

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/2367776>

An Introduction to Compilers

Article · May 2001

Source: CiteSeer

CITATION

1

READS

611

1 author:



[Dirk Vermeir](#)

Vrije Universiteit Brussel

155 PUBLICATIONS 1,449 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Syllabus [View project](#)

An introduction to compilers

DRAFT

D. Vermeir
Dept. of Computer Science
Vrij Universiteit Brussel, VUB
dvermeir@vub.ac.be



July 11, 2008

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 6 |
| 1.1 | Compilers and languages | 6 |
| 1.2 | Applications of compilers | 7 |
| 1.3 | Overview of the compilation process | 9 |
| 1.3.1 | Micro | 9 |
| 1.3.2 | x86 code | 10 |
| 1.3.3 | Lexical analysis | 12 |
| 1.3.4 | Syntax analysis | 13 |
| 1.3.5 | Semantic analysis | 14 |
| 1.3.6 | Intermediate code generation | 15 |
| 1.3.7 | Optimization | 16 |
| 1.3.8 | Code generation | 17 |
| 2 | Lexical analysis | 18 |
| 2.1 | Introduction | 18 |
| 2.2 | Regular expressions | 24 |
| 2.3 | Finite state automata | 26 |
| 2.3.1 | Deterministic finite automata | 26 |
| 2.3.2 | Nondeterministic finite automata | 28 |
| 2.4 | Regular expressions vs finite state automata | 31 |
| 2.5 | A scanner generator | 32 |

| | | |
|----------|--|-----------|
| 3 | Parsing | 35 |
| 3.1 | Context-free grammars | 35 |
| 3.2 | Top-down parsing | 38 |
| 3.2.1 | Introduction | 38 |
| 3.2.2 | Eliminating left recursion in a grammar | 41 |
| 3.2.3 | Avoiding backtracking: LL(1) grammars | 43 |
| 3.2.4 | Predictive parsers | 44 |
| 3.2.5 | Construction of first and follow | 48 |
| 3.3 | Bottom-up parsing | 50 |
| 3.3.1 | Shift-reduce parsers | 50 |
| 3.3.2 | LR(1) parsing | 54 |
| 3.3.3 | LALR parsers and yacc/bison | 62 |
| 4 | Checking static semantics | 65 |
| 4.1 | Attribute grammars and syntax-directed translation | 65 |
| 4.2 | Symbol tables | 68 |
| 4.2.1 | String pool | 69 |
| 4.2.2 | Symbol tables and scope rules | 69 |
| 4.3 | Type checking | 71 |
| 5 | Intermediate code generation | 74 |
| 5.1 | Postfix notation | 75 |
| 5.2 | Abstract syntax trees | 76 |
| 5.3 | Three-address code | 78 |
| 5.4 | Translating assignment statements | 79 |
| 5.5 | Translating boolean expressions | 81 |
| 5.6 | Translating control flow statements | 85 |
| 5.7 | Translating procedure calls | 86 |
| 5.8 | Translating array references | 88 |

| | | |
|----------|---|------------|
| 6 | Optimization of intermediate code | 92 |
| 6.1 | Introduction | 92 |
| 6.2 | Local optimization of basic blocks | 94 |
| 6.2.1 | DAG representation of basic blocks | 95 |
| 6.2.2 | Code simplification | 99 |
| 6.2.3 | Array and pointer assignments | 100 |
| 6.2.4 | Algebraic identities | 101 |
| 6.3 | Global flow graph information | 101 |
| 6.3.1 | Reaching definitions | 103 |
| 6.3.2 | Reaching definitions using datalog | 105 |
| 6.3.3 | Available expressions | 106 |
| 6.3.4 | Available expressions using datalog | 109 |
| 6.3.5 | Live variable analysis | 110 |
| 6.3.6 | Definition-use chaining | 112 |
| 6.3.7 | Application: uninitialized variables | 113 |
| 6.4 | Global optimization | 113 |
| 6.4.1 | Elimination of global common subexpressions | 113 |
| 6.4.2 | Copy propagation | 114 |
| 6.4.3 | Constant folding and elimination of useless variables | 116 |
| 6.4.4 | Loops | 116 |
| 6.4.5 | Moving loop invariants | 120 |
| 6.4.6 | Loop induction variables | 123 |
| 6.5 | Aliasing: pointers and procedure calls | 126 |
| 6.5.1 | Pointers | 128 |
| 6.5.2 | Procedures | 128 |
| 7 | Code generation | 130 |
| 7.1 | Run-time storage management | 131 |
| 7.1.1 | Global data | 131 |
| 7.1.2 | Stack-based local data | 132 |
| 7.2 | Instruction selection | 134 |

| | | |
|----------|--|------------|
| 7.3 | Register allocation | 136 |
| 7.4 | Peephole optimization | 137 |
| A | A Short Introduction to x86 Assembler Programming under Linux | 139 |
| A.1 | Architecture | 139 |
| A.2 | Instructions | 140 |
| A.2.1 | Operands | 140 |
| A.2.2 | Addressing Modes | 140 |
| A.2.3 | Moving Data | 141 |
| A.2.4 | Integer Arithmetic | 142 |
| A.2.5 | Logical Operations | 142 |
| A.2.6 | Control Flow Instructions | 142 |
| A.3 | Assembler Directives | 143 |
| A.4 | Calling a function | 144 |
| A.5 | System calls | 144 |
| A.6 | Example | 146 |
| B | Mc: the Micro-x86 Compiler | 149 |
| B.1 | Lexical analyzer | 149 |
| B.2 | Symbol table management | 151 |
| B.3 | Parser | 152 |
| B.4 | Driver script | 155 |
| B.5 | Makefile | 156 |
| B.6 | Example | 157 |
| B.6.1 | Source program | 157 |
| B.6.2 | Assembly language program | 157 |
| C | Minic parser and type checker | 159 |
| C.1 | Lexical analyzer | 159 |
| C.2 | String pool management | 161 |
| C.3 | Symbol table management | 163 |

| | | |
|---------------------|--|------------|
| C.4 | Types library | 166 |
| C.5 | Type checking routines | 172 |
| C.6 | Parser with semantic actions | 175 |
| C.7 | Utilities | 178 |
| C.8 | Driver script | 179 |
| C.9 | Makefile | 180 |
| Index | | 181 |
| Bibliography | | 186 |

Chapter 1

Introduction

1.1 Compilers and languages

A *compiler* is a program that translates a *source language* text into an equivalent *target language* text.

E.g. for a C compiler, the source language is C while the target language may be Sparc assembly language.

Of course, one expects a compiler to do a faithful translation, i.e. the *meaning* of the translated text should be the same as the meaning of the source text.

One would not be pleased to see the C program in Figure 1.1

```
1 #include          <stdio.h>
2
3 int
4 main(int,char**)
5 {
6 int x = 34;
7 x = x*24;
8 printf("%d\n",x);
9 }
```

Figure 1.1: A source text in the C language

translated to an assembler program that, when executed, printed “Goodbye world” on the standard output.

So we want the translation performed by a compiler to be *semantics preserving*. This implies that the compiler is able to “understand” (compute the semantics of)

the source text. The compiler must also “understand” the target language in order to be able to generate a semantically equivalent target text.

Thus, in order to develop a compiler, we need a precise definition of both the source and the target language. This means that both source and target language must be *formal*.

A language has two aspects: a *syntax* and a *semantics*. The syntax prescribes which texts are grammatically correct and the semantics specifies how to derive the meaning from a syntactically correct text. For the C language, the syntax specifies e.g. that

“the body of a function must be enclosed between matching braces (“{ }”)”.

The semantics says that the meaning of the second statement in Figure 1.1 is that

“the value of the variable x is multiplied by 24 and the result becomes the new value of the variable x ”

It turns out that there exist excellent formalisms and tools to describe the syntax of a formal language. For the description of the semantics, the situation is less clear in that existing semantics specification formalisms are not nearly as simple and easy to use as syntax specifications.

1.2 Applications of compilers

Traditionally, a compiler is thought of as translating a so-called “high level language” such as C¹ or Modula2 into assembly language. Since assembly language cannot be directly executed, a further translation between assembly language and (relocatable) machine language is necessary. Such programs are usually called *assemblers* but it is clear that an assembler is just a special (easier) case of a compiler.

Sometimes, a compiler translates between high level languages. E.g. the first C++ implementations used a compiler called “cfront” which translated C++ code to C code. Such a compiler is often called a “cross-compiler”.

On the other hand, a compiler need not target a real assembly (or machine) language. E.g. Java compilers generate code for a virtual machine called the “Java

¹If you want to call C a high-level language

Virtual Machine” (JVM). The JVM interpreter then interprets JVM instructions without any further translation.

In general, an interpreter needs to understand only the source language. Instead of translating the source text, an interpreter immediately executes the instructions in the source text. Many languages are usually “interpreted”, either directly, or after a compilation to some virtual machine code: Lisp, Smalltalk, Prolog, SQL are among those. The advantages of using an interpreter are that is easy to port a language to a new machine: all one has to do is to implement the virtual machine on the new hardware. Also, since instructions are evaluated and examined at run-time, it becomes possible to implement very flexible languages. E.g. for an interpreter it is not a problem to support variables that have a dynamic type, something which is hard to do in a traditional compiler. Interpreters can even construct “programs” at run time and interpret those without difficulties, a capability that is available e.g. for Lisp or Prolog.

Finally, compilers (and interpreters) have wider applications than just translating programming languages. Conceivably any large and complex application might define its own “command language” which can be translated to a virtual machine associated with the application. Using compiler generating tools, defining and implementing such a language need not be difficult. Hence SQL can be regarded as such a language associated with a database management system. Other so-called “little languages” provide a convenient interface to specialized libraries. E.g. the language (n)awk is a language that is very convenient to do powerful pattern matching and extraction operations on large text files.

1.3 Overview of the compilation process

In this section we will illustrate the main phases of the compilation process through a simple compiler for a toy programming language. The source for an implementation of this compiler can be found in Appendix B and on the web site of the course.

| | | |
|------------------|---|---------------------------------|
| program | : | declaration_list statement_list |
| | : | |
| declaration_list | : | declaration ; declaration_list |
| | | ε |
| | : | |
| declaration | : | declare var |
| | : | |
| statement_list | : | statement ; statement_list |
| | | ε |
| | : | |
| statement | : | assignment |
| | | read_statement |
| | | write_statement |
| | : | |
| assignment | : | var = expression |
| | : | |
| read_statement | : | read var |
| | : | |
| write_statement | : | write expression |
| | : | |
| expression | : | term |
| | | term + term |
| | | term − term |
| | : | |
| term | : | <i>NUMBER</i> |
| | | var |
| | | (expression) |
| | : | |
| var | : | <i>NAME</i> |
| | : | |

Figure 1.2: The syntax of the Micro language

1.3.1 Micro

The source language “Micro” is very simple. It is based on the toy language described in [FL91].

The syntax of Micro is described by the rules in Figure 1.2. We will see in Chapter 3 that such rules can be formalized into what is called a *grammar*.

Note that *NUMBER* and *NAME* have not been further defined. The idea is, of course, that *NUMBER* represents a sequence of digits and that *NAME* represents a string of letters and digits, starting with a letter.

A simple Micro program is shown in Figure 1.3

```
{  
  declare xyz;  
  xyz = (33+3)-35;  
  write xyz;  
}
```

Figure 1.3: A Micro program

The semantics of Micro should be clear²: a Micro program consists of a sequence of read/write or assignment statements. There are integer-valued variables (which need to be declared before they are used) and expressions are restricted to addition and subtraction.

1.3.2 x86 code

The target language will be code for the x86 processor family. Figure 1.4 shows part of the output of the compiler for the program of Figure 1.3. The full output can be found in Section B.6.2, page 157.

X86 processors have a number of registers, some of which are special purpose, such as the **esp** register which always points to the top of the stack (which grows downwards). More information on x86 assembler programming can be found in the Appendix, Section A, page 139.

- line 1 The code is divided into data and text sections where the latter contains the actual instructions.
- line 2 This defines a data area of 4 bytes wide which can be referenced using the name *xyz*. This definition is the translation of a Micro *declaration*.

²The output of the program in Figure 1.3 is, of course, 1.

```

1      .section .data
2      .lcomm xyz, 4
3      .section .text
...
44     .globl main
45     .type main, @function
46     main:
47         pushl %ebp
48         movl %esp, %ebp
49         pushl $33
50         pushl $3
51         popl %eax
52         addl %eax, (%esp)
53         pushl $35
54         popl %eax
55         subl %eax, (%esp)
56         popl xyz
57         pushl xyz
58         call print
59         movl %ebp, %esp
60         popl %ebp
61         ret
...

```

Figure 1.4: X86 assembly code generated for the program in Figure 1.3

- line 44 This defines *main* as a globally available name. It will be used to refer to the single function (line 45) that contains the instructions corresponding to the Micro program. The function starts at line 46 where the location corresponding to the label '*main*' is defined.
- line 47 Together with line 48, this starts off the function according to the C calling conventions: the current top of stack contains the return address. This address has been pushed on the stack by the (function) **call** instruction. It will eventually be used (and popped) by a subsequent **ret** (return from function call) instruction. Parameters are passed by pushing them on the stack just before the **call** instruction. Line 47 saves the caller's "base pointer" in the **ebp** register by pushing it on the stack before setting (line 48) the current top of the stack as a new base pointer **ebp**. When returning from the function call, the original stack is restored by copying the saved value from **ebp** to **esp** (line 59) and popping the saved base pointer (line 60).
- line 49 The evaluation of the subexpression $33 + 3$ is initiated by pushing both constants (indicated by the use of a '\$' prefix) on the stack (lines 49, 50)

- line 51 Once both operands are on the top of the stack, the operation (corresponding to the $33 + 3$ expression) is executed by popping the second argument to the **eax** register (line 51) which is then added to the first argument on the top of the stack (line 52). The net result is that the two arguments of the operation are replaced on the top of the stack by the (single) result of the operation.
- line 53 To subtract 35 from the result, this second argument of the subtraction is pushed on the stack (line 53), after which the subtraction is executed on the two operands on the stack, replacing them on the top of the stack by the result (lines 54,55).
- line 56 The result of evaluating the expression is assigned to the variable *xyz* by popping it from the stack to the appropriate address.
- line 57 In order to print the value at address *xyz*, it is first pushed on the stack as a parameter for a subsequent call (line 58) to a *print* function (the code of which can be found in Section B.6.2).

1.3.3 Lexical analysis

The raw input to a compiler consists of a string of bytes or characters. Some of those characters, e.g. the “{” character in Micro, may have a meaning by themselves. Other characters only have meaning as part of a larger unit. E.g. the “y” in the example program from Figure 1.3, is just a part of the *NAME* “xyz”. Still others, such as “ ”, “\n” serve as separators to distinguish one meaningful string from another.

The first job of a compiler is then to group sequences of raw characters into meaningful *tokens*. The lexical analyzer module is responsible for this. Conceptually, the lexical analyzer (often called *scanner*) transforms a sequence of characters into a sequence of tokens. In addition, a lexical analyzer will typically access the *symbol table* to store and/or retrieve information on certain source language concepts such as variables, functions, types.

For the example program from Figure 1.3, the lexical analyzer will transform the character sequence

```
{ declare xyz; xyz = (33+3)-35; write xyz; }
```

into the token sequence shown in Figure 1.5.

Note that some tokens have “properties”, e.g. a $\langle \text{NUMBER} \rangle$ token has a *value* property while a $\langle \text{NAME} \rangle$ token has a symbol table reference as a property.

```

⟨LBRACE⟩
⟨DECLARE symbol_table_ref=0⟩
⟨NAME symbol_table_ref=3⟩
⟨SEMICOLON⟩
⟨NAME symbol_table_ref=3⟩
⟨ASSIGN⟩
⟨LPAREN⟩
⟨NUMBER value=33⟩
⟨PLUS⟩
⟨NUMBER value=3⟩
⟨RPAREN⟩
⟨MINUS⟩
⟨NUMBER value=35⟩
⟨SEMICOLON⟩
⟨WRITE symbol_table_ref=2⟩
⟨NAME symbol_table_ref=3⟩
⟨SEMICOLON⟩
⟨RBRACE⟩

```

Figure 1.5: Result of lexical analysis of program in Figure 1.3

After the scanner finishes, the symbol table in the example could look like

| | | |
|---|-----------|---------|
| 0 | "declare" | DECLARE |
| 1 | "read" | READ |
| 2 | "write" | WRITE |
| 3 | "xyz" | NAME |

where the third column indicates the type of symbol.

Clearly, the main difficulty in writing a lexical analyzer will be to decide, while reading characters one by one, when a token of which type is finished. We will see in Chapter 2 that regular expressions and finite automata provide a powerful and convenient method to automate this job.

1.3.4 Syntax analysis

Once lexical analysis is finished, the *parser* takes over to check whether the sequence of tokens is grammatically correct, according to the rules that define the syntax of the source language.

Looking at the grammar rules for Micro (Figure 1.2), it seems clear that a program is syntactically correct if the structure of the tokens matches the structure of a *⟨program⟩* as defined by these rules.

Such matching can conveniently be represented as a *parse tree*. The parse tree corresponding to the token sequence of Figure 1.5 is shown in Figure 1.6.

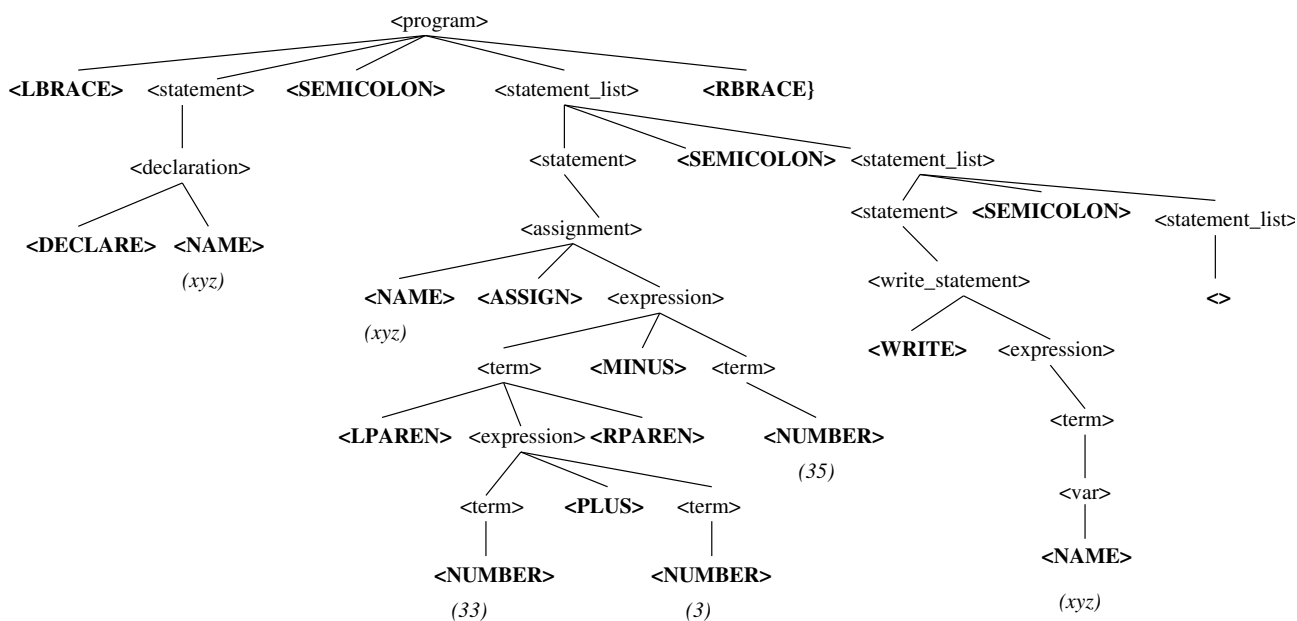


Figure 1.6: Parse tree of program in Figure 1.3

Note that in the parse tree, a node and its children correspond to a rule in the syntax specification of Micro: the parent node corresponds to the left hand side of the rule while the children correspond to the right hand side. Furthermore, the yield³ of the parse tree is exactly the sequence of tokens that resulted from the lexical analysis of the source text.

Hence the job of the parser is to construct a parse tree that fits, according to the syntax specification, the token sequence that was generated by the lexical analyzer.

In Chapter 3, we'll see how context-free grammars can be used to specify the syntax of a programming language and how it is possible to automatically generate parser programs from such a context-free grammar.

1.3.5 Semantic analysis

Having established that the source text is syntactically correct, the compiler may now perform additional checks such as determining the type of expressions and

³The *yield* of a tree is the sequence of leafs of the tree in lexicographical (left-to-right) order

checking that all statements are correct with respect to the typing rules, that variables have been properly declared before they are used, that functions are called with the proper number of parameters etc.

This phase is carried out using information from the parse tree and the symbol table. In our example, very little needs to be checked, due to the extreme simplicity of the language. The only check that is performed verifies that a variable has been declared before it is used.

1.3.6 Intermediate code generation

In this phase, the compiler translates the source text into an simple intermediate language. There are several possible choices for an intermediate language. but in this example we will use the popular “three-address code” format. Essentially, three-address code consists of assignments where the right-hand side must be a single variable or constant or the result of a binary or unary operation. Hence an assignment involves at most three variables (addresses), which explains the name. In addition, three-address code supports primitive control flow statements such as *goto*, *branch-if-positive* etc. Finally, retrieval from and storing into a one-dimensional array is also possible.

The translation process is *syntax-directed*. This means that

- Nodes in the parse tree have a set of *attributes* that contain information pertaining to that node. The set of attributes of a node depends on the kind of syntactical concept it represents. E.g. in Micro, an attribute of an $\langle expression \rangle$ could be the sequence of x86 instructions that leave the result of the evaluation of the expression on the top of the stack. Similarly, both $\langle var \rangle$ and $\langle expression \rangle$ nodes have a *name* attribute holding the name of the variable containing the current value of the $\langle var \rangle$ or $\langle expression \rangle$

We use $n.a$ to refer to the value of the attribute a for the node n .

- A number of *semantic rules* are associated with each syntactic rule of the grammar. These semantic rules determine the values of the attributes of the nodes in the parse tree (a parent node and its children) that correspond to such a syntactic rule. E.g. in Micro, there is a semantic rule that says that the code associated with an $\langle assignment \rangle$ in the rule

$$assignment : var = expression$$

consists of the code associated with $\langle expression \rangle$ followed by a three-address code statement of the form

$$var.name = expression.name$$

More formally, such a semantic rule might be written as

$$assignment.code = expression.code \parallel \text{“var.name = expression.name”}$$

- The translation of the source text then consists of the value of a particular attribute for the root of the parse tree.

Thus intermediate code generation can be performed by computing, using the semantic rules, the attribute values of all nodes in the parse tree. The result is then the value of a specific (e.g. “code”) attribute of the root of the parse tree.

For the example program from Figure 1.3, we could obtain the three-address code in Figure 1.7.

```
T0 = 33 + 3
T1 = T0 - 35
XYZ = T1
WRITE XYZ
```

Figure 1.7: three-address code corresponding to the program of Figure 1.3

Note the introduction of several temporary variables, due to the restrictions inherent in three-address code. The last statement before the *WRITE* may seem wasteful but this sort of inefficiency is easily taken care of by the next optimization phase.

1.3.7 Optimization

In this phase, the compiler tries several optimization methods to replace fragments of the intermediate code text with equivalent but faster (and usually also shorter) fragments.

Techniques that can be employed include common subexpression elimination, loop invariant motion, constant folding etc. Most of these techniques need extra information such as a flow graph, live variable status etc.

In our example, the compiler could perform constant folding and code reordering resulting in the optimized code of Figure 1.8.

XYZ = 1
WRITE XYZ

Figure 1.8: Optimized three-address code corresponding to the program of Figure 1.3

1.3.8 Code generation

The final phase of the compilation consists of the generation of target code from the intermediate code. When the target code corresponds to a register machine, a major problem is the efficient allocation of scarce but fast registers to variables. This problem may be compared with the paging strategy employed by virtual memory management systems. The goal is in both cases to minimize traffic between fast (the registers for a compiler, the page frames for an operating system) and slow (the addresses of variables for a compiler, the pages on disk for an operating system) memory. A significant difference between the two problems is that a compiler has more (but not perfect) knowledge about future references to variables, so more optimization opportunities exist.

Chapter 2

Lexical analysis

2.1 Introduction

As seen in Chapter 1, the lexical analyzer must transform a sequence of “raw” characters into a sequence of *tokens*. Often a token has a structure as in Figure 2.1.

```
1 #ifndef LEX_H
2 #define LEX_H
3 //      %M% (%I%)      %U%      %E%
4
5 typedef enum { NAME, NUMBER, LBRACE, RBRACE, LPAREN, RPAREN, ASSIGN,
6               SEMICOLON, PLUS, MINUS, ERROR } TOKEN;
7
8 typedef struct
9 {
10     TOKEN  type;
11     union {
12         int    value; /* type == NUMBER */
13         char   *name; /* type == NAME */
14     } info;
15     }      TOKEN;
16
17 extern TOKEN *lex();
18 #endif LEX_H
```

Figure 2.1: A declaration for TOKEN and lex()

Actually, the above declaration is not realistic. Usually, more “complex” tokens such as *NAMEs* will refer to a symbol table entry rather than simply their string representation.

Clearly, we can split up the scanner using a function *lex()* as in Figure 2.1 which returns the next token from the source text.

It is not impossible¹ to write such a function by hand. A simple implementation of a hand-made scanner for Micro (see Chapter 1 for a definition of “Micro”) is shown below.

```

1  //      %M% (%I%)      %U%      %E%
2
3  #include      <stdio.h>      /* for getchar() and friends */
4  #include      <ctype.h>      /* for isalpha(), isdigit() and friends */
5  #include      <stdlib.h>      /* for atoi() */
6  #include      <string.h>      /* for strdup() */
7
8  #include      "lex.h"
9
10 static int      state      = 0;
11
12 #define MAXBUF      256
13 static char      buf[MAXBUF];
14 static char*      pbuf;
15
16 static char*      token_name[]      =
17 {
18     "NAME",      "NUMBER",      "LBRACE",      "RBRACE",
19     "LPAREN",      "RPAREN",      "ASSIGN",      "SEMICOLON",
20     "PLUS",      "MINUS",      "ERROR"
21 };
22
23 static TOKEN      token;
24 /*
25  *      This code is not robust: no checking on buffer overflow, ...
26  *      Nor is it complete: keywords are not checked but lumped into
27  *      the 'NAME' token type, no installation in symbol table, ...
28  */
29 TOKEN*
30 lex()
31 {
32     char      c;
33
34     while (1)
35         switch(state)
36         {
37             case 0: /* stands for one of 1,4,6,8,10,13,15,17,19,21,23 */
38                 pbuf      = buf;
39                 c = getchar();
40                 if (isspace(c))
41                     state = 11;

```

¹Some people actually enjoy this.

```

42         else if (isdigit(c))
43             {
44                 *pbuf++ = c; state = 2;
45             }
46         else if (isalpha(c))
47             {
48                 *pbuf++ = c; state = 24;
49             }
50         else switch(c)
51             {
52                 case '{': state = 5; break;
53                 case '}': state = 7; break;
54                 case '(': state = 9; break;
55                 case ')': state = 14; break;
56                 case '+': state = 16; break;
57                 case '-': state = 18; break;
58                 case '=': state = 20; break;
59                 case ';': state = 22; break;
60                 default:
61                     state = 99; break;
62             }
63         break;
64     case 2:
65         c = getchar();
66         if (isdigit(c))
67             *pbuf++ = c;
68         else
69             state = 3;
70         break;
71     case 3:
72         token.info.value= atoi(buf);
73         token.type      = NUMBER;
74         ungetc(c,stdin);
75         state = 0; return &token;
76         break;
77     case 5:
78         token.type      = LBRACE;
79         state = 0; return &token;
80         break;
81     case 7:
82         token.type      = RBRACE;
83         state = 0; return &token;
84         break;
85     case 9:
86         token.type      = LPAREN;
87         state = 0; return &token;
88         break;
89     case 11:
90         c = getchar();

```

```
91         if (isspace(c))
92             ;
93         else
94             state = 12;
95         break;
96     case 12:
97         ungetc(c, stdin);
98         state = 0;
99         break;
100    case 14:
101        token.type      = RPAREN;
102        state = 0; return &token;
103        break;
104    case 16:
105        token.type      = PLUS;
106        state = 0; return &token;
107        break;
108    case 18:
109        token.type      = MINUS;
110        state = 0; return &token;
111        break;
112    case 20:
113        token.type      = ASSIGN;
114        state = 0; return &token;
115        break;
116    case 22:
117        token.type      = SEMICOLON;
118        state = 0; return &token;
119        break;
120    case 24:
121        c = getchar();
122        if (isalpha(c) || isdigit(c))
123            *pbuf++ = c;
124        else
125            state = 25;
126        break;
127    case 25:
128        *pbuf = (char)0;
129        token.info.name = strdup(buf);
130        token.type      = NAME;
131        ungetc(c, stdin);
132        state = 0; return &token;
133        break;
134    case 99:
135        if (c==EOF)
136            return 0;
137        fprintf(stderr, "Illegal character: \'%c\'\n", c);
138        token.type      = ERROR;
139        state = 0; return &token;
```

```

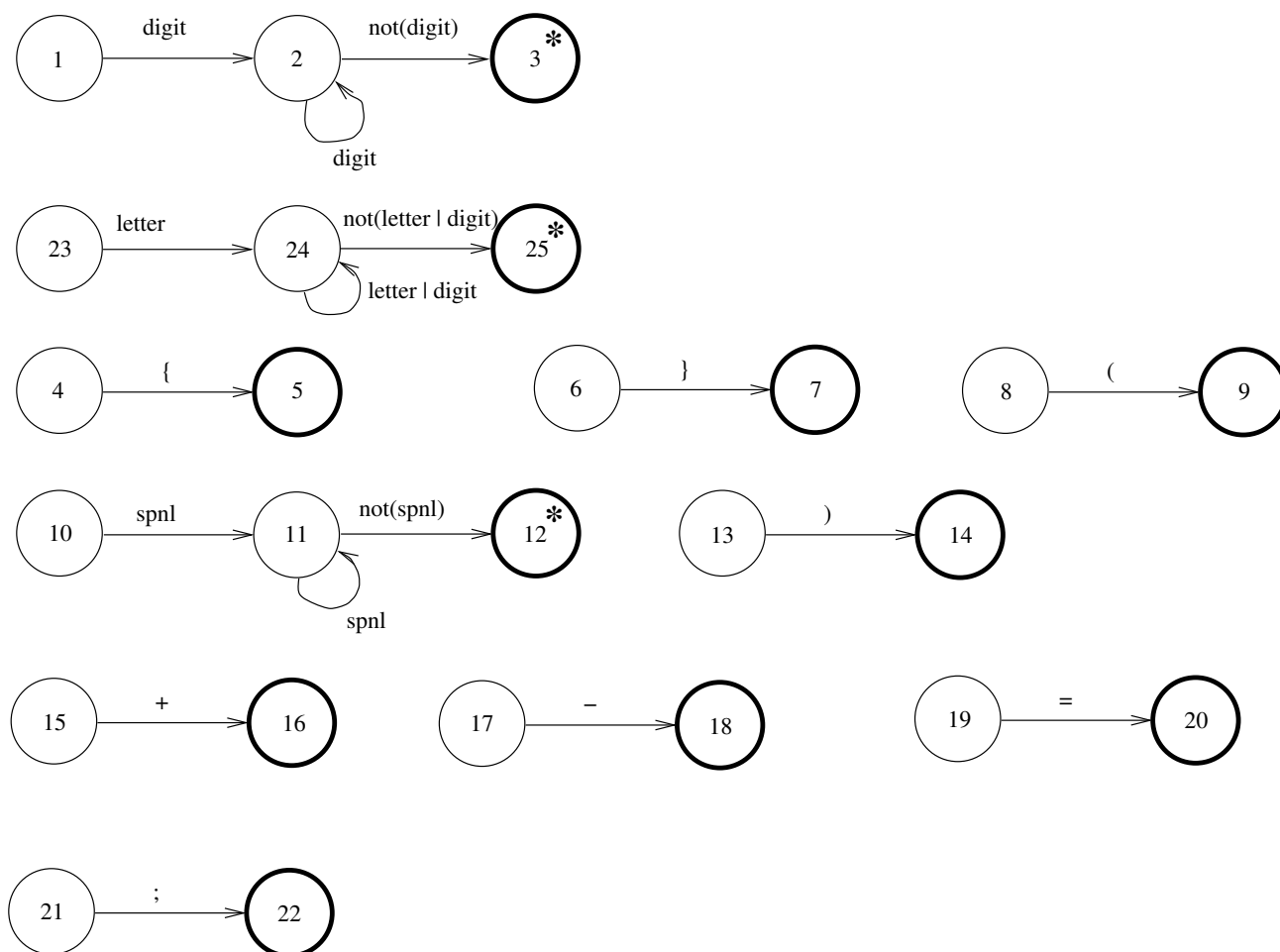
140             break;
141         default:
142             break; /* Cannot happen */
143     }
144 }
145
146 int
147 main()
148 {
149     TOKEN    *t;
150
151     while ((t=lex()))
152     {
153         printf("%s", token_name[t->type]);
154         switch (t->type)
155         {
156             case NAME:
157                 printf(": %s\n", t->info.name);
158                 break;
159             case NUMBER:
160                 printf(": %d\n", t->info.value);
161                 break;
162             default:
163                 printf("\n");
164                 break;
165         }
166     }
167     return 0;
168 }

```

The control flow in the above *lex()* procedure can be represented by a combination of so-called transition diagrams which are shown in Figure 2.2.

