *f* saved in its own prologue. This makes tail-call optimisation a bit more complicated, but it is normally possible (with a bit of care) to move the restoring of callee-saves registers to *before* the call to *g* instead of waiting until after *g* returns. This allows *g* to overwrite *f*'s frame, which no longer holds any useful information (except the return address, which we explicitly avoid overwriting).

Exercise 11.4 asks you to do such tail-call optimisation for a mixed caller-saves/callee-saves strategy.

11.8 Specialisation

Modern programs consist mainly of calls to predefined library functions (or procedures or methods). Calls to such library functions often fix some of the parameters to constants, as the function is written more generally than needed for the particular call. For example, a function that raises a number to a power is often called with a constant as the power, say, `power(x, 5)` which raises *x* to its fifth power. In such cases, the generality of the library function is wasted, and speed can be gained by using a more specific function, say, one that raises its argument to the fifth power.

A possible implementation of the general power function (in C) is:


```

double power(double x, int n)
{
    double p=1.0;
    while (n>0)
        if (n%2 == 0) {
            x = x*x;
            n = n/2;
        } else {
            p = p*x;
            n = n-1;
        }
    return(p);
}

```

If we have a call `power(x,5)`, we can replace this by a call `power5(x)` to a specialised function. We now need to add a definition of this specialised function to the program. The most obvious idea would be to take the above code for the `power` function and replace all occurrences of `n` by 5, but this will not work, as `n` changes value inside the body of `power`. What we do instead is to observe the following:

1. The loop condition depends only on `n`.
2. Every change to `n` depends only on `n`.

So it is safe to unroll the loop at compile-time, doing all the computations on `n`, but leaving the computations on `p` and `x` for run-time. This yields the following specialised definition:

```

double power5(double x)
{
    double p=1.0;
    p = p*x;
    x = x*x;
    x = x*x;
    p = p*x;
    return(p);
}

```

Executing `power5(x)` is, obviously, a lot faster than executing `power(x,5)`. Since `power5` is fairly small, we can additionally inline the call, as described in section 11.7.1.

This kind of specialisation may not always be applicable, even if a function call has constant parameters, for example if the call was

$\text{power}(3.14159, p)$, where p is not a constant, but when the method is applicable, the speedup can be dramatic.

Similar specialisation techniques are used in C++ compilers for compiling templates: When a call specifies template parameters, the definition of the template is specialised with respect to the actual template parameters. Since templates in C++ can be recursively defined, an infinite number of specialised versions might be required. Most C++ compilers put a limit on the recursion level of template instantiations and stop with an error message when this limit is exceeded.

Suggested exercises: 11.5.

11.9 Further reading

We have covered only a small portion of the optimisations that are found in optimising compilers. More examples of optimisations (including value numbering) can be found in advanced compiler textbooks, such as [4, 7, 9, 35].

A detailed treatment of program analysis can be found in [36]. Specialisation techniques like those mentioned in section 11.8 are covered in more detail in [17, 21]. The book [34] has good articles on both program analysis and transformation.

Additionally, the conferences “Compiler Construction” (CC), “Programming Language Design and Implementation” (PLDI) and other programming-language-oriented conferences often present new optimisation techniques, so past proceedings from these is a good source for advanced optimisation methods.

Exercises

Exercise 11.1

In the program in figure 11.2, replace instructions 13 and 14 by

```
13:  $h := n * 3$ 
14: IF  $i < h$  THEN loop ELSE end
```

- a) Repeat common subexpression elimination on this modified program.
- b) Repeat, again, common subexpression elimination on the modified program, but, prior to the fixed-point iteration, initialise all sets to the empty set instead of the set of all assignments.

What differences does this make to the final result of fixed-point iteration, and what consequences do these differences have for the optimisation?

Exercise 11.2

Write a program that has jumps to jumps and perform jump-to-jump optimisation of it as described in section 11.3. Try to make the program cover all the three optimisation cases described at the end of section 11.3.

Exercise 11.3

- a) As described in the beginning of section 11.4, add extra labels and gotos for each IF-THEN-ELSE in the program in figure 11.7.
- b) Do the fixed-point iteration for index-check elimination on the result.
- c) Eliminate the redundant tests.
- d) Do jump-to-jump elimination as described in section 11.3 on the result to remove the extra labels and gotos introduced in question a.

Exercise 11.4

Section 11.7.2 describes tail-call optimisation for a pure caller-saves strategy. Things become somewhat more complicated when you use a mixed caller-saves/callee-saves strategy.

Using the call sequence from figure 10.10 and the epilogue from figure 10.9, describe how the combined sequence can be rewritten to get some degree of tail-call optimisation. Your main focus should be on reusing stack space, and secondarily on saving time.

Exercise 11.5

Specialise the power function in section 11.8 to $n = 12$.

Chapter 12

Memory management

12.1 Introduction

In chapter 7, we mentioned that arrays, records and other multi-word objects could be allocated either statically, on the stack or in the heap. We will now look into more detail of how these three kinds of allocation can be implemented and what their relative merits are.

12.2 Static allocation

Static allocation means that the data is allocated at a place in memory that has both known size and address at compile time. Furthermore, the allocated memory stays allocated throughout the execution of the program.

Most modern computers divide their logical address space into a text section (used for code) and a data section (used for data). Assemblers (programs that convert symbolic machine code into binary machine code) usually maintain “current address” pointers to both the text area and the data area. They also have pseudo-instructions (directives) that can place labels at these addresses and move them. So you can allocate space for, say, an array in the data space by placing a label at the current-address pointer in the data space and then move the current-address pointer up by the size of the array. The code can use the label to access the array. Allocation of space for an array A of 1000 32-bit integers (*i.e.*, 4000 bytes) can look like this in symbolic code:

```
.data          # go to data area for allocation
baseofA:       # label for array A
.space 4000    # move current-address pointer up 4000 bytes
.text         # go back to text area for code generation
```

The base address of the array A is at the label `baseofA`.

The assembler (possibly assisted by a linker) translates the symbolic code to binary code where references to labels are replaced by references to numeric addresses.

In the programming language C, all global variables, regardless of size, are statically allocated.

12.2.1 Limitations

The size of an array must be known at compile time if it is statically allocated, and the size can not change during execution. Furthermore, the space is allocated throughout the entire execution of the program, even if the array is only in use for a small fraction of this time. In essence, there is little reuse of statically allocated space within one program execution. The programming language Fortran allows the programmer to specify that several statically allocated arrays can share the same space, but this feature is rarely seen in modern languages. It is, in theory, possible to analyse the liveness of arrays in the same way that we analysed liveness of local variables in chapter 9, and use this to define interference between arrays in such a way that two arrays that do not interfere can share the same space. This is, however, rarely done, as it typically requires an expensive analysis of the entire program.

12.3 Stack allocation

As mentioned in chapter 10, the call stack can also be used to allocate arrays and other data structures. This is done by making room in the current (topmost) frame on the call stack. This is done by moving the frame pointer and/or the stack pointer.

Stack allocation has both advantages and disadvantages:

- The space for the array is freed up when the function that allocated the array returns (as the frame in which the array is allocated is taken off the stack). This allows easy reuse of the memory for later stack allocations in the same program execution.
- Allocation is fairly quick: You just move the stack pointer or the frame pointer. Releasing the memory again is also fast – again, you only move a pointer.
- The size of the array need not be known at compile time: The frame is at runtime created to be large enough to hold the array. If the size of the topmost frame can be extended after it is created, you can even stack-allocate arrays during execution of the function body.
- An unbounded number of arrays can be allocated at runtime, since recursive functions can allocate an array in each invocation of the function.

- The array will not survive return from the function in which it is allocated, so it can be used only locally inside this function and functions that it calls.
- Once allocated, the array can not be extended, as you can not (easily) “squeeze” more space into the middle of the stack, where the array is allocated, nor can you reallocate it at the top of the stack (unless it is the same frame), as the reallocated array will be released earlier than the original array.

In C, arrays that are declared locally in a function are stack allocated (unless declared `static`, in which case they are statically allocated). C allows you to return pointers to stack-allocated arrays and it is up to the programmer to make sure he never follows a pointer to an array that is no longer allocated. This often goes wrong. Other languages (like Pascal) avoid the problem by not allowing pointers to stack-allocated data to be returned from a function.

12.4 Heap allocation

The limitations of static allocation and stack allocation are often too restricting: You might want arrays that can be resized or which survive the function invocation in which they are allocated. Hence, most modern languages allow heap allocation, which is also called dynamic memory allocation.

Data that is heap allocated stays allocated until the program execution ends, until the data is explicitly deallocated by a program statement or until the data is deemed dead by a run-time memory-management system, which then deallocates it. The size of the array (or other data) need not be known until it is allocated, and if the size needs to be increased later, a new array can be allocated and references to the old array can be redirected to point to the new array. The latter requires that all references to the old array can be tracked down and updated, but that can be arranged, for example, by letting all references to an array go through an indirection node that is never moved.

Languages that use heap allocation can be classified by how data is deallocated: Explicitly by commands in the program or automatically by the run-time memory-management system. The first is called “manual memory management” and the latter “automatic memory management” or “automatic garbage collection”. We will look into both of these in more detail below.

12.5 Manual memory management

Most operating systems allow a program to allocate and free (deallocates) chunks of memory while the program runs. These chunks are, typically, fairly large and operating-system calls to allocate and free memory can be slow. So, it is normal

for programs that use heap allocation to allocate a large chunk of memory from the operating system and then manage this itself when the program allocates and deallocates data on the heap. This is normally done by library functions.

In C, these are called `malloc()` and `free()`. `malloc(n)` allocates a block of at least *n* bytes on the heap and returns a pointer to this block. If there is not enough memory available to allocate a new block of size *n*, `malloc(n)` returns a null pointer. `free(p)` takes a pointer to a block that was previously allocated by `malloc()` and deallocates this. If *p* is a pointer to a block that was previously allocated by `malloc()` and not already deallocated, `free(p)` will always succeed, otherwise the result is undefined. The behaviour of following a pointer to a deallocated block is also undefined and is the source of many errors in C programs.

Object-oriented languages often allocate and free memory with object constructors and destructors, but these typically work the same way as `malloc()` and `free()`, except that object constructors and destructors may do more than just allocation and deallocation, *e.g.*, initialising fields or opening and closing files.

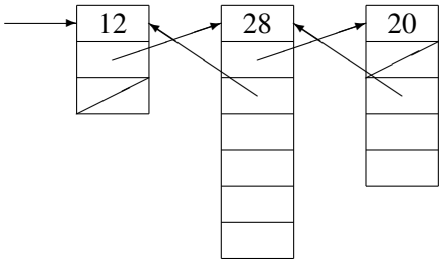
12.5.1 A simple implementation of `malloc()` and `free()`

Initially, the allocation library will allocate a large chunk of memory from the operating system. If this turns out to be too small to satisfy `malloc()` calls, the library will allocate another chunk from the operating system and so on, until the operating system refuses to allocate more, in which case the call to `malloc()` will fail (and return a null pointer).

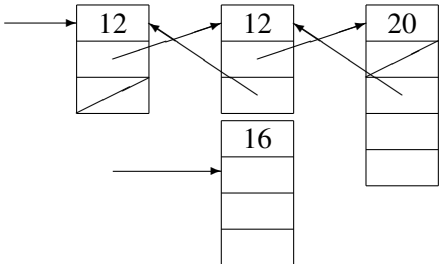
All the chunks that are allocated from the operating system and all blocks that have been freed by the program (by calling `free()`) are linked in a list called *the free list*: The first word of each block (or chunk) contains the size *s* (in bytes) of the block, the second word contains a pointer to the next block and the third word a pointer to the previous block, making the free-list a doubly-linked list. In the last block in the free list, the next-pointer is a null pointer and in the first block, the previous-pointer is a null pointer. Note that a block must be at least three words in size, so it can hold the size field and the two pointer fields. Figure 12.1(a) shows an example of a free list containing three small blocks. The number of bytes in a word (*wordsize*) is assumed to be 4. A null pointer is shown as a slash.

A call to `malloc(n)` with $n \geq 2 \cdot \text{wordsize}$ will now do the following:

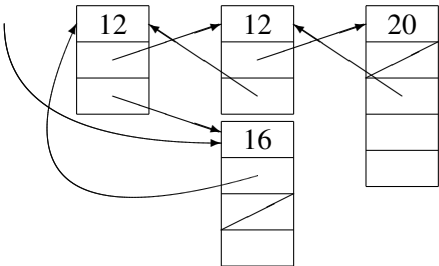
1. Search through the free list for a block of size at least $n + \text{wordsize}$. If none is found, ask the operating system for a new chunk of memory (that is at least large enough to accommodate the request) and put this at the end of the free list. If this fails, return a null pointer (indicating failure to allocate).
- 2a. If the found block is just large enough (at most $n + 3 \cdot \text{wordsize}$ bytes), remove it from the free list (adjusting the pointers of the previous and next blocks)



(a) The initial free list.



(b) After allocating 12 bytes.



(c) After freeing the same 12 bytes.

Figure 12.1: Operations on a free list

and return a pointer to the second word. The first word still contains the size of the block.

- 2b. If the block is more than $n + 3 \cdot \text{wordsize}$ bytes, it is split in two: The last $n + \text{wordsize}$ bytes (rounded up to a multiple of the word size) is made into a new block, the first word of which holds its size (*i.e.*, $n + \text{wordsize}$). A pointer to the word after the size field is returned. $n + \text{wordsize}$ is subtracted from the size field of the original block, which stays in the free list.

We assume n is a multiple of the word size and at least $2 \cdot \text{wordsize}$. Otherwise, `malloc()` will round n up to the nearest acceptable value.

Figure 12.1(b) shows the free list from figure 12.1(a) after a call `malloc(12)`. The second block in the free list has been split in two, and a pointer to the second word of the second part has been returned. Note that the split-up block has size 16, since it needs to hold the requested 12 bytes and four bytes for the size field.

A call to `free(p)` adds the block at p to the front of the free list: The second word of the block is updated to point to the first block in the free list, the previous-pointer of the first element in the free-list and the free-list pointer are updated to point to $p - \text{wordsize}$, *i.e.*, to the size field of the released block. Figure 12.1(c) shows the free list from figure 12.1(b) after freeing the previously-allocated 12-byte block.

As allocating memory involves a search for a block that is large enough, the time is, in the worst case, linear in the number of blocks in the free list, which is proportional to the number of calls to `free()` that the program has made. Freeing a block is done in constant time.

Another problem with this implementation is *fragmentation*: We split blocks into smaller blocks but never join released blocks, so we will, over time, accumulate a large number of small blocks in the free list. This will increase the search time when calling `malloc()` and, more seriously, a call to `malloc(n)` may fail even though there is sufficient free memory – it is just divided into a lot of blocks that are all smaller than n bytes. For example, if we with the free list in figure 12.1(c) try to allocate 20 bytes, there is no block that is large enough, so if the operating system will not give more memory to the memory allocator, the call to `malloc()` will fail.

It is possible to resize heap-allocated arrays if all accesses to the array happen through an indirection node: A one-word node that contains a pointer to the array. When the array is resized, a new block is allocated, the elements of the old array are copied to the new array and the old array is freed. Finally, the indirection node is made to point to the new array. Using an indirection node makes array accesses slower, but the actual address can be loaded into a register if repeated accesses to the array are made.

Suggested exercises: 12.1.

12.5.2 Joining freed blocks

We can try to remedy the fragmentation problem by joining a freed block to a block already in the free list: When we free a block, we don't just add it to the front of the free list, but search through the list to see if it contains blocks that are adjacent to the freed block. If this is the case, we join these blocks into one large block. Otherwise, we add the freed block as a new block to the free list.

With such joining, freeing the 12-byte block in figure 12.1(b) will join this to the second block in the free list, restoring the free list to the state shown in figure 12.1(a).

Now `malloc()` and `free()` both use time that is linear in the length of the free list (which may, however, be shorter). In the worst case, no blocks are ever joined, so we could just be wasting time in trying to join blocks. But if, for example, n same-sized objects are allocated with no other calls to `malloc()` or `free()` in between, and they are later freed, also with no calls to `malloc()` or `free()` in between, all blocks that were split by the allocation will be joined again by the deallocation. Overall, joining blocks will reduce fragmentation, but can not eliminate it, as there can still be freed blocks that can not be joined with previously freed blocks.

To make joining of blocks faster, we can allow a freed block to be joined with the next adjacent block only: Since the size of the freed block is known, it is easy to find the address of the next block. We now just need a faster way than searching the free list to see if the neighbouring is free or allocated. We can use a bit in the size field for this. Sizes are always a full number of machine words, so the size will be an even number. Adding one to the size (making it an odd number) when a block is allocated can indicate that it is in use. When freeing a block, we subtract one from the size, look that many bytes ahead to see if the size field stored in this word is even. If so, we add the sizes and store the sum in the size field of the released block, remove its neighbour from the free list and add the joined block to it.

Note that the order of freeing blocks can determine if blocks are joined: Two neighbours are only joined if the one at the highest address is freed before its lower-address neighbour.

With this modification, `free()` is again constant time, but we merge fewer blocks than otherwise. For example, freeing the 12-byte block in figure 12.1(b) will not make it join with the second block in the free list, as this neighbours the freed block on the wrong side.

Suggested exercises: 12.2.

12.5.3 Sorting by block size

To reduce the time used for searching for a sufficiently large block when calling `malloc()`, we can keep free blocks sorted by size. A common strategy for this is limiting block sizes (including size fields) to powers of two and keeping a free list for each size. Given the exponential growth in sizes, the number of free lists is modest, even on a system with large memory.

When `malloc(n)` is called, we find the nearest power of two that is at least $n + \text{wordsize}$, remove the first block in the free list for that size and return this. If the relevant free list is empty, we take the first block in the next free list and split this into two equal-size blocks, put one of these into the previous free list and return the other. If the next free list is also empty, we go to the next again and so on. If all free lists are empty, we allocate a new chunk from the operating system. The worst-case time for a single call to `malloc()` is, hence, logarithmic in the size of the total heap memory.

To reduce fragmentation and allow fast joining of freed blocks, some memory managers restrict blocks to be aligned to block size. So a 16-word block would have an address that is 16-word aligned, a 32-word block would be 32-word aligned, and so on. When joining blocks, we can only do so if the resulting joined block is aligned to its size, so some blocks can only be joined with the next (higher-address) block of the same size and some only with the previous (lower-address) block of the same size. Two blocks that can be joined are called “buddies”, and each block has a unique buddy. Like in section 12.5.2, we use a bit in the size field to indicate that a block is in use. When we free a block, we clear this bit and look at the size field of the buddy. If the buddy is not in use (indicated by having an even size field) and is the same size as the freed block, we merge the two blocks simply by doubling the size in the size field of the buddy with the lowest address. We remove the already-free buddy from its free list and add the joined block to the next free list. Once we have joined two buddies, the newly joined block might be joined with its own buddy in the same way. This can cause a chain of joinings, making `free()` take time logarithmic in the size of memory. The order of freeing blocks does not affect which blocks are joined, unlike in the system described in section 12.5.2, where a newly freed block could only be joined with the next adjacent block.

With the buddy system, we get both `malloc()` and `free()` down to logarithmic time, but we use extra space by rounding allocation sizes up to powers of two. This waste is sometimes called *internal fragmentation*, where having many small free blocks is called *external fragmentation*.

Even when joining blocks as described above, we can still get external fragmentation. For example, we might get a heap that alternates between four-word blocks that are in use and four-word blocks that are free. The free blocks can have arbitrarily large combined size, but, since no two of these are adjacent, we can not

join them and, hence, not allocate a block larger than four words.

A variant of the above method uses block sizes that are generalised Fibonacci sequences of numbers. In the original Fibonacci sequence, a number is the sum of the two immediately preceding numbers in the sequence, so the sequence is 0, 1, 1, 3, 5, 8, 13, 21, 34 and so on.

In a generalised Fibonacci sequence, the next number in the sequence is the sum of two previous numbers in the sequence, but not necessarily the immediately preceding numbers. For example, we can define a sequence where the first four numbers are 0, 1, 2, 3 and where the number at position $n + 4$ is the sum of the numbers at position n and $n + 3$. Hence, the sequence continues 3, 4, 6, 9, 12 and so on.

If block sizes are numbers in a generalised Fibonacci sequence, we can split a block into two that have sizes corresponding to earlier numbers in the sequence. In the example above, we can split a block of size 12 into blocks of size 3 and 9. The advantage over using blocks of two is that the difference in size is smaller, so there is less waste (internal fragmentation) when you have to round a requested size up to the nearest available block size. For example, in the sequence above the next number is approximately 38% larger than the preceding number, so internal fragmentation can be no more than 38% overhead, where it for power-of-two block sizes can be up to 100%. By choosing other generalised Fibonacci sequences, you can get internal fragmentation even lower. The disadvantage is that you need more free lists and that joining of blocks becomes more complicated.

Suggested exercises: 12.3, 12.4.

12.5.4 Summary of manual memory management

Manual memory management means that both allocation and deallocation of heap-allocated objects (arrays etc.) is explicit in the code. In C, this is done using the library functions `malloc()` and `free()`.

It is possible to get allocation (`malloc()`) time down to logarithmic in the size of memory. Freeing (`free()`) can be made constant time, but external fragmentation can be very bad, so it is more common to make `free()` join blocks to reduce external fragmentation. With such joining, the time for a call to `free()` also becomes logarithmic in the size of memory.

Manual memory management has two serious problems:

- Manually inserting calls to `free()` in the program is error-prone and a common source of errors. If memory is used after it is erroneously freed, the results can be unpredictable and compromise the security of a software system. If memory is not freed or freed too late, the amount of in-use memory can increase, which is called a space leak.