There is a transition diagram for each token type and another one for white space (blank, tab, newline). The code for *lex()* simply implements those diagrams. The only complications are

- When starting a new token (i.e., upon entry to *lex()*), we use a “special” state 0 to represent the fact that we didn’t decide yet which diagram to follow. The choice here is made on the basis of the next input character.
- In Figure 2.2, bold circles represent states where we are sure which token has been recognized. Sometimes (e.g. for the *LBRACE* token type) we know this immediately after scanning the last character making up the token. However, for other types, like *NUMBER*, we only know the full extent of the token after reading an extra character that will not belong to the to-

Figure 2.2: The transition diagram for *lex()*

ken. In such a case we must push the extra character back onto the input before returning. Such states have been marked with a * in Figure 2.2.

- If we read a character that doesn't fit any transition diagram, we return a special *ERROR* token type.

Clearly, writing a scanner by hand seems to be easy, once you have a set of transition diagrams such as the ones in Figure 2.2. It is however also boring, and error-prone, especially if there are a large number of states.

Fortunately, the generation of such code can be automated. We will describe how a specification of the various token types can be automatically converted in code that implements a scanner for such token types.

First we will design a suitable formalism to specify token types.

2.2 Regular expressions

In Micro, a *NUMBER* token represents a digit followed by 0 or more digits. A *NAME* consists of a letter followed by 0 or more alphanumeric characters. A *LBRACE* token consists of exactly one “{” character, etc.

Such specifications are easily formalized using *regular expressions*. Before defining regular expressions we should recall the notion of *alphabet* (a finite set of abstract symbols, e.g. ASCII characters), and (*formal*) *language* (a set of strings containing symbols from some alphabet).

The length of a string w , denoted $|w|$ is defined as the number of symbols occurring in w . The prefix of length l of a string w , denoted $\text{pref}_l(w)$ is defined as the longest string x such that $|x| \leq l$ and $w = xy$ for some string y . The *empty string* (of length 0) is denoted ϵ . The *product* $L_1.L_2$ of two languages is the language

$$L_1.L_2 = \{xy \mid x \in L_1 \wedge y \in L_2\}$$

The *closure* L^* of a language L is defined by

$$L^* = \bigcup_{i \in \mathbb{N}} L^i$$

(where, of course, $L^0 = \{\epsilon\}$ and $L^{i+1} = L.L^i$).

Definition 1 *The following table, where r and s denote arbitrary regular expressions, recursively defines all regular expressions over a given alphabet Σ , together with the language L_x each expression x represents.*

| Regular expression | Language |
|--------------------|-----------------|
| \emptyset | \emptyset |
| ϵ | $\{\epsilon\}$ |
| a | $\{a\}$ |
| $(r + s)$ | $L_r \cup L_s$ |
| (rs) | $L_r \cdot L_s$ |
| (r^*) | L_r^* |

In the table, r and s denote arbitrary regular expressions, and $a \in \Sigma$ is an arbitrary symbol from Σ .

A language L for which there exists a regular r expression such that $L_r = L$ is called a **regular language**.

We assume that the operators $+$, concatenation and $*$ have increasing precedence, allowing us to drop many parentheses without risking confusion. Thus, $((0(1^*)) + 0)$ may be written as $01^* + 0$.

From Figure 2.2 we can deduce regular expressions for each token type, as shown in Figure 2.1. We assume that

$$\Sigma = \{a, \dots, z, A, \dots, Z, 0, \dots, 9, \text{SP}, \text{NL}, (,), +, =, \{, \}, ;, -\}$$

| Token type or abbreviation | Regular expression |
|----------------------------|--|
| letter | $a + \dots + z + A + \dots + Z$ |
| digit | $0 + 1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + 9$ |
| NUMBER | $\text{digit}(\text{digit})^*$ |
| NAME | $\text{letter}(\text{letter} + \text{digit})^*$ |
| space | $(\text{SP} + \text{NL})(\text{SP} + \text{NL})^*$ |
| LBRACE | $\{$ |
| .. | .. |

Table 2.1: Regular expressions describing Micro tokens

A full specification, such as the one in Section B.1, page 149, then consists of a set of (extended) regular expressions, plus C code for each expression. The idea is that the generated scanner will

- Process input characters, trying to find a longest string that matches any of the regular expressions².
- Execute the code associated with the selected regular expression. This code can, e.g. install something in the symbol table, return a token type or whatever.

In the next section we will see how a regular expression can be converted to a so-called *deterministic finite automaton* that can be regarded as an abstract machine to recognize strings described by regular expressions. Automatic translation of such an automaton to actual code will turn out to be straightforward.

²If two expressions match the same longest string, the one that was declared first is chosen.

2.3 Finite state automata

2.3.1 Deterministic finite automata

Definition 2 A *deterministic finite automaton (DFA)* is a tuple

$$(Q, \Sigma, \delta, q_0, F)$$

where

- Q is a finite set of *states*,
- Σ is a finite *input alphabet*
- $\delta : Q \times \Sigma \rightarrow Q$ is a (total) transition function
- $q_0 \in Q$ is the *initial state*
- $F \subseteq Q$ is the set of *final states*

Definition 3 Let $M = (Q, \Sigma, \delta, q_0, F)$ be a DFA. A **configuration** of M is a pair $(q, w) \in Q \times \Sigma^*$. For a configuration (q, aw) (where $a \in \Sigma$), we write

$$(q, aw) \vdash_M (q', w)$$

just when $\delta(q, a) = q'$ ³. The reflexive and transitive closure of the binary relation \vdash_M is denoted as \vdash_M^* . A sequence

$$c_0 \vdash_M c_1 \vdash_M \dots \vdash_M c_n$$

is called a **computation** of $n \geq 0$ steps by M .

The **language accepted by** M is defined by

$$L(M) = \{w \mid \exists q \in F \cdot (q_0, w) \vdash_M^* (q, \epsilon)\}$$

We will often write $\delta^*(q, w)$ to denote the unique $q' \in Q$ such that $(q, w) \vdash_M^* (q', \epsilon)$.

³We will drop the subscript M in \vdash_M if M is clear from the context.

Example 1 Assuming an alphabet $\Sigma = \{l, d, o\}$ (where “*l*” stands for “letter”, “*d*” stands for “digit” and “*o*” stands for “other”), a DFA recognizing Micro *NAMEs* can be defined as follows:

$$M = (\{q_0, q_e, q_1\}, \{l, d, o\}, \delta, q_0, \{q_1\})$$

where δ is defined by

$$\begin{aligned} \delta(q_0, l) &= q_1 \\ \delta(q_0, d) &= q_e \\ \delta(q_0, o) &= q_e \\ \delta(q_1, l) &= q_1 \\ \delta(q_1, d) &= q_1 \\ \delta(q_1, o) &= q_e \\ \delta(q_e, l) &= q_e \\ \delta(q_e, d) &= q_e \\ \delta(q_e, o) &= q_e \end{aligned}$$

M is shown in Figure 2.3 (the initial state has a small incoming arrow, final states are in bold):

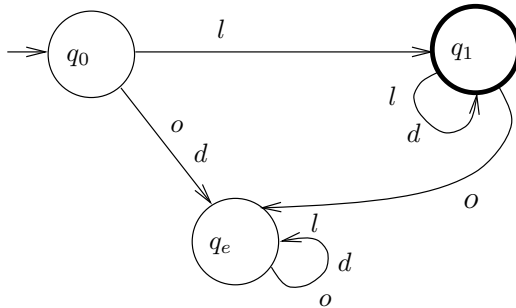


Figure 2.3: A DFA for *NAME*

Clearly, a DFA can be efficiently implemented, e.g. by encoding the states as numbers and using an array to represent the transition function. This is illustrated in Figure 2.4. The `next_state` array can be automatically generated from the DFA description.

What is not clear is how to translate regular expressions to DFA's. To show how this can be done, we need the more general concept of a nondeterministic finite automaton (NFA).

```

1  typedef int      STATE;
2  typedef char     SYMBOL;
3  typedef enum {false,true} BOOL;
4
5  STATE  next_state[SYMBOL][STATE];
6  BOOL   final[STATE];
7
8  BOOL
9  dfa(SYMBOL *input, STATE q)
10 {
11     SYMBOL  c;
12
13     while (c=*input++)
14         q = next_state[c,q];
15     return final[q];
16 }

```

Figure 2.4: DFA implementation

2.3.2 Nondeterministic finite automata

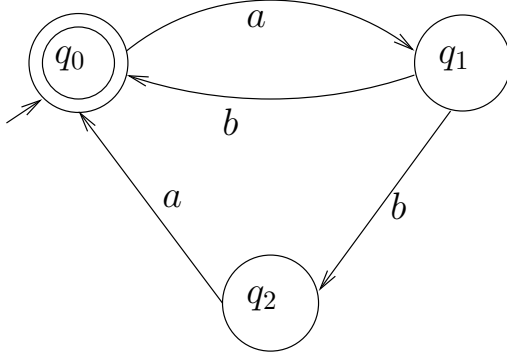
A *nondeterministic finite automaton* is much like a deterministic one except that we now allow several possibilities for a transition on the same symbol from a given state. The idea is that the automaton can arbitrarily (nondeterministically) choose one of the possibilities. In addition, we will also allow ϵ -**moves** where the automaton makes a state transition (labeled by ϵ) without reading an input symbol.

Definition 4 A *nondeterministic finite automaton (NFA)* is a tuple

$$(Q, \Sigma, \delta, q_0, F)$$

where

- Q is a finite set of *states*,
- Σ is a finite *input alphabet*
- $\delta : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$ is a (total) transition function ⁴
- $q_0 \in Q$ is the *initial state*
- $F \subseteq Q$ is the set of *final states*

Figure 2.5: M_1

It should be noted that $\emptyset \in 2^Q$ and thus Definition 4 sanctions the possibility of there not being any transition from a state q on a given symbol a .

Example 2 Consider $M_1 = (\{q_0, q_1, q_2\}, \delta_1, q_0, \{q_0\})$ as depicted in Figure 2.5. The table below defines δ_1 :

| $q \in Q$ | $\sigma \in \Sigma$ | $\delta_1(q, \sigma)$ |
|-----------|---------------------|-----------------------|
| q_0 | a | $\{q_1\}$ |
| q_0 | b | \emptyset |
| q_1 | a | \emptyset |
| q_1 | b | $\{q_0, q_2\}$ |
| q_2 | a | $\{q_0\}$ |
| q_2 | b | \emptyset |

The following definition formalizes our intuition about the behavior of nondeterministic finite automata.

Definition 5 Let $M = (Q, \Sigma, \delta, q_0, F)$ be a NFA. A **configuration** of M is a pair $(q, w) \in Q \times \Sigma^*$. For a configuration (q, aw) (where $a \in \Sigma \cup \{\epsilon\}$), we write

$$(q, aw) \vdash_M (q', w)$$

just when $q' \in \delta(q, a)$ ⁵. The reflexive and transitive closure of the binary relation \vdash_M is denoted as \vdash_M^* . The **language accepted by** M is defined by

$$L(M) = \{w \mid \exists q \in F \cdot (q_0, w) \vdash_M^* (q, \epsilon)\}$$

⁴For any set X , we use 2^X to denote its power set, i.e. the set of all subsets of X .

⁵We will drop the subscript M in \vdash_M if M is clear from the context.

Example 3 The following sequence shows how M_1 from Example 2 can accept the string $abaab$:

$$\begin{aligned}
 (q_0, abaab) &\vdash_{M_1} (q_1, baab) \\
 &\vdash_{M_1} (q_2, aab) \\
 &\vdash_{M_1} (q_0, ab) \\
 &\vdash_{M_1} (q_1, b) \\
 &\vdash_{M_1} (q_0, \epsilon)
 \end{aligned}$$

$$L(M_1) = \{w_0w_1 \dots w_n \mid n \in \mathbb{N} \wedge \forall 0 \leq i \leq n \cdot w_i \in \{ab, aba\}\}$$

Although nondeterministic finite automata are more general than deterministic ones, it turns out that they are not more powerful in the sense that any NFA can be simulated by a DFA.

Theorem 1 *Let M be a NFA. There exists a DFA M' such that $L(M') = L(M)$.*

Proof: (sketch) Let $M = (Q, \Sigma, \delta, q_0, F)$ be a NFA. We will construct a DFA M' that simulates M . This is achieved by letting M' be “in all possible states” that M could be in (after reading the same symbols). Note that “all possible states” is always an element of 2^Q , which is finite since Q is.

To deal with ϵ -moves, we note that, if M is in a state q , it could also be in any state q' to which there is an ϵ -transition from q . This motivates the definition of the ϵ -closure $\mathcal{C}_\epsilon(S)$ of a set of states S :

$$\mathcal{C}_\epsilon(S) = \{p \in Q \mid \exists q \in S \cdot (q, \epsilon) \vdash_M^* (p, \epsilon)\} \quad (2.1)$$

Now we define

$$M' = (2^Q, \Sigma, \delta', s_0, F')$$

where

- δ' is defined by

$$\forall s \in 2^Q, a \in \Sigma \cdot \delta'(s, a) = \cup_{q \in s} \mathcal{C}_\epsilon(\delta(q, a)) \quad (2.2)$$

- $s_0 = \mathcal{C}_\epsilon(q_0)$, i.e. M' starts in all possible states where M could go to from q_0 without reading any input.
- $F' = \{s \in 2^Q \mid s \cap F \neq \emptyset\}$, i.e. if M could end up in a final state, then M' will do so.

It can then be shown that $L(M') = L(M)$. □

2.4 Regular expressions vs finite state automata

In this section we show how a regular expression can be translated to a nondeterministic finite automata that defines the same language. Using Theorem 1, we can then translate regular expressions to DFA's and hence to a program that accepts exactly the strings conforming to the regular expression.

Theorem 2 *Let r be a regular expression. Then there exists a NFA M_r such that $L(M_r) = L_r$.*

Proof:

We show by induction on the number of operators used in a regular expression r that L_r is accepted by an NFA

$$M_r = (Q, \Sigma, \delta, q_0, \{q_f\})$$

(where Σ is the alphabet of L_r) which has exactly one final state q_f satisfying

$$\forall a \in \Sigma \cup \{\epsilon\} \cdot \delta(q_f, a) = \emptyset \quad (2.3)$$

Base case

Assume that r does not contain any operator. Then r is one of \emptyset , ϵ or $a \in \Sigma$.

We then define M_\emptyset , M_ϵ and M_a as shown in Figure 2.6.

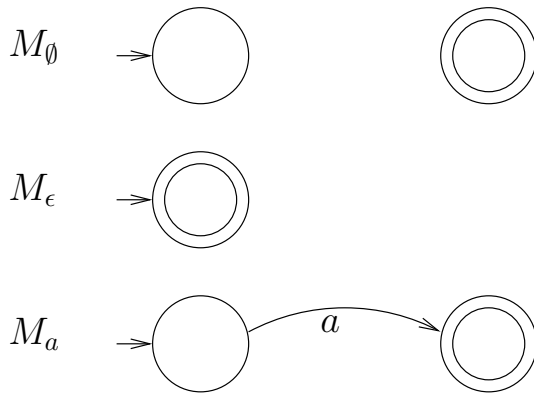
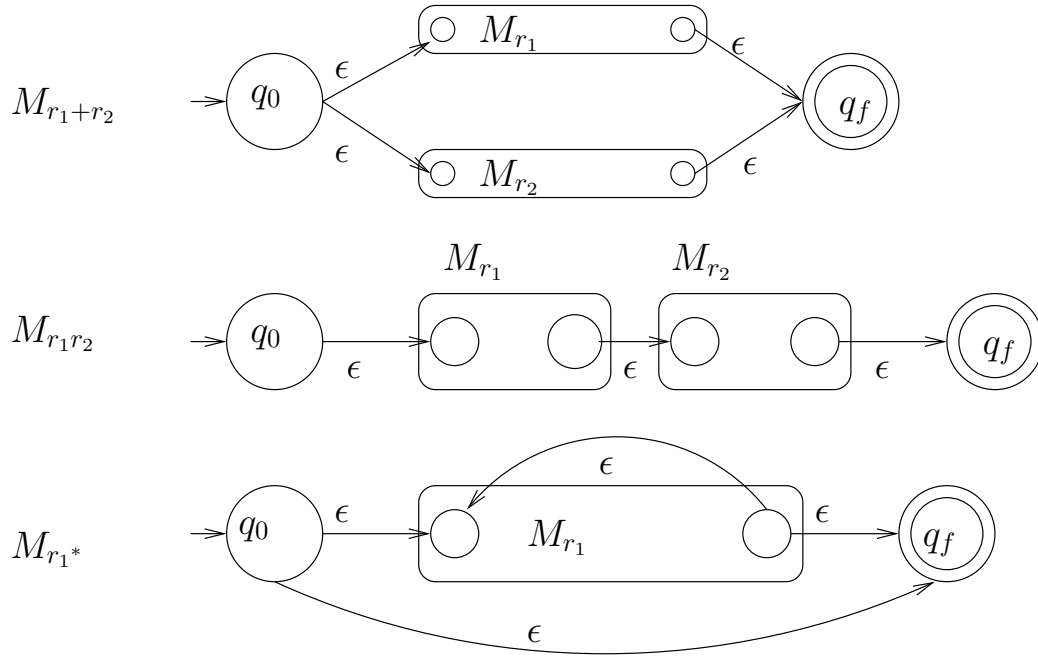


Figure 2.6: M_\emptyset , M_ϵ and M_a

Induction step

More complex regular expressions must be of one of the forms $r_1 + r_2$, $r_1 r_2$ or r_1^* . In each case, we can construct a NFA $M_{r_1 + r_2}$, $M_{r_1 r_2}$ or $M_{r_1^*}$, based on M_{r_1} and M_{r_2} , as shown in Figure 2.7.

Figure 2.7: $M_{r_1+r_2}$, $M_{r_1r_2}$ and $M_{r_1}^*$

It can then be shown that

$$\begin{aligned}
 L(M_{r_1+r_2}) &= L(M_{r_1}) \cup L(M_{r_2}) \\
 L(M_{r_1r_2}) &= L(M_{r_1}) L(M_{r_2}) \\
 L(M_{r_1}^*) &= L(M_{r_1})^*
 \end{aligned}$$

□

2.5 A scanner generator

We can now be more specific on the design and operation of a scanner generator such as *lex(1)* or *flex(1L)*, which was sketched on page 25.

First we introduce the concept of a “dead” state in a DFA.

Definition 6 Let $M = (Q, \Sigma, \delta, q_0, F)$ be a DFA. A state $q \in Q$ is called **dead** if there does not exist a string $w \in \Sigma^*$ such that $(q, w) \vdash_M^* (q_f, \epsilon)$ for some $q_f \in F$.

Example 4 The state q_e in Example 1 is dead.

It is easy to determine the set of dead states for a DFA, e.g. using a marking algorithm which initially marks all states as “dead” and then recursively works backwards from the final states, unmarking any states reached.

The generator takes as input a set of regular expressions, $R = \{r_1, \dots, r_n\}$ each of which is associated with some code c_{r_i} to be executed when a token corresponding to r_i is recognized.

The generator will convert the regular expression

$$r_1 + r_2 + \dots + r_n$$

to a DFA $M = (Q, \Sigma, \delta, q_0, F)$, as shown in Section 2.4, with one addition: when constructing M , it will remember which final state of the DFA corresponds with which regular expression. This can easily be done by remembering the final states in the NFA’s corresponding to each of the r_i while constructing the combined DFA M . It may be that a final state in the DFA corresponds to several patterns (regular expressions). In this case, we select the one that was defined first.

Thus we have a mapping

$$pattern : F \rightarrow R$$

which associates the first (in the order of definition) pattern to which a certain final state corresponds. We also compute the set of dead states of M .

The code in Figure 2.8 illustrates the operation of the generated scanner.

The scanner reads input characters, remembering the last final state seen and the associated regular expression, until it hits a dead state from where it is impossible to reach a final state. It then backs up to the last final state and executes the code associated with that pattern. Clearly, this will find the longest possible token on the input.

```

1  typedef int      STATE;
2  typedef char     SYMBOL;
3  typedef enum {false,true} BOOL;
4
5  typedef struct { /* what we need to know about a user defined pattern */
6      TOKEN*      (*code)(); /* user-defined action */
7      BOOL        do_return; /* whether action returns from lex() or not */
8      }           PATTERN;
9
10 static STATE     next_state[SYMBOL][STATE];
11 static BOOL      dead[STATE];
12 static BOOL      final[STATE];
13 static PATTERN*  pattern[STATE]; /* first regexp for this final state */
14 static SYMBOL    *last_input = 0; /* input pointer at last final state */
15 static STATE     last_state, q = 0; /* assuming 0 is initial state */
16 static SYMBOL    *input; /* source text */
17
18 TOKEN*
19 lex()
20 {
21     SYMBOL    c;
22     PATTERN *last_pattern    = 0;
23
24     while (c=*input++) {
25         q = next_state[c,q];
26         if (final[q]) {
27             last_pattern = pattern[q];
28             last_input    = input;
29             last_state    = q;
30         }
31         if (dead[q]) {
32             if (last_pattern) {
33                 input    = last_input;
34                 q        = 0;
35                 if (last_pattern->do_return)
36                     return pattern->code();
37                 else
38                     pattern->code();
39             }
40             else /* error */
41                 ;
42         }
43     }
44
45     return (TOKEN*) 0;
46 }

```

Figure 2.8: A generated scanner

Chapter 3

Parsing

3.1 Context-free grammars

As mentioned in Section 1.3.1, page 9, the rules (or railroad diagrams) used to specify the syntax of a programming language can be formalized using the concept of context-free grammar.

Definition 7 A *context-free grammar* (cfg) is a tuple

$$G = (V, \Sigma, P, S)$$

where

- V is a finite set of **nonterminal** symbols
- Σ is a finite set of **terminal** symbols, disjoint from V : $\Sigma \cap V = \emptyset$.
- P is a finite set of **productions** of the form $A \rightarrow \alpha$ where $A \in V$ and $\alpha \in (V \cup \Sigma)^*$
- $S \in V$ is a nonterminal **start symbol**

Note that terminal symbols correspond to token types as delivered by the lexical analyzer.

Example 5 The following context-free grammar defines the syntax of simple arithmetic expressions:

$$G_0 = (\{E\}, \{+, \times, (,), \text{id}\}, P, E)$$

where P consists of

$$\begin{aligned} E &\rightarrow E + E \\ E &\rightarrow E \times E \\ E &\rightarrow (E) \\ E &\rightarrow \text{id} \end{aligned}$$

We shall often use a shorter notation for a set of productions where several right-hand sides for the same nonterminal are written together, separated by “|”. Using this notation, the set of rules of G_0 can be written as

$$E \rightarrow E + E \mid E \times E \mid (E) \mid \text{id}$$

Definition 8 Let $G = (V, \Sigma, P, S)$ be a context-free grammar. For strings $x, y \in (V \cup \Sigma)^*$, we say that x **derives** y in one step, denoted $x \Rightarrow_G y$ iff $x = x_1 A x_2$, $y = x_1 \alpha x_2$ and $A \rightarrow \alpha \in P$. Thus \Rightarrow_G is a binary relation on $(V \cup \Sigma)^*$. The relation \Rightarrow_G^* is the reflexive and transitive closure of \Rightarrow_G . The **language** $L(G)$ **generated by** G is defined by

$$L(G) = \{w \in \Sigma^* \mid S \Rightarrow_G^* w\}$$

A language is called **context-free** if it is generated by some context-free grammar. A **derivation** in G of w_n from w_0 is any sequence of the form

$$w_0 \Rightarrow_G w_1 \Rightarrow_G \dots \Rightarrow_G w_n$$

where $n \geq 0$ (we say that the derivation has n **steps**) and $\forall 1 \leq i \leq n \cdot w_i \in (V \cup \Sigma)^*$. We write $v \Rightarrow_G^n w$ ($n \geq 0$) when w can be derived from v in n steps.

Thus a context-free grammar specifies precisely which sequences of tokens are valid sentences (programs) in the language.

Example 6 Consider the grammar G_0 from Example 5. The following is a derivation in G where at each step, the symbol to be rewritten is underlined.

$$\begin{aligned} \underline{S} &\Rightarrow_{G_0} \underline{E} \times E \\ &\Rightarrow_{G_0} (\underline{E}) \times E \\ &\Rightarrow_{G_0} (E + \underline{E}) \times E \\ &\Rightarrow_{G_0} (E + \text{id}) \times \underline{E} \\ &\Rightarrow_{G_0} (\underline{E} + \text{id}) \times \text{id} \\ &\Rightarrow_{G_0} (\text{id} + \text{id}) \times \text{id} \end{aligned}$$

A derivation in a context-free grammar is conveniently represented by a *parse tree*.

Definition 9 Let $G = (V, \Sigma, P, S)$ be a context-free grammar. A **parse tree** corresponding to G is a labeled tree where each node is labeled by a symbol from $V \cup \Sigma$ in such a way that, if A is the label of a node and $A_1 A_2 \dots A_n$ ($n > 0$) are the labels of its children (in left-to-right order), then

$$A \rightarrow A_1 A_2 \dots A_n$$

is a rule in P . Note that a rule $A \rightarrow \epsilon$ gives rise to a leaf node labeled ϵ .

As mentioned in Section 1.3.4, it is the job of the parser to convert a string of tokens into a parse tree that has precisely this string as yield. The idea is that the parse tree describes the syntactical structure of the source text.

However, sometimes, there are several parse trees possible for a single string of tokens, as can be seen in Figure 3.1.

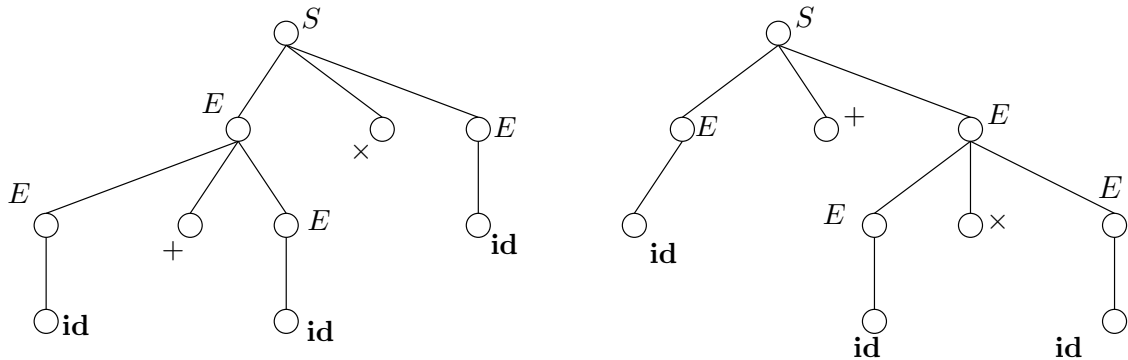


Figure 3.1: Parse trees in the ambiguous context-free grammar from Example 5

Note that the two parse trees intuitively correspond to two evaluation strategies for the expression. Clearly, we do not want a source language that is specified using an ambiguous grammar (that is, a grammar where a legal string of tokens may have different parse trees).

Example 7 Fortunately, we can fix the grammar from Example 5 to avoid such ambiguities.

$$G_1 = (\{E, T, F\}, \{+, \times, (,), \text{id}\}, P', E)$$

where P' consists of

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T \times F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

is an unambiguous grammar generating the same language as the grammar from Example 5.

Still, there are context-free languages such as $\{a^i b^j c^k \mid i = j \vee j = k\}$ for which only ambiguous grammars can be given. Such languages are called *inherently ambiguous*. Worse still, checking whether an arbitrary context-free grammar allows ambiguity is an unsolvable problem[HU69].

3.2 Top-down parsing

3.2.1 Introduction

When using a top-down (also called *predictive*) parsing method, the parser tries to find a *leftmost derivation* (and associated parse tree) of the source text. A leftmost derivation is a derivation where, during each step, the leftmost nonterminal symbol is rewritten.