- The available memory can over time be split into a large number of small blocks. This is called (external) fragmentation.

Space leaks and fragmentation can be quite serious in systems that run for a long time, such as telephone exchange systems.

12.6 Automatic memory management

Because manual memory management is error prone, many modern languages support automatic memory management, also called *garbage collection*.

With automatic memory management, heap allocation is still done in the same way as with manual memory management: By calling `malloc()` or invoking an object constructor. But the programmer does not need to call `free()` or object destructors: It is up to the compiler to generate code that frees memory. This is typically not done only by making the compiler insert calls to `free()` in the generated code, as finding the right places to do so requires very complicated analysis of the whole program.

Instead, a block is freed when the program can no longer access it. If no pointer to a block exists, the program can no longer access the block, so this is typically the property that is used when freeing blocks. Note that this can cause blocks to stay allocated even if the program will never access them: Just because the program *can* access a block does not mean that it *will* ever do so. To reduce space leaks caused by keeping pointers to blocks that will never be accessed, it is often advised that programmers overwrite (with null) pointers that will definitely never be followed again. The benefit of automatic memory management is, however, reduced if programmers spend much time considering when they can safely overwrite pointers with null, and such modifications can make programs harder to maintain. So overwriting pointers for the purpose of freeing space should only be considered if actual space leaks are observed.

Automatic memory management usually require pointers only to point to the start of blocks, as it can be difficult to find the start of a block that a pointer points to if the pointer can point anywhere within the block.

We will look at two types of automatic memory management: Reference counting and tracing garbage collection.

12.7 Reference counting

Reference counting uses the property that if no pointer to a block of memory exists, then the block can not be accessed by the program and can, hence, safely be freed.

To detect when there are no pointers to a block, the block keeps a count of incoming pointers. The counter is an extra field in the block in addition to the

size field. When a block is allocated, its counter field is set to 1 (representing the pointer that is returned by `malloc()`). When the program adds or removes pointers to the block, the counter is incremented and decremented. If, after decrementing a counter, the counter becomes 0, the block is freed by calling `free()`.

Reference counting usually uses free lists, possibly organised using the buddy system described in section 12.5.3.

The cost of maintaining the counter is substantial: If `p` and `q` are both pointers, an assignment `p := q` requires the following operations:

1. The counter field of the block B_1 that `p` points to is loaded into a register variable c .
2. If $c = 1$, the assignment will remove the last pointer to B_1 , so B_1 can be freed. Otherwise, $c - 1$ is written back to the counter field of B_1 .
3. The counter field of the block B_2 that `q` points to is incremented. This requires loading the counter into c , incrementing it and writing the modified counter back to memory.

In total, four memory accesses are required for an operation that could otherwise just be a register-to-register move.

A further complication happens if a block can contain pointers to other blocks. When a block is freed, these pointers are no longer accessible by the program, so the blocks these point to can have their counters decreased. Hence, we must go through a freed block and decrement the counters of all targets of pointers from the block, freeing those targets that get their counters decremented to zero.

This requires that we can see which fields of a block are heap pointers and which are other values (such as integers) that might just look like heap pointers.

In a strongly typed language, the type will typically have information about which fields are pointers. In a dynamically typed language, type information is available at runtime, so, when freeing a block, this type information is inspected and used to determine which pointers need to be followed. Basically, freeing a block is done by calling `free()` with the type information as argument in addition to the pointer. In a statically typed language, the compiler knows the type, so it can generate a specialised `free()` procedure for each type and insert a call to the relevant specialised procedure when freeing an object.

If full type information is not available (which can happen in weakly typed languages and polymorphically typed languages), the compiler must ensure that sufficient information is available to distinguish pointers from non-pointers. This is often done by observing that heap pointers point to word boundaries, so these will be even machine numbers. By forcing all other values to be represented as odd machine numbers, pointers are easily identifiable at runtime. An integer value

n will, in this representation, be represented as the machine number $2n + 1$. This means that integers have one bit less than a machine number and care must be done when doing arithmetic on integers, so the result is represented in the correct way. For example, when adding integers m and n (represented as $2m + 1$ and $2n + 1$), we must represent the result as $2(m + n) + 1$.

The requirement to distinguish pointers from other values also means that the fields of an allocated block must be initialised, so they don't hold any spurious pointers. This is not normally done by `malloc()` itself, but by the object constructors that call `malloc()`.

If a list or tree structure loses the last pointer to its root, the entire data structure is freed recursively. This takes time proportional to the size of the data structure, so if the freed data structure is very large, there can be a noticeable pause in program execution.

Another complication is circular data structures such as doubly-linked lists. Even if the last pointer to a doubly-linked list disappears, the elements point to each other, so their reference counts will not become zero, and the list is not freed. This can be handled by not counting back-references (*i.e.*, pointers that create cycles) in the reference counts. Pointers that are not counted are called *weak pointers*. For example, in a doubly-linked list, the backwards pointers are weak pointers while the forward pointers are normal pointers.

It is not always easy for the compiler to determine which pointer fields should be weak pointers, so reference counting is rarely used for languages that allow construction of circular data structures.

Suggested exercises: 12.6.

12.8 Tracing garbage collectors

A tracing garbage collector determines which blocks are reachable from variables in the program and frees all other blocks, as these can never be accessed by the program. Reachability can handle cycles, so tracing collectors (unlike reference counting) have no problem with circular data structures.

The variables in the program include register-allocated variables, stack-allocated variables, global variables and variables that have been spilled on the stack or saved on the stack during function calls. The collection of all these variables is called *the root set* of the heap, as all reachable heap-allocated tree or graph structures will be rooted in one of these variables.

Like with reference counting, we need to distinguish heap pointers from non-pointers. This can be done in the ways described above, but as an additional complication, we need such distinctions also for the root set. So all stack frames and

global variables must have descriptors, or we must use distinct representations for pointers and non-pointers also for values in the root set.

A tracing collector maintains an invariant during the reachability analysis by classifying all nodes (blocks and roots) in three categories represented by colours:

White: The node is not (yet) found to be reachable from the root set.

Grey: The node itself has been classified as reachable, but some of its children might not have been classified yet.

Black: Both the node itself and all its immediate children are classified as reachable.

By this invariant, a black node can not point to a white node. A grey node represents an unfinished reachability analysis, as it is clear that its children *are* reachable – they just haven’t been classified as such yet.

Initially, the root set is classified as grey and all heap-allocated blocks as white. When the reachability analysis is complete, all nodes are classified as either black or white. This means that all reachable nodes are now black and all white nodes are definitely unreachable, so white nodes can safely be freed.

While manual memory management and reference counting frees memory blocks one at a time using relatively short time for each (unless a large data structure is freed), tracing collectors are not called every time a pointer is moved or a block is freed – the overhead would be far too high. Instead, a tracing collector will typically be called when there is no free block that can accommodate a `malloc()` request. The tracing collector will then free all unreachable blocks before returning to `malloc()`, which can then return one of these freed blocks (or ask the operating system for more space, if no block is large enough). This can take long time if the heap is big, so noticeable pauses in execution are common when using tracing collection. We will look at ways to reduce these pauses in section 12.8.3.

12.8.1 Scan-sweep collection

Scan-sweep collection (also called *mark-sweep collection*) marks all reachable blocks using a depth-first scan starting from each root in the root set. When scanning is complete, the collector goes through (sweeps) the heap and frees all unmarked nodes. Freeing (as well as allocation) is done using free lists as described in section 12.5. We can sketch the two phases with the following pseudo-code:

```

scan(p)
  if the block at p is marked (* classified as reachable *)
  then return
  else
    mark it          (* classify it as grey *)
    for all fields q in the block at p do
      if q is a pointer then scan(q)
    (* the node at p is now black *)

sweep(b)
  if the block at b is unmarked (* classified as white *)
  then add it to the free list
  else clear the mark bit
  b := b + b.size
  if b > end-of-heap
  then return
  else sweep(b)

```

Scan-sweep collection requires a bit in every block for marking. This bit can be stored in the same machine word as the size field of the block. There is no explicit mark to distinguish grey and black nodes. Instead, for every grey node p , there is an uncompleted call to $\text{scan}(p)$. In short:

White: Mark bit is clear.

Grey: Mark bit is set and there is an uncompleted call to $\text{scan}(p)$.

Black: Mark bit is set and there is no uncompleted call to $\text{scan}(p)$.

Both scan and sweep are described above as recursive functions, but implementing them as such can make them use a lot of stack space. So scan is usually implemented using an explicit stack (containing pointers to all grey nodes) and sweep as a loop (since it is tail recursive). The stack used in the scan phase can, in the worst case, hold a pointer to every heap-allocated object. If blocks are at least four machine words, this overhead is at most 25% extra memory on top of the heap space.

It is possible to implement scan without using a stack by replacing the mark bit of a block with a counter that can count up to $1 + \text{the size of the block}$. In white nodes, the counter is 0, in grey nodes it is between 1 and the size of the block and in black nodes it is equal to $1 + \text{the size of the block}$. Part of the invariant is that if the counter is between 1 and the size of the block, the corresponding word of the block points to a parent node instead of to the child node. When the counter is incremented, the pointer is restored and the next word (if any) is made to point to the parent. When the last word is restored, $\text{scan}()$ follows the parent-pointer and continues scanning there.

This technique, called *pointer reversal*, requires more accesses to memory than scanning using an explicit stack, so what you save in space is paid for in time. Also, it is impossible to make pointer reversal concurrent (see section 12.8.3), as data structures are changed during the scan phase.

Suggested exercises: 12.7.

12.8.2 Two-space collection

Both reference counting and scan-sweep collection use free lists, so they suffer from the same fragmentation issues as manual memory management. Two-space collection avoids both internal and external fragmentation completely at the cost of having to move blocks during collection. Also, while scan-sweep collection uses time both for live (reachable) and dead (unreachable) nodes, two-space collection uses time only for live nodes. Functional and object-oriented languages often allocate a lot of small, short-lived objects, so the number of live objects will often be a lot smaller than the number of dead objects. So having to use time only for live objects can be a big saving, even if the cost of handling each dead object is small.

In two-space collection, allocation is done from a single, large contiguous chunk of memory called the *allocation area*. The memory manager maintains a pointer *next* to the first unused address in this block and a pointer *last* to the end of this block. Allocation is simple:

```
malloc(n)
  n := n + wordsize;    if last - next < n
  then do garbage collection
  store n at next      (* set size field *)
  next := next + n
  return next - n      (* old value of next *)
```

If garbage collection is not required, allocation is done in constant time. After returning from a garbage collection, it is assumed that there is sufficient memory for the allocation. We will later, briefly, return to the case where this is not true.

In addition to the allocation area, The garbage collector needs an extra chunk of memory (called the *to-space*) at least as large as the allocation area (which during collection is called the *from-space*). Collection will copy all live nodes from the from-space to one end of the to-space and then free the entire from-space. The to-space becomes the new allocation area, and in the next collection, the rôles of the two spaces are reversed, so the previous from-space becomes the new to-space and the allocation space becomes the new from-space. Once collection starts, the *next* pointer is made to point to the first free address in to-space and the *last* pointer is made to point to the end of to-space.