Definition 10 Let $G = (V, \Sigma, P, S)$ be a context-free grammar. For strings $x, y \in (V \cup \Sigma)^*$, we say that x derives y in a leftmost fashion and in one step, denoted

$$x \xRightarrow{L}_G y$$

iff $x = x_1 A x_2$, $y = x_1 \alpha x_2$, $A \rightarrow \alpha$ is a rule in P and $x_1 \in \Sigma^*$ (i.e. the leftmost occurrence of a nonterminal symbol is rewritten).

The relation \xRightarrow{L}_G^* is the reflexive and transitive closure of \xRightarrow{L}_G . A derivation

$$y_0 \xRightarrow{L}_G y_1 \xRightarrow{L}_G \dots \xRightarrow{L}_G y_n$$

is called a **leftmost** derivation. If $y_0 = S$ (the start symbol) then we call each y_i in such a derivation a **left sentential form**.

It is not hard to see that restricting to leftmost derivations does not alter the language of a context-free grammar.

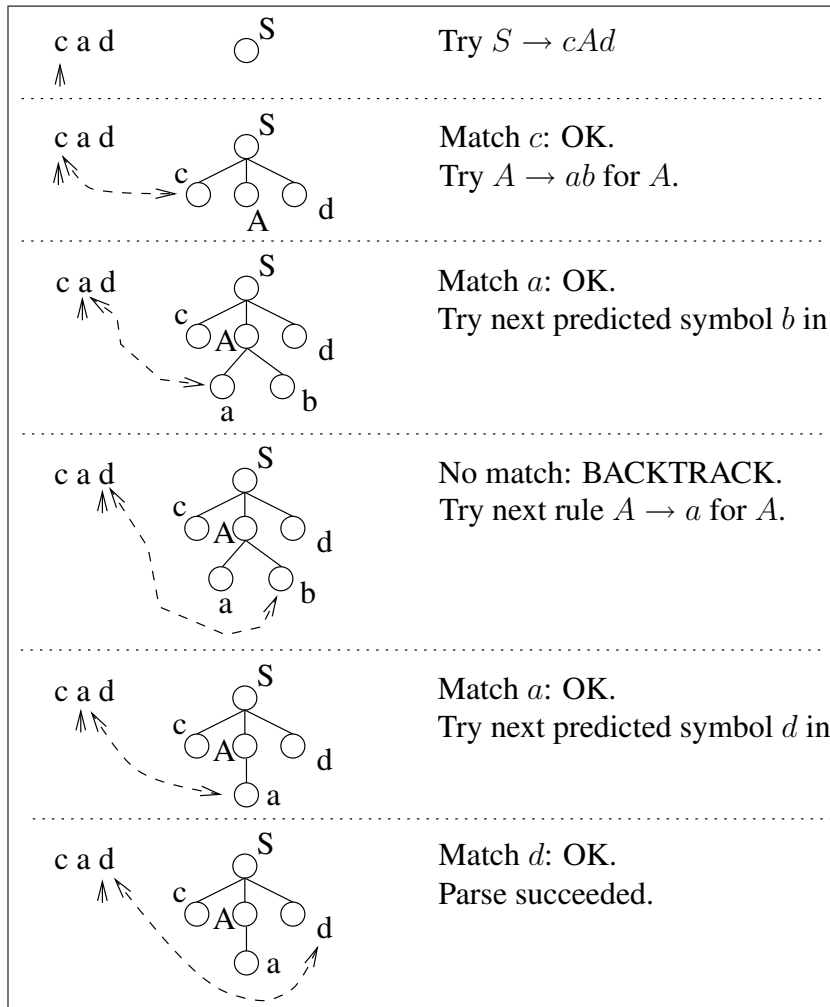


Figure 3.2: A simple top-down parse

Theorem 3 Let $G = (V, \Sigma, P, S)$ be a context-free grammar. If $A \in V \cup \Sigma$ then

$$A \Rightarrow_G^* w \in \Sigma^* \text{ iff } A \xRightarrow{L}_G^* w \in \Sigma^*$$

Example 8 Consider the trivial grammar

$$G = (\{S, A\}, \{a, b, c, d\}, P, S)$$

where P contains the rules

$$\begin{aligned} S &\rightarrow cAd \\ A &\rightarrow ab \mid a \end{aligned}$$

Let $w = cad$ be the source text. Figure 3.2 shows how a top-down parse could proceed.

The reasoning in Example 8 can be encoded as shown below.

```

1  typedef enum {false,true} BOOL;
2
3  TOKEN* input; /* output array from scanner */
4  TOKEN* token; /* current token from input */
5
6  BOOL
7  parse_S() { /* Parse something derived from S */
8      /* Try rule S --> c A d */
9      if (*token=='c') {
10         ++token;
11         if (parse_A()) {
12             if (*token=='d') {
13                 ++token;
14                 return true;
15             }
16         }
17     }
18     return false;
19 }
20
21 BOOL
22 parse_A() { /* Parse stuff derived from A */
23     TOKEN* save; /* for backtracking */
24
25     save = token;
26
27     /* Try rule A --> a b */
28     if (*token=='a') {
29         ++token;
30         if (*token=='b') {
31             ++token;
32             return true;
33         }
34     }
35
36     token = save; /* didn't work: backtrack */
37
38     /* Try rule A --> a */
39     if (*token=='a') {
40         ++token;
41         return true;
42     }
43
44     token = save; /* didn't work: backtrack */

```

```

45
46  /* no more rules: give up */
47  return false;
48  }

```

Note that the above strategy may need recursive functions. E.g. if the grammar contains a rule such as

$$E \rightarrow (E)$$

the code for *parse*_E() will contain a call to *parse*_E().

The method illustrated above has two important drawbacks:

- It cannot be applied if the grammar G is *left-recursive*, i.e. $A \xRightarrow{L}_G^* Ax$ for some $x \in (V \cup \Sigma)^*$.

Indeed, even for a “natural” rule such as

$$E \rightarrow E + E$$

it is clear that the *parse*_E() function would start by calling itself (in an infinite recursion), before reading any input.

We will see that it is possible to eliminate left recursion from a grammar without changing the generated language.

- Using backtracking is both expensive and difficult. It is difficult because it usually does not suffice to simply restore the input pointer: all actions taken by the compiler since the backtrack point (e.g. symbol table updates) must also be undone.

The solution will be to apply top-down parsing only for a certain class of restricted grammars for which it can be shown that backtracking will never be necessary.

3.2.2 Eliminating left recursion in a grammar

First we look at the elimination of *immediate* left recursion where we have rules of the form

$$A \rightarrow A\alpha \mid \beta$$

where β does not start with A .

The idea is to reorganize the rules in such a way that derivations are simulated as shown in Figure 3.3

This leads to the following algorithm to remove immediate left recursion.

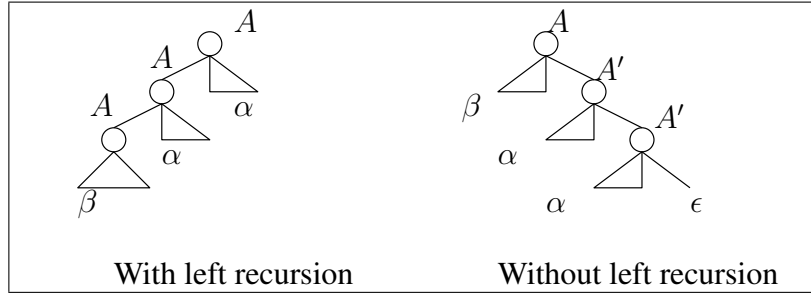


Figure 3.3: Eliminating immediate left recursion

Algorithm 1 [Removing immediate left recursion for a nonterminal A from a grammar $G = (V, \Sigma, P, S)$]

Let

$$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \dots \mid \beta_n$$

be all the rules with A as the left hand side. Note that $m > 0$ and $n \geq 0$ and $\forall 0 \leq i \leq n \cdot \beta_i \notin A(V \cup \Sigma)^*$

1. If $n = 0$, no terminal string can ever be derived from A , so we may as well remove all A -rules.
2. Otherwise, define a new nonterminal A' , and replace the A rules by

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \dots \mid \beta_n A' \\ A' &\rightarrow \alpha_1 A' \mid \dots \mid \alpha_m A' \mid \epsilon \end{aligned}$$

□

Example 9 Consider the grammar G_1 from Example 7. Applying Algorithm 1 results in the grammar $G_2 = (\{E, E', T, T', F\}, \{+, \times, (,), \text{id}\}, P_{G_2}, E)$ where P_{G_2} contains

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow \times FT' \mid \epsilon \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

It is also possible (see e.g. [ASU86], page 177) to eliminate general (also indirect) left recursion from a grammar.

3.2.3 Avoiding backtracking: LL(1) grammars

When inspecting Example 8, it appears that in order to avoid backtracking, one should be able to predict with certainty, given the nonterminal X of the current function $parse_X()$ and the next input token a , which production $X \rightarrow x$ was used as a first step to derive a string of tokens that starts with the input token a from X .

The problem is graphically represented in Figure 3.4

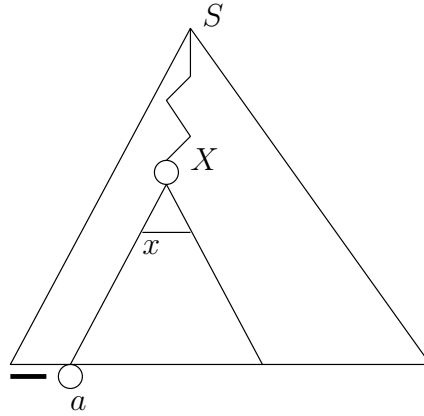


Figure 3.4: The prediction problem

The following definition defines exactly the class of context-free grammars for which this is possible.

Definition 11 An *LL(1) grammar* is a context-free grammar $G = (V, \Sigma, P, S)$ such that if

$$S \xRightarrow{*}_G^L w_1 X x_3 \xRightarrow{L}_G w_1 x_2 x_3 \xRightarrow{*}_G^L w_1 a w_4$$

and

$$S \xRightarrow{*}_G^L w_1 X \hat{x}_3 \xRightarrow{L}_G w_1 \hat{x}_2 \hat{x}_3 \xRightarrow{*}_G^L w_1 a \hat{w}_4$$

where $a \in \Sigma$, $X \in V$, $w_1 \in \Sigma^*$, $w_4 \in \Sigma^*$, $\hat{w}_4 \in \Sigma^*$ then

$$x_2 = \hat{x}_2$$

A language L is called an *LL(1) language* if there exists a *LL(1) grammar* G such that $L(G) = L$.

Intuitively, Definition 11 just says that if there are two possible choices for a production, these choices are identical.

Thus for LL(1) grammars¹, we know for sure that we can write functions as we did in Example 8, without having to backtrack. In the next section, we will see how to automatically generate parsers based on the same principle but without the overhead of (possibly recursive) function calls.

3.2.4 Predictive parsers

Predictive parsers use a stack (representing strings in $(V \cup \Sigma)^*$) and a parse table as shown in Figure 3.5.

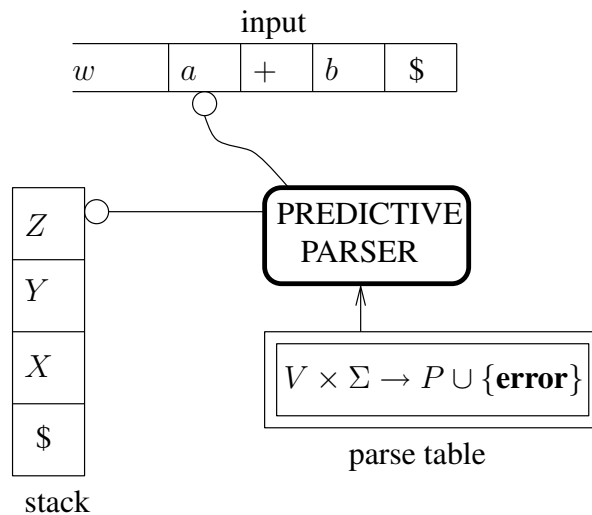


Figure 3.5: A predictive parser

The figure shows the parser simulating a leftmost derivation

$$S \xRightarrow{L}_G \dots \xRightarrow{L}_G \underbrace{wZYX}_{\text{depicted position of parser}} \xRightarrow{L}_G \dots \xRightarrow{L}_G wa + b$$

Note that the string consisting of the input read so far, followed by the contents of the stack (from the top down) constitutes a left sentential form. End markers (depicted as $\$$ symbols) are used to mark both the bottom of the stack and the end of the input.

¹It is possible to define LL(k) grammars and languages by replacing a in Definition 11 by a terminal string of (up to) k symbols.

The parse table M represents the production to be chosen: when the parser has X on the top of the stack and a as the next input symbol, $M[X, a]$ determines which production was used as a first step to derive a string starting with a from X .

Algorithm 2 [Operation of an LL(1) parser]

The operation of the parser is shown in Figure 3.6. □

Intuitively, a predictive parser simulates a leftmost derivation of its input, using the stack to store the part of the left sentential form that has not yet been processed (this part includes all nonterminals of the left sentential form). It is not difficult to see that if a predictive parser successfully processes an input string, then this string is indeed in the language generated by the grammar consisting of all the rules in the parse table.

To see the reverse, we need to know just how the parse table is constructed.

First we need some auxiliary concepts:

Definition 12 Let $G = (V, \Sigma, P, S)$ be a context-free grammar.

- The function **first** : $(V \cup \Sigma)^* \rightarrow 2^{(\Sigma \cup \{\epsilon\})}$ is defined by

$$\text{first}(\alpha) = \{a \mid \alpha \Rightarrow_G^* aw \in \Sigma^*\} \cup X_\alpha$$

where

$$X_\alpha = \begin{cases} \{\epsilon\} & \text{if } \alpha \Rightarrow_G^* \epsilon \\ \emptyset & \text{otherwise} \end{cases}$$

- The function **follow** : $V \rightarrow 2^{(\Sigma \cup \{\$ \})}$ is defined by

$$\text{follow}(A) = \{a \in \Sigma \mid S \Rightarrow_G^* \alpha A a \beta\} \cup Y_\alpha$$

where

$$Y_\alpha = \begin{cases} \{\$ \} & \text{if } S \Rightarrow_G^* \alpha A \\ \emptyset & \text{otherwise} \end{cases}$$

Intuitively, $\text{first}(\alpha)$ contains the set of terminal symbols that can appear at the start of a (terminal) string derived from α , including ϵ if α can derive the empty string.

On the other hand, $\text{follow}(A)$ consists of those terminal symbols that may follow a string derived from A in a terminal string from $L(G)$.

The construction of the LL(1) parse table can then be done using the following algorithm.

```

1  PRODUCTION* parse_table[NONTERMINAL,TOKEN];
2  SYMBOL End, S; /* marker and start symbol */
3  SYMBOL* stack;
4  SYMBOL* top_of_stack;
5
6  TOKEN   *input;
7
8  BOOL
9  parse() { /* LL(1) predictive parser */
10     SYMBOL X;
11     push(End); push(S);
12
13     while (*top_of_stack!=End) {
14         X = *top_of_stack;
15         if (is_terminal(X)) {
16             if (X==*input) { /* predicted outcome */
17                 ++input; /* advance input */
18                 pop(); /* pop X from stack */
19             }
20             else
21                 error("Expected %s, got %s", X, *input);
22         }
23         else { /* X is nonterminal */
24             PRODUCTION* p = parse_table[X,*input];
25             if (p) {
26                 pop(); /* pop X from stack */
27                 for (i=length_rhs(p)-1; (i>=0);--i)
28                     push(p->rhs[i]); /* push symbols of rhs, last first */
29             }
30             else
31                 error("Unexpected %s", *input);
32         }
33     }
34     if (*input==End)
35         return true;
36     else
37         error("Redundant input: %s", *input);
38 }

```

Figure 3.6: Predictive parser operation

Algorithm 3 [Construction of LL(1) parse table]

Let $G = (V, \Sigma, P, S)$ be a context-free grammar.

1. Initialize the table: $\forall A \in V, b \in \Sigma \cdot M[A, b] = \emptyset$
2. For each production $A \rightarrow \alpha$ from P

- (a) For each $b \in \text{first}(\alpha) \cap \Sigma$, add $A \rightarrow \alpha$ to $M[A, b]$.
- (b) If $\epsilon \in \text{first}(\alpha)$ then
 - i. For each $c \in \text{follow}(A) \cap \Sigma$, add $A \rightarrow \alpha$ to $M[A, c]$.
 - ii. If $\$ \in \text{follow}(A)$ then add $A \rightarrow \alpha$ to $M[A, \$]$.
- 3. If each entry in M contains at most a single production then return success, else return failure.

□

It can be shown that, if the parse table for a grammar G was successfully constructed using Algorithm 3, then the parser of Algorithm 2 accepts exactly the strings of $L(G)$. Also, if Algorithm 3 fails, then G is not an LL(1) grammar.

Example 10 Consider grammar G_2 from Example 9. The *first* and *follow* functions are computed in Examples 11 and 12. The results are summarized below.

| | E | E' | T | T' | F |
|--------|-------|---------------|----------|-----------------------|---------------------|
| first | (, id | +, ϵ | (, id | \times , ϵ | (, id |
| follow | \$,) | \$,) | +, \$,) | \$,), + | \times , \$,), + |

Applying Algorithm 3 yields the following LL(1) parse table.

| | E | E' | T | T' | F |
|----------|---------------------|---------------------------|---------------------|-----------------------------|---------------------------|
| id | $E \rightarrow TE'$ | | $T \rightarrow FT'$ | | $F \rightarrow \text{id}$ |
| + | | $E' \rightarrow +TE'$ | | $T' \rightarrow \epsilon$ | |
| \times | | | | $T' \rightarrow \times FT'$ | |
| (| $E \rightarrow TE'$ | | $T \rightarrow FT'$ | | $F \rightarrow (E)$ |
|) | | $E' \rightarrow \epsilon$ | | $T' \rightarrow \epsilon$ | |
| \$ | | $E' \rightarrow \epsilon$ | | $T' \rightarrow \epsilon$ | |

The operation of a predictive parser using the above table is illustrated below for

the input string $\text{id} + \text{id} \times \text{id}$.

| stack | input | rule |
|-------------------|--|-----------------------------|
| $\$E$ | $\text{id} + \text{id} \times \text{id}\$$ | |
| $\$E'T$ | $\text{id} + \text{id} \times \text{id}\$$ | $E \rightarrow TE'$ |
| $\$E'T'F$ | $\text{id} + \text{id} \times \text{id}\$$ | $T \rightarrow FT'$ |
| $\$E'T'\text{id}$ | $\text{id} + \text{id} \times \text{id}\$$ | $F \rightarrow \text{id}$ |
| $\$E'T'$ | $+\text{id} \times \text{id}\$$ | |
| $\$E'$ | $+\text{id} \times \text{id}\$$ | $T' \rightarrow \epsilon$ |
| $\$E'T+$ | $+\text{id} \times \text{id}\$$ | $E' \rightarrow +TE'$ |
| $\$E'T$ | $\text{id} \times \text{id}\$$ | |
| $\$E'T'F$ | $\text{id} \times \text{id}\$$ | $T \rightarrow FT'$ |
| $\$E'T'\text{id}$ | $\text{id} \times \text{id}\$$ | $F \rightarrow \text{id}$ |
| $\$E'T'$ | $\times \text{id}\$$ | |
| $\$E'T'F \times$ | $\times \text{id}\$$ | $T' \rightarrow \times FT'$ |
| $\$E'T'F$ | $\text{id}\$$ | |
| $\$E'T'\text{id}$ | $\text{id}\$$ | $F \rightarrow \text{id}$ |
| $\$E'T'$ | $\$$ | |
| $\$E'$ | $\$$ | $T' \rightarrow \epsilon$ |
| $\$$ | $\$$ | $E' \rightarrow \epsilon$ |

3.2.5 Construction of first and follow

Algorithm 4 [Construction of first]

Let $G = (V, \Sigma, P, S)$ be a context-free grammar. We construct an array $F : V \cup \Sigma \rightarrow 2^{(\Sigma \cup \{\epsilon\})}$ and then show a function to compute $\text{first}(\alpha)$ for arbitrary $\alpha \in (V \cup \Sigma)^*$.

1. F is constructed via a fixpoint computation:

- (a) for each $X \in V$, initialize $F[X] \leftarrow \{a \mid X \rightarrow a\alpha, a \in \Sigma\}$
- (b) for each $a \in \Sigma$, initialize $F[a] \leftarrow \{a\}$

2. repeat the following steps

- (a) $F' \leftarrow F$
- (b) for each rule $X \rightarrow \epsilon$, add ϵ to $F[X]$.
- (c) for each rule $X \rightarrow Y_1 \dots Y_k$ ($k > 0$) do
 - i. for each $1 \leq i \leq k$ such that $Y_1 \dots Y_{i-1} \in V^*$ and $\forall 1 \leq j < i \cdot \epsilon \in F[Y_j]$ do $F[X] \leftarrow F[X] \cup (F[Y_i] \setminus \{\epsilon\})$
 - ii. if $\forall 1 \leq j \leq k \cdot \epsilon \in F[Y_j]$ then $F[X] \leftarrow F[X] \cup \{\epsilon\}$

until $F' = F$

Define

$$\text{first}(X_1 \dots X_n) = \begin{cases} \bigcup_{1 \leq j \leq (k+1)} (F[X_j] \setminus \{\epsilon\}) & \text{if } k(X_1 \dots X_n) < n \\ \bigcup_{1 \leq j \leq n} F[X_j] \cup \{\epsilon\} & \text{if } k(X_1 \dots X_n) = n \end{cases}$$

where $k(X_1 \dots X_n)$ is the largest index k such that $X_1 \dots X_k \Rightarrow_G^* \epsilon$, i.e.

$$k(X_1 \dots X_n) = \max\{1 \leq i \leq n \mid \forall 1 \leq j \leq i \cdot \epsilon \in F[X_j]\}$$

□

Example 11 Consider grammar G_2 from Example 9.

The construction of F is illustrated in the following table.

| E | E' | T | T' | F | |
|-------|---------------|-------|--------------------|-------|--|
| | + | | \times | (, id | initialization |
| | ϵ | (, id | ϵ | | rule $E' \rightarrow \epsilon$ rule $T \rightarrow FT'$ rule $T' \rightarrow \epsilon$ |
| (, id | | | | | rule $E \rightarrow TE'$ |
| (, id | +, ϵ | (, id | \times, ϵ | (, id | |

Algorithm 5 [Construction of follow]

Let $G = (V, \Sigma, P, S)$ be a context-free grammar. The follow set of all symbols in V can be computed as follows:

1. $\text{follow}(S) \leftarrow \$$
2. Apply the following rules until nothing can be added to $\text{follow}(A)$ for any $A \in V$.

(a) If there is a production $A \rightarrow \alpha B \beta$ in P then

$$\text{follow}(B) \leftarrow \text{follow}(B) \cup (\text{first}(\beta) \setminus \{\epsilon\})$$

(b) If there is a production $A \rightarrow \alpha B$ or a production $A \rightarrow \alpha B \beta$ where $\epsilon \in \text{first}(\beta)$ then $\text{follow}(B) \leftarrow \text{follow}(B) \cup \text{follow}(A)$.

□

Example 12 Consider again grammar G_2 from Example 9.

The construction of *follow* is illustrated in the following table.

| E | E' | T | T' | F | |
|--------|--------|-----------|-----------|--------------------|--|
| \$ | | | | | initialization |
|) | | + | | | rule $E \rightarrow TE'$ rule $F \rightarrow (E)$ rule $T \rightarrow FT'$ |
| | \$(,) | \$(,) | | | rule $E \rightarrow TE'$ rule $E \rightarrow TE'$ rule $T \rightarrow FT'$ rule $T \rightarrow FT'$ |
| | | | \$(,), + | | |
| | | | | \$(,), + | |
| \$(,) | \$(,) | +, \$(,) | \$(,), + | $\times, $(,), +$ | |

3.3 Bottom-up parsing

3.3.1 Shift-reduce parsers

As described in Section 3.2, a top-down parser simulates a leftmost derivation in a top-down fashion. This simulation predicts the next production to be used, based on the nonterminal to be rewritten and the next input symbol.

A bottom-up parser simulates a rightmost derivation in a bottom-up fashion, where a rightmost derivation is defined as follows.

Definition 13 Let $G = (V, \Sigma, P, S)$ be a context-free grammar. For strings $x, y \in (V \cup \Sigma)^*$, we say that x derives y in a rightmost fashion and in one step, denoted

$$x \xRightarrow{R}_G y$$

iff $x = x_1 A x_2$, $y = x_1 \alpha x_2$, $A \rightarrow \alpha$ is a rule in P and $x_2 \in \Sigma^*$ (i.e. the rightmost occurrence of a nonterminal symbol is rewritten).

The relation \xRightarrow{R}_G^* is the reflexive and transitive closure of \xRightarrow{R}_G . A derivation

$$x_0 \xRightarrow{R}_G x_1 \xRightarrow{R}_G \dots \xRightarrow{R}_G x_n$$

is called a **rightmost** derivation. If $x_0 = S$ (the start symbol) then we call each x_i in such a derivation a **right sentential form**.

A **phrase** of a right sentential form is a sequence of symbols that have been derived from a single nonterminal symbol occurrence (in an earlier right sentential form of the derivation). A **simple phrase** is a phrase that contains no smaller phrase. The **handle** of a right sentential form is its leftmost simple phrase (which is unique). A prefix of a right sentential form that does not extend past its handle is called a **viable prefix** of G .

Example 13 Consider the grammar

$$G_3 = (\{S, E, T, F\}, \{+, \times, (,), \text{id}\}, P_{G_3}, S)$$

where P_{G_3} consists of

$$\begin{aligned} S &\rightarrow E \\ E &\rightarrow E + T \mid T \\ T &\rightarrow T \times F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

G_3 is simply G_1 from Example 7, augmented by a new start symbol whose only use is at the left hand side of a single production that has the old start symbol as its right hand side.

The following table shows a rightmost derivation of $\text{id} + \text{id} \times \text{id}$; where the handle in each right sentential form is underlined.

$$\begin{aligned} &S \\ &\underline{E} \\ &\underline{E + T} \\ &\underline{E + T \times F} \\ &\underline{E + T \times \text{id}} \\ &\underline{E + \underline{F} \times \text{id}} \\ &\underline{E + \text{id} \times \text{id}} \\ &\underline{T + \text{id} \times \text{id}} \\ &\underline{F + \text{id} \times \text{id}} \\ &\underline{\text{id} + \text{id} \times \text{id}} \end{aligned}$$

When doing a bottom-up parse, the parser will use a stack to store the right sentential form up to the handle, using **shift** operations to push symbols from the input onto the stack. When the handle is completely on the stack, it will be popped and replaced by the left hand side of the production, using a **reduce** operation.

Example 14 The table below illustrates the operation of a shift-reduce parser to simulate the rightmost derivation from Example 13. Note that the parsing has to be read from the bottom to the top of the figure and that both the stack and the input contain the end marker “\$”.

| derivation | stack | input | operation |
|--|----------------------------|--|-----------------------------------|
| S | $\$S$ | $\$$ | accept |
| \underline{E} | $\$E$ | $\$$ | reduce $S \rightarrow E$ |
| $\underline{E} + \underline{T}$ | $\$E + T$ | $\$$ | reduce $E \rightarrow E + T$ |
| $E + \underline{T \times F}$ | $\$E + T \times F$ | $\$$ | reduce $T \rightarrow T \times F$ |
| | $\$E + T \times \text{id}$ | $\$$ | reduce $F \rightarrow \text{id}$ |
| | $\$E + T \times$ | $\text{id}\$$ | shift |
| $E + T \times \underline{\text{id}}$ | $\$E + T$ | $\times \text{id}\$$ | shift |
| $E + \underline{F} \times \text{id}$ | $\$E + F$ | $\times \text{id}\$$ | reduce $T \rightarrow F$ |
| | $\$E + \text{id}$ | $\times \text{id}\$$ | reduce $F \rightarrow \text{id}$ |
| | $\$E +$ | $\text{id} \times \text{id}\$$ | shift |
| $E + \underline{\text{id}} \times \text{id}$ | $\$E$ | $+ \text{id} \times \text{id}\$$ | shift |
| $\underline{T} + \text{id} \times \text{id}$ | $\$T$ | $+ \text{id} \times \text{id}\$$ | reduce $E \rightarrow T$ |
| $\underline{F} + \text{id} \times \text{id}$ | $\$F$ | $+ \text{id} \times \text{id}\$$ | reduce $T \rightarrow F$ |
| | $\$\text{id}$ | $+ \text{id} \times \text{id}\$$ | reduce $F \rightarrow \text{id}$ |
| $\underline{\text{id}} + \text{id} \times \text{id}$ | $\$$ | $\text{id} + \text{id} \times \text{id}\$$ | shift |

Clearly, the main problem when designing a shift-reduce parser is to decide when to shift and when to reduce. Specifically, the parser should reduce only when the handle is on the top of the stack. Put in another way, the stack should always contain a **viable prefix**.

We will see in Section 3.3.2 that the language of all viable prefixes of a context-free grammar G is regular, and thus there exists a DFA $M_G = (Q, V \cup \Sigma, \delta, q_0, F)$ accepting exactly the viable prefixes of G . The operation of a shift-reduce parser is outlined in Figure 3.7².

However, rather than checking each time the contents of the stack (and the next input symbol) vs. the DFA M_G , it is more efficient to store the states together with the symbols on the stack in such a way that for each prefix xX of (the contents of) the stack, we store the state $\delta^*(q_0, xX)$ together with X on the stack.

In this way, we can simply check $\delta_{M_G}(q, a)$ where $q \in Q$ is the state on the top of the stack and a is the current input symbol, to find out whether shifting the input would still yield a viable prefix on the stack. Similarly, upon a $\text{reduce}(A \rightarrow \alpha)$

²We assume, without losing generality, that the start symbol of G does not occur on the right hand side of any production

```

1  while (true) {
2      if ((the top of the stack contains the start symbol) &&
3          (the current input symbol is the endmarker '$'))
4          accept;
5      else if (the contents of the stack concatenated with the next
6                input symbol is a viable prefix)
7          shift the input onto the top of the stack;
8      else
9          if (there is an appropriate production)
10             reduce by this production;
11         else
12             error;
13 }

```

Figure 3.7: Naive shift-reduce parsing procedure

operation, both A and $\delta(q, A)$ are pushed where q is the state on the top of the stack after popping α .

In practice, a shift-reduce parser is driven by two tables:

- An **action** table that associates a pair consisting of a state $q \in Q$ and a symbol $a \in \Sigma$ (the input symbol) with an action of one of the following types:
 - *accept*, indicating that the parse finished successfully
 - *shift*(q'), where $q \in Q$ is a final state of M_G . This action tells the parser to push both the input symbol and q on top of the stack. In this case $\delta_{M_G}(q, a) = q'$.
 - *reduce*($A \rightarrow \alpha$) where $A \rightarrow \alpha$ is a production from G . This instructs the parser to pop α (and associated states) from the top of the stack, then pushing A and the state q' given by $q' = \text{goto}(q'', A)$ where q'' is the state left on the top of the stack after popping α . Here $\text{goto}(q'', A) = \delta_{M_G}(q'', A)$.
 - **error**, indicating that $\delta_{M_G}(q, a)$ yields a dead state.
- A **goto** table:

$$\text{goto} : Q \times V \rightarrow Q$$

which is used during a reduce operation.