When a node is copied from from-space to to-space, all nodes that point to this node have to have their pointers updated, so they point to the new copy instead of the old. The colours of the garbage-collection invariant indicate the status of copying of nodes and updates of pointers:

White: The node has not been copied to to-space.

Grey: The node has been copied to to-space, but some of the pointer fields in the new copy may still point to nodes in from-space. Pointer fields in the old copy are unchanged.

Black: The node has been copied to to-space and all of the pointer fields in the new copy have been updated to point to nodes in to-space. Pointer fields in the old copy are still unchanged.

Once all nodes in to-space are black, all live blocks have been copied to to-space and no pointers exist from to-space to from-space. So the entire from-space can be safely freed. Freeing requires no explicit action – the from-space is just re-used as to-space in the next collection.

The basic operation of two-space collection is *forwarding* a pointer: The block pointed to by the pointer is copied to to-space (unless this has already happened) and the address of the new copy is returned. To show that the node has already been copied, the first word of the old copy (*i.e.*, its size field) is overwritten with a pointer to the new copy. To distinguish sizes from pointers, sizes can have their least significant bit set (just like we in section 12.5.2 used this bit to mark a block in use). In pseudo code, the forwarding function can be written as

```
forward(p)
  if p points to from-space
  then
    q := the content of the word p points to
    if q is even    (* a forwarding pointer *)
    then return q
    else          (* q is 1 + the size of the block *)
      q := q - 1
      copy q bytes from p to next
      overwrite the word at p with the value of next
      next := next + q
      return next - q
  else return p
```

Forwarding a pointer copies the node to to-space if necessary, but it does not update the pointer-fields of the copy to point to to-space, so the node is grey. The necessary updating is done by the function `update`, that takes a pointer to a grey node (located in to-space) and forwards all its pointer-fields, making the node black:

```

update(p)
  for all heap-pointer fields q of p do
    q := forward(q)

```

The entire garbage collection process can be written in pseudo code as

```

next := first word of to-space
last := last word of to-space
scan := next
for all heap-pointer variables p in the root set do
  p := forward(p)
while scan < next do
  update(scan)
  scan := scan + the size of the node at scan

```

By the invariant, nodes between the start of to-space and scan are black and nodes between scan and next are grey. So when scan=next, all nodes in to-space are black and all reachable nodes have been copied to to-space.

If this garbage collection does not free up sufficient space to allocate the requested block, the memory manager will ask the operating system for a chunk of memory at least large enough to hold all the live blocks plus the requested block and immediately make a new garbage collection to the new chunk (as the new to-space). Both the old spaces are freed and a new to-space is allocated from the operating system at the next garbage collection.

Two-space collection changes the value of pointers while preserving only that a given pointer field or variable points to the new copy of the node that it pointed to before collection. So the language must allow pointer values to change in this way. This means that it is not possible to cast a pointer to an integer and cast the integer back to a pointer, since the new pointer might not be valid. Also, since the order of blocks is changed by garbage collection, pointers to different heap blocks can not be compared to see which has a lower address, as this can change without warning. So you can only compare pointers for equality.

Since collection time is proportional only to the live (reachable) nodes, it is opportune to do collection when there is little live data. Often, a programmer can identify place in the program where this is true and force a collection to happen there. Not all memory managers offer facilities to force collections, though, and when they do, it is rarely portable between different implementations.

Suggested exercises: 12.8.

12.8.3 Generational and concurrent collectors

Tracing collectors can be refined in various ways. We will look at generational collectors and concurrent collectors.

A problem with copying collectors is that long-lived blocks are repeatedly copied back and forth between the two spaces. To reduce this problem, we can use two observations: Larger blocks tend to live longer than smaller blocks, and blocks that have already lived long are likely to continue living for a long time. So we aim to not copy such blocks as often as small, young blocks. We do this by not having only two spaces, but n spaces, where each space is larger than the previous (typically 4–16 times larger). These spaces are called *generations*, and the small spaces are called the young generations. Small blocks are allocated in the smaller (younger) generations and large blocks in the larger (older) generations. We keep as an invariant that generation $g + 1$ should have at least as much free space as the total size of generation g , so when collecting generation g , generation $g + 1$ can be used as to-space. All collections start with collecting generation 0 (the smallest) using generation 1 as to-space. If generation 1 after the collection does not have enough free space, we collect this, using generation 2 as to-space and so on. After this, generation $1 - g$ will be empty and generation $g + 1$ will have free space at least the size of g .

If allocation of a new block would cause its intended generation to have too little free space, we collect all generations up to and including this before allocating the block.

The last generation does not have a bigger generation to use as to-space, so we collect this using a non-generational collection method such as scan-sweep or two-space copying collection.

Generally, there will be 10–100 collections of generation g between collections in generation $g + 1$, so blocks in the older generations are not copied very often. Also, collection of a small generation is very fast, so pauses (though more common) are typically much shorter. When all generations need to be collected, though, the pause is the same as with a normal two-space collector. The smallest generations will typically fit in the cache of the system, so allocation of small objects and collection in the youngest generation can be done quite quickly.

Remembered sets Generational collectors have a problem with pointers from older generations to younger generations: When a block is copied to an older space (generation), all pointers to this block must be updated to point to the new copy. In two-space collection, this is done by calling `update()` on all blocks, but doing so on all blocks in all generations is much too expensive. So we need another mechanism for updating pointers from older generations to younger generations. The traditional solution is let each generation have a list of pointers to those blocks in older generations that may point into the generation. This list, called *the remembered set*, can be allocated in the generation itself (starting from the opposite end of normal allocations). When a generation is collected, we need only call `update()` on the forwarded blocks and the blocks in the remembered set. The blocks pointed

to by the remembered set are also used as extra roots when collecting the generation. This way, we don't have to trace reachability through the older generations: Only root pointers that point directly into the collected generation and pointers in the remembered set of the generation are treated as roots.

Remembered sets adds an extra burden on the program, though: Whenever the program updates a pointer field in a block, it has to check if the pointer points into a younger generation than the generation in which the block itself is allocated and, if so, add to the remembered set of the younger generation a pointer to the updated block. Adding an element to the remembered set of a generation might cause the free space in the generation to fall below the collection threshold, so this must be checked at any pointer-field update and a collection started if this is the case.

Studies have shown that adding pointers from old generations to younger generations is relatively rare (especially in functional languages), so remembered sets are typically small.

Concurrent and incremental collection While generational collectors can reduce the average length of pauses caused by garbage collection, there can still be long pauses (up to a couple of seconds) when older (larger) generations are collected. For highly interactive programs (such as computer games) or programs that need to react quickly to outside changes, such pauses can be unacceptable. Hence, we want to be able to execute the program concurrently with the garbage collector in such a way that the long pauses caused by garbage collection are replaced by a larger number of much shorter pauses. So, at regular intervals, the program gives time to the garbage collector to do some work. This can be done by letting the garbage collector be a separate thread which synchronises with the program in the usual way, or the program can simply call the collector regularly and expect the garbage collector to return after a short time, so the program and the garbage collector work like co-routines. This is often called *incremental garbage collection*. The garbage collector can rarely complete a collection before returning control to the program, so it is important to keep both the program state and the garbage-collection state consistent when control passes between the program and the collector.

Reference counting is easily made incremental: It is only when freeing large data structure that pauses can be long, but the collector can return to the program before this is complete and resume the process when next called. Making tracing collectors concurrent is more difficult, though.

A tracing collector uses the invariant that a black node can never point to a white node. If the program updates a pointer field in a root or heap-allocated block, it might violate this invariant. To avoid this, the program can make sure that whenever it stores a pointer in a root or heap-allocated block, that either the node at the end of the pointer or the root or the heap-allocated block that contain the pointer is grey.

It can do that by forcing either of these to become grey.

In the scan phase of a scan-sweep collector that uses an explicit stack of grey objects, we can mark the node at the end of the pointer and put the pointer value on this stack. This makes the node that is pointed to grey, so whatever colour the updated root or heap-allocated block has, the invariant is preserved. In the sweep phase, no action is needed: The program can only access black nodes and the sweep phase concerns only white nodes, so the invariant can not be violated. Blocks that are allocated during collection are immediately marked. If fields are initialised with pointers, these are handled like pointer updates: If it happens during scan, we add the pointers to the stack, and if it happens during sweep, we do nothing.

In a copying collector, things are more complicated, as blocks are moved by the collector, so the program state is made inconsistent during collection. At the start of the collection, all roots are forwarded, so all roots will point to blocks that are already moved. In other words, all roots are black at this time. But if a pointer is read from memory to a program variable (*i.e.*, a root), the pointer might point to a block in from-space, *i.e.*, a white node. Forwarding the pointer when it is read will make sure the program variable points to a grey or black node, so the invariant is preserved. In a generational collector, you only need to forward pointers that point into the generation that is currently being collected.

Reads of heap pointers are much more common than writes to heap pointers, so adding work at every pointer read is a high overhead. To reduce the overhead, some concurrent copying collectors use more complicated methods that require extra work only at pointer updates.

Concurrent collectors risk running out of space even when there is free memory, since this might just not have been collected yet. If this happens, the collector can continue collection until enough free space is available, even if this takes time, or it can ask the operating system for more space. A two-space copying collector can not easily increase the size of the to-space in the middle of a collection, so concurrent copying collectors are normally also generational – it is easy enough to add an extra generation or, if the last generation uses a free-list, add more blocks to this.

12.9 Summary of automatic memory management

Automatic memory management needs to distinguish heap-pointers from other values. This can be done by using type information or by having non-overlapping representations of pointers and non-pointers. Newly allocated blocks must be initialised so they contain no spurious pointers.

Reference counting keeps in every block a count of incoming pointers, so adding or removing pointers to a block requires an update of its counter. Reference counting can not (without assistance) collect circular data structures, but it

mostly avoids the long pauses that (non-concurrent) tracing collectors can cause.

Scan-sweep is a simple tracing collector that uses a free list.

Copying collectors don't use free lists, so fragmentation is avoided and allocation is constant time (except when garbage collection is needed). Furthermore, the time to do a collection is independent of the amount of unreachable (dead) nodes, so it can be faster than scan-sweep collection. There is, however, an extra cost in copying live blocks during garbage collection. Also, the language must allow pointers to change value.

Generational collection can speed up copying collectors by reducing copying of old and large blocks at the cost of increasing the cost of updating pointer fields.

Concurrent (incremental) collection can reduce the pauses of tracing collectors, but adds overhead to ensure that the program does not invalidate the invariant of the collector and *vice versa*.

12.10 Further reading

Techniques for manual memory management including the buddy system mentioned in section 12.5.3 is described in detail in [26].

Automatic memory management, including generational and concurrent collection, is described in detail in [44], [22] and [9].