The architecture of a shift-reduce parser is shown in Figure 3.8.

Using the *action* and *goto* tables, we can refine the shift-reduce parsing algorithm as follows.

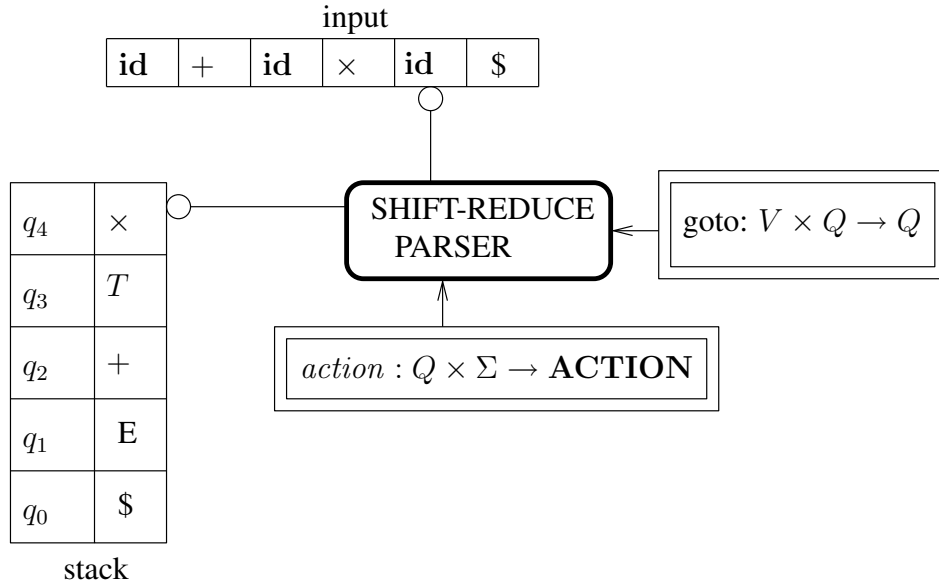


Figure 3.8: A shift-reduce parser

Algorithm 6 [Shift-reduce parsing]

See Figure 3.9 on page 55.

□

3.3.2 LR(1) parsing

Definition 14 Let $G = (V, \Sigma, P, S')$ be a context-free grammar such that $S' \rightarrow S$ is the only production for S' and S' does not appear on the right hand side of any production.

- An **LR(1) item** is a pair $[A \rightarrow \alpha \bullet \beta, a]$ where $A \rightarrow \alpha\beta$ is a production from P , \bullet is a new symbol and $a \in \Sigma \cup \{\$\}$ is a terminal or end marker symbol. We use items_G to denote the set of all items of G .
- The **LR(1) viable prefix NFA** of G is

$$N_G = (\text{items}_G, V \cup \Sigma, \delta, [S' \rightarrow \bullet S], \text{items}_G)$$

where δ is defined by

$$\begin{aligned} [A \rightarrow \alpha X \bullet \beta, a] &\in \delta([A \rightarrow \alpha \bullet X \beta, a], X) && \text{for all } X \in V \cup \Sigma \\ [X \rightarrow \bullet \gamma, b] &\in \delta([A \rightarrow \alpha \bullet X \beta, a], \epsilon) && \text{if } X \rightarrow \gamma \text{ is in } P \text{ and } b \in \text{first}(\beta a) \end{aligned}$$


```

1  TOKEN* input; /* array of input tokens */
2  STATE* states; /* (top of) stack of states */
3  SYMBOL* symbols; /* (top of) stack of symbols */
4
5  BOOL
6  parse() {
7      while (true)
8          switch (action[*states,*input]) {
9              case shift(s):
10                 ++symbols = *input++; /* push symbol */
11                 ++states = s; /* push state */
12                 break;
13              case reduce(X->x):
14                 symbols -= length(x); /* pop rsh */
15                 states -= length(x); /* also for states */
16                 STATE s = goto[*states,X];
17                 ++states = s; ++symbols = X;
18                 break;
19              case accept:
20                 return true;
21              case error:
22                 return false;
23          }
24 }

```

Figure 3.9: Shift-reduce parser algorithm

Note that all states in N_G are final. This does not mean that N_G accepts any string since there may be no transitions possible on certain inputs from certain states.

Intuitively, an item $[A \rightarrow \alpha \bullet \beta, a]$ can be regarded as describing the state where the parser has α on the top of the stack and next expects to process the result of a derivation from β , followed by an input symbol a .

To see that $L(N_G)$ accepts all viable prefixes. Consider a (partial)³ rightmost derivation

$$S = X_0 \xRightarrow{R}_G x_0 X_1 y_0 \xRightarrow{R}_G^* x_0 x_1 X_2 y_1 \xRightarrow{R}_G^* x_0 \dots x_{n-1} X_n y_n \xRightarrow{R}_G x_0 \dots x_{n-1} x_n y_n$$

where $y_i \in \Sigma^*$ for all $0 \leq i \leq n$, as depicted in Figure 3.10 where we show only the steps involving the ancestors of the left hand side (X_n) of the last step (note that none of the x_i changes after being produced because we consider a rightmost derivation).

³By a partial rightmost derivation we mean a rightmost derivation that can be extended to a successful full rightmost derivation.

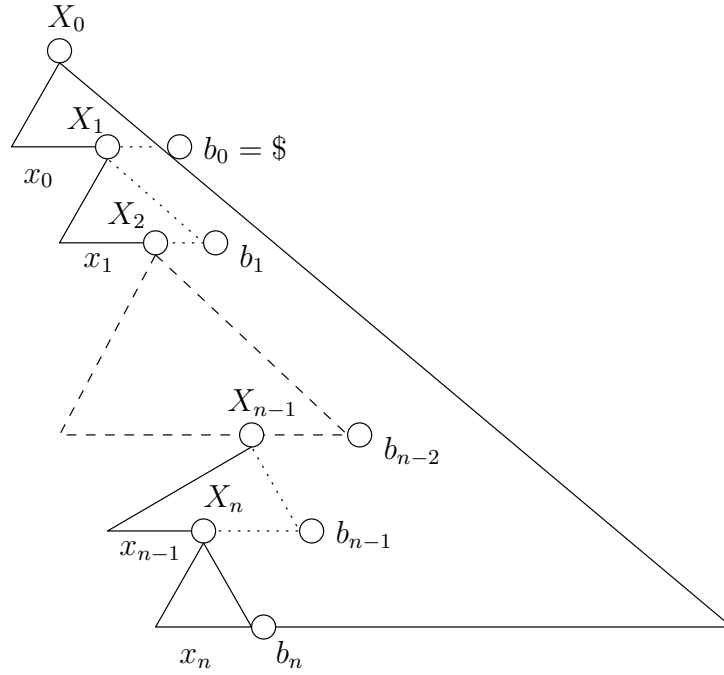


Figure 3.10: A (partial) rightmost derivation

N_G can accept any viable prefix of the right sentential form $x_0 \dots x_{n-1}x_n y_n$, i.e. any prefix of $x_0 \dots x_{n-1}x_n$ as follows:

| | |
|--|--|
| $[S' \rightarrow \bullet X_0, \$]$ | the initial state |
| $[X_0 \rightarrow \bullet x_0 X_1 z_0, b_0]$ | using an ϵ -move, $b_0 \in \text{first}(\epsilon \$) = \{\$ \}$ |
| $[X_0 \rightarrow x_0 \bullet X_1 z_0, b_0]$ | reading x_0 |
| $[X_1 \rightarrow \bullet x_1 X_2 z_1, b_1]$ | using an ϵ -move, $b_1 \in \text{first}(z_0 b_0)$ |
| $[X_1 \rightarrow x_1 \bullet X_2 z_1, b_1]$ | reading x_1 |
| \dots | |
| $[X_{n-1} \rightarrow \bullet x_{n-1} X_n z_{n-1}, b_{n-1}]$ | |
| $[X_{n-1} \rightarrow x_{n-1} \bullet X_n z_{n-1}, b_{n-1}]$ | reading x_{n-1} |
| $[X_n \rightarrow \bullet x_n, b_n]$ | using an ϵ -move, $b_n \in \text{first}(z_{n-1} b_{n-1})$ |
| \dots | any prefix of x_n can be read |

Note that all b_i exist since each z_i derives a terminal string.

To see the reverse, it suffices to note that each accepting sequence in M_G is of the form sketched above and thus, a partial rightmost derivation like the one in Figure 3.10 can be inferred from it.

Theorem 4 Let $G = (V, \Sigma, P, S')$ be a context-free grammar and let N_G be its viable prefix NFA. Then $L(N_G)$ contains exactly the viable prefixes of G .

Because of Theorem 1, we can convert N_G to an equivalent DFA M_G that can be used to build the *action* and *goto* tables of a shift-reduce parser as follows.

Algorithm 7 [Construction of LR(1) tables]

Let

$$M_G = (2^{\text{items}_G}, V \cup \Sigma, \delta, C_\epsilon([S' \rightarrow \bullet S, \$]), 2^{\text{items}_G} \setminus \{\emptyset\})$$

be the DFA corresponding to the viable prefix automaton N_G of G .

- Add an action *accept* to $\text{action}[s, \$]$ whenever $[S' \rightarrow S\bullet, \$] \in s$.
- Add an action *shift*(s') to $\text{action}[s, a]$ whenever $\delta(s, a) = s'$.
- Add an action *reduce*($A \rightarrow \alpha$) to $\text{action}[s, a]$ whenever s contains an item $[A \rightarrow \alpha\bullet, a]$ where $A \neq S'$.
- For any $X \in V$, s a state of M_G , define $\text{goto}(s, X) = \delta(s, X)$.

After the above rules have been exhaustively applied, add **error** to all entries in *action* that remain empty.

The algorithm succeeds if $\text{action}[s, a]$ contains exactly one action for each reachable $s \in 2^{\text{items}_G}$, $a \in \Sigma$. □

Note that it follows from Theorem 4 and Algorithm 7 that the resulting parser will announce an error at the latest possible moment, that is only when shifting would result in a stack that does not represent a viable prefix.

If Algorithm 7 does not succeed, there exists either a **shift-reduce conflict**, when an entry in the action table contains both a shift and a reduce operation, or a **reduce-reduce conflict**, when an entry in the action table contains two or more reduce actions.

Example 15 Consider G_3 from Example 13. The DFA accepting the viable prefixes of G is shown in Figures 3.11 and 3.12. Note that $[A \rightarrow \alpha, abc]$ is used as shorthand for $\{[A \rightarrow \alpha, a], [A \rightarrow \alpha, b], [A \rightarrow \alpha, c]\}$.

Theorem 5 Let $G = (V, \Sigma, P, S')$ be a context-free grammar such that Algorithm 7 succeeds. Then Algorithm 6 accepts exactly $L(G)$.

| State | Items | State | Items |
|-------|---|-------|--|
| 0 | $S' \rightarrow \bullet S$ \$ $S \rightarrow \bullet E$ \$ $E \rightarrow \bullet E + T$ +\$ $E \rightarrow \bullet T$ +\$ $T \rightarrow \bullet T \times F$ $\times + \$$ $T \rightarrow \bullet F$ $\times + \$$ $F \rightarrow \bullet id$ $\times + \$$ $F \rightarrow \bullet (E)$ $\times + \$$ | 11 | $E \rightarrow E + T \bullet$ +\$ $T \rightarrow T \bullet \times F$ $\times + \$$ |
| | | 12 | $F \rightarrow (E \bullet)$ $\times + \$$ $E \rightarrow E \bullet + T$)+\$ |
| | | 13 | $F \rightarrow (E) \bullet$ $\times + \$$ |
| | | 14 | $T \rightarrow F \bullet$) $\times +$ |
| | | 15 | $F \rightarrow (E) \bullet$) $\times +$ |
| | | 16 | $F \rightarrow (E \bullet)$) $\times +$ $E \rightarrow E \bullet + T$)+\$ |
| 1 | $S \rightarrow E \bullet$ \$ $E \rightarrow E \bullet + T$ +\$ | 17 | $E \rightarrow E + T \bullet$)+\$ $T \rightarrow \bullet T \times F$) $\times +$ $T \rightarrow \bullet F$) $\times +$ $F \rightarrow \bullet id$) $\times +$ $F \rightarrow \bullet (E)$) $\times +$ |
| 2 | $S' \rightarrow S \bullet$ \$ | | |
| 3 | $E \rightarrow E + \bullet T$ +\$ $T \rightarrow \bullet T \times F$ $\times + \$$ $T \rightarrow \bullet F$ $\times + \$$ $F \rightarrow \bullet id$ $\times + \$$ $F \rightarrow \bullet (E)$ $\times + \$$ | 18 | $F \rightarrow (\bullet E)$) $\times +$ $E \rightarrow \bullet E + T$)+\$ $E \rightarrow \bullet T$)+\$ $T \rightarrow \bullet T \times F$) $\times +$ $T \rightarrow \bullet F$) $\times +$ $F \rightarrow \bullet id$) $\times +$ $F \rightarrow \bullet (E)$) $\times +$ |
| 4 | $T \rightarrow F \bullet$ $\times + \$$ | 19 | $E \rightarrow T \bullet$)+\$ $T \rightarrow T \bullet \times F$) $\times +$ |
| 5 | $F \rightarrow id \bullet$ $\times + \$$ | 20 | $E \rightarrow E + T \bullet$)+\$ $T \rightarrow T \bullet \times F$) $\times +$ |
| 6 | $F \rightarrow (\bullet E)$ $\times + \$$ $E \rightarrow \bullet E + T$)+\$ $E \rightarrow \bullet T$)+\$ $T \rightarrow \bullet T \times F$) $\times +$ $T \rightarrow \bullet F$) $\times +$ $F \rightarrow \bullet id$) $\times +$ $F \rightarrow \bullet (E)$) $\times +$ | 21 | $T \rightarrow T \times \bullet F$) $\times +$ $F \rightarrow \bullet id$) $\times +$ $F \rightarrow \bullet (E)$) $\times +$ |
| 7 | $E \rightarrow T \bullet$ +\$ $T \rightarrow T \bullet \times F$ $\times + \$$ | 22 | $T \rightarrow T \times F \bullet$) $\times +$ |
| 8 | $T \rightarrow T \times \bullet F$ $\times + \$$ $F \rightarrow \bullet id$ $\times + \$$ $F \rightarrow \bullet (E)$ $\times + \$$ | | |
| 9 | $T \rightarrow T \times F \bullet$ $\times + \$$ | | |
| 10 | $F \rightarrow id \bullet$) $\times +$ | | |

Figure 3.11: States of the viable prefix DFA of G_3

Example 16 The action and goto tables of the LR(1) parser for G_3 are shown in Figure 3.13. The rules of G_3 have been numbered as follows:

| | | | |
|---|----------------------------|---|---------------------|
| 1 | $S \rightarrow E$ | 5 | $T \rightarrow F$ |
| 2 | $E \rightarrow E + T$ | 6 | $F \rightarrow id$ |
| 3 | $E \rightarrow T$ | 7 | $F \rightarrow (E)$ |
| 4 | $T \rightarrow T \times F$ | | |

A parse of

$id + id \times id$

| $\delta_{M_{G_3}}$ | S | E | T | F | \mathbf{id} | $($ | $)$ | $+$ | \times |
|--------------------|-----|-----|-----|-----|---------------|-----|-----|-----|----------|
| 0 | 2 | 1 | 7 | 4 | 5 | 6 | | | |
| 1 | | | | | | | | 3 | |
| 2 | | | | | | | | | |
| 3 | | | 11 | 4 | 5 | 6 | | | |
| 4 | | | | | | | | | |
| 5 | | | | | | | | | |
| 6 | | | 19 | 14 | 10 | 18 | | | |
| 7 | | | | | | | | | 8 |
| 8 | | | | 9 | 5 | 6 | | | |
| 9 | | | | | | | | | |
| 10 | | | | | | | | | |
| 11 | | | | | | | | | 8 |
| 12 | | | | | | | 13 | 17 | |
| 13 | | | | | | | | | |
| 14 | | | | | | | | | |
| 15 | | | | | | | | | |
| 16 | | | | | | | | 17 | |
| 17 | | | 20 | 14 | 10 | 18 | | | |
| 18 | | | 19 | 14 | 10 | 18 | | | |
| 19 | | | | | | | | | 21 |
| 20 | | | | | | | | | 21 |
| 21 | | | | | | | | | |
| 22 | | | | 22 | 10 | | | | |

Figure 3.12: Transition function of the viable prefix DFA of G_3

is shown below (states on the stack are between “[]”).

| stack | input | operation |
|---|--|------------------------------------|
| $\$[0]$ | $\mathbf{id} + \mathbf{id} \times \mathbf{id}\$$ | shift |
| $\$[0]\mathbf{id}[5]$ | $+\mathbf{id} \times \mathbf{id}\$$ | reduce $F \rightarrow \mathbf{id}$ |
| $\$[0]F[4]$ | $+\mathbf{id} \times \mathbf{id}\$$ | reduce $T \rightarrow F$ |
| $\$[0]T[7]$ | $+\mathbf{id} \times \mathbf{id}\$$ | reduce $E \rightarrow T$ |
| $\$[0]E[1]$ | $+\mathbf{id} \times \mathbf{id}\$$ | shift |
| $\$[0]E[1] + [3]$ | $\mathbf{id} \times \mathbf{id}\$$ | shift |
| $\$[0]E[1] + [3]\mathbf{id}[5]$ | $\times \mathbf{id}\$$ | reduce $F \rightarrow \mathbf{id}$ |
| $\$[0]E[1] + [3]F[4]$ | $\times \mathbf{id}\$$ | reduce $T \rightarrow F$ |
| $\$[0]E[1] + [3]T[11]$ | $\times \mathbf{id}\$$ | shift |
| $\$[0]E[1] + [3]T[11] \times [8]$ | $\mathbf{id}\$$ | shift |
| $\$[0]E[1] + [3]T[11] \times [8]\mathbf{id}[5]$ | $\$$ | reduce $F \rightarrow \mathbf{id}$ |
| $\$[0]E[1] + [3]T[11] \times [8]F[9]$ | $\$$ | reduce $T \rightarrow T \times F$ |
| $\$[0]E[1] + [3]T[11]$ | $\$$ | reduce $E \rightarrow E + T$ |
| $\$[0]E[1]$ | $\$$ | reduce $S \rightarrow E$ |
| $\$[0]S[2]$ | $\$$ | accept |

It will turn out that Algorithm 7 succeeds exactly for so-called LR(1) grammars.

| | goto | | | | action | | | | | |
|-------|------|---|----|----|--------|-----|-----|-----|-----|----|
| state | S | E | T | F | id | (|) | + | × | \$ |
| 0 | 2 | 1 | 7 | 4 | s5 | s6 | | | | |
| 1 | | | | | | | | s3 | | r1 |
| 2 | | | | | | | | | | a |
| 3 | | | 11 | 4 | s5 | s6 | | | | |
| 4 | | | | | | | | r5 | r5 | r5 |
| 5 | | | | | | | | r6 | r6 | r6 |
| 6 | | | 19 | 14 | s10 | s18 | | | | |
| 7 | | | | | | | | r3 | s8 | r3 |
| 8 | | | | 9 | s5 | s6 | | | | |
| 9 | | | | | | | | r4 | r4 | r4 |
| 10 | | | | | | | r6 | r6 | r6 | |
| 11 | | | | | | | | r2 | s8 | r2 |
| 12 | | | | | | | s13 | s17 | | |
| 13 | | | | | | | | r7 | r7 | r7 |
| 14 | | | | | | | r5 | r5 | r5 | |
| 15 | | | | | | | r7 | r7 | r7 | |
| 16 | | | | | | | | s17 | | |
| 17 | | | 20 | 14 | s10 | s18 | | | | |
| 18 | | | 19 | 14 | s10 | s18 | | | | |
| 19 | | | | | | | r3 | r3 | s21 | |
| 20 | | | | | | | r2 | r2 | s21 | |
| 21 | | | | 22 | s10 | | | | | |
| 22 | | | | | | | r4 | r4 | r4 | |

Figure 3.13: LR(1) tables for G_3

Definition 15 An $LR(1)$ **grammar** is a context-free grammar $G = (V, \Sigma, \delta, P, S)$ such that S does not appear at the right hand side of any production. Moreover, if

$$S \xRightarrow*_G^R x_1 X w_3 \xRightarrow*_G^R x_1 x_2 w_3$$

and

$$S \xRightarrow*_G^R \hat{x}_1 \hat{X} \hat{w}_3 \xRightarrow*_G^R \hat{x}_1 \hat{x}_2 \hat{w}_3$$

and

$$\text{pref}_{l+1}(x_1 x_2 w_3) = \text{pref}_{l+1}(\hat{x}_1 \hat{x}_2 \hat{w}_3)$$

where $X \in V$, $w_3 \hat{w}_3 \in \Sigma^*$, $l = |x_1 x_2|$, then

$$x_1 = \hat{x}_1 \wedge \hat{X} = X$$

A language L is called $LR(1)$ if there exists an $LR(1)$ context-free grammar G such that $L(G) = L$.

Intuitively, in the above definition, $x_1 x_2$ represents the stack just before a reduce operation in a shift-reduce parser. The definition then says that there is only a single choice for when to reduce, based only on the next input symbol⁴.

Theorem 6 A context-free grammar $G = (V, \Sigma, \delta, P, S)$, where S does not appear on the right hand side of any production, is $LR(1)$ iff Algorithm 7 succeeds.

One may wonder about the relationship between $LL(k)$ and $LR(k)$ grammars and languages. It turns out (see e.g. [Sal73]) that every $LR(k)$ ($k \geq 0$) language can be generated by an $LR(1)$ grammar. On the other hand, there are context free languages, e.g. the inherently ambiguous language of Section 3.1, page 38, that are not $LR(1)$. In fact, it can be shown that all deterministic context-free languages⁵ are $LR(1)$ (and, obviously, also the reverse). As for $LL(k)$ languages, they form a proper hierarchy in that for every $k > 1$ there is an $LL(k)$ language L_k for which no $LL(k-1)$ grammar exists. Still, every $LL(k)$ language is $LR(1)$.

Summarizing, as far as parsing power is concerned, bottom-up ($LR(1)$) parsers are provably superior to top-down ($LL(1)$) parsers.

⁴One can also define $LR(k)$ grammars and languages for any $k \geq 0$. It suffices to replace $l+1$ by $l+k$ in Definition 15

⁵A context-free language is deterministic if it can be accepted by a deterministic pushdown automaton.

3.3.3 LALR parsers and yacc/bison

LR(1) parsing tables can become quite large due to the large number of states in the viable prefix DFA. A DFA for a language like Pascal may have thousands of states.

In an LALR parser, one tries to reduce the number of states by merging states that have the same (extended) productions in their items. E.g. in the DFA of Example 15, one may merge, among others, state 3 and state 17.

It is then not hard to see that, due to the way states are constructed using ϵ -closures, if s_1 and s_2 can be merged, then so can $\delta(s_1, X)$ and $\delta(s_2, X)$. This yields a more compact DFA which is “almost” equivalent to the original one in the sense that compaction cannot produce shift-reduce conflicts although new reduce-reduce conflicts are possible.

It can be shown that if the LALR compaction succeeds (i.e. does not lead to new reduce-reduce conflicts), the resulting parser accepts the same language as the original LR(1) parser.

The parser-generator yacc (and bison) generates an LALR parser based on an input grammar. In addition, yacc (and bison) provide facilities for resolving shift-reduce conflicts based on the declared precedence of operators.

Bison will resolve a shift reduce conflict on a state like

$$\begin{aligned} E &\rightarrow E o_1 E \bullet, o_2 \\ E &\rightarrow E \bullet o_2 E \end{aligned}$$

by comparing the declared precedences of o_1 and o_2 : if o_1 has the higher precedence, then reduce will be selected. If o_2 has a higher precedence than o_1 then shift will be selected. If o_1 and o_2 have the same precedence, associativity information on o_1 and o_2 will be used: if both symbols (operators) are left-associative, reduce will be selected; if both symbols are right-associative, shift will be selected; if both symbols are non-associative an error will be reported⁶. If a symbol (such as MINUS in Example 17 below) is used with two precedences, bison allows you to define the precedence of the rule using a `%prec` directive⁷

The default conflict resolution strategy of bison (and yacc) is as follows:

⁶Note that operators with the same precedence must have the same associativity.

⁷The above description is made more concrete in the bison manual:

The first effect of the precedence declarations is to assign precedence levels to the terminal symbols declared. The second effect is to assign precedence levels to certain rules: each rule gets its precedence from the last terminal symbol mentioned in the components (o_1 in the example). (You can also specify explicitly the precedence of a rule.)

- Choose shift when resolving a shift-reduce conflict.
- Choose reduce by the first production (in the order of declaration) when resolving a reduce-reduce conflict.

Example 17 The following is the bison input file for a highly ambiguous grammar, where ambiguities are resolved using the precedence/associativity defining facilities of bison/yacc.

```

1
2 %token NAME  NUMBER LPAREN RPAREN EQUAL PLUS MINUS
3 %token TIMES DIVIDE IF THEN ELSE
4
5 /*      associativity and precedence: in order of increasing precedence */
6
7 %nonassoc      LOW  /* dummy token to suggest shift on ELSE */
8 %nonassoc      ELSE /* higher than LOW */
9
10 %nonassoc      EQUAL
11 %left          PLUS    MINUS
12 %left          TIMES   DIVIDE
13 %left          UMINUS  /* dummy token to use as precedence marker */
14
15 %%
16
17 stmt           : IF exp THEN stmt           %prec LOW
18                | IF exp THEN stmt ELSE stmt /* shift will be selected */
19                | /* other types of statements would come here */
20                ;
21
22 exp            : exp PLUS exp
23                | exp MINUS exp
24                | exp TIMES exp
25                | exp DIVIDE exp
26                | exp EQUAL exp
27                | LPAREN exp RPAREN
28                | MINUS exp %prec UMINUS /* this will force a reduce */
29                | NAME
30                | NUMBER

```

Finally, the resolution of conflicts works by comparing the precedence of the rule (of o_1 in the example, unless that rule's precedence has been explicitly defined) being considered (for reduce) with that of the look-ahead token (o_2 in the example). If the token's precedence is higher, the choice is to shift. If the rule's precedence is higher, the choice is to reduce. If they have equal precedence, the choice is made based on the associativity of that precedence level. The verbose output file made by '-v' (see section Invoking Bison) says how each conflict was resolved.

Not all rules and not all tokens have precedence. If either the rule or the look-ahead token has no precedence, then the default is to shift.

31 ;
32 %%

Chapter 4

Checking static semantics

4.1 Attribute grammars and syntax-directed translation

Attribute grammars extend context-free grammars by associating a set of **attributes** with each grammar symbol. Such attributes can be used to present information that cannot easily (if at all) be expressed using only syntactical constructs (e.g. productions). Thus each node in a parse tree has values for each of the attributes associated with the corresponding symbol.

For example, we could associate an attribute *type* with the nonterminal *expression* to represent the type of the expression.

Attributes for nodes are computed using functions that take attribute values of other nodes as input. The computation is specified by *semantic actions* which are associated with each production. Since such actions can only reference the symbols in the production, the value of an attribute of a node depends only on the attribute values of its parent, children and siblings.

Example 18 Consider the grammar

$$G_4 = (\{E\}, \{+, \times, -, /, (,), \text{number}\}, P, E)$$

where P consists of

$$E \rightarrow \text{number} \mid (E) \mid E + E \mid E \times E \mid E / E \mid E - E \mid -E$$

We define an attribute *value* for E and **number**. The semantic actions are as follow:

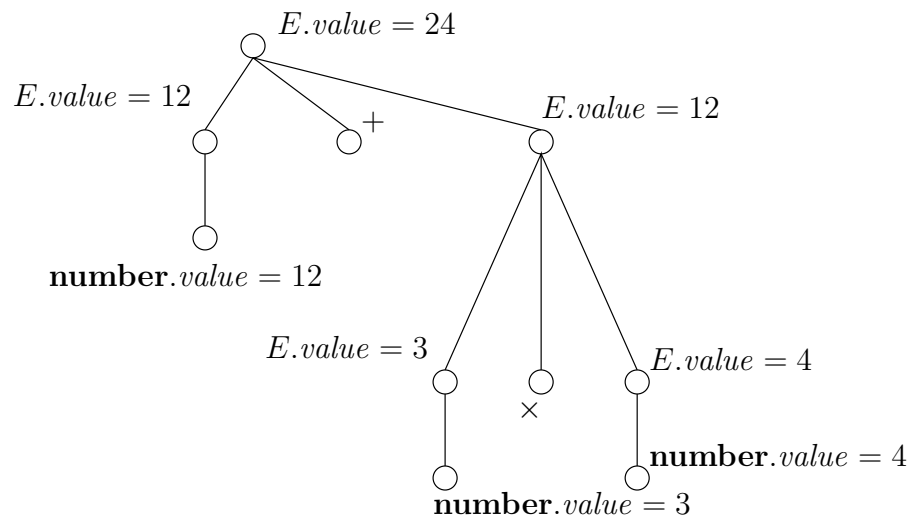


Figure 4.1: Computing synthesized attribute values

| | |
|-------------------------------|--|
| $E \rightarrow \text{number}$ | $E_0.value = \text{number}_1.value$ |
| $E \rightarrow (E)$ | $E_0.value = E_1.value$ |
| $E \rightarrow E + E$ | $E_0.value = E_1.value + E_3.value$ |
| $E \rightarrow E \times E$ | $E_0.value = E_1.value \times E_3.value$ |
| $E \rightarrow E - E$ | $E_0.value = E_1.value - E_3.value$ |
| $E \rightarrow E / E$ | $E_0.value = E_1.value / E_3.value$ |
| $E \rightarrow -E$ | $E_0.value = -E_1.value$ |

Note that we use subscripts to refer to a particular occurrence of a symbol in the rule (subscript 0 is used for the left-hand side of the rule).

The attribute values for a parse tree of $12 + 3 * 4$ are shown in Figure 4.1. Note that the **number.value** is supposed to have been provided by the lexical analyzer.

The actions in Example 18 all compute an attribute value for a node based on the attribute values of its children. Attributes that are always computed in this way are called *synthesized attributes*. Attributes whose value depends on values of its parent and/or siblings are called *inherited attributes*.

The following example shows a grammar with both synthesized and inherited attributes.

Example 19 Consider the grammar

$$G_5 = (\{D, T, L\}, \{', ', ;, \text{id}, \text{int}, \text{double}\}, P, D)$$

Here D represents a C-style declaration consisting of a type followed by a comma-separated list of identifiers. The symbols T , L , and id have a *type* attribute which is inherited for L and id .