In a weakly-typed low-level language like C, there is no way to distinguish pointers and integers. Also, pointers can point into the middle of blocks. So automatic memory management with the restrictions we have described above can not be used for C. Even so, garbage collectors for C have been made [10]. The idea is that any value that *looks like* a pointer into the heap is treated as if it *is* a pointer into the heap. If unlimited type casting is allowed (as it is in C), any value can be cast into a pointer, so this is not an unreasonable strategy. Any heap-allocated block that is pointed (in)to by something that looks like a pointer is preserved during garbage collection. This idea is called *conservative garbage collection*. A conservative collector can not change values of pointers (as it might accidentally change the value of a non-pointer that just looks like a pointer), so copying collectors can not be used.

Automatic memory management based on compile-time analysis of lifetimes is discussed in [43].

Exercises

Exercise 12.1

Show the free list from figure 12.1(b) after another call to `malloc(12)`.

Exercise 12.2

In section 12.5.2, a fast method for joining blocks is described, where a newly freed block can be joined with the adjacent block at higher address (if this is free), but not with the adjacent block at lower address. If a memory manager uses this method (and the allocation method described in section 12.5.1), what is the best strategy: Freeing blocks in the same order that they are allocated or freeing blocks in the opposite order of the allocation order?

Exercise 12.3

In the buddy system described in section 12.5.3, block sizes are always powers of two, so instead of storing the actual size of the block in its size field, we can store the base-2 logarithm of the size. *i.e.*, if the size is 2^n , we store n in the size field.

Consider advantages and disadvantages of storing the logarithm of the size instead of the size itself.

Exercise 12.4

Given the buddy system described in section 12.5.3, the worst-case execution time for `malloc()` and `free()` is logarithmic in the size of the total heap memory. Describe a sequence of allocations and deallocations that cause this worst case to happen every time. You can assume that the initial heap memory is a single block of size 2^n for some n .

Exercise 12.5

In section 12.5.3, it is suggested that block sizes can be generalised Fibonacci numbers instead of powers of two. Describe how you would join freed blocks in this case.

Exercise 12.6

In section 12.7, it is suggested that heap pointers can be distinguished from other values by representing heap pointers by even machine numbers and all other values by odd machine numbers.

An integer n would, hence, be represented by a machine number with the value $2n + 1$, which is guaranteed to be odd.

Consider how you would effectively (using the shortest possible sequence of machine-code instructions) add two integers m and n that are both represented this way (*i.e.*, as machine numbers $2m + 1$ and $2n + 1$) in such a way that the result is also represented the same way (*i.e.* as a machine number $2(m + n) + 1$).

Then do the same for multiplication of m and n .

Exercise 12.7

In section 12.8.1, both `scan` and `sweep` are described as recursive procedures, but it is suggested that an explicit stack can be used to avoid recursion in `scan` and that `sweep` can be rewritten to use a loop.

Write pseudo-code for non-recursive versions of `scan` and `sweep` using these ideas.

Exercise 12.8

In section 12.8.2, the space between `scan` and `next` works like a queue, so copying live nodes to to-space is done as a breadth-first traversal of the live nodes. This means that adjacent elements of a linked list can be copied to far-apart locations in to-space, which is bad for locality of reference. Suggest a modification of the `forward()` function that will make adjacent elements of a singly-linked list be copied to adjacent blocks in to-space. Try to avoid using extra space.

Chapter 13

Bootstrapping a compiler

13.1 Introduction

When writing a compiler, one will usually prefer to write it in a high-level language. A possible choice is to use a language that is already available on the machine where the compiler should eventually run. It is, however, quite common to be in the following situation:

You have a completely new processor for which no compilers exist yet. Nevertheless, you want to have a compiler that not only targets this processor, but also runs on it. In other words, you want to write a compiler for a language A, targeting language B (the machine language) and written in language B.

The most obvious approach is to write the compiler in language B. But if B is machine language, it is a horrible job to write any non-trivial compiler in this language. Instead, it is customary to use a process called “bootstrapping”, referring to the seemingly impossible task of pulling oneself up by the bootstraps.

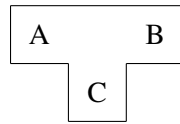
The idea of bootstrapping is simple: You write your compiler in language A (but still let it target B) and then let it compile itself. The result is a compiler from A to B written in B.

It may sound a bit paradoxical to let the compiler compile itself: In order to use the compiler to compile a program, we must already have compiled it, and to do this we must use the compiler. In a way, it is a bit like the chicken-and-egg paradox. We shall shortly see how this apparent paradox is resolved, but first we will introduce some useful notation.

13.2 Notation

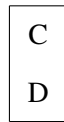
We will use a notation designed by H. Bratman [11]. The notation is hence called “Bratman diagrams” or, because of their shape, “T-diagrams”.

In this notation, a compiler written in language C, compiling from the language A and targeting the language B is represented by the diagram



In order to use this compiler, it must “stand” on a solid foundation, *i.e.*, something capable of executing programs written in the language C. This “something” can be a machine that executes C as machine-code or an interpreter for C running on some other machine or interpreter. Any number of interpreters can be put on top of each other, but at the bottom of it all, we need a “real” machine.

An interpreter written in the language D and interpreting the language C, is represented by the diagram

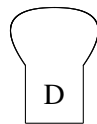


A machine that directly executes language D is written as

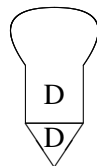


The pointed bottom indicates that a machine need not stand on anything; it is itself the foundation that other things must stand on.

When we want to represent an unspecified program (which can be a compiler, an interpreter or something else entirely) written in language D, we write it as

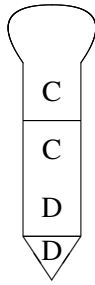


These figures can be combined to represent executions of programs. For example, running a program on a machine D is written as



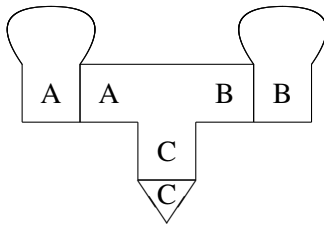
Note that the languages must match: The program must be written in the language that the machine executes.

We can insert an interpreter into this picture:



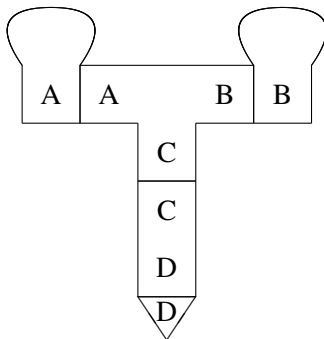
Note that, also here, the languages must match: The interpreter can only interpret programs written in the language that it interprets.

We can run a compiler and use this to compile a program:



The input to the compiler (*i.e.*, the source program) is shown at the left of the compiler, and the resulting output (*i.e.*, the target program) is shown on the right. Note that the languages match at every connection and that the source and target program are not “standing” on anything, as they are not executed in this diagram.

We can insert an interpreter in the above diagram:



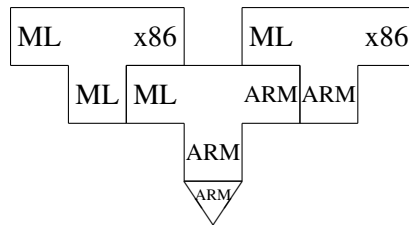
13.3 Compiling compilers

The basic idea in bootstrapping is to use compilers to compile themselves or other compilers. We do, however, need a solid foundation in form of a machine to run the compilers on.

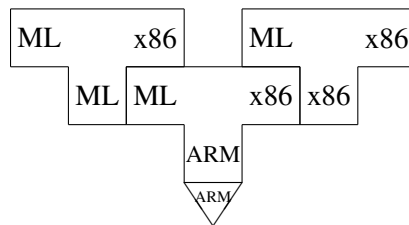
It often happens that a compiler does exist for the desired source language, it just does not run on the desired machine. Let us, for example, assume we want a compiler for ML to x86 machine code and want this to run on an x86. We have access to an ML-compiler that generates ARM machine code and runs on an ARM machine, which we also have access to. One way of obtaining the desired compiler would be to do *binary translation*, i.e., to write a compiler from ARM machine code to x86 machine code. This will allow the translated compiler to run on an x86, but it will still generate ARM code. We can use the ARM-to-x86 compiler to translate this into x86 code afterwards, but this introduces several problems:

- Adding an extra pass makes the compilation process take longer.
- Some efficiency will be lost in the translation.
- We still need to make the ARM-to-x86 compiler run on the x86 machine.

A better solution is to write an ML-to-x86 compiler in ML. We can compile this using the ML compiler on the ARM:



Now, we can run the ML-to-x86 compiler on the ARM and let it compile itself¹:



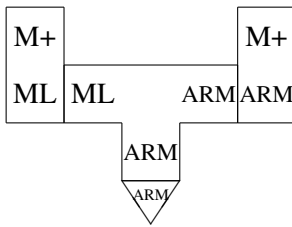
We have now obtained the desired compiler. Note that the compiler can now be used to compile itself directly on the x86 platform. This can be useful if the compiler is later extended or, simply, as a partial test of correctness: If the compiler, when compiling itself, yields a different object code than the one obtained with the above process, it must contain an error. The converse is not true: Even if the same target is obtained, there may still be errors in the compiler.

¹We regard a compiled version of a program as the same program as its source-code version.

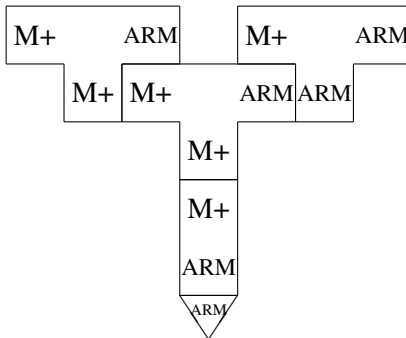
required to get the compiler to run on the same machine that it targets. We chose the target language to make a point: Bootstrapping might not be complete even if a compiler with the right functionality has been obtained.

Using an interpreter

Instead of writing a QAD compiler, we can write a QAD interpreter. In our example, we could write an M+ interpreter in ML. We would first need to compile this:



We can then use this to run the M+ compiler directly:



Since the “real” compiler has been used to do the compilation, nothing will be gained by using the generated compiler to compile itself, though this step can still be used as a test and for extensions.

Though using an interpreter requires fewer steps, this should not really be a consideration, as the computer(s) will do all the work in these steps. What is important is the amount of code that needs to be written by hand. For some languages, a QAD compiler will be easier to write than an interpreter, and for other languages an interpreter is easier. The relative ease/difficulty may also depend on the language used to implement the QAD interpreter/compiler.

Incremental bootstrapping

It is also possible to build the new language and its compiler incrementally. The first step is to write a compiler for a small subset of the language, using that same

subset to write it. This first compiler must be bootstrapped in one of the ways described earlier, but thereafter the following process is done repeatedly:

- 1) Extend the language subset slightly.
- 2) Extend the compiler so it compiles the extended subset, but without using the new features.
- 3) Use the previous compiler to compile the new.

In each step, the features introduced in the previous step can be used in the compiler. Even when the full language is compiled, the process can be continued to improve the quality of the compiler.