The productions and actions of P are

$$\begin{array}{ll} D \rightarrow TL; & L_2.type = T_1.type \\ T \rightarrow \text{int} & T_0.type = \text{int} \\ T \rightarrow \text{double} & T_0.type = \text{double} \\ L \rightarrow L, \text{id} & L_1.type = L_0.type \\ & \text{id}_3.type = L_0.type \\ L \rightarrow \text{id} & \text{id}_1.type = L_0.type \end{array}$$

The attribute values for the parse tree of

int $x, y, z;$

are shown in Figure 4.2. The dependencies between values are shown by dashed arrows.

In general, one can define a **dependency graph** for a parse tree where the nodes are pairs (n, a) where n is a node in the parse tree and a is an attribute for (the symbol of) n . There is an edge from (n_1, a_1) to (n_2, a_2) if the value of $n_2.a_2$ depends on the value of (n_1, a_1) , i.e. an action computing $n_2.a_2$ uses (n_1, a_1) .

Attribute values can then be computed in any order that corresponds to a topological sort¹ of the nodes of the dependency graph. Clearly, for this to work, the dependency graph may not contain a cycle.

Thus, the syntax-directed translation $G(\alpha)$ of an input string α using an attribute grammar G can be computed by executing the following steps:

1. Parse α , according to the context-free grammar underlying G , yielding a parse tree t_α .
2. Construct the dependency graph g_{t_α} , and compute the attribute values of the nodes of t_α as described above.

¹A topological sort of a directed (acyclic) graph is any total ordering $x_1 < x_2 < \dots < x_n < \dots$ of the nodes of the graph that satisfies $x_i < x_j$ whenever there is an edge from x_i to x_j .

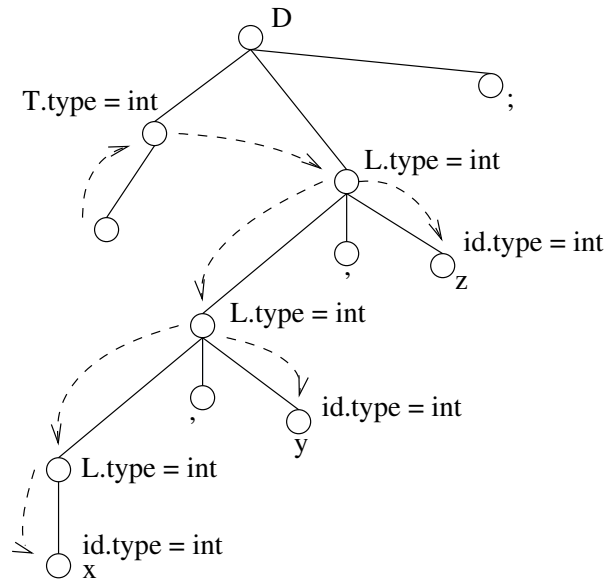


Figure 4.2: Computing inherited attribute values

3. $G(\alpha)$ is the value corresponding to some selected attribute of the root of t_α .

In practice, the full parse tree for an input program is often not explicitly constructed. Rather, attributes are computed and stored on a stack that mirrors the parsing stack (just like the stack containing the states of the viable prefix DFA) and actions are performed upon a reduce operations. This allows computation of synthesized attributes but inherited attributes are not easily supported although in certain cases it is possible to also evaluate inherited attributes (see [ASU86], pages 310-311).

In yacc and bison, each symbol (nonterminal or token) can have only one attribute, not larger than a single word. The type of an attribute is fixed per symbol using `%type` declarations, as can be seen in the example parser in Appendix C.6, page 175. The set of used types must also be declared in a `%union` declaration. The restriction to a single word is not limiting: usually the attribute is declared to be a pointer to a structure.

4.2 Symbol tables

The restriction of shift-reduce parsers to use only synthesized attributes can be overcome by using a global symbol table to share information between various

parts of the parse tree. In general, a symbol table maps names (of variables, functions, types, ...) to information on this name which has been accumulated by the parser and semantic actions.

We discuss two aspects of symbol tables:

- The management of names (the domain of the mapping).
- The organization of the symbol table to take into account the scope rules of the language.

4.2.1 String pool

Since a symbol table maps names to information, one could implement a symbol table entry as a structure where the first field is the name and the other fields contain the information associated with the name

```
1 struct symtab_entry {
2     char    name[MAX_NAME_LENGTH];
3     INFO    info;
4     };
```

However, because modern programming languages do not impose a small limit on the length of a name, it is wasteful or impossible to keep this organization. Hence one typically stores names contiguously in a so-called “string pool” array, declaring a symbol table entry as follows:

```
1 struct symtab_entry {
2     char*   name;    /* points into string pool */
3     INFO    info;
4     };
```

For fast retrieval, entries may be organized into a hash table, as illustrated in Figure 4.3 where separate chaining is used to link entries that correspond to the same hash value.

4.2.2 Symbol tables and scope rules

Many programming languages impose *scope rules* that determine which definition of a name corresponds to a particular occurrence of that name in a program text. E.g. in Pascal, and in C, scopes may be nested to an arbitrary depth, as is illustrated in the following fragment.

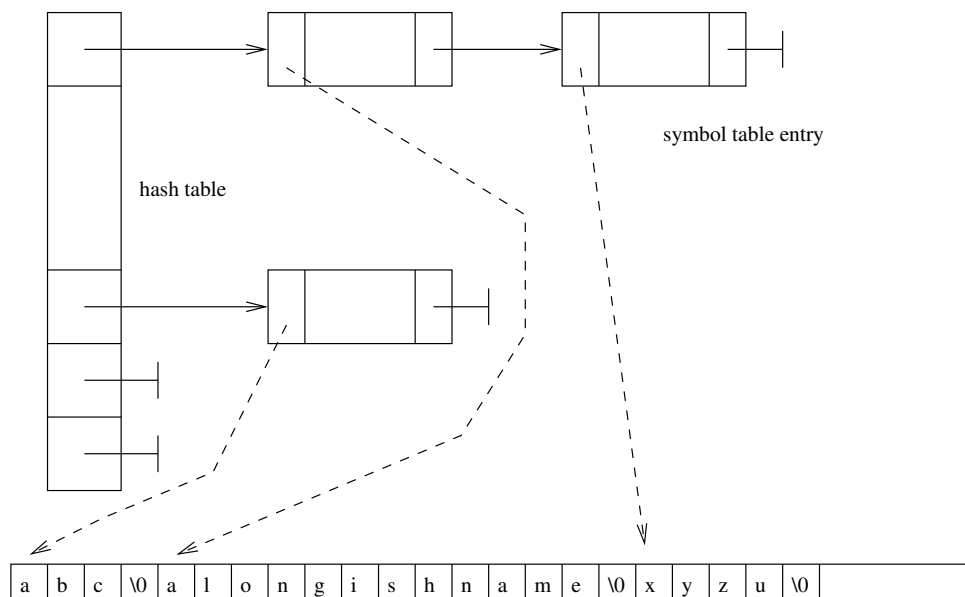


Figure 4.3: A simple symbol table using a string pool and hashing

```

1 float   x; /* def 1 */
2
3 int
4 f(int a)
5 {
6 while (a>0) {
7     double x = a*a,i; /* def 2 */
8
9     for (i=0; (i<a); ++i) {
10         int x = a+i,b; /* def 3 */
11         b = x+2; /* refers to def 3 */
12     }
13 }
14 }
15 }

```

In the above examples, the symbol table will contain three entries (definitions) for x at the time the inner loop is processed. Which definition is meant is determined using the “most closely nested scope” rule.

In a compiler for a language with scope rules, when a semantic action wants to retrieve the information associated with a certain occurrence of a name, it has to know which definition of the name corresponds to the given occurrence. If the language allows only nested scopes for which the “most closely nested scope” rule holds, this can be solved by replacing the symbol table by a stack of simple

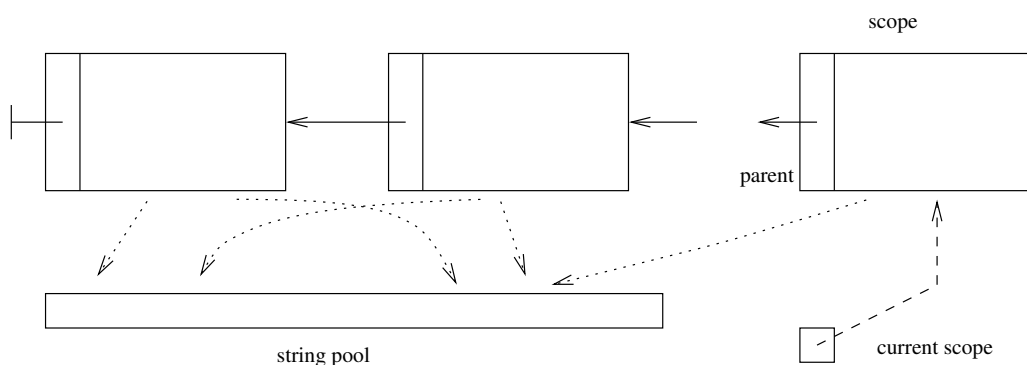


Figure 4.4: A stacked symbol table

symbol tables, one for each scope. Moreover, each simple table refers to its parent in the stack which contains symbol definitions from the enclosing scope. When looking up a name, it suffices to search the simple table at the top of the stack. If the name is not found, the parent table (corresponding to the enclosing scope) is searched and so on.

When the parser finishes a scope it pops the current scope from the stack and continues with the parent scope as the current one. Of course, if the scope is still needed, as is the case e.g. in C++ with scopes that correspond to class declarations, the scope will not be thrown away but a reference to it will be kept, e.g. in the symbol table entry corresponding to the class name.

This organization is depicted in Figure 4.4.

An simple implementation of this scheme, where hashing is replaced by linear search (first in the current scope, then in the parent scope), may be found in Appendix C.3, page 163.

4.3 Type checking

Informally, the so-called *static semantics* of a programming language describes the properties of a program that can be verified without actually executing it. Typically, for statically-typed languages, this consists of verifying the “well-typed” property where a program is well-typed if all its constructs (statements, expressions, ...) obey the type rules of the language.

The *type system* of a programming language consists of

- A definition of the set of all possible types.

- Type rules that determine the type of a language construct.

The language from Appendix C, Minic, has two primitive types: integers and floats and a few type constructors, as shown by the following domain² definition³:

$$\begin{aligned}
 \text{type} &:: \text{int} \\
 &:: \text{float} \\
 &:: \text{array}(\text{type}) \\
 &:: \text{struct}(\text{tenv}) \\
 &:: \text{type}^* \rightarrow \text{type} \\
 \text{tenv} &:: \text{name} \rightarrow \text{type}
 \end{aligned}$$

Note that the above definition is more general than is actually needed in Minic, which has no higher order functions. E.g. it is not possible to define a function with type

$$\text{int} \rightarrow (\text{int} \rightarrow \text{int})$$

in Minic.

The type rules of Minic can be symbolically represented as follows (we use $x : t$ to denote that x has type t). In the table below we use $+$ to use any of the arithmetic operators $+$, \times , $/$, $-$.

$$\begin{aligned}
 t \ x; &\Rightarrow x : t \quad // \text{declaration of } x \\
 x : \text{int} \wedge y : \text{int} &\Rightarrow (x + y) : \text{int} \\
 x : \text{float} \wedge y : \text{float} &\Rightarrow (x + y) : \text{float} \\
 x : t \wedge y : t &\Rightarrow (x == y) : \text{int} \quad // \text{support} = \text{test on any type} \\
 x : \text{array}(t) \wedge i : \text{int} &\Rightarrow x[i] : t \\
 x : \text{struct}(s) \wedge n \in \text{dom}(s) &\Rightarrow x.n : s(n) \\
 x = e \wedge x : t \wedge e : t' \wedge t \neq t' &\Rightarrow \text{error} \\
 \forall 1 \leq i \leq n \cdot x_i : t_i \wedge f : t_1 t_2 \dots t_n \rightarrow t &\Rightarrow f(x_1, \dots, x_n) : t \\
 \text{return } x \wedge x : t' \wedge f : t_1 t_2 \dots t_n \rightarrow t \wedge t' \neq t &\Rightarrow \text{error} \quad // \text{within the body of } f
 \end{aligned}$$

Here **error** is used to indicate that a certain construct is not well-typed, i.e. no type can be assigned to it without violating the type rules.

²See e.g. the “Foundations of Computer Science II” course.

³Note that there is no boolean type: like C, Minic uses integers instead of booleans.

Type checking then consists of using the rules to assign a type to each expression, variable etc. If no type, or more than one type, can be assigned to a construct, or if the rules generate **error**, the program is not well-typed.

It should be obvious that with the help of attributes, semantic actions and a symbol table, such type checking is possible. An example of a type-checking parser for Minic can be found in Appendix C. In the example implementation, types are represented by structures where each structure contains a tag corresponding to the type constructor. Furthermore, types are shared, which saves space and facilitates later checking for equality (only pointers must be compared).

Note that there are statically typed languages (e.g. ML or Haskell) that do not require, as Minic does, that each variable is declared to have a specific type. Instead, the type of a variable is inferred from its use. In addition, such languages usually have a more sophisticated type system where types can have variables that may be instantiated. In this way, a language can allow for polymorphism where a function has several (possibly an infinite number of) types.

A similar feature is supported by C++ through the template facility. E.g. a function

```

1      template <class T>
2      int      length(list<T>) {
3          // ...
4      }
```

has a type $\text{list}(\alpha) \rightarrow \text{int}$ where α is a type variable. The compiler will instantiate α to **int** when processing the expression `length(l)` in

```

1      list<int>      l;
2
3      length(l);
```

Chapter 5

Intermediate code generation

While it is possible to construct a compiler that translates the source text directly into the target language, especially if sophisticated optimization is not required, there are important advantages of using an intermediate representation:

- It facilitates retargeting. A compiler employing an intermediate representation can be viewed as consisting of a “front end” (everything up to intermediate code generation¹) and a “back end” (optimization, code generation). Note that all dependencies on the source language are located in the front end while dependencies on the target language are confined to the back end. Thus, to retarget the compiler (i.e. to generate code for another processor), it suffices to adapt the back end only. Similarly, if the intermediate code is sufficiently general, the back end could be reused in a compiler for another language.
- Many machine-independent optimizations, especially those that require code movement etc., can be done more easily on intermediate code.

Several forms of intermediate code have been proposed (and used):

- Postfix notation.
- Abstract syntax trees.
- Three-address code.

¹Any machine-independent optimization that can be performed on the intermediate code should also be done in the front end.

5.1 Postfix notation

Postfix notation, also called “postfix Polish” after the originator J. Lukasiewicz, represents expressions and flow of control statements in postfix form, which is free of parentheses, as in the following example:

| infix | postfix |
|--------------------------|--------------------|
| $a + b$ | $ab+$ |
| $(a + b) \times (c - d)$ | $ab + cd - \times$ |

In general, an expression consisting of a k -ary operator θ , applied on k (sub)expressions e_1, \dots, e_k is represented as

$$e'_1 e'_2 \dots e'_k \theta$$

where e'_i , $1 \leq i \leq k$, is the postfix form of e_i .

Evaluation of postfix expressions is trivial, provided one uses an operand stack, as shown below.

```

1 operand
2 eval_operator(OPERATOR o,item args[]);
3
4 operand
5 eval_postfix(item code[],int code_len)
6 {
7   operand stack[MAX_STACK];
8   int    tos    = 0; // top of stack
9
10  for (int i=0; i<code_len;++i)
11      if (is_operand(code[i]))
12          stack[tos++] = code[i]; // push on stack
13      else
14          { // operator: pop args, eval, push result
15            k = arity(code[i]);
16            stack[tos-k] = eval_operator(code[i].oprtr, stack-k);
17            tos -= (k-1);
18          }
19  return stack[0];
20 }
```

By introducing labels which refer to a position in the postfix expression, arbitrary control structures can be implemented using instructions like

label jump
 $e_1 e_2$ *label jumplt*

In order to interpret labels and jump instructions, the above algorithm must of course be extended to use a program counter (an index in the `code` array),

Postfix notation is attractive because it can be easily generated using syntax-directed translation. Below, we use an synthesized attribute *code* representing the postfix notation of the corresponding symbol occurrence. The statements associated with a rule define the attribute of the left hand side of the rule (denoted `$$`), based on the attribute values of the symbols in its right hand side (`$i` refers to the *i*'th symbol in the rules right hand side).

```

1 exp      : exp binop exp {
2              $$code = concat ($1.code, $3.code, operator($2));
3              }
4 exp      : ( exp) {
5              $$code = $1.code;
6              }
7 exp      : var { $$code = $1; /* probably a symtab ref */ }
```

This scheme has the so-called *simple postfix* property which means that for a production

$$A \rightarrow A_1 A_2 \dots A_n$$

the generated code is such that²

$$\text{code}(A) = \text{code}(A_1) + \dots + \text{code}(A_n) + \text{tail}$$

which implies that code can be emitted as soon as a reduction occurs.

Postfix code is often used for stack-oriented virtual machines without general purpose registers, such as the Java virtual machine (see e.g. the micro compiler in Appendix B. It is not particularly suitable for optimization (code movement is hard).

5.2 Abstract syntax trees

Abstract syntax trees are “pruned” parse trees where all nodes that contain redundant (e.g. with respect to the tree structure) information have been removed. Indeed, many nodes in the parse tree result from syntactic sugar in the grammar which is needed to ensure that the language can be efficiently parsed. Once we have parsed the text, there is no need to keep these superfluous nodes³. E.g. the abstract syntax tree corresponding to the parse tree of Figure 1.6 on page 14 is shown in Figure 5.1.

²We use $+$ to denote concatenation.

³It should be noted that a formal definition of a language’s semantics is often based on a so-called *abstract syntax* which only specifies the structure of the language’s constructs, rather than

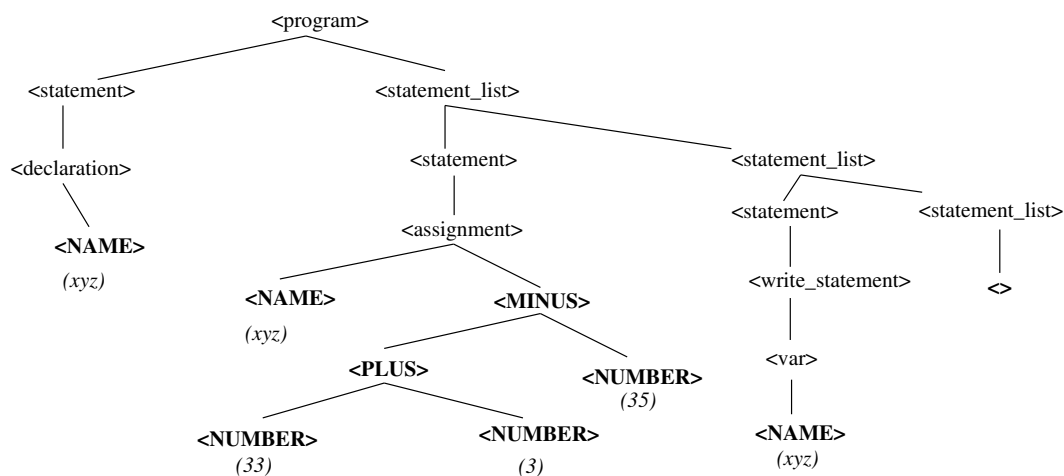


Figure 5.1: Abstract syntax tree corresponding to the program in Figure 1.3

Abstract syntax trees (and directed acyclic graphs), particularly of expressions, are useful for certain optimizations such as common subexpression elimination and optimal register allocation (for expression computation).

An abstract syntax tree can be generated using synthesized attributes.

```

1  exp      : exp binop exp {
2              $$ . tree = new_node_2($2,$1.tree,$3,tree);
3              }
4  %*
5  exp      : ( exp ) {
6              $.tree = $2.tree;
7              }
8  %*
9  exp      : unop exp {
10             $$ . tree = new_node_1($1,$2.tree);
11             }
12 %*
13 exp      : var {
14             $$ . tree = new_node_0($1);
15             }

```

Abstract syntax trees can also be constructed for complete source texts. They then

the concrete syntax. E.g. in the abstract syntax specification, one would have a rule

$$\text{assignment} \rightarrow \text{var expression}$$

rather than

$$\text{assignment} \rightarrow \text{var} = \text{expression}$$

form the input for code generators that are based on pattern matching (in trees) and tree transformations. This strategy is employed by so-called *code-generator generators* (see e.g. Section 9.12 in [ASU86]) which take as input *tree-rewriting rules* of the form

$$\text{replacement} \leftarrow \text{template} \{ \text{action} \}$$

where the template specifies a tree pattern that must be matched, the replacement is a tree (usually a single node) that is to replace the matched template and the action is a code fragment, as in syntax-directed translation. The code fragment may check additional constraints and generates target instructions. Figure 5.2 illustrates such tree-rewriting rules. The idea is that to retarget the compiler to

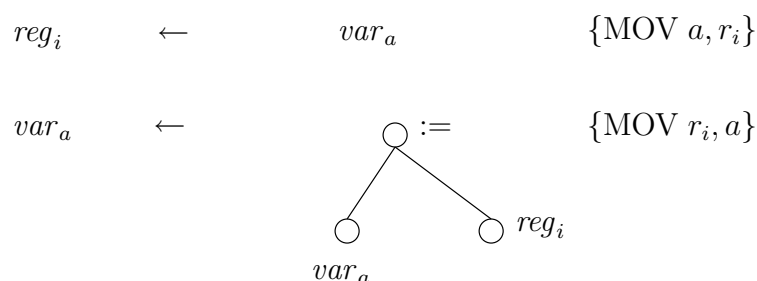


Figure 5.2: Tree-rewriting rules

a processor with different instruction sets (and costs), one simply replaces the tree-rewriting rules.

5.3 Three-address code

Three address code is a convenient and uniform intermediate representation. It consists of a number of simple instructions which involve at most three addresses, hence the name. The addresses refer to user-defined names (variables, functions, ...) or constants, or to temporary variables generated by the compiler.

In this and the following chapter, we will use three-address code as an intermediate representation. The instructions are listed in Figure 5.3.

An important advantage of three-address code is that it can be easily reordered (at least within a function body) for optimization. A convenient representation for three-address code instructions is as a quadruple.

```

1  typedef enum {A2PLUS, A2TIMES, A2MINUS, A2DIV, A1MINUS,
2                A1FTOI, A1ITOF, A0, GOTO, IFEQ, ..., PARAM,
3                CALL, AAC, AAS, ADDR, DEREf, DEREFA } OPCODE;
```


| opcode | instruction | meaning |
|-----------|-------------------------|--|
| A2(binop) | A = B binop C | assignment, binop can be +, ×, ... |
| A1(unop) | A = unop B | assignment, unop can be −, float2int, ... |
| A0 | A = B | simple assignment |
| GOTO | goto label | |
| IF(relop) | if A relop B goto label | relop can be <, >, ≤, ==, ... |
| PARAM | param A | push parameter before function call |
| CALL | call F, n | function call, n is number of parameters |
| AAC | A = B[I] | array access, B contains base address, I is offset in memory units |
| AAS | A[I] = B | array modification, A contains base address |
| ADDR | A = address B | take address of B |
| DEREF | A = * B | as in C |
| DEREFA | *A = B | as in C |
| RETURN | RETURN A | as in C |

Figure 5.3: Three-address code instructions

```

4 typedef struct {
5     OPCODE          op_code;
6     SYM_INFO*       args[2];
7     SYM_INFO*       result;
8     } INSTRUCTION;

```

5.4 Translating assignment statements

The actions below use a simplified grammar for expressions. Also, we assume that all expressions are of type integer.

We use two synthesized attributes:

- *exp.place* represents a reference to the location (variable, symbol table reference) containing the value of the expression *exp*.
- *exp.code* represents the sequence of three-address statements computing the value of *exp*. Code concatenation is written as “+”.

```

1 extern INSTRUCTION /* create quadruple representing instruction */
2 gen3ai(short op_code, SYM_INFO* arg1, SYM_INFO* arg2, SYM_INFO* result);
3
4 assign : var = exp {
5     $$ .code = $3.code + gen3ai(A0, $3.place, 0, $1.place);
6     }

```

```

7  exp      : exp + exp {
8              $$place = newvar();
9              $$code = $1.code + $3.code
10                 + gen3ai(A2PLUS, $1.place, $3.place, $$place);
11          }
12  exp      : - exp {
13              $$place = newvar();
14              $$code = $2.code
15                 + gen3ai(A1MINUS, $2.place, 0, $$place);
16          }
17  exp      : ( exp ) {
18              $$place = $2.place;
19              $$code = $2.code;
20          }
21  exp      : var {
22              $$place = $1.place;
23              $$code = 0;
24          }

```

Note that the actions for the *code* attribute have the “simple postfix” property. Hence we can simplify the actions.

```

1  extern void
2  emit(INSTRUCTION i);
3
4  assign    : var = exp {
5              emit(gen3ai(A0, $3.place, 0, $1.place));
6          }
7  exp      : exp + exp {
8              $$place = newvar();
9              emit(gen3ai(A2PLUS, $1.place, $3.place, $$place));
10         }
11  exp      : - exp {
12              $$place = newvar();
13              emit(gen3ai(A1MINUS, $2.place, 0, $$place));
14         }
15  exp      : ( exp ) {
16              $$place = $2.place;
17         }
18  exp      : var {
19              $$place = $1.place;
20         }

```

As an example of the integration of static (notably type) checking with intermediate code generation, we show below how appropriate conversions can be generated for languages that support them. Note that we use an extra attribute *type* representing the type of an expression, which can be either `t_int` (integer) or `t_float` (floating point). We also assume that there are two instructions for

addition: A2PLUSI and A2PLUSF for integer and floating point addition, respectively.

```

1  exp      : exp + exp {
2      $$place = newvar();
3      if (($1.type==t_int)&&($3.type==t_int)) {
4          $.type = t_int; symtab_set_type($$.place,t_int);
5          emit(gen3ai(A2PLUSI,$1.place,$3.place,$$.place));
6      }
7      else if (($1.type==t_float)&&($3.type==t_float)) {
8          $.type = t_float; symtab_set_type($$.place,t_float);
9          emit(gen3ai(A2PLUSF,$1.place,$3.place,$$.place));
10     }
11     else if (($1.type==t_int)&&($3.type==t_float)) {
12         SYM_INFO* tmpf = newvar(t_float);
13         emit(gen3ai(A1ITOF,$1.place,0,tmpf));
14         $.type = t_float; symtab_set_type($$.place,t_float);
15         emit(gen3ai(A2PLUSF,$1.place,tmpf,$$.place));
16     }
17     else if (($1.type==t_float)&&($3.type==t_int)) {
18         SYM_INFO* tmpf = newvar(t_float);
19         emit(gen3ai(A1ITOF,$3.place,0,tmpf));
20         $.type = t_float; symtab_set_type($$.place,t_float);
21         emit(gen3ai(A2PLUSF,tmpf,$3.place,$$.place));
22     }
23 }
```

5.5 Translating boolean expressions

First we consider the translation of boolean expressions (conditions). An obvious possibility is to translate such expressions *by value*, i.e. by creating, if necessary, a temporary variable to hold its value. In this case, we agree to use 0 to represent “false” and 1 to represent “true”.

```

1  extern int
2  next3ai(); /* returns number of next (free) location in code sequence */
3
4  bexp      : bexp OR bexp {
5      $$place = newvar(t_bool);
6      emit(gen3ai(A2OR,$1.place,$3.place,$$.place));
7      }
8
9  bexp      : bexp > bexp {
10     $$place = newvar(t_bool);
11     emit(gen3ai(IFGT,$1.place,$3.place,next3ai()+3));
12     emit(gen3ai(A0C,0,0,$$.place));

```

```

13         emit (gen3ai (GO, next3ai () + 2, 0, 0));
14         emit (gen3ai (A0C, 1, 0, $$ . place));
15     }

```

Note the use of the function `next3ai ()` which is used to generate forward jump instructions.

As an example, one can verify that the C expression

```

1     a < b || c

```

will be translated to

```

1   if a < b goto 4
2   t1 = 0
3   goto 5
4   t1 = 1
5   t2 = t1 || c

```

Translation by value is appropriate e.g. for code computing the value of a boolean variable. However, for conditions in while- or if-statements, translating by value may be wasteful. Indeed, suppose that $a < b$ in the following statement.

```

1   if ( (a < b) || (c < d) )
2       x = y + z;

```

When translating by value, the value of $c < d$ will be computed, even though it is already known (because $a < b$) that the then-branch will be taken. A better translation scheme would branch to the code for $x = y + z$ as soon as possible, saving on superfluous computations⁴.

The alternative translation strategy for boolean expressions is *by flow of control*. This scheme is used for conditions in the context of control-flow statements. Intuitively the value of a condition is not represented directly as (the value of) a variable but rather indirectly as the location to which control transfers after evaluating (part of) the condition. This is illustrated in Figure 5.4.

The example statement above should then be translated as

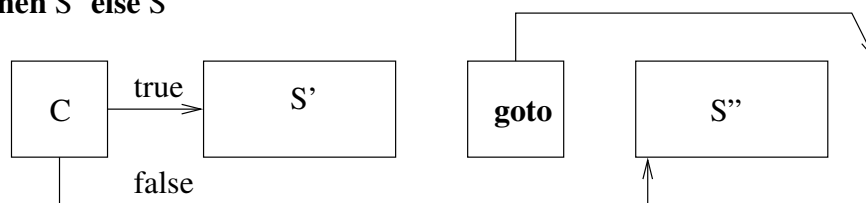
```

1   if a < b goto 4
2   if c < d goto 4
3   goto 6
4   t1 = y + z      /* true exit */
5   x = t1
6                   /* false exit */

```

⁴Note that it depends on the semantics of the source language whether such short-circuit translation of conditions is allowed. E.g. in C, translating `if ((a < b) || ((u = z) < d)) x = y + z;` by value will always execute `u = z` while translating “by flow of control” may not.

if C then S' else S''



while C S

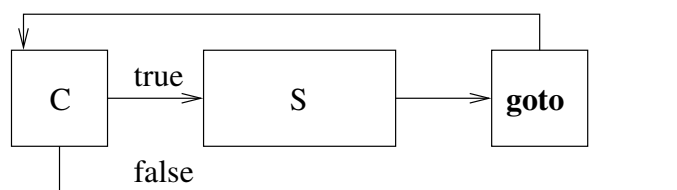


Figure 5.4: Conditional control flow statements

Thus a condition is associated with two locations in the code: a true and false exit. The generated code is such that control will transfer to the true or false exit as soon as possible. However, as can be seen in Figure 5.4, the exits are usually unknown when the condition is parsed, since they refer to locations further in the code array. The solution is to generate a dummy “**goto ?**” instruction wherever a jump to an exit should come. Here “?” is a placeholder for the, as yet unknown, location of the exit. The locations of the dummy jump instructions are kept in two synthesized attributes associated with the condition:

- *bexp.true* contains the set of all “**goto ?**” instructions generated so far where “?” should eventually be replaced by the address of the true exit.
- *bexp.false* contains the set of all “**goto ?**” instructions generated so far where “?” should eventually be replaced by the address of the false exit.

When the exact location *l* of e.g. the true exit becomes known, one can *backpatch* the instructions in *bexp.true*, replacing “?” by *l* in each of them.

The semantics of boolean operations such as disjunction, conjunction and negation can easily be expressed by manipulation of the *true* and *false* sets, as illustrated in Figure 5.5.

The actions are shown below.

```

1
2 cond      : cond AND marker cond {
3             $$.true = $4.true;
```

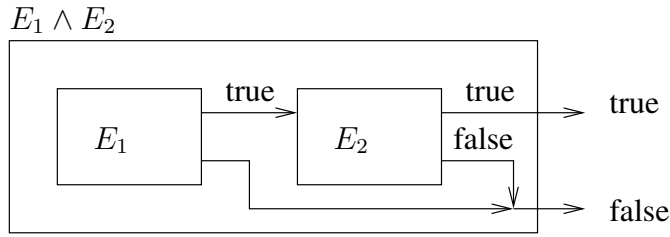


Figure 5.5: Exit sets after “and”

```

4         backpatch($1.true,$3.location);
5         $$false = locations_union($1.false,$4.false);
6     }
7 marker : /* empty right hand side */ {
8         $$location = next3ai();
9     }
10 cond   : NOT cond {
11         $$true = $2.false;
12         $$false = $2.true;
13     }
14 cond   : var {
15         $$true = locations_add(locations_create_set(),next3ai());
16         emit(gen3ai(IFNEQ,$1.place,0,0)); /* true, backpatch later */
17         $$false = locations_add(locations_create_set(),next3ai());
18         emit(gen3ai(GOTO,0,0,0)); /* false, backpatch later */
19     }
20 cond   : exp < exp {
21         $$true = locations_add(locations_create_set(),next3ai());
22         emit(gen3ai(IFLT,$1.place,$3.place,0)); /* true, backpatch later */
23         $$false = locations_add(locations_create_set(),next3ai());
24         emit(gen3ai(GOTO,0,0,0)); /* false, backpatch later */
25     }

```

Note the dummy *marker* $\rightarrow \epsilon$ production: since it is reduced just before code for the second *bexp* will be generated, it can be used to store the location of the start of this code in an attribute *marker.location*.

The following example illustrates the translation of

```
1      p<q || (r<s && t<u)
```

The generated intermediate code is shown below:

```

100  if p < q goto ?t    /* true exit */
101  goto 102
102  if r < s goto 104
103  goto ?f             /* false exit */
104  if y < u goto ?t    /* true exit */
105  goto ?f             /* false exit */

```

5.6 Translating control flow statements

We will use the following example grammar:

```

stmt  →  if cond stmt else stmt
        |  while cond stmt
        |  { slist }
        |  assign
slist →  slist ; stmt
        |  stmt

```

Observe that the code corresponding to a *stmt*, e.g. a **while** statement, may contain forward references to a “next” location (meaning, after the code of the present statement) which is unknown at the time code is generated for the current statement. Therefore we use an attribute *stmt.next* which contains the locations containing instructions that need to be backpatched with the location of the first statement *after* the current one. We will also use dummy nonterminals to mark certain locations in the code. The actions for the above grammar are shown below.

```

1  marker  : /* empty */ {
2           $$ .location = next3ai();
3       }
4
5  goend   : /* empty */ {
6           $$ .next = locations_add(locations_create_set(), next3ai());
7           emit(gen3ai(GOTO, 0, 0, 0)); /* goto end, backpatch later */
8       }
9
10 stmt    : IF cond marker stmt goend ELSE marker stmt {
11         backpatch($2.true, $3.location);
12         backpatch($2.false, $7.location);
13         $$ .next = locations_union(
14             locations_union($4.next, $5.next),
15             $8.next);
16     }
17

```

```

18 stmt      : WHILE marker cond marker stmt {
19             backpatch($3.true,$4.location);
20             backpatch($5.next,$2.location);
21             $$next = $3.false;
22             emit(gen3ai(GOTO,$2.location,0,0));
23             }
24
25 stmt      : assign {
26             $$next = locations_create_set();
27             }
28
29 stmt      : LBRACE slist RBRACE {
30             $$next = $2.next;
31             }
32
33 slist     : stmt {
34             $$next = $1.next;
35             }
36
37 slist     : slist SEMICOLON marker stmt {
38             backpatch($1.next,$3.location);
39             $$next = $4.next;
40             }

```

The following example illustrates the translation of

```

1      while (a<b)
2          if (c<d)
3              x = y + z;

```

The generated intermediate code is shown below:

```

100  if a < b goto 102
101  goto 107
102  if c < d goto 104
103  goto 100          /* false exit */
104  t = y + z
105  x = t
106  goto 100

```

5.7 Translating procedure calls

We will use the following example grammar:

$$\begin{array}{ll}
 \text{stmt} & \rightarrow \text{id}(\text{arg_list}) \\
 \text{arg_list} & \rightarrow \text{arg_list}, \text{exp} \\
 & | \text{exp}
 \end{array}$$

Note that, as in C, we assume that all parameters are passed “by value”. If, as e.g. in Pascal or C++, the source language also supports passing parameters “by address”, extra instructions that take the address of an operand, may need to be inserted before the PARAM instruction.

A statement such as

```
1          f(e1, e2)
```

will be translated as

```

          .. ..          /* code for e1, place t1 */
          .. ..          /* code for e2, place t2 */
102  param t1
103  param t2
104  call f

```

The actions are:

```

1  typedef struct place_list {
2      ...
3      } PLACE_LIST;
4
5  typedef struct {
6      ..
7      } PLACE_LIST_CURSOR;
8
9  /*
10 *      list of places, return new list after successful operation
11 */
12 PLACE_LIST*      place_list_create();
13 void             place_list_destroy(PLACE_LIST*);
14 PLACE_LIST*      place_list_append(PLACE_LIST*, SYM_INFO*);
15
16 /*
17 *      cursor into list of places, return null pointer upon failure.
18 */
19 PLACE_LIST_CURSOR*      place_list_first(PLACE_LIST*);
20 void                     place_list_cursor_destroy(PLACE_LIST_CURSOR*);
21 PLACE_LIST_CURSOR*      place_list_next(PLACE_LIST_CURSOR*);
22 SYM_INFO*               place_list_current(PLACE_LIST_CURSOR*);
23
24 stmt      : ID ( arg_list) {

```

```

25         PLACE_LIST_CURSOR*      c;
26         int                      n = 0; /* number of arguments */
27
28         while (c = place_list.first($3.places)) {
29             ++n;
30             emit(gen3ai(PARAM, place_list_current(c), 0, 0));
31             c = place_list_next(c);
32         }
33         emit(gen3ai(CALL, ID.place, n, 0));
34         /* clean up */
35         place_list_cursor_destroy(c);
36         place_list_destroy($3.places);
37     }
38
39 arg_list: arg_list , exp {
40     $$places = place_list_append($1.places, $3.place);
41 }
42 arg_list: exp {
43     $$places = place_list_append(place_list_create(), $1.place);
44 }

```

5.8 Translating array references

Consider the declaration

```

1 int    a[10,20];
2 int    i,j;
3
4 ..
5 .. a[i,j]

```

We will show how to generate intermediate code for array references such as $a[i, j]$.

We assume that the size of an integer is 4 bytes, that a is the address of the array and that the array is stored in row-major order.

The address of $a[i, j]$ is then

$$a + (20 \times (i - 1) + (j - 1)) \times 4$$

where $20 \times (i - 1)$ is the number of integers in rows appearing before the i 'th row and $(j - 1)$ is the number of integers on the same row but before $a[i, j]$.

The address calculation can be rewritten as

$$a - 21 \times 4 + (20 \times i + j) \times 4$$

Note that $a - 21 \times 4$ is a fixed “base” quantity that only depends on the (declaration of) the array. The remainder $(20 \times i + j) \times 4$ is called the “offset”.

The corresponding intermediate code is shown below.

```

.. ..      /* compute i, place t1 */
.. t2 = t1 * 20 /* 20 × i */
.. ..      /* compute j, place t3 */
102 t4 = t2 + t3 /* (20 × i + j) */
103 t5 = addr a
104 t5 = t5 - 84 /* (a - 21 × 4) */
105 t = t5[t4] /* a[i, j] */

```

In general, we consider an array

$$\text{type } a[d_1, \dots, d_k];$$

where $s = \text{sizeof}(\text{type})$, d_i is the number of elements corresponding to the i 'th dimension ($1 \leq i \leq k$), and a is the address of a .

The address of an element $a[i_1, \dots, i_k]$ is

$$\begin{aligned}
 \text{address}(a[i_1, \dots, i_k]) &= \\
 & a + [(i_1 - 1) \times d_2 \times \dots \times d_k + (i_2 - 1) \times d_3 \times \dots \times d_k + \dots \\
 & \quad + (i_{k-1} - 1) \times d_k + (i_k - 1)] \times s \\
 &= a + [(i_1 \times d_2 \times \dots \times d_k + i_2 \times d_3 \times \dots \times d_k + \dots + i_{k-1} \times d_k + i_k \\
 & \quad - (d_2 \times \dots \times d_k + d_3 \times \dots \times d_k + \dots + d_k + 1)] \times s \\
 &= a - C + [(\dots ((i_1 \times d_2) + i_2) \times d_3 + i_3) \dots] \times d_k + i_k] \times s
 \end{aligned}$$

where

$$C = (d_2 \times \dots \times d_k + d_3 \times \dots \times d_k + \dots + d_k + 1) \times s$$

depends only on (the declaration of) the array.

We will use the following example grammar (an *lvalue* is an expression that can be assigned to):

```

exp      → lvalue
lvalue   → ID /* simple variable */
          | elist ] /* array reference */
elist    → elist , exp
          | ID [ exp

```

For *lvalue*, we use two attributes: *lvalue.place* and *lvalue.offset* representing the base address and the offset, respectively (the offset is only meaningful for array references).

For *elist*, we have the attributes *elist.place*, containing the offset computed so far, *elist.array* containing a reference to the symbol table entry for the array variable and *elist.ndim* representing the number of dimensions seen so far.

The offset will be computed by successively computing (in *elist.place*)

$$\begin{aligned}
 o_1 &= i_1 \\
 o_2 &= o_1 \times d_2 + i_2 \\
 o_3 &= o_2 \times d_3 + i_3 \\
 &\dots \\
 o_k &= o_{k-1} \times d_k + i_k \\
 \text{offset} &= o_k \times s
 \end{aligned}$$

The actions are shown below.

```

1  #define NIL      0x0FFFFFFF
2
3  int      array_limit(SYM_INFO*,int); /* return size of dimension i of array */
4  int      array_base(SYM_INFO*); /* return constant C associated with array */
5  int      array_el_size(SYM_INFO*); /* return size of element of array */
6
7  elist    : ID [ exp {
8              $.place = exp.place; /* first index */
9              $.ndim = 1;
10             $.array = $1.place;
11             }
12  elist    : elist , exp {
13             int limit = array_limit($1.array,$1.ndim+1);
14             $.place = newvar();
15             /* offset(next) = offset(prev)*limit(prev)+index(next) */
16             emit(gen3ai(A2TIMES,$1.place,limit,$$.place));
17             emit(gen3ai(A2PLUS,$$.place,$3.place,$$.place));
18             $.array = $1.array;
19             $.ndim = $1.ndim+1;
20             }
21  lvalue    : elist ] {
22             $.place = newvar();
23             $.offset = newvar();
24             /* base = addr a - array_base(a) */
25             emit(gen3ai(ADDR,$1.array,0,$$.place));

```

```
26         emit(gen3ai(A2MINUS,$$.place,array_base($1.array),$$place);
27         /* offset = elist.offset * sizeof(element) */
28         emit(gen3ai(A2TIMES,$1.place,array_el_size($1.array),$$offset);
29         /* now $$place[$$.offset] references array element */
30     }
31 lvalue   : ID {
32         $$place = $1.place;
33         $$offset = NIL;
34     }
35 exp      : lvalue {
36         if ($1.offset==NIL) { /* note: NIL!=0 */
37             $$place = $1.place
38         }
39         else
40         {
41             $$place = newvar();
42             emit(gen3ai(AAC,$1.place,$1.offset,$$.place));
43         }
44     }
```

Chapter 6

Optimization of intermediate code

6.1 Introduction

In this chapter we look at techniques for improving intermediate (machine-independent) code with respect to two criteria: *time*, i.e. execution time, and *space*, i.e. the space taken up by the intermediate code instructions and variables. Note that optimization is actually a misnomer since there is no guarantee that the resulting code is the best possible.

Usually, improvements in time will also result in more compact code since instructions and variables will be eliminated. This is due to the “technical” nature of the improvements. On a more abstract, e.g. the algorithm, level, the opposite is usually the case in that one can obtain faster algorithms at the cost of more space for data structures. As a simple example, consider the problem of designing a function that takes any ASCII character as input and converts upper case to lower case letters.

A naive “small” algorithm will look something like

```
1 char
2 tolower(char c) /* assume 0<=c<=127 */
3 {
4     if ((c>='A') && (c<='Z'))
5         return c-'A'+'a';
6     else
7         return c;
8 }
```

By storing the function as a table, hence taking up more space, one can obtain a faster algorithm as shown below.

```

1  static char TOLOWER[] = {
2      '\0',    /* 0 */
3      ...
4      'a',     /* 65 ('A') */
5      'b',     /* 66 ('B') */
6      ..
7      '~',     /* 126 ('~') */
8      '\0177' /* 127 */
9  };
10
11 char
12 tolower(char c) /* assume 0 <= c <= 127 */
13 {
14     return TOLOWER[c];
15 }

```

In general, it is worth noting that, if speed is an issue, optimization by the compiler can make a fast algorithm run faster but it cannot make a slow algorithm run fast enough.

In the example below, the running time of a bad algorithm (bubblesort, B) is compared with its optimized version B^* and with a good algorithm (quicksort, Q). The number n is a measure for the size of the problem (i.e. the number of records to be sorted).

| | $n = 10^4$ | $n = 10^5$ |
|-------|------------|------------|
| B | 1 | 100 |
| B^* | .1 | 10 |
| Q | .1 | 1.5 |

As the table shows, while compiler optimization can result in a tenfold speedup, as the size of the problem grows, an unoptimized superior algorithm easily catches up with the optimized version of a slow algorithm. Hence, if speed is an issue, there is no excuse for not designing an efficient algorithm before relying on the compiler for further optimization.

Finally, one should be aware that machine-independent code improvement techniques, such as the ones discussed in this chapter, are only one source of optimization that can be performed by the compiler. Indeed, during the code generation phase, there are additional opportunities for the compiler to dramatically improve program running times, e.g. by judicious allocation of variables to registers.

6.2 Local optimization of basic blocks

A sequence of intermediate code instructions is converted to a *flow graph* where the nodes consist of *basic blocks* as follows.

Definition 16 Let s be a sequence of intermediate code instructions. A **leader instruction** is any instruction $s[i]$ such that

- $i = 0$, i.e. $s[i]$ is the first instruction; or
- $s[i]$ is the target of a (conditional or unconditional) jump instruction; or
- $s[i - 1]$ is a (conditional or unconditional) jump instruction

A **basic block** is a maximal subsequence $s[i], \dots, s[i + k]$ such that $s[i]$ is a leader and none of the instructions $s[i + j]$, $i < j \leq k$ is a leader.

Thus, a basic block is a sequence of consecutive instructions in which flow of control enters at the beginning and leaves at the end.

Definition 17 Let s be a sequence of intermediate code instructions. The **flow graph** of s is a directed graph where the nodes are the basic blocks in s and there is an edge from B_1 to B_2 iff

- The last instruction of B_1 is a (conditional or unconditional) jump instruction to the first instruction of B_2 ; or
- B_2 follows B_1 in s and the last instruction of B_1 is not an unconditional jump instruction.

As an example, consider the following fragment which computes the inner product of two vectors.

```

1 {
2 product = 0;
3 for (i=0; i<20; i=i+1)
4     product = product + a[i]*b[i];
5 }
```

The intermediate code is shown below. It contains two basic blocks, B_1 and B_2 , comprising the instructions 1-2 and 3-16, respectively. The flow graph is shown in Figure 6.1


```

1  product = 0
2  i = 1
3  t1 = 4*i
4  t2 = addr a
5  t3 = t2 - 4
6  t4 = t3[t1]
7  t5 = addr b
8  t6 = t5 - 4
9  t7 = 4*i
10 t8 = t6[t7]
11 t9 = t4 * t8
12 t10 = product + t9
13 product = t10
14 t11 = i+1
15 i = t11
16 if (i<20) goto 3

```

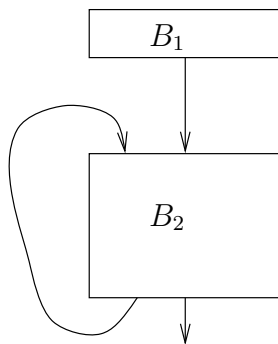


Figure 6.1: Flow graph of example code

6.2.1 DAG representation of basic blocks

A basic block can be represented as a *directed acyclic graph* (dag) which abstracts away from the non-essential ordering of instructions in the block. On the other hand, the dag will be constructed such that expressions are evaluated only once, since they correspond to single nodes in the graph.

Pseudo-code for the algorithm to construct the dag for a basic block is shown below.

Algorithm 8 [Constructing a DAG from a basic block]

```

1  /*
2  *  node(SYM_INFO*) and set_node(SYM_INFO*,NODE*) maintain a mapping:
3  *
4  *      node:  SYM_INFO* -> NODE*
5  *
6  *  which associates a (at most 1) node with a symbol (the node containing
7  *  the "current value" of the symbol).
8  */
9  extern NODE* node(SYM_INFO*); /* get node associated with symbol */
10 extern void set_node(SYM_INFO*,NODE*); /* (re)set node associated with symbol */
11
12 /*
13 *  leaf(SYM_INFO*) and set_leaf(SYM_INFO*,NODE*) maintain a mapping:
14 *
15 *      leaf:  SYM_INFO* -> LEAF*
16 *
17 *  which associates a (at most 1) leaf node with a symbol (the node
18 *  containing the "initial value" of the symbol)
19 *
20 *  set_leaf() is called only once for a symbol, by the leaf
21 *  creating routine new_leaf(SYM_INFO*)
22 */
23 extern NODE* leaf(SYM_INFO*); /* get leaf associated with symbol */
24 extern void set_leaf(SYM_INFO*,NODE*); /* set leaf associated with symbol */
25
26 /*
27 *  node creation functions. Each node has an opcode and a set
28 *  of symbols associated with it. A node may have up to 2 children.
29 */
30 extern new_0_node(SYM_INFO*);
31 extern new_1_node(OPCODE,SYM_INFO*);
32 extern new_2_node(OPCODE,SYM_INFO*, SYM_INFO*);
33 extern node_add_sym(SYM_INFO*,NODE*);
34 /*
35 *  node finding function: returns node with given opcode and node arguments
36 */
37 extern NODE* match(OPCODE,NODE*,NODE*)
38
39 NODE*
40 new_leaf(SYM_INFO* s) {
41     NODE* n;
42     n = new_0_node(s);
43     set_leaf(s,n);
44     set_node(s,n);
45     node_add_sym(s,n);
46     return n;
47 }
48
49 void

```

```

50 basic_block_2_dag(INSTRUCTION bb[],int len) {
51     for (i=0;(i<len);++i) { /* for each instruction in the basic block */
52         switch (type_of(bb[i]) {
53             case BINOP: /* A = B op C */
54                 nb = node(B);
55                 if (!nb) /* B not yet in dag */
56                     nb = new_leaf(B);
57                 nc = node(C);
58                 if (!nc) /* C not yet in dag */
59                     nc = new_leaf(C);
60                 /* find op-node wich children nb,nc */
61                 n = match(op,nb,nc);
62                 if (!n)
63                     n = new_2_node(op,nb,nc);
64                 node_add_sym(A,n); set_node(A,n);
65                 break;
66             case UNOP: /* A = op B */
67                 nb = node(B);
68                 if (!nb) /* B not yet in dag */
69                     nb = new_leaf(B);
70                 /* find op-node wich children nb,nc */
71                 n = match(op,nb,nc);
72                 if (!n)
73                     n = new_1_node(op,nb);
74                 node_add_sym(A,n); set_node(A,n);
75                 break;
76             case ZOP: /* A = B */
77                 n = node(B);
78                 if (!n) /* B not yet in dag */
79                     n = new_leaf(B);
80                 node_add_sym(A,n); set_node(A,n);
81                 break;
82             case CJUMP: /* if B relop C GOTO L */
83                 /* this must be the last instruction in the block! */
84                 nb = node(B);
85                 if (!nb) /* B not yet in dag */
86                     nb = new_leaf(B);
87                 nc = node(C);
88                 if (!nc) /* C not yet in dag */
89                     nc = new_leaf(C);
90                 n = new_2_node(relop,nb,nc);
91                 /* set_node(L,n); */
92                 break;
93             default:
94                 break;
95         }
96     }
97 }

```

□

The result of executing the algorithm on the basic block B_2 from Figure 6.1 is shown in Figure 6.2.

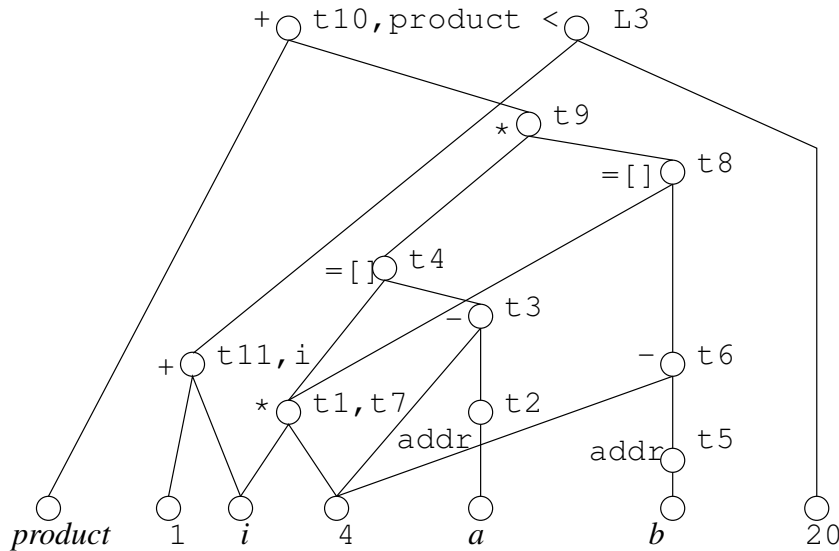


Figure 6.2: Dag of example code

The DAG representation of a basic block provides useful information:

- The range of the *leaf* map is the set of all symbols that are used in the basic block.
- The range of the *node* map contains the set of symbols that are available for use outside the basic block. In the absence of more precise global information, we must assume that they *will* be used and therefore we assume that these variables are *live*¹ at the end of the basic block. However, temporary variables created during intermediate code generation may be assumed not to be live.

Note that a variable available at the end of a basic block is not necessarily “live”. Section 6.3.5 (page 110), contains a global analysis that obtains a more precise smaller estimate of which variables are really “live”.

¹See Definition 24, page 110, for a precise definition of “live variable”.

6.2.2 Code simplification

Algorithm 9 [Code simplification in a basic block]

The DAG can also be used to perform code simplification on the basic block. Basically, intermediate code is generated for each node in the block, subject to the following rules and constraints (we assume that the code does not contain array assignments or pointer manipulations):

- *Code for children should precede the code for the parent. The code for the unique jump instruction should be generated last.*
- *When performing the operation corresponding to a node n we prefer a live target variable from the set $node(n)$.*
- *If $node(n)$ contains several live variables, add additional simple assignments of the form $A = B$ to ensure that all live variables are assigned to.*
- *If $node(n)$ is empty, use a new temporary variable as the target.*
- *Do not assign to a variable v if its current value is still needed (e.g. because its leaf node still has an unevaluated parent).*

□

The result of applying the algorithm sketched above is shown below. Note that, compared with the original version, both the number of instructions and the number of temporary variables has decreased.

```

1
2
3  t1 = 4*i
4  t2 = addr a
5  t3 = t1 -4
6  t4 = t3[t1]
7  t5 = addr b
8  t6 = t5 - 4
9  t8 = t6[t1]
10 t9 = t4 * t8
11 product = product + t9
12 i = i+1
13 if (i<20) goto 3

```

It should be noted that there are heuristics (see e.g Section 9.8 in [ASU86]) to schedule node evaluation of a dag which tend to have a beneficial effect on the number of registers that will be needed when machine code is generated².

²For trees there is even an optimal, w.r.t. register usage, algorithm.

6.2.3 Array and pointer assignments

The presence of array references and pointer manipulations complicates the code simplification algorithm.

Consider the intermediate code

```

1  x = a[i]
2  a[j] = y
3  z = a[i]

```

The corresponding dag is shown in Figure 6.3.

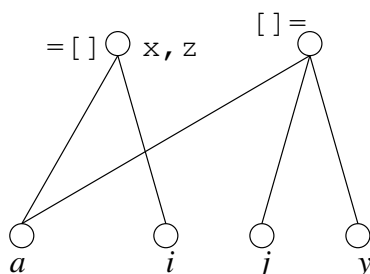


Figure 6.3: Dag of code with array manipulation

A possible simplified instruction sequence would be

```

1  x = a[i]
2  z = x
3  a[j] = y

```

which is clearly not equivalent to the original (consider e.g. what would happen if $i == j$ but $y \neq a[i]$).

The effect of an array assignment ($a[x] = y$) node should be that all nodes referring to an element of a should be *killed* in the sense that no further identifiers can be added to such nodes. Moreover, when generating simplified code, care should be taken that all array reference nodes that existed before the creation of the array assignment node should be evaluated before the assignment node. Similarly, all array reference nodes created after the array assignment node should be evaluated after that node. One way to represent such extra constraints is by adding extra edges to the dag from the “old” array reference nodes to the array assignment node and from that node to the “new” array reference nodes, as shown in Figure 6.4 (the extra constraining edges are shown as dashed arrows).

The effect of a pointer assignment $*a = b$ is even more drastic. Since a can point anywhere, the node corresponding to such an instruction kills every node in the dag so far. Moreover, when generating simplified code, all “old” nodes will

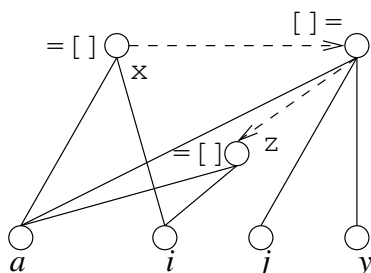


Figure 6.4: Corrected dag of code with array manipulation

have to be evaluated before the pointer assignment and all “new” nodes after it. Similarly, if we assume that a procedure can have arbitrary side effects, the effect of a `CALL` instruction is similar to the one of a pointer assignment.

6.2.4 Algebraic identities

The dag creation algorithm can be made more sophisticated by taking into account certain algebraic identities such as commutativity ($a * b = b * a$) and associativity.

Another possibility when processing a conditional jump instruction on $a > b$ is to check whether there is an identifier x holding $a - b$ and to replace the condition by $x > 0$ which is probably cheaper.

As an example, consider the intermediate code below.

```

1  t1 = b + c
2  a = t1
3  t2 = c + d
4  t3 = t2 + b
5  e = t3

```

The corresponding dags (with and without algebraic optimization) are shown in Figure 6.5. In the example, the first version of the dag is optimized by applying associativity to recognize $b + c$ as a common subexpression, eliminating the need for the temporary variable $t2$.

6.3 Global flow graph information

In this section we discuss several kinds of information that can be derived from the flow graph (and the instructions in it). In the next section, we will design code improving transformations that rely on this information.

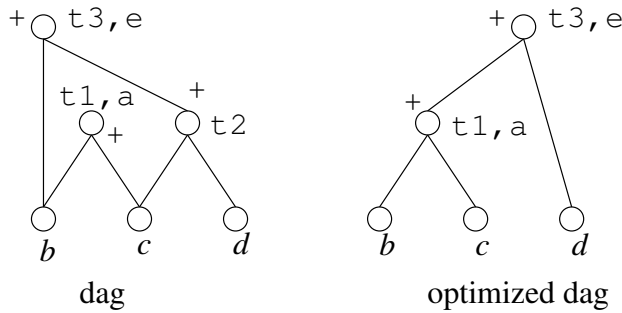


Figure 6.5: Dag using algebraic identities

Note that we assume that a flow graph corresponds to a single procedure. Also, we will not consider problems caused by *aliasing*, where two variables refer to the same location (aliasing may be caused by pointer manipulation and/or reference parameters in procedure calls). The impact of aliasing is briefly discussed in Section 6.5.

In the discussion below we use the following concepts related to flow graphs of basic blocks, which themselves consist of intermediate code instruction sequences.

- A basic block B' is a **successor** of a basic block B , denoted $B < B'$, if there is an edge from B to B' in the flow graph.
- We assume the existence of an **initial basic block** which represents the entry to the procedure.
- There are $n + 1$ **points** in a basic block with n instructions: one between each pair of consecutive instructions, one before the first instruction and one just after the last instruction.
- Consider all points in all basic blocks of a flow graph. A **path** from a point p_0 to a point p_n is a sequence of points $p_0, p_1, \dots, p_{n-1}, p_n$ such that for each $0 \leq i < n$, either
 - p_i is the point immediately before an instruction s and p_{i+1} is the point immediately after s ; or
 - p_i is (after) the end of some block and p_{i+1} is the point before the first instruction in a successor block.

Often, we will confuse an instruction s with the point immediately before it.

6.3.1 Reaching definitions

Definition 18 A *definition* of a variable is an instruction that assigns, or may assign, a value to it. In the first case, we say that a definition is **unambiguous**, otherwise it is called **ambiguous**. The set of all definitions of a variable a is denoted by $Def(a)$, while $UDef(a)$ represents the set of unambiguous definitions of a .

A definition d **reaches** a point p if there is a path from the point immediately following d to p such that d is not **killed** by an unambiguous definition of the same variable along that path.

The phrase “may assign” is necessary since, e.g. a call to a procedure with x as a formal (reference) parameter may or may not result in an assignment to x . A similar reasoning holds for an assignment through a pointer or to an array element.

Note that $Def(a)$ is easy to compute.

Thus, if a definition d for x reaches a point p , the value of x at p may be given by the value assigned at d .