Suggested exercises: 13.1.

13.4 Further reading

Bratman's original article, [11], only describes the T-shaped diagrams. The notation for interpreters, machines and unspecified programs was added later in [15].

An early bootstrapped compiler was LISP 1.5 [30].

The first Pascal compiler [45] was made using incremental bootstrapping.

Though we in section 13.3 dismissed binary translation as unsuitable for porting a compiler to a new machine, it is occasionally used. The advantage of this approach is that a single binary translator can port any number of programs, not just compilers. It was used by Digital Equipment Corporation in their FX!32 software [18] to enable programs compiled for Windows on an x86-platform to run on their Alpha RISC processor.

Exercises

Exercise 13.1

You have a machine that can execute *Alpha* machine code and the following programs:

- 1: A compiler from C to *Alpha* machine code written in *Alpha* machine code.
- 2: An interpreter for ML written in C.
- 3: A compiler from ML to C written in ML.

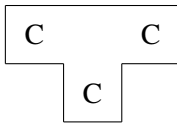
Now do the following:

- a) Describe the above programs and machine as diagrams.
- b) Show how a compiler from ML to C written in *Alpha* machine code can be generated from the above components. The generated program must be stand-alone, *i.e.*, it may not consist of an interpreter and an interpreted program.
- c) Show how the compiler generated in question b can be used in a process that compiles ML programs to *Alpha* machine code.

Exercise 13.2

A source-code optimiser is a program that can optimise programs at source-code level, *i.e.*, a program O that reads a program P and outputs another program P' , which is equivalent to P , but may be faster.

A source-code optimiser is like a compiler, except the source and target languages are the same. Hence, we can describe a source-code optimizer for C written in C with the diagram



Assume that you additionally have the following components:

- A compiler, written in ARM code, from C to ARM code.
- A machine that can execute ARM code.
- Some unspecified program P written in C.

Now do the following:

- a) Describe the above components as diagrams.
- b) Show, using Bratman diagrams, the steps required to optimise P to P' and then execute P' .

Appendix A

Set notation and concepts

A.1 Basic concepts and notation

A set is a collection of items. You can write a set by listing its elements (the items it contains) inside curly braces. For example, the set that contains the numbers 1, 2 and 3 can be written as $\{1, 2, 3\}$. The order of elements do not matter in a set, so the same set can be written as $\{2, 1, 3\}$, $\{2, 3, 1\}$ or using any permutation of the elements. The number of occurrences also does not matter, so we could also write the set as $\{2, 1, 2, 3, 1, 1\}$ or an infinity of other ways. All of these describe the same set. We will normally write sets without repetition, but the fact that repetitions do not matter is important to understand the operations on sets.

We will typically use uppercase letters to denote sets and lowercase letters to denote elements in a set, so we could write $M = \{2, 1, 3\}$ and $x = 2$ as an element of M . The empty set can be written either as an empty list of elements ($\{\}$) or using the special symbol \emptyset . The latter is more common in mathematical texts.

A.1.1 Operations and predicates

We will often need to check if an element belongs to a set or select an element from a set. We use the same notation for both of these: $x \in M$ is read as “ x is an element of M ” or “ x is a member of M ”. The negation is written as $x \notin M$, which is read as “ x is not an element of M ” or “ x is not a member of M ”.

We can use these in conditional statements like “if $3 \in M$ then ...”, for asserting a fact “since $x \notin M$, we can conclude that ...” or for selecting an element from a set: “select $x \in M$ ”, which will select an arbitrary element from M and let x be equal to this element.

We can combine two sets M and N into a single set that contains all elements from both sets. We write this as $M \cup N$, which is read as “ M union N ” or “the union of M and N ”. For example, $\{1, 2\} \cup \{5, 1\} = \{1, 2, 5, 1\} = \{1, 2, 5\}$. The

following statement holds for membership and union:

$$x \in (M \cup N) \Leftrightarrow x \in M \vee x \in N$$

where \Leftrightarrow is bi-implication (“if and only if”) and \vee is logical disjunction (“or”).

We can also combine two sets M and N into a set that contains only the elements that occur in both sets. We write this as $M \cap N$, which is read as “ M intersect N ” or “the intersection of M and N ”. For example, $\{1, 2\} \cap \{5, 1\} = \{1\}$. The following statement holds for membership and intersection:

$$x \in (M \cap N) \Leftrightarrow x \in M \wedge x \in N$$

where \wedge is logical conjunction (“and”).

We can also talk about set difference (or set subtraction), which is written as $M \setminus N$, which is read as “ M minus N ” or “ M except N ”. $M \setminus N$ contains all the elements that are members of M but not members of N . For example, $\{1, 2\} \setminus \{5, 1\} = \{2\}$. The following statement holds for membership and set difference:

$$x \in (M \setminus N) \Leftrightarrow x \in M \wedge x \notin N$$

Just like arithmetic operators, set operators have precedence rules: \cap binds more tightly than \cup (just like multiplication binds tighter than addition). So writing $A \cup B \cap C$ is the same as writing $A \cup (B \cap C)$. Set difference has the same precedence as union (just like subtraction has the same precedence as addition).

If all the elements of a set M are also elements of a set N , we call M a subset of N , which is written as $M \subseteq N$. This can be defined by

$$M \subseteq N \Leftrightarrow (x \in M \Rightarrow x \in N)$$

where \Rightarrow is logical implication (“only if”).

The converse of subset is superset: $M \supseteq N \Leftrightarrow N \subseteq M$.

A.1.2 Properties of set operations

Just like we have mathematical laws saying that, for example $x + y = y + x$, there are also similar laws for set operations. Here is a selection of the most commonly used laws:

$A \cup A = A$	union is idempotent
$A \cap A = A$	intersection is idempotent
$A \cup B = B \cup A$	union is commutative
$A \cap B = B \cap A$	intersection is commutative
$A \cup (B \cup C) = (A \cup B) \cup C$	union is associative
$A \cap (B \cap C) = (A \cap B) \cap C$	intersection is associative
$A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$	union distributes over intersection
$A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$	intersection distributes over union
$A \cup \emptyset = A$	the empty set is a unit element of union
$A \cap \emptyset = \emptyset$	the empty set is a zero element of intersection
$A \subseteq B \Leftrightarrow A \cup B = B$	subset related to union
$A \subseteq B \Leftrightarrow A \cap B = A$	subset related to intersection
$A \subseteq B \Leftrightarrow A \setminus B = \emptyset$	subset related to set difference
$A \subseteq B \wedge B \subseteq A \Leftrightarrow A = B$	subset is antisymmetric
$A \subseteq B \wedge B \subseteq C \Leftrightarrow A \subseteq C$	subset is transitive
$A \setminus (B \cup C) = (A \setminus B) \setminus C$	corresponds to $x - (y + z) = (x - y) - z$

Since \cup and \cap are associative, we will often omit parentheses and write, *e.g.*, $A \cup B \cup C$ or $A \cap B \cap C$.

A.2 Set-builder notation

We will often build a new set by selecting elements from other sets and doing operations on these elements. We use the very flexible set-builder notation for this. A set builder has the form $\{e \mid p\}$, where e is an expression and p is a list of predicates separated by commas. Typically, p will contain predicates of the form $x \in M$, which defines x to be any element of M . The set builder will evaluate the expression e for all elements x of M that fulfills the other predicates in p and build a set of the results. We read $\{e \mid p\}$ as “the set of all elements of the form e where p holds”, or just “ e where p ”. Some mathematical texts use a colon instead of a bar, *i.e.*, writing $\{e : p\}$ instead of $\{e \mid p\}$.

A simple example is

$$\{x^3 \mid x \in \{1, 2, 3, 4\}, x < 3\}$$

which builds the set $\{1^3, 2^3\} = \{1, 8\}$, as only the elements 1 and 2 from the set $\{1, 2, 3, 4\}$ fulfill the predicate $x < 3$.

We can take elements from more than one set, for example

$$\{x + y \mid x \in \{1, 2, 3\}, y \in \{1, 2, 3\}, x < y\}$$

which builds the set $\{1 + 2, 1 + 3, 2 + 3\} = \{3, 4, 5\}$. All combinations of elements from the two sets that fulfill the predicate are used.

We can separate the predicates in a set builder by \wedge or “and” instead of commas. So the example above can, equivalently, be written as

$$\{x + y \mid x \in \{1, 2, 3\}, y \in \{1, 2, 3\} \text{ and } x < y\}$$

A.3 Sets of sets

The elements of a set can be other sets, so we can, for example, have the set $\{\{1, 2\}, \{2, 3\}\}$ which is a set that has the two sets $\{1, 2\}$ and $\{2, 3\}$ as elements. We can “flatten” a set of sets to a single set which is the union of the element sets using the “big union” operator:

$$\bigcup \{\{1, 2\}, \{2, 3\}\} = \{1, 2, 3\}$$

Similarly, we can take the intersection of the element sets using the “big intersection” operator:

$$\bigcap \{\{1, 2\}, \{2, 3\}\} = \{2\}$$

We can use these “big” operators together with set builders, for example

$$\bigcap \{\{x^n \mid n \in \{0, 1, 2\}\} \mid x \in \{1, 2, 3\}\}$$

which evaluates to $\bigcap \{\{1\}, \{1, 2, 4\}, \{1, 3, 9\}\} = \{1\}$.

When a big operator is used in combination with a set builder, a special abbreviated notation can be used: $\bigcup \{e \mid p\}$ and $\bigcap \{e \mid p\}$ can be written, respectively, as

$$\bigcup_p e \quad \text{and} \quad \bigcap_p e$$

For example,

$$\bigcap \{\{x^n \mid n \in \{0, 1, 2\}\} \mid x \in \{1, 2, 3\}\}$$

can be written as

$$\bigcap_{x \in \{1, 2, 3\}} \{x^n \mid n \in \{0, 1, 2\}\}$$

A.4 Set equations

Just like we can have equations where the variables represent numbers, we can have equations where the variables represent sets. For example, we can write the equation

$$X = \{x^2 \mid x \in X\}$$

This particular equation has several solutions, including $X = \{0\}$, $X = \emptyset$ and $X = \{0, 1\}$ or even $X = [0, 1]$, where $[0, 1]$ represents the interval of real numbers between 0 and 1. Usually, we have an implied universe of elements that the sets can draw from. For example, we might only want sets of integers as solutions, so we don't consider intervals of real numbers as valid solutions.

When there are more solutions, we are often interested in a solution that has the minimum or maximum possible number of elements. In the above example (assuming we want sets of integers), there is a unique minimal (in terms of number of elements) solution, which is $X = \emptyset$ and a unique maximal solution $X = \{0, 1\}$. Not all equations have unique minimal or maximal solutions. For example, the equation

$$X = \{1, 2, 3\} \setminus X$$

has no solution at all, and the equation

$$X = \{1, 2, 3\} \setminus \{6/x \mid x \in X\}$$

has exactly two solutions: $X = \{1, 2\}$ and $X = \{1, 3\}$, so there are no unique minimal or maximal solutions.