The following basic block properties will be useful later on:

Definition 19 Let B be a basic block.

$$\begin{aligned} GEN(B) &= \{d \in B \mid d \text{ is a definition that reaches the end of } B\} \\ KILL(B) &= \{d \notin B \mid d \in Def(x) \text{ and } B \cap UDef(x) \neq \emptyset\} \\ IN(B) &= \{d \mid d \text{ is a definition and } d \text{ reaches the start of } B\} \\ OUT(B) &= \{d \mid d \text{ is a definition and } d \text{ reaches past the end of } B\} \end{aligned}$$

Observe that $GEN(B)$, which depends only on B , is straightforward to construct using a backward scan of B (after seeing an unambiguous definition for a variable x , no more definitions for x should be added to $GEN(B)$). Similarly, given the set DEF of all definitions in the graph, the set $KILL(B)$ can easily be computed by keeping only those definitions from DEF that define a variable for which an unambiguous definition in B exists.

It is easy to see that $IN(B)$ and $OUT(B)$ must satisfy the following data-flow equations:

$$\begin{aligned} OUT(B) &= (IN(B) \setminus KILL(B)) \cup GEN(B) \\ IN(B) &= \cup_{C < B} OUT(C) \end{aligned}$$

If there are n blocks in the graph, we have $2n$ equations with $2n$ unknowns ($GEN(B)$ and $KILL(B)$ can be determined independently). Due to the possible presence of cycles in the flow graph, these equations usually do not have a unique solution.

E.g. consider the graph in Figure 6.6.

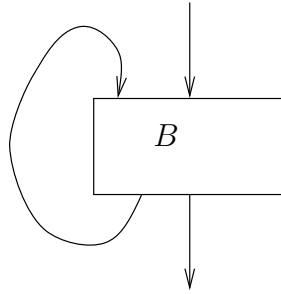


Figure 6.6: Flow graph with cycle

Clearly, if $IN(B) = I_0$ and $OUT(B) = O_0$ form a solution and d is a definition such that $d \notin I_0 \cup O_0 \cup KILL(B)$, then $IN(B) = I_0 \cup \{d\}$ is also a solution. It turns out that the smallest solution is the desired one³. It can be easily computed using the following (smallest) fixpoint computation.

Algorithm 10 [Computing definitions available at a basic block]

```

1  typedef  set<Instruction*>  ISET;
2  typedef  ..          BLOCK;
3
4  ISET  in[BLOCK], out[BLOCK], kill[BLOCK], gen[BLOCK];
5
6  void
7  compute_in_out() {
8      for all b in BLOCK {
9          in[b] = emptyset;
10         out[b] = gen[b];
11     }
12
13     bool change = false;
14
15     do {
16         change = false;
17         foreach b in BLOCK {
```

³The solutions are closed under intersection

```

18     newin = set_union { out[c] | c < b };
19     change = change || (newin != in[b]);
20     in[b] = newin;
21     out[b] = (in[b] - kill[b]) set_union gen[b]
22 }
23 while change;
24 }
25 }

```

□

Definition 20 Let a be a variable and let u be a point in the flow graph. The set of definitions of a that reach u is defined by

$$UD(a, u) = \{d \in Def(a) \mid d \text{ reaches } u\}$$

Such a set is often called a **use-definition chain**.

$UD(a, u)$ can be easily computed using $IN(B)$ where u is the basic block containing u .

Algorithm 11 [Computing use-definition chains]

```

1  iset
2  defu(SYM_INFO* a, POINT u) {
3      BLOCK b = block_of(u);
4
5      if (Def(a) intersection B != emptyset)
6          /* return the singleton containing the last
7           * definition in b of a before u
8           */
9          for (p=u; (p>0); --p)
10             if (p in UDEF(a))
11                 return {p}
12     return Def(a) intersect in[b]
13 }

```

□

6.3.2 Reaching definitions using datalog

An alternative definition uses logic programming (in particular, datalog) to define reaching definitions, see [ALSU07], page 925.

We will use a term $B.N$, $N \geq 0$ to denote the point in block B us before the $N + 1$ 'th statement (in block B). The latter statement is said to be at point $B.N$.

We start with the following base predicates.

```

1 assign(B.N, X, E) % the statement following B.N has the
2 % form  $X = E$ ,  $E$  an expression (with at most 2 operands)
3 use(B.N, X) % the statement at B.N uses the variable  $X$ 
4 type(X, T) % type type of  $X$  is  $T$ 
5 size(B, N) % block  $B$  has  $N$  statements
6 succ(B, C) % true iff block  $B$  has a successor block  $C$ 
7 initial(B) :- not( succ(C, B) ). % initial block has no predecessors

```

We can then define what it means for a statement to (ambiguously or unambiguously) define a variable X (here we assume a strongly typed language without casting).

```

1 % unambig_def(B.N, X) -- statement  $N$  in block  $B$  unambiguously defines variable  $X$ 
2 unambig_def(B.N, X) :- assign(B.N, X, _).
3
4 % def(B.N, X) -- statement  $N$  in block  $B$  may define variable  $X$ 
5 define(B.N, X) :- unambig_def(B.N, X).
6 define(B.N, X) :- assign(B.N, *P, _), type(X, T), type(P, T*).

```

Then the following program computes the reaching definitions.

```

1 % rd(B.N, C.M, X) -- definition of  $X$  at  $C.M$  reaches  $B.N$ 
2 rd(B.N, B.N, X) :- define(B.N, X).
3 rd(B.N, C.M, X) :- rd(B.N-1, C.M, X), not(unambig_def(B.N-1, X)).
4 rd(B.0, C.M, X) :- rd(D.N, C.M, X), succ(D, B), size(D, N).

```

The first rule simply states that a definition reaches itself while the second rule indicates that if a definition reaches a point $N - 1$ in a block B and the next statement does not unambiguously (re)define X , then the definition also reaches point N (after the next statement) in B . The third rule says that a definition reaches the beginning of a block B if it reaches the end of a predecessor D of B .

A use-definition predicate be defined as follows, where the use u is represented by a point N in a block B , and d by its block C and point M):

```

ud(A, u(B.N), d(C.M)) :- rd(B.N, C.M, A), use(B.N, A).

```

6.3.3 Available expressions

Definition 21 An expression $e = x + y$ ⁴ is **available** at a point p , denoted $AVAIL(e, p)$ if every path from the initial basic block to p evaluates $z = x + y$, and after the

⁴We use '+' as an example operator. The same definition applies to other binary or unary operators.

last such evaluation prior to p , there are no subsequent definitions for x , y or z . A block **kills** an expression $x + y$ if it defines x or y and does not subsequently (re)compute $x + y$. A block **generates** an expression $x + y$ if it unambiguously evaluates $x + y$ and does not subsequently redefine x or y .

In order to be able to determine the availability of an expression at some point, we define

Definition 22 Let B be a basic block.

$$\begin{aligned} GEN_e(B) &= \{e \mid B \text{ generates } e\} \\ KILL_e(B) &= \{e \mid B \text{ kills } e\} \\ IN_e(B) &= \{e \mid e \text{ is available at the start of } B\} \\ OUT_e(B) &= \{e \mid e \text{ is available past the end of } B\} \end{aligned}$$

Observe that $GEN_e(B)$, which depends only on B , is straightforward to construct using a forward scan of B , keeping track of available expressions by appropriately processing definitions. E.g. an instruction $z = x + y$ makes $x + y$ available but kills any available expressions involving z (note that $z = x$ is possible). The set $KILL_e(B)$ consists of all expressions $y + z$ such that either y or z is defined in the block and $y + z$ is not generated by the block.

It is easy to see that $IN_e(B)$ and $OUT_e(B)$ must satisfy the following data-flow equations:

$$\begin{aligned} OUT_e(B) &= (IN_e(B) \setminus KILL_e(B)) \cup GEN_e(B) \\ IN_e(B) &= \begin{cases} \cap_{C < B} OUT_e(C) & \text{if } B \text{ is not initial} \\ \emptyset & \text{if } B \text{ is initial} \end{cases} \end{aligned}$$

Note the similarity of these equations with the ones for reaching definitions. The main difference is the use of \cap instead of \cup in the equation for $IN_e(B)$ (\cap is necessary since e must be available on *all* paths into B in order to be sure that e is available at the start of B).

Consider the flow graph in Figure 6.7 and assume that neither x nor y are defined in block B after d which is itself the final definition of z in block B . Assume furthermore that block C does not contain any definitions for x , y or z . Hence $x + y$ should be available at the beginning of block C , i.e. $x + y \in IN_e(C)$, which would not be found if the minimal solution for the above equations were computed.

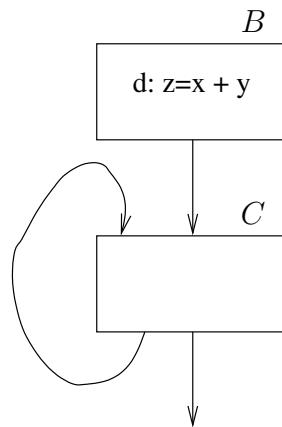


Figure 6.7: Example flow graph for available expressions

Since, clearly, the solutions for $IN_e(C)$ should be closed under union, it turns out that the intended solution is the maximal one. Hence, in the following fixpoint procedure, we compute the greatest fixpoint.

Algorithm 12 [Computing expressions available at a basic block]

```

1  typedef .. EXPRESSION;
2  typedef set<EXPRESSION> ESET;
3  typedef .. BLOCK;
4
5  extern BLOCK* initial_block();
6  extern ESET all_expressions();
7
8  ESET ine[BLOCK], oute[BLOCK], kille[BLOCK], gene[BLOCK];
9
10 void
11 compute_avail_exp() {
12     BLOCK b1 = initial_block();
13
14     /* initialization */
15
16     ine[b1] = empyset; /* never changes */
17     oute[b1] = gene[b1]; /* never changes */
18
19     for all b in (BLOCK - {b1}) {
20         oute[b] = all_expressions() - kille[b]
21     }
22
23     bool change = false;

```

```

24
25  do {
26    change = false;
27    for all b in (BLOCK - {b1}) {
28      ine[b] = intersection { oute[c] | c < b }
29      oldout = oute[b];
30      oute[b] = gene[b] set_union (ine[b] - kille[b]);
31      change = change || (oute[b] != oldout);
32    }
33  }
34  while (change);
35
36 }

```

□

Note that, based on $IN_e(B)$, it is easy to determine $AVAIL(e = x + y, p)$, for $p \in B$.

Algorithm 13 [Computing expression availability at a point]

Let $p \in B$ be a point where we want to compute $AVAIL(e, p)$.

- Check whether $e \in IN_e(B)$.
- If so, check whether x or y are redefined before p in B . If not, e is available; else not.
- If $e \notin IN_e(B)$, check whether e is generated before p in B . If so, e is available; else not.

□

6.3.4 Available expressions using datalog

In datalog, the formulation is slightly more involved than might be expected because the usual datalog semantics computes minimal solutions. The approach is to define a predicate `unavail_exp/2` asserting that an expression is not available at some point. The availability of an expression at a point then holds if `unavail` does not hold at that point.

```

1  define_exp(B.N, X+Y) :- assign(B.N, _, X+Y).
2  % a possibled definition of any operand kills the expression
3  kill_exp(B.N, X+Y) :- define(B.N, X, _).
4  kill_exp(B.N, X+Y) :- define(B.N, Y, _).

```

```

5
6 unavail_exp(B.N, X+Y) :- kill_exp(B.N
7 unavail_exp(B.N+1, X+Y) :-
8   unavail_exp(B.N, X+Y), not(define_exp(B.N+1, X+Y)).
9 % if X+Y is unavaible at the end of a predecessor of B,
10 % it is also unavailable at the beginning of B
11 unavail_exp(B.0, X+Y) :- succ(C, B), size(C,N), unavail_exp(C.N, X+Y).
12 % at the start of the initial block, all expressions are unavailable
13 unavail_exp(B.0, X+Y) :- initial(B).
14 % an expression is available if it is not unavailable
15 avail(B.N, X+Y) :- not(unavail_exp(B.N, X+Y)).

```

6.3.5 Live variable analysis

Definition 23 A variable x is *used* at instruction s if its value may be required by s .

Definition 24 Let x be a variable and let p be a point in the flow graph. We say that x is *live* at p if there is some path from p on which the value of x at p can be used. If x is not live at p , we say that it is *dead* at p .

In datalog this becomes:

```

1 live(B.N, X) :- assign(B.N, _, X + Y).
2 live(B.N, X) :- assign(B.N, _, Y + X).
3 live(B.N-1, X) :- live(B.N, X),
4   not(unambig_def(B.N-1, X)).
5 live(B.N, X) :- size(B, N), succ(B, C), live(C.0, X).

```

In order to be able to determine whether a variable is live at some point, we define

Definition 25 Let B be a basic block.

$$\begin{aligned}
 KILL_{live}(B) &= \{x \mid x \text{ is unambiguously defined before used in } B\} \\
 GEN_{live}(B) &= \{x \mid x \text{ may be used in } B \text{ before it is unambiguously defined in } B\} \\
 IN_{live}(B) &= \{x \mid x \text{ is live at the start of } B\} \\
 OUT_{live}(B) &= \{x \mid x \text{ is live past the end of } B\}
 \end{aligned}$$

Note that, in the definition of $KILL_{live}(B)$, it is not necessary that x is actually used in B , as long as there is an unambiguous definition of x before any use. Similarly, in the definition of $IN_{live}(B)$, it may well be that $x \in IN_{live}(B)$ although x is not (unambiguously) defined in B .

Observe that both $KILL_{live}(B)$ and $GEN_{live}(B)$, which depend only on B , are straightforward to construct.

To compute $IN_{live}(B)$ and $OUT_{live}(B)$, we note that they satisfy the following data-flow equations:

$$\begin{aligned} IN_{live}(B) &= GEN_{live}(B) \cup (OUT_{live}(B) - KILL_{live}(B)) \\ OUT_{live}(B) &= \cup_{B < C} IN_{live}(C) \end{aligned}$$

Again we will use a fixpoint computation to determine the smallest solution of the above equations. Note that intuitively, the algorithm works “backwards” in the graph, as can be seen from the use of successors instead of predecessors in the equation for $OUT_{live}(B)$.

Algorithm 14 [Computing live variables for a basic block]

```

1  typedef set<SYM_INFO*>  VSET;
2  typedef ..               BLOCK;
3
4  VSET    in_live[BLOCK], out_live[BLOCK], use[BLOCK], gen_live[BLOCK];
5
6  void
7  compute_live()
8  {
9    for all  $b$  in BLOCK
10       in_live[b] = emptyset;
11
12    bool change = true;
13
14    while (change) {
15       change = false;
16       foreach  $b$  in BLOCK {
17          oldin = in_live[b];
18          out_live[b] = union { in[s] |  $b < s$  };
19          in_live[b] = use[b] union (out_live[b] - gen_live[b]);
20          change = change || (oldin != in_live[b]);
21       }
22    }
23 }
```

□

Note that the information from $OUT_{live}(B)$ can be profitably used to refine Algorithm 9.

6.3.6 Definition-use chaining

Definition 26 *The definition-use chain of a definition d is defined as*

$$DU(d, x) = \{s \mid s \text{ uses } x \text{ and there is a path from } d \text{ to } s \text{ that does not redefine } x\}$$

The following sets, defined for basic blocks, facilitate the computation of $DU(d, x)$.

Definition 27 *Let B be a basic block.*

$$\begin{aligned} KILL_{use}(B) &= \{(s, x) \mid s \notin B \text{ uses } x \text{ and } B \text{ defines } x\} \\ GEN_{use}(B) &= \{(s, x) \mid s \in B \text{ uses } x \text{ and } x \text{ is not defined prior to } s \text{ in } B\} \\ IN_{use}(B) &= \{(s, x) \mid s \text{ uses } x \text{ and } s \text{ reachable from the beginning of } B\} \\ OUT_{use}(B) &= \{(s, x) \mid s \notin B \text{ uses } x \text{ and } s \text{ reachable from the end of } B\} \end{aligned}$$

Here, “reachable” means without any intervening definitions of x .

Observe that both $KILL_{use}(B)$ and $GEN_{use}(B)$ are straightforward to construct.

To compute $IN_{use}(B)$ and $OUT_{use}(B)$, we note that they satisfy the following data-flow equations:

$$\begin{aligned} IN_{use}(B) &= GEN_{use}(B) \cup (OUT_{use}(B) - KILL_{use}(B)) \\ OUT_{use}(B) &= \cup_{B < C} IN_{use}(C) \end{aligned}$$

Since the above equations are isomorphic to the ones in Section 6.3.5, an algorithm similar to Algorithm 14 can be used to compute $OUT_{use}(B)$. $DU(d, x)$ can then be computed by.

$$DU(d, x) = \begin{cases} OUT_{use}(B) \cup \{(s, x) \mid s \in B \text{ uses } x\} & \text{if } B \text{ does not contain a definition of } x \text{ after } d \\ \{(s, x) \mid s \in B \text{ uses } x \text{ and } s \text{ comes before } d'\} & \text{if } d' \in B \text{ is the first definition of } x \text{ after } d \end{cases}$$

6.3.7 Application: uninitialized variables

To illustrate the use of global flow information, we show how to detect uninitialized variables.

It suffices to add a dummy assignment $d_x:x = 0$ for each variable x in front of the initial block. Now if for any such d_x , it turns out that $DU(d_x, x) \neq \emptyset$, we can deduce that x is possibly used before it is assigned to.

6.4 Global optimization

6.4.1 Elimination of global common subexpressions

Algorithm 15 [Elimination of global common subexpressions]

For every instruction $s:x = y + z$ ⁵ such that $AVAIL(y + z, p)$ where p is the point before s , do the following (assume $s \in B$):

1. *Find the evaluations of $y + z$ that reach B by searching the graph (and each block) backwards from B , not continuing after a block that evaluates $y + z$.*
2. *Create a new temporary variable u .*
3. *Replace each instruction $w = y + z$ found in step 1 by $u = y + z; w = u$.*
4. *Replace s by $x = u$.*

□

Note that the algorithm will miss the fact that $a * z$ is the same as $c * z$ in

```

1  a = x + y
2  b = a * z
3  ...
4  c = x + y
5  d = c * z

```

Figure 6.8: Code before first subexpression elimination

Algorithm 15 will produce the code in Figure 6.9. However, repeated application of Algorithm 15, combined with copy propagation (Algorithm 16), will eventually optimize the code further.

⁵Here '+' stands for any operator.

```

1  u = x + y
2  a = u
3  b = a * z
4  ...
5  c = u
6  d = c * z

```

Figure 6.9: Code after first subexpression elimination

6.4.2 Copy propagation

Copy instructions of the form $x = y$ are generated e.g. by intermediate code generation and by Algorithm 15. Copy propagation can be used to eliminate many such instructions by replacing further uses of x by y .

One can eliminate a copy instruction $s:x = y$ and substitute y for x in all uses of x in $s' \in DU(s, x)$ provided that:

- $UD(s', x) = \{s\}$, i.e. s is the only definition of x to reach s' .
- No paths from s to s' , including paths that go through s' several times, contains definitions of y .

In datalog:

```

1  % copy(B.N, X, Y) -- X is copy of Y at B.N
2  copy(B.N, X, Y) :- assign(B.N, X, Y) .
3  copy(B.N, X, Y) :- copy(B.N-1, X, Y) ,
4      not ( def(B.N-1, X) ), not ( def(B.N-1, Y) ) .
5  copy(B.0, X, Y) :- not ( initial(B) ) ,
6      not (
7          succ(C, B) , size(C, M) , not ( copy(C.M, X, Y) )
8      ) .

```

To solve the latter question, we solve yet another data-flow problem where $IN_{copy}(B)$ ($OUT_{copy}(B)$) is the set of copies $s : x = y$ such that *every* path from the initial node to the beginning (end) of B contains s , and subsequent to the last occurrence of s , there are no definitions of y . We also consider $GEN_{copy}(B)$, the set of copies $s:x = y$ generated by B where $s \in B$ and there are no subsequent definitions of x or y within B . $KILL_{copy}(B)$ contains the set of copies $s:x = y$ where either x or y is defined in B and $s \notin B$. Note that $IN_{copy}(B)$ can contain only one copy instruction with x on the left since, otherwise, both instructions would kill each other.

Then the following equations hold:

$$\begin{aligned} OUT_{copy}(B) &= GEN_{copy}(B) \cup (IN_{copy}(B) \setminus KILL_{copy}(B)) \\ IN_{copy}(B) &= \begin{cases} \cap_{C < B} OUT_{copy}(C) & \text{if } B \text{ is not initial} \\ \emptyset & \text{if } B \text{ is initial} \end{cases} \end{aligned}$$

Note that these equations are isomorphic to the ones used for Algorithm 12 and therefore we can use a similar algorithm to solve them.

Given $IN_{copy}(B)$, i.e. the set of copies $x = y$ that reach B along every path, with no definition of x or y following the last occurrence of $x = y$ on the path, we can propagate copies as follows.

Algorithm 16 [Copy propagation]

For each copy instruction $s:x = y$ do the following:

1. *Determine $DU(s, x)$, i.e. the set of uses s' that are reached by s .*
2. *Check whether for every $s' \in DU(s, x)$, $s \in IN_{copy}(B_{s'})$ ⁶, i.e. s reaches $B_{s'}$ on every path and, moreover, no definitions of x or y occur prior to s' in B .*
3. *If the previous step succeeds, remove s and replace x by y in each $s' \in DU(s, x)$.*

□

Applying Algorithm 16 on the code of Figure 6.9 yields

```

1  u = x + y
2  b = u * z
3  ...
4  d = u * z
```

Now, running Algorithm 15 again will result in

```

1  u = x + y
2  v = u * z
3  b = v
4  ...
5  d = v
```

⁶ $B_{s'}$ is the basic block containing s' .

6.4.3 Constant folding and elimination of useless variables

Another transformation concerns the propagation of constants. As an example, consider the code

```

1  t0 = 33
2  t1 = 3
3  t2 = t0 + t1
4  t3 = t2 - 35
5  x = t3;
```

Algorithm 17 [Constant folding]

Call a definition $s:x = a + b$ ⁷ “constant” if all its operands a and b are constant.

- *We can then replace s by $x = c$ where $c = a + b$ is a new constant.*
- *For each constant definition s and for each $s' \in DU(x, s)$, if $UD(x, s') = \{s\}$, i.e. s is the only x -definition reaching s' , we can replace x in s' by c , possibly making s' constant. If after this, $DU(x, s) = \emptyset$, we can eliminate s .*
- *Repeat the above steps until there are no further changes.*

□

Note that the second step above can be seen as a special simple case of copy propagation since the right hand side c of the copy instruction $x = c$ is constant and therefore there are no c -assignments to worry about.

In the example, application of the algorithm sketched above yields

```

1  x = 1
```

6.4.4 Loops

It is a well-known piece of computer science folklore that programs tend to spend 90% of their execution time in 10% of their code. One can use a profiling tool such as *gprof* to collect statistics on a program run from which it can be deduced which routines and instructions are most often executed. These “hot spots” will almost always turn out to be loops. So it is natural to concentrate optimization efforts on loops.

Before considering loop optimization, we need to define what constitutes a loop in a flow graph.

⁷Here $+$ may be any operator.

Definition 28 Let G be a flow graph with an initial basic block node b_0 .

- A basic block b **dominates** another block b' iff every path from b_0 to b' passes through b .
- A **back edge** is an edge $n \rightarrow b$ such that b dominates n .
- The **(natural) loop** L of a back edge $n \rightarrow h$ consists of h and all basic blocks x such that there is a path from x to n not going through h ; h is called the **header** of the loop.
- An **inner loop** is a loop that contains no other loops.
- An **exit node** of a loop L is a block $e \in L$ that has a successor outside L .

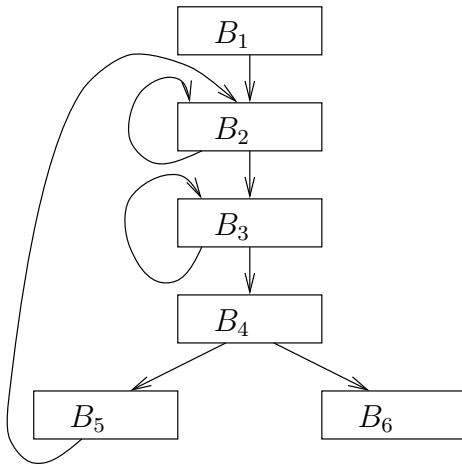


Figure 6.10: Flow graph with loops

In Figure 6.10, there are three loops: $L_0 = \{B_2\}$, $L_1 = \{B_3\}$, $L_2 = \{B_2, B_3, B_4, B_5\}$. B_2 is the header of L_2 , while B_4 is L_2 's single exit node.

One crude way to compute the dominators of a block b is to take the intersection of all acyclic paths from the initial node to b .

A smarter algorithm relies on the fact that

$$\text{dom}(b) = \{b\} \cup \bigcap_{c < b} \text{dom}(c)$$

i.e. a node that dominates all predecessors of b also dominates b .

Algorithm 18 [Computing dominators]

```

1  typedef set<BLOCK*>      BLOCKS;
2
3  BLOCK*          b0; /* initial basic block */
4  set<BLOCK*>      B; /* set of all blocks */
5
6  void
7  compute_dominates()
8  {
9  /* initialize */
10
11  dom(b0) = { b0 };
12  for b in (B - {b0})
13      dom(b) = B;
14
15  /* compute fixpoint */
16
17  do
18      {
19      change = false;
20      for b in (B-{b0})
21          {
22              old_dom = dom(b);
23              dom(b) = {b} union intersection { dom(c) | c < b }
24              change = change || (dom(b) != old_dom);
25          }
26      }
27  while (change)
28  }
```

□

The table below shows the application of Algorithm 19 to the flow graph of Figure 6.10.

| block | 1 | 2 | 3 | 4 | 5 | 6 |
|------------------------|---|-------------|-------------|-------------|-------------|-------------|
| predecessors | | 1,2,5 | 2,3 | 3 | 4 | 4 |
| dominators initialize | 1 | 1,2,3,4,5,6 | 1,2,3,4,5,6 | 1,2,3,4,5,6 | 1,2,3,4,5,6 | 1,2,3,4,5,6 |
| dominators iteration 1 | 1 | 1,2 | 1,2,3 | 1,2,3,4 | 1,2,3,4,5 | 1,2,3,4,6 |
| dominators iteration 2 | 1 | 1,2 | 1,2,3 | 1,2,3,4 | 1,2,3,4,5 | 1,2,3,4,6 |

The following algorithm can be used to compute the natural loop corresponding to a back edge. It corresponds to a backward scan from the header.

Algorithm 19 [Computing the natural loop of a back edge]


```

1  type      set<BLOCK*>      BLOCKS;
2
3  typedef struct {
4      BLOCKS  nodes;
5      BLOCK*  header;
6      } LOOP;
7
8  BLOCKS
9  predecessors(BLOCK *b); /* return predecessors of b in flow graph */
10
11 LOOP
12 natural_loop(BLOCK* src,BLOCK* trg)
13 {
14  /* compute natural loop corresponding to back edge src-->trg */
15
16  BLOCK*  stack[];
17  BLOCK*  tos = stack;
18  LOOP    loop = { { trg }, trg }; /* trg is header */
19
20  if (trg!=src) {
21      loop.nodes = loop.nodes union { src }
22      stack[tos++] = src; /* push src */
23  }
24
25  while (tos!=stack) { /* stack not empty */
26      BLOCK*  b = stack[--tos]; /* pop b */
27      BLOCKS  P = predecessors(b);
28
29      for each p in P
30          if (!(p in loop)) {
31              loop.nodes = loop.nodes union {p};
32              stack[tos++] = p;
33          }
34
35  }
36 }

```

□

Applying Algorithm 19 to the back edge $B_5 \rightarrow B_2$ in the flow graph of Figure 6.10 is illustrated in the table below.

| faze | loop.nodes | stack |
|-------------|--------------|-------|
| initialize | {2} | |
| adding src | {2, 5} | 5 |
| iteration 1 | {2, 5, 4} | 4 |
| iteration 2 | {2, 5, 4, 3} | 3 |
| iteration 3 | {2, 5, 4, 3} | |

6.4.5 Moving loop invariants

A *loop invariant* is an instruction of the form $s:x = y + z$ ⁸ such that x does not change as long as control stays in the loop. Such instructions can, under certain conditions, be moved outside of the loop so that the value x is computed only once.

Example 20 Consider a C function to compute the scalar product of two vectors (represented as arrays).

```

1  int
2  innerproduct(int a[], int b[], int s) {
3  // s is size of arrays
4  int r = 0;
5  int i = 0;
6  while (i < s) {
7      r = r + a[i]*b[i];
8      i = i+1;
9  }
10 return r;
11 }
```

The flow graph, after code simplification, is shown in Figure 6.11 on page 121. Note that the size of an `int` is assumed to be 4 bytes, as witnessed by the offset computation $t1 = 4 * i$.

Applying Algorithm 19 on the graph of Figure 6.11 yields that, among others, C dominates A , B and D and that, consequently, $D \rightarrow C$ is a back edge. It then follows from Algorithm 19 that $L = \{D, C\}$ is an inner loop with exit node C .

Algorithm 20 [Computing loop invariants]

Given a loop L , we find the list of loop invariants from L as follows (note that the order of the list is important, as it reflects the order of the instructions after having been moved):

1. Mark “invariant” (and append to the list of invariants) all instructions $s:x = y + z$ such that each operand y (or z) is either constant or is such that $UD(y, s) \cap L = \emptyset$, i.e. all its reaching definitions are outside L .
2. Repeat the following step until no more instructions are marked.
3. Mark (and append to the list) “invariant” all instructions $s:x = y + z$ such that each operand y (or z) is

⁸Here $+$ stands for any operator.

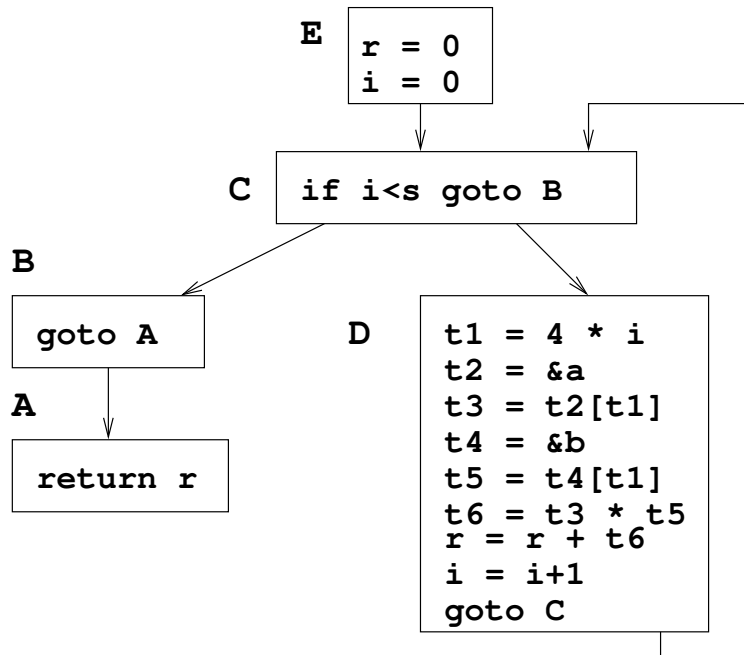


Figure 6.11: Flow graph of Example 20

- *constant; or*
- *all its reaching definitions are outside L ; or*
- *$UD(y, s) = \{q\}$ where $q \in L$ is marked invariant; i.e. the only reaching definition of y is an invariant in the loop.*

□

Applying Algorithm 20 for Example 20 yields two loop invariants: $t2 = \&a$ and $t4 = \&b$, both in block D .