A.4.1 Monotonic set functions

The set equations we have seen so far are of the form $X = F(X)$, where F is a function from sets to sets. A solution to such an equation is called a *fixed-point* for F .

As we have seen, not all such equations have solutions, and when they do, there are not always unique minimal or maximal solutions. We can, however, define a property of the function F that guarantees a unique minimal and a unique maximal solution to the equation $X = F(X)$.

We say that a set function F is *monotonic* if $X \subset Y \Rightarrow F(X) \subseteq F(Y)$.

Theorem A.1 *If we draw elements from a finite universe U and F is a monotonic function over sets of elements from U , then there exist natural numbers m and n , so the unique minimal solution to the equation $X = F(X)$ is equal to $F^m(\emptyset)$ and the unique maximal solution to the equation $X = F(X)$ is equal to $F^n(U)$.*

where $F^i(A)$ is F applied i times to A . For example $F^3(A) = F(F(F(A)))$.

Proof: It is trivially true that $\emptyset \subseteq F(\emptyset)$. Since F is monotonic, this implies $F(\emptyset) \subseteq F(F(\emptyset))$. This again implies $F(F(\emptyset)) \subseteq F(F(F(\emptyset)))$ and, by induction, $F^i(\emptyset) \subseteq F^{i+1}(\emptyset)$. So we have a chain

$$\emptyset \subseteq F(\emptyset) \subseteq F(F(\emptyset)) \subseteq F(F(F(\emptyset))) \subseteq \dots$$

Since the universe U is finite, the sets $F^i(\emptyset)$ can not all be different. Hence, there exist an m such that $F^m(\emptyset) = F^{m+1}(\emptyset)$, which means $X = F^m(\emptyset)$ is a solution to the equation $X = F(X)$. To prove that it is the unique minimal solution, assume that another solution A exist. Since $A = F(A)$, we have $A = F^m(A)$. Since $\emptyset \subseteq A$ and F is monotonic, we have $F^m(\emptyset) \subseteq F^m(A) = A$. This implies that $F^m(\emptyset)$ is a subset of all solutions to the equation $X = F(X)$, so there can not be a minimal solution different from $F^m(\emptyset)$.

The proof for the maximal solution is left as an exercise.

fixed-point iteration

The proof provides an algorithm for finding minimal solutions to set equations of the form $X = F(X)$, where F is monotonic and the universe is finite: Simply compute $F(\emptyset)$, $F^2(\emptyset)$, $F^3(\emptyset)$ and so on until $F^{m+1}(\emptyset) = F^m(\emptyset)$. This is easy to implement on a computer:

```

X :=  $\emptyset$ ;
repeat
  Y := X;
  X := F(X)
until X = Y;
return X

```

A.4.2 Distributive functions

A function can have a stronger property than being monotonic: A function F is *distributive* if $F(X \cup Y) = F(X) \cup F(Y)$ for all sets X and Y . This clearly implies monotonicity, as $Y \supseteq X \Leftrightarrow Y = X \cup Y \Rightarrow F(Y) = F(X \cup Y) = F(X) \cup F(Y) \supseteq F(X)$.

We also solve set equations over distributive functions with fixed-point iteration, but we exploit the distributivity to reduce the amount of computation we must do: If we need to compute $F(A \cup B)$ and we have already computed $F(A)$, then we need only compute $F(B)$ and add the elements from this to $F(A)$. We can implement an algorithm for finding the minimal solution that exploits this:

```

X := ∅;
W := F(∅);
while W ≠ ∅ do
  pick x ∈ W;
  W := W \ {x};
  X := X ∪ {x};
  W := W ∪ (F({x}) \ X);
return X

```

We keep a work set W that by invariant is equal to $F(X) \setminus X$. A solution must include any $x \in W$, so we move this from W to X while keeping the invariant by adding $F(x) \setminus X$ to W . When W becomes empty, we have $F(X) = X$ and, hence, a solution. While the algorithm is more complex than the simple fixed-point algorithm, we can compute F one element at a time and we avoid computing F twice for the same element.

A.4.3 Simultaneous equations

We sometimes need to solve several simultaneous set equations:

$$\begin{aligned}
 X_1 &= F_1(X_1, \dots, X_n) \\
 &\vdots \\
 X_n &= F_n(X_1, \dots, X_n)
 \end{aligned}$$

If all the F_i are monotonic in all arguments, we can solve these equations using fixed-point iteration. To find the unique minimal solution, start with $X_i = \emptyset$ and then iterate applying all F_i until a fixed-point is reached. The order in which we do this doesn't change the solution we find (it will always be the unique minimal solution), but it might affect how fast we find the solution. Generally, we need only recompute X_i if a variable used by F_i changes.

If all F_i are distributive in all arguments, we can use a work-set algorithm similar to the algorithm for a single distributive function.

Exercises

Exercise A.1

What set is built by the set builder

$$\{x^2 + y^2 \mid x \in \{1, 2, 3, 4\}, y \in \{1, 2, 3, 4\}, x < y^2\}?$$

Exercise A.2

What set is built by the set expression

$$\bigcup_{x \in \{1, 2, 3\}} \{x^n \mid n \in \{0, 1, 2\}\}?$$

Exercise A.3

Find all solutions to the equation

$$X = \{1, 2, 3\} \setminus \{x + 1 \mid x \in X\}$$

Hint: Any solution must be a subset of $\{1, 2, 3\}$.

Exercise A.4

Prove that if elements are drawn from a finite universe U and F is a monotonic function over sets of elements from U , then there exists an n such that $X = F^n(U)$ is the unique maximal solution to the set equation $X = F(X)$.

Bibliography

- [1] A. Aasa. Precedences in specification and implementations of programming languages. In J. Maluszyński and M. Wirsing, editors, *Proceedings of the Third International Symposium on Programming Language Implementation and Logic Programming*, number 528 in LNCS, pages 183–194. Springer Verlag, 1991.
- [2] Harold Abelson, Gerald Jay Sussman, and Julie Sussman. *Structure and Interpretation of Computer Programs*. MIT Press, 1996. Also downloadable from <http://mitpress.mit.edu/sicp/full-text/sicp/book/>.
- [3] Alfred V. Aho, John E. Hopcroft, and Jeffrey D. Ullman. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, 1974.
- [4] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers; Principles, Techniques and Tools*. Addison-Wesley, 2007. Newer edition of [5].
- [5] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers; Principles, Techniques and Tools*. Addison-Wesley, 1986. Rereleased in extended form as [4].
- [6] Hassan Aït-Kaci. *Warren’s Abstract Machine – A Tutorial Reconstruction*. MIT Press, 1991.
- [7] John R. Allen and Ken Kennedy. *Optimizing compilers for modern architectures: a dependence-based approach*. Morgan Kaufmann, 2001.
- [8] Andrew W. Appel. *Compiling with Continuations*. Cambridge University Press, 1992.
- [9] Andrew W. Appel. *Modern Compiler Implementation in ML*. Cambridge University Press, 1998.

- [10] H. Boehm and M. Weiser. Garbage collection in an uncooperative environment. *Software Practice and Experience*, 18(9):807–820, 1988.
- [11] H. Bratman. An alternative form of the ‘UNCOL’ diagram. *Communications of the ACM*, 4(3):142, 1961.
- [12] Preston Briggs. *Register Allocation via Graph Coloring, Tech. Rept. CPC-TR94517-S*. PhD thesis, Rice University, Center for Research on Parallel Computation, Apr. 1992.
- [13] J. A. Brzozowski. Derivatives of regular expressions. *Journal of the ACM*, 1(4):481–494, 1964.
- [14] Noam Chomsky. Three models for the description of language. *IRE Transactions on Information Theory*, IT-2(3):113–124, 1956.
- [15] J. Earley and H. Sturgis. A formalism for translator interactions. *Communications of the ACM*, 13:607–617, 1970.
- [16] Peter Naur (ed.). Revised report on the algorithmic language Algol 60. *Communications of the ACM*, 6(1):1–17, 1963.
- [17] John Hatcliff, Torben Mogensen, and Peter Thiemann (Eds.). *Partial Evaluation: Practice and Theory*, volume 1706 of *Lecture Notes in Computer Science*. Springer Verlag, 1999.
- [18] Raymond J. Hookway and Mark A. Herdeg. Digital fx!32: Combining emulation and binary translation.
<http://www.cs.tufts.edu/comp/150PAT/optimization/DTJP01PF.pdf>, 1997.
- [19] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. *Introduction to Automata Theory, Languages and Computation*, 2nd ed. Addison-Wesley, 2001.
- [20] Kathleen Jensen and Niklaus Wirth. *Pascal User Manual and Report (2nd ed.)*. Springer-Verlag, 1975.
- [21] Neil D. Peyton Jones, Carsten Gomard, and Peter Sestoft. *Partial Evaluation and Automatic Program Generation*. Prentice Hall, 1993.
- [22] Richard E. Jones and Rafael Dueire Lins. *Garbage Collection: Algorithms for Automatic Dynamic Memory Management*. John Wiley, 1996.
- [23] Simon L. Peyton Jones and David Lester. *Implementing Functional Languages – A Tutorial*. Prentice Hall, 1992.

- [24] J. P. Keller and R. Paige. Program derivation with verified transformations – a case study. *Communications in Pure and Applied Mathematics*, 48(9–10), 1996.
- [25] B. W. Kernighan and D. M. Ritchie. *The C Programming Language*. Prentice-Hall, 1978.
- [26] Donald Knuth. *The Art of Computer Programming, Volume 1: Fundamental Algorithms*. Addison-Wesley, 1997.
- [27] James Larus. Assembler, linkers and the spim simulator.
http://pages.cs.wisc.edu/~larus/HP_AppA.pdf, 1998.
- [28] M. E. Lesk. Lex: a Lexical Analyzer Generator. Technical Report 39, AT&T Bell Laboratories, Murray Hill, N. J., 1975.
- [29] T. Lindholm and F. Yellin. *The Java Virtual Machine Specification, 2nd ed.* Addison-Wesley, Reading, Massachusetts, 1999.
- [30] John McCarthy, Paul W. Abrahams, Daniel J. Edwards, Timothy P. Hart, and Michael I. Levin. *LISP 1.5 Programmer's Manual*. The M.I.T. Press, 1962.
- [31] R. McNaughton and H. Yamada. Regular expressions and state graphs for automata. *IEEE Transactions on Electronic Computers*, 9(1):39–47, 1960.
- [32] Robin Milner. A theory of type polymorphism in programming. *Journal of Computational Systems Science*, 17(3):348–375, 1978.
- [33] Robin Milner. *Communication and Concurrency*. Prentice-Hall, 1989.
- [34] Torben Æ. Mogensen, David A. Schmidt, and I. Hal Sudborough, editors. *The essence of computation: complexity, analysis, transformation*. Springer-Verlag New York, Inc., New York, NY, USA, 2002.
- [35] Steven S. Muchnick. *Advanced Compiler Design and Implementation*. Morgan Kaufmann, 1997.
- [36] Flemming Nielson, Hanne R. Nielson, and Chris Hankin. *Principles of Program Analysis*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 1999.
- [37] Chris Okasaki. *Purely Functional Data Structures*. Cambridge University Press, 1998.
- [38] Scott Owens, John Reppy, and Aaron Turon. Regular-expression derivatives re-examined. *J. Funct. Program.*, 19(2):173–190, 2009.

- [39] David A. Patterson and John L. Hennessy. *Computer Organization & Design, the Hardware/Software Interface*. Morgan Kaufmann, 1998.
- [40] Vern Paxson. Flex, version 2.5, a fast scanner generator.
http://www.gnu.org/software/flex/manual/html_mono/flex.html, 1995.
- [41] G. L. Steele and G. J. Sussman. The Art of the Interpreter or, The Modularity Complex. Technical Report AIM-453, Massachusetts Institute of Technology, Cambridge, MA, USA, 1978.
- [42] Mikkel Thorup. All structured programs have small tree-width and good register allocation. *Information and Computation*, 142(2):159–181, 1998.
- [43] Mads Tofte and Jean-Pierre Talpin. Region-based memory management, 1997.
- [44] Paul R. Wilson. Uniprocessor garbage collection techniques. In *IWMM '92: Proceedings of the International Workshop on Memory Management*, pages 1–42, London, UK, 1992. Springer-Verlag.
- [45] Niklaus Wirth. The design of a Pascal compiler. *Software - Practice and Experience*, 1(4):309–333, 1971.

Index

- abstract syntax, **99**, 122
- accept, 89, 93, 94
- action, 41, 99, 100
- activation record, 210
- alias, 222, 223
- allocation, 166, 226
 - dynamic, 259
 - heap, 259
 - stack, 258
 - static, 257
- Alpha, 180, 288
- alphabet, 10
- ARM, 181, 284
- array, 165, 246
- assembly, 3
- assignment, 150
- associative, 64, 65
- attribute, 135
 - inherited, 135
 - synthesised, 135
- available assignments, 233

- back-end, 147
- biased colouring, 205
- binary translation, 288
- binding
 - dynamic, 114
 - static, 114
- bootstrapping, 281, 283
 - full, 285
 - half, 285
 - incremental, 287
- Bratman diagram, 281

- C, 4, 40, 64, 66, 100, 102, 105, 118, 153,
159, 160, 164, 167, 221, 222,
227, 245, 252, 258–260
- C++, 254
- cache, 246
- cache line, 246
- call sequence, 248
- call stack, 209
- call-by-reference, 222
- call-by-value, 209
- call-sequence, 212
- callee-saves, 213, 215
- caller-saves, 213, 215
- caller/callee, 209
- calling convention, 210
- CISC, 181
- closure, 117, 131, 144, 228
- coalescing, 205
- code generator, 180, 183
- code hoisting, 187, **245**
- column-major, 167
- comments
 - nested, 42
- common subexpression elimination, 186,
233, **237**, 246
- compilation, 129
- compile-time, 152
- compiling compilers, 283
- conflict, 82, 87, 94, 97, 102
 - reduce-reduce, 94, 97
 - shift-reduce, 94, 97
- consistent, 31
- constant in operand, 181

- constant propagation, 187
- context-free, 133
 - grammar, 53, **54**, 58
 - language, 104
- dangling-else, 66, 95, 97
- data-flow analysis, 232, 244
- dead code elimination, 241
- dead variable, 182, 192
- declaration, 113
 - global, 113
 - local, 113
- derivation, 58, **58**, 60, 68, 82
 - left, 60, 80
 - leftmost, 60
 - right, 60
 - rightmost, 60, 88
- DFA, 16, **22**, 44, 88, 90
 - combined, 37
 - converting NFA to, 23, 26
 - equivalence of, 30
 - minimisation, **30**, 31, 37
 - unique minimal, 30
- Digital Vax, 228
- distributive, 25
- domain specific language, 4
- dynamic programming, 183
- environment, 114, 135
- epilogue, 211, 248
- epsilon transition, 16
- epsilon-closure, **23**
- execution, 121
- FA, 16
- finite automaton
 - graphical notation, 17
- finite automaton, 10, **16**
 - deterministic, **22**
 - nondeterministic, **16**
- FIRST*, 69, 73
- fixed-point, 24, 71, 72, 194, 195
- flag, 180
 - arithmetic, 181
- floating-point constant, 14
- floating-point numbers, 151
- FOLLOW*, 74
- FORTTRAN, 40
- Fortran, 258
- frame, 210
- frame pointer, 211
- free list, 260
- front-end, 147
- function call, 151, 248
- function calls, 179, 209
- functional, 115
- garbage collection, 266
 - incremental, 275
- garbage collection
 - concurrent, 275
 - scan-sweep, 269
 - tracing, 268
 - two-space, 271
- gen and kill sets, 193
- generic types, 142
- global variable, 221
- go, 89, 92
- grammar, 68
 - ambiguous, 60, 62, 63, 65, 69, 73, 82, 94
 - equivalent, 62
- graph colouring, 199, **200**
- greedy algorithm, 183
- hashing, 117
- Haskell, 102, 115
- heuristics, 199, **202**
- IA-32, 180, 205
- IA-64, 180
- IBM System/370, 227
- imperative, 115
- implicit types, 143

- in and out sets, 193
- index check
 - translation of, 170
- index check, 170
- index-check
 - elimination, 188
- index-check elimination, 241
- inlining, 249
- instruction set description, 183
- integer, 14, 151
- interference, 196
- interference graph, 198
- intermediate code, 2, 147, 191
- intermediate language, 2, 148, 179, 186
 - tree-structured, 188
- interpretation, 121
 - of expressions, 124
 - of function calls, 126
 - of programs, 128
- interpreter, 3, 121, 147, 149, 282
- Java, 40, 100, 148
- jump, 151
 - conditional, 151, 180
- jump-to-jump optimisation, 175, 240
- just-in-time compilation, 148
- keyword, 14
- label, 150
- LALR(1), 98, 105
- language, 10, 58
 - context-free, 104
 - high-level, 147, 281
- left-associative, 64, 97
- left-derivation, 68
- left-factorisation, 86
- left-recursion, 64, 65, 86, 100
 - elimination of, 84
 - indirect, 86
- lexer, 9, **35**, 68
- lexer generator, 36, 41
- lexical, 9
 - analysis, 9
 - error, 40
- lexical analysis, 2
- lexing, 133
- linking, 3
- LISP, 227
- live variable, 192, 209
 - at end of procedure, 194
- live-range splitting, 205
- liveness, **192**
- liveness analysis, 193
- LL(1), 53, 80, **81**, 82, 88, 95, 100, 105
- local variables, 209
- longest prefix, 40
- lookahead, 80
- LR, 88
- machine code, 3, 147, 149, 179
- machine language, 191
- malloc(), 260
- memory management
 - automatic, 266
 - manual, 259
- memory transfer, 151
- MIPS, 180–182, **183**, 188, 229
- monotonic, 24
- name space, 118, 122
- nested scopes, 223, 225
- NFA, 16, 90, 92, 103
 - combined, 36, 37
 - converting to DFA, 23, 26
 - fragment, 18
- non-associative, 64, 97
- non-local variable, 221
- non-recursive, 65
- nonterminal, 54
- Nullable*, 70, 73
- operator, 150
- operator hierarchy, 63, 64

- optimisations, 186
- overloading, 140
- PA-RISC, 181
- parser, 62
 - generator, 63, 98, 102
 - predictive, 68, 73
 - shift-reduce, 88
 - table-driven, 88
 - top-down, 68
- parsing, 53, 60, 133
 - bottom-up, 68
 - predictive, 73, 79, 80
 - table-driven, 81
- Pascal, 4, 64, 66, 100, 105, 118, 222, 223, 259
- pattern, 182
- persistent, 115
- pointer, 222
- polymorphism, 142
- PowerPC, 180
- precedence, 56, 63, 65, 66, 88, 95
 - declaration, 96, 97, 104
 - rules, 63
- prefetch, 246
- processor, 281
- production, 54, 55
 - empty, 55, 73
 - nullable, 70, 73
- prologue, 211, 248
- recursive descent, **80**
- reduce, 88, 89, 93
- reference counting, 266
- register, 191
 - for passing function parameters, 215
- register allocation, 191
 - global, 198
- register allocation, 2, 179
 - by graph colouring, 199
- register allocator, 219
- regular expression, **10**, 41
 - converting to NFA, 18
 - equivalence of, 30
- regular language, 30, 42
- return address, 210, 216
- right-associative, 64, 97
- right-recursion, 65
- RISC, 179, 181, 216
- row-major, 167
- run-time, 152
- Scheme, 102, 115
- scope, 113
 - nested, 223, 225
- select, 200
- sequential logical operators, 159, 160
- set constraints, 75
- set equation, 23, **23**
- shift, 88, 89, 92
- simplify, 200
- SLR, 53, 88, 94, 95
 - algorithm, 91
 - construction of table, 90, 95
- SML, 4, 40, 64, 102, 115, 118, 223
- source program, 283
- Sparc, 180
- spill, 211
- spill code, 202
- spilling, 191, **200**
- stack automaton, 53
- stack automaton, 104
- stack pointer, 226
- start symbol, 54, 68
- starting state, 16
- state, 16, 17
 - accepting, 16–18, 27, 31, 36
 - dead, 34
 - final, 16
 - initial, 16, 17
 - starting, 16, 18, 27, 36
- static links, 225
- subset construction, 26

- symbol table, 114, **114**, 124, 135
 - implemented as list, 115
 - implemented as function, 116
 - implemented as stack, 117
- syntactic category, 57, 122, 135
- syntax analysis, 2, 9, 53, 58, 60, **68**
- syntax tree, 53, **60**, 68, 85
- T-diagram, 281
- tail call, 250
- tail-call optimisation, 250
- target program, 283
- templates, 254
- terminal, 54
- thunk, 228
- token, 9, 36, 37, 41, 68
- transition, 16, 17, 27, 31
 - epsilon, 16, 93
- translation
 - of arrays, 165
 - of case-statements, 164
 - of declarations, 172
 - of expressions, 152
 - of function, 220
 - of index checks, 170
 - of logical operators, 159, 160
 - of multi-dimensional arrays, 167
 - of non-zero-based arrays, 170
 - of records/structs, 171
 - of statements, 156
 - of strings, 171
 - of break/exit/continue, 164
 - of goto, 164
- type checking, 129
- type checking, 2, 133, 135
 - of assignments, 140
 - of data structures, 140
 - of expressions, 136
 - of function declarations, 138
 - of programs, 139
- type conversion, 142
- type error, 136
- undecidable, 62
- value numbering, 239
- value numbering, 254
- variable
 - global, 221
 - non-local, 221
- variable name, 14
- white-space, 9, 41
- word length, 166
- work-list algorithm, 26
- x86, 180, 284, 288