Having found the invariants in a loop, we can move some of them outside the loop. They will be put in a new node, called a *preheader* which has as its single successor the header of the loop. The flow graph is further updated such that all edges that pointed from outside the loop to the header, will be made to point to the new preheader node.

Algorithm 21 [Moving invariants]

Let I be the list of invariants of the loop L .

1. For each $s: x = y + z$ from I , check that:

- (a) s is in a block b that dominates all exits of L or $x \notin IN_{live}(e)$ for any successor $e \notin L$ of an exit block of L ; and
 - (b) x is not defined anywhere else in L ; and
 - (c) For each use s' of x in L , $UD(x, s') \cap L = \{s\}$, i.e. all uses of x in L can only be reached by the definition s of x
2. Move instructions from I that satisfy the above test, and that do not depend on definitions in I that do not pass the test, to the preheader, in the same order as they appear in I .

□

Figure 6.12 shows the result of applying Algorithm 21 to Example 20. Note that the block (D) containing $t2 = \&a$ does not dominate all exits of L . However, since $t2 \notin IN_{live}(B)$, condition (1a) above is still satisfied.

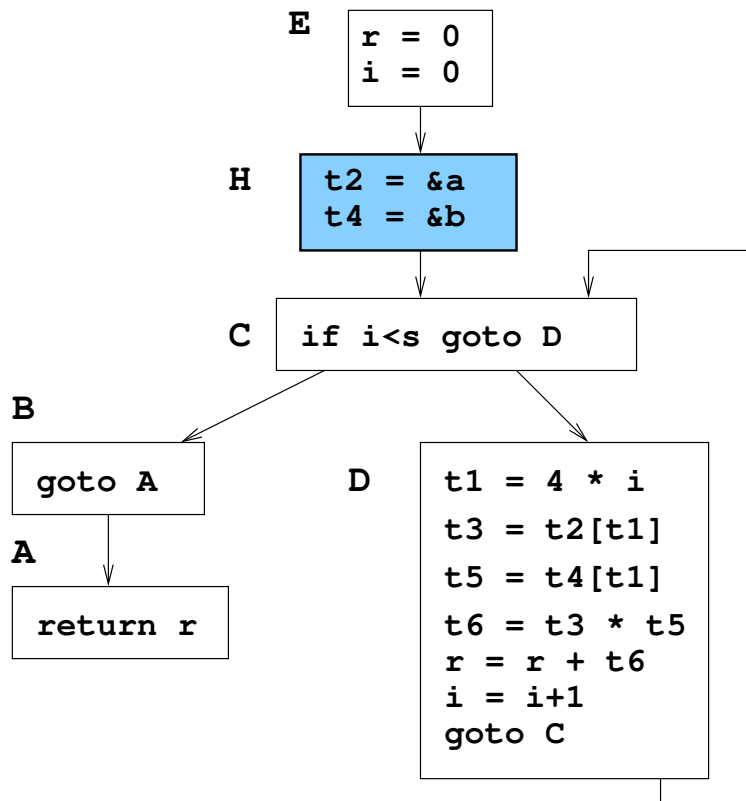


Figure 6.12: Flow graph of Example 20 after moving loop invariants

6.4.6 Loop induction variables

A variable x is called an *induction variable* of a loop L if every time it changes value, it is incremented or decremented by some constant. Examples include array indexes that are incremented on each pass through the loop, counters corresponding to *for* instructions etc.

Often, one induction variable depends on another one.

E.g. in Example 20, both i and $t1$ are induction variables where $t1$ depends on i .

Definition 29 Let L be a loop. A **basic induction variable** is a variable i such that the only assignments to i in L are of the form $i = i \pm c$ where c is a constant. An **induction variable** is a variable j such that j is defined only once in L and its value is a linear function of some basic induction variable i , i.e. $j = a \times i + b$ for some constants a, b . We say that j belongs to the **family** of i , and associate the triple (i, a, b) with j .

The following algorithm determines the induction variables of a loop.

Algorithm 22 [Finding the induction variables of a loop]

Given a loop L , we proceed as follows:

1. Find all basic induction variables i . Such variables are associated with the triple $(i, 1, 0)$.
2. Find variables k such that $\text{Def}(k) \cap L$ is a singleton s which has one of the forms $k = j \times b$, $k = b \times j$, $k = j/b$, $k = b/j$, $k = j \pm b$, or $k = b \pm j$, where b is a constant and j is an induction variable, basic or otherwise.

- If j is basic, then k is in the family of j and we associate k with the appropriate triple. E.g. if $s: k = j \times b$ then the triple of k is $(j, b, 0)$.
- If j is not basic, then j belongs to the family of some basic induction variable i . We first check the following additional conditions:
 - There is no definition of i between the (single) definition of j in L and the (single) definition of k ; and
 - $(\text{UD}(j, s) \setminus L) = \emptyset$, i.e. no definition of j outside L reaches s .

Provided these conditions are met, we add k to the family of i and associate it with the triple $(i, b \times c, b \times d)$, where we assume that s is of the form $k = j \times b$ and the triple of j is (i, c, d) .

□

Applying Algorithm 22 to the loop $L = \{C, D\}$ in Figure 6.12 yields:

- a basic induction variable \mathfrak{i} which is associated with the tuple $(i, 1, 0)$.
- an induction variable $\mathfrak{t}1$ in the family of \mathfrak{i} . $\mathfrak{t}1$ is associated with the tuple $(i, 4, 0)$.

Once we have found suitable induction variables, we can perform strength reduction, where expensive multiplications are replaced by cheaper additions, as in the following algorithm

Algorithm 23 [Strength reduction]

Given a loop L and a set of induction variables, each with associated triple, we consider each basic induction variable i in turn.

For each non-basic induction variable j with associated triple (i, c, d) we do the following:

1. *Create a new variable t (but if two variables j_1 and j_2 have the same triples, create only one variable for both of them).*
2. *Replace the unique assignment to j in L by $j = t$.*
3. *Immediately after each assignment $i = i + a$ (a constant) to the basic induction variable i in L , add an instruction $t = t + c \times a$ and place t in the family of i with triple (i, c, d) .*
4. *To properly initialize t , put*

$$\begin{aligned} t &= c \times i \\ t &= t + d \end{aligned}$$

at the end of the preheader of L .

□

The result of applying Algorithm 23 to the loop $L = \{C, D\}$ in Figure 6.12 is shown in Figure 6.13.

Block D can be further simplified using Algorithm 9 (code simplification): because $\mathfrak{t}1 \notin OUT_{live}(D)$, $\mathfrak{t}1$ will be eliminated. On the other hand, block B can be simplified using copy propagation and constant folding (Algorithms 16 and 17). The result is shown in Figure 6.14.

After strength reduction, it may be that the only use of some induction variable i is in tests such as

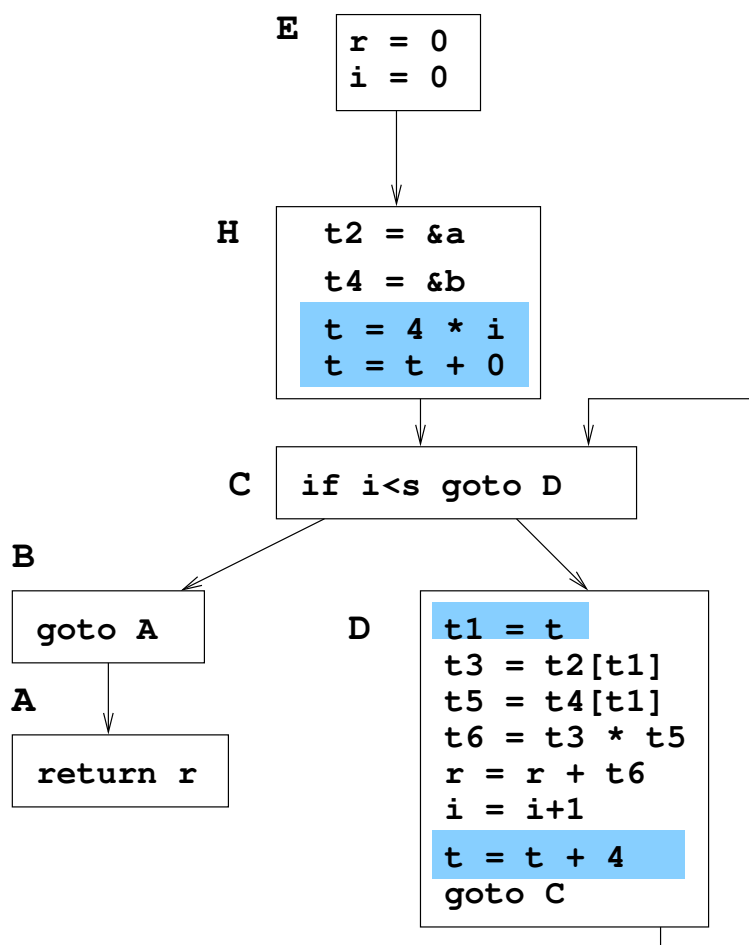


Figure 6.13: Flow graph of Example 20 after strength reduction

```
1 if i > C goto XXX
```

(and in the computation of other induction variables). In such a case, it may be profitable to replace the variable in the test by another induction variable j from the same family, and modify the test accordingly.

E.g., in our example (Figure 6.14), the triple of t is $(i, 4, 0)$, and we can replace the code in block C by

```
1 s1 = 4 * s
2 if t < s1 goto B
```

After this it may become possible to eliminate i since it is only used to compute itself. It may also be that the definition of the new temporary corresponds to a loop invariant which may be moved outside the loop.

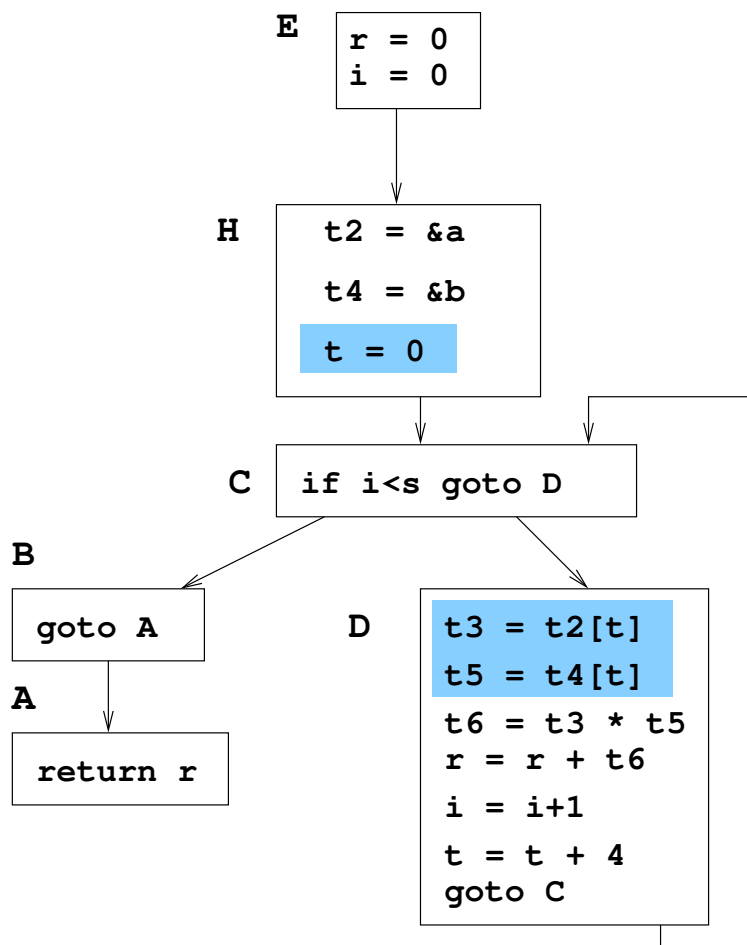


Figure 6.14: Flow graph of Figure 6.13 after copy propagation and constant folding

A more detailed description of this algorithm can be found in [ASU86].

Applying the above transformations to Example 20 results in the final graph depicted in Figure 6.15.

6.5 Aliasing: pointers and procedure calls

Two expressions are *aliases* of each other if they denote the same memory address. Both pointers and procedure calls introduce aliases. Aliases complicate data-flow analysis and code improvement, locally (see e.g. Section 6.2.3, page 100) as well as globally.

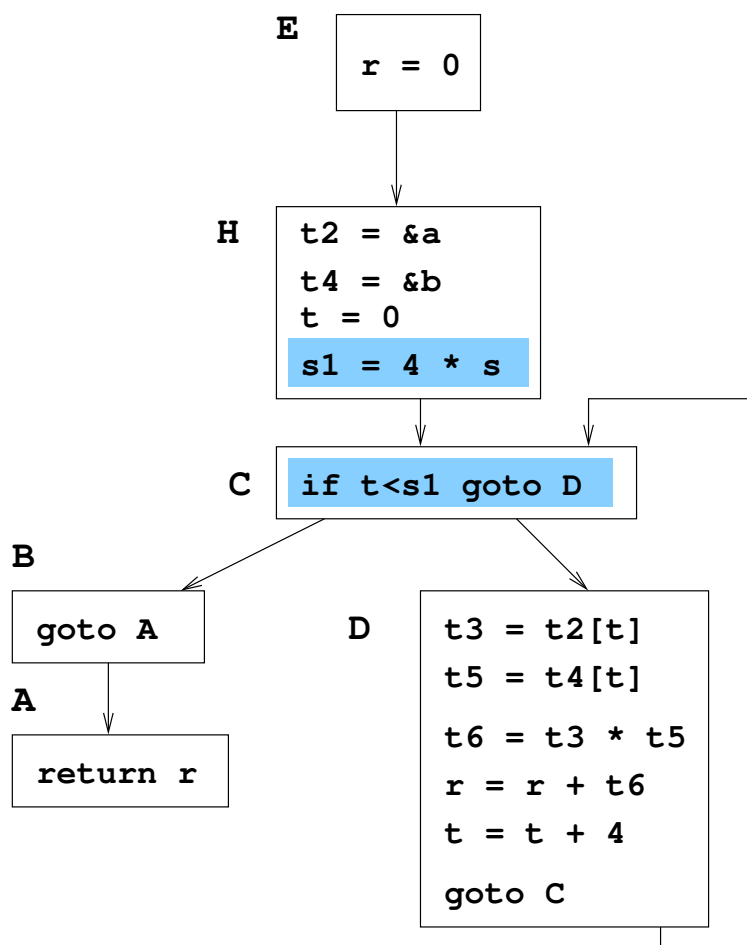


Figure 6.15: optimized flow graph for Example 20

If nothing is known about where a pointer points to, as is e.g. the case in a liberal language such as C, we must assume that any pointer assignment $*p = e$ can potentially define any variable. Likewise, dereferencing a pointer, as in $x = *p$ must be counted as a use of any variable. This obviously increases the number of live variables and reaching definition, while decreasing the number of available expressions. The net result is fewer code improvement opportunities.

A similar reasoning holds for procedure calls: if nothing is known about which variables the procedure may change, we must assume that it potentially defines all variables.

If, e.g. because of the rules of the source language, we know more about the possible usage of pointers and/or procedure calls, we may improve matters.

6.5.1 Pointers

We deal with pointers first. For each basic block B , we define two sets $IN_{ptr}(B)$ and $OUT_{ptr}(B)$ containing pairs (p, a) where p is a pointer variable and a is another variable. Intuitively, $(p, a) \in IN_{ptr}(B)$ means that, on entry to B , p may point to a (note that there may be several a such that $(p, a) \in IN_{ptr}(B)$), while $OUT_{ptr}(B)$ encodes where pointers may point upon exit of B .

The knowledge about pointer rules is encoded in a set of functions $trans_s$, one for each type of instruction s , which transforms a set of pairs (p, a) into another such set which is such that, if p may point to any a such that $(p, a) \in S$, then, after s , p may point to any b such that $(p, b) \in trans_s(S)$. $trans_s$ can be extended to $trans_B$ for any basic block B by composing, in order, the functions corresponding to the instructions of B . $IN_{ptr}(B)$ and $OUT_{ptr}(B)$ then satisfy the data-flow equations.

$$\begin{aligned} OUT_{ptr}(B) &= trans_B(IN_{ptr}(B)) \\ IN(B) &= \cup_{C < B} OUT_{ptr}(B) \end{aligned}$$

Note that these equations are isomorphic to the ones used to compute reaching definitions (Section 6.3.1, page 103, and hence one can use Algorithm 10 (page 104) to compute the smallest fixpoint.

$IN_{ptr}(B)$ can then be used to derive a precise set of variables that may be defined by an instruction such as $*p = e$. Similarly, the set of variables that may be used by an instruction such as $a = *p$ can be derived from $IN_{ptr}(B)$, and the data-flow algorithms can be suitably refined to take this information into account (note that, e.g. for use-definition chains, unless p can point to only one variable, we must not consider $*p = a$ as a definition that may hide earlier ones).

6.5.2 Procedures

If we have knowledge about procedures, we can apply a similar reasoning as for pointers in order to better estimate the effect of `CALL` instructions. In this section, we will use $p(\&x, y)$ to denote a procedure calling sequence in three address code where x is passed by reference while y is not.

As an example, we consider the case of C where procedures may only modify local and global variables. Note that procedures may also pass global variables as reference parameters⁹. Clearly, passing a variable x for a reference formal

⁹In C, reference parameters are simulated by passing a pointer by value.

parameter y makes x and y aliases, denoted $x \equiv y$ (we assume that the names of all local variables, including formal parameters, in all procedures and the names of global variables are all distinct).

The following algorithm computes the set of possible aliases amongst all variables.

Algorithm 24 [Computing aliases]

We show only reference formal parameters for the procedures.

1. *For each procedure $p(\&x_1, \dots, \&x_n)$ and each call $p(\&y_1, \dots, \&y_n)$, add $x_1 \equiv y_1, \dots, x_n \equiv y_n$ to the set of aliases.*
2. *Take the reflexive and transitive closure of the relation \equiv resulting from step 1.*

□

Now we can determine the set $change(p)$, where p is a procedure, of global variables and formal parameters¹⁰ that may possibly be changed by a call to p .

Define

- $def(p)$ as the set of formal parameters and globals of p that have definitions in the flow graph of p .
- A_p as the set of global variables or formal parameters a of p such that p contains a call $q(\dots, \&a, \dots)$ and the corresponding formal parameter is in $change(q)$.
- G_p as the set of global variables in $change(q)$ where p calls q .

Then

$$change(p) = def(p) \cup A_p \cup G_p$$

and this equation can be solved by a smallest fixpoint computation, starting with $change(p) = def(p)$ as an underestimate.

$change(p)$ can then be used as an estimate of possible definitions of variables, e.g. when computing available expressions (see also Section 10.8 in [ASU86]). A similar analysis can be made for the possible uses represented by a procedure call (for local variables this is easy: they are at most the parameters, for globals, one would need to analyze the call graph of the called procedure).

¹⁰Note that formal parameters of p are also local variables of p .

Chapter 7

Code generation

Code generation concerns the generation of machine code from (relatively) machine-independent intermediate code.

Actually, there are a few possibilities as to exactly what kind of “machine code” is generated:

- One could generate so-called *absolute* machine code in which addresses are fixed (“absolute”) as much as possible. This strategy predates multiprocessing operating systems and was of little use in modern operating systems. It may however be relevant again in the context of so-called “just-in-time” compilers.
- It is also feasible to generate *relocatable* code where most addresses are relative to some unknown base addresses. The code then has to be processed by a linker-loader before it can be executed. Since this option duplicates the work done by an assembler, we prefer the next option.
- We generate *assembler source code* where addresses are symbolic and leave it to a separate program to translate further to relocatable code.

Code generation deals mainly with the following issues:

- *Run-time storage management* is concerned with the mapping of names to addresses. We will look mainly at stack-based allocation of local variables into so-called *activation records*.
- *Instruction selection*. Depending on the characteristics of the hardware and the context, the same instruction may be translated in several ways, each with different costs in time and/or space.

- Closely related to the previous concern is the issue of *register allocation*. Many processors have very efficient instructions that take their operands from a limited set of registers. The aim is to allocate registers to variables in order to optimize global performance.

7.1 Run-time storage management

7.1.1 Global data

The global memory layout for a traditional language such as C is shown in Figure 7.1. Of course, in modern operating systems, these areas may be in different memory segments, with different permissions, e.g. the code segment may be read-only and shared.

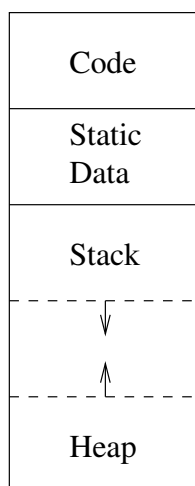


Figure 7.1: Memory layout in C

For static global variables, there is usually no problem: we simply use assembler pseudo-operations like

```
1 GLOBAL BYTE 4,0 ; allocate & initialize 4 byte quantity GLOBAL
```

For dynamic global variables (e.g. the **new** operator in C++), memory is (de)allocated at run-time from a so-called *heap*. Efficient algorithms for heap management, including e.g. garbage collection, are outside the scope of this text.

7.1.2 Stack-based local data

If we are dealing with a language such as C, we also have to provide storage for *local variables*. If the language supports recursion, we are naturally lead to a stack-based organization where, upon a function/procedure call, a so-called *activation record* is pushed onto the stack. Upon return from the call, the activation record is popped.

A possible layout for an activation record is shown in Figure 7.2.

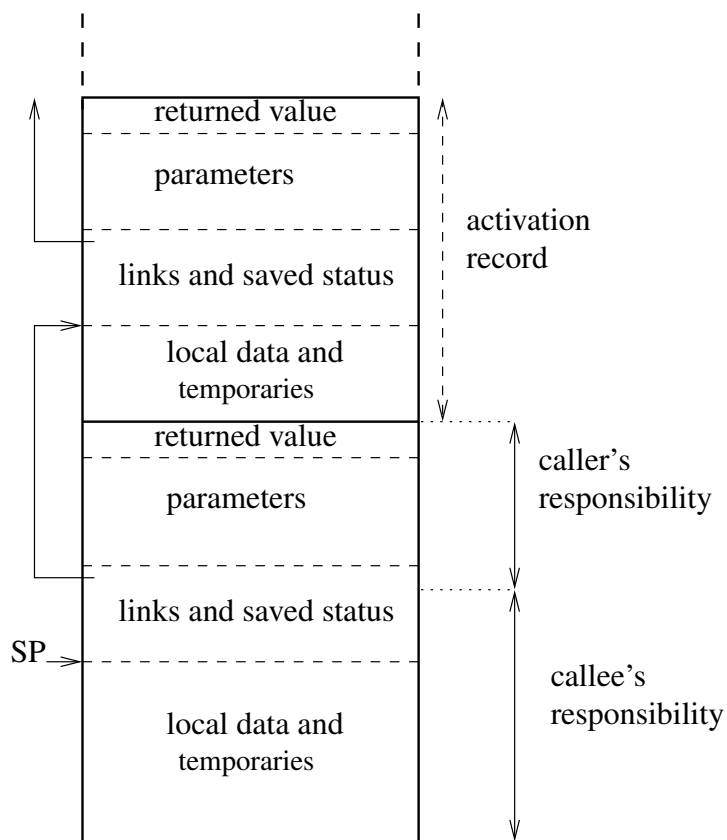


Figure 7.2: Activation records

Such a record contains, in growing stack order,

- Room for the value to be returned by the procedure (if any).
- The actual parameters.
- Links and saved status. The first item here is the stack pointer to be restored upon exit of the procedure. The next one might be the return address, i.e. the

address to which control should be transferred upon exit of the procedure. There may be additional hardware status information that is save here for restoration upon exit of the call.

- Room for local and temporary variables. By local variables, we mean the ones explicitly declared by the programmer. Temporary variables have been generated by the compiler during intermediate code generation and/or optimization.

The above organization is convenient because it allows an easy division of the work between caller and callee. The scenario for a procedure call is:

1. The caller evaluates the actual parameters and pushes them on the stack (note that the caller knows where is the top of the stack, since it knows *SP* and the size of its local and temporary data).
2. The caller pushes the value of *SP*, followed by the return address. It then puts *SP* to its new value and transfers control to the callee, e.g. by a `jump` instruction.
3. The callee saves further status information on the stack (using *SP*). The callee may further update *SP* to point just past the saved information (as in the figure).
4. The callee initializes its local data and goes about its further business.

A possible return (from a call) sequence is:

1. The callee places the return value in the appropriate place (note that it knows how many and which parameters it received, so the return value will be stored at $SP[-c]$ for some known c).
2. The callee restores the saved information, including the *SP* and jumps to the saved return address.
3. The caller has access to the returned value (right at the top of the stack), which it may copy into its local data.

It should be clear that, at run time, a register will be used to store *SP*. Access to local data will then be achieved using indirect addressing through the contents of *SP* and some constant offset which is known at compile time. This offset can easily be computed and stored in the symbol table. An interesting aspect is that local variables or temporaries may share a slot in the activation record provided they are never live at the same time. This occurs, e.g. when compiling the following C code.

```

1  int f(int x)
2  {
3  int a, b; /* slot 0, 1 */
4
5  ...
6
7  if (x>1) {
8      int y; /* slot 2 */
9
10     ...
11     }
12 else {
13     int z; /* slot 2 */
14     }
15 }

```

In the example, y and z will never be alive at the same time since they are in different (non-nested) scopes. Therefore, both can share slot 2.

Note also that the organization above leaves the possibility for a procedure to have local data whose size is only known at run-time. E.g. variable-sized arrays can be allocated on top of the stack, with a pointer to them reserved in the fixed-size local data area.

7.2 Instruction selection

Even abstracting from the register allocation program (Section 7.3, page 136), there may be several possibilities and constraints involved in the translation of an intermediate code instruction.

As a simple example, consider the intermediate code instruction

```
1  a = a + 1
```

This may be translated as (note the use of indirect addressing to access local names):

```

1  load [sp]+a,r0 ; copy a to register 0
2  add #1
3  store r0,[sp]+a ; copy register 0 back to a

```

but if the processor supports an *inc* (increment) instruction, it could also be translated as

```
1  inc [sp]+a ; increment a
```


Some operations may only be possible with certain registers. E.g. on the IBM/370 architecture, multiplication needs to be performed as shown in Figure 7.3, i.e. one operand must be available in the odd register of an even-odd register pair which will be filled with the result.

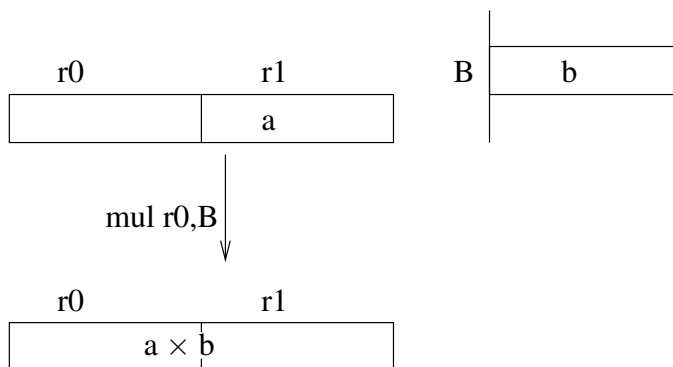


Figure 7.3: Multiplication on IBM/370

Addressing modes also influence the cost of an instruction. As an example, suppose the following table gives addressing modes and (relative) costs.

| | | |
|---------------------------------|---|---------------------------------|
| <code>move r0,r1</code> | 1 | register-to-register |
| <code>move [r0],[r1]</code> | 1 | register-to-register (indirect) |
| <code>move r5,x</code> | 2 | register-to-memory |
| <code>move #1,r3</code> | 2 | constant-to-register |
| <code>move [r0]+4,[r1]+5</code> | 3 | indexed-to-indexed |

We can then consider the following alternative code sequences for $a = b + c$ (assuming they are all global):

- We can use a scratch register (total cost 6):

```
1 move b,r0
2 add c,r0      ; result in r0
3 move r0,a
```

- We can perform addition in “memory” (total cost 6):

```
1 move b,a
2 add c,a
```

- If the *values* of b and c are in registers $r1$ and $r2$, we can reduce the cost to 3, provided that b is not live after the instruction:

```

1  add r2,r1
2  move r1,a

```

As the example illustrates, register allocation can have a dramatic effect on performance.

7.3 Register allocation

We describe how register allocation can be reduced to a graph coloring problem.

The idea is that, for a given flow graph (procedure), every variable corresponds to a *symbolic register*¹. The problem then is to assign each of a limited number of real registers to symbolic ones without introducing errors.

Clearly, we cannot use the same real register for symbolic registers (variables) r_1 and r_2 if e.g. r_1 is defined at a time when r_2 is live.

Definition 30 *The **register inference graph** G corresponding to a flow graph is an undirected graph where the nodes are the symbolic registers and there is an edge between r_1 and r_2 if r_1 is defined at a point where r_2 is live.*

If there are k real registers available, every symbolic register can be mapped to a real one without causing errors if G can be colored using k colors (no connected nodes may have the same color). Graph coloring with k colors is known to be NP-complete if $k > 2$ but the following heuristic procedure gives good results in practice.

Algorithm 25 [Register allocation using graph coloring]

Let G be the graph to be colored.

1. *While there is a node n with fewer than k connected edges, remove the n (and remember it).
The reason for this is that n can always receive a color, independent of the color chosen for its neighbors.*
2. *If we are left with an empty graph, G was colorable and we can assign colors by coloring the nodes in the reverse order (of removal).*

¹The algorithm can be adapted to hardware peculiarities where, e.g., a special register needs to be used for certain instructions. In this case, we associate a symbolic register with these instructions; the register is defined at the time the instruction is used.

3. *If we are left with a nonempty graph, it is not k -colorable since it has a node with at least k neighbors. We then choose a node n to be spilled, remove the node from the graph and start again with step 1. Spilling a node means that the corresponding variable will be loaded (into a register) on each use and stored back on each definition.*

Several heuristics have been proposed to choose a spilled node:

- *Choose a node with a maximal number of attached edges, maximizing the chance of ending up with a k -colorable graph.*
- *Estimate the cost of spilling a node, and choose the one with the lowest cost. A simple measure would count the number of extra load/store instructions, taking care to multiply with a suitable factor any such instructions that appear in a loop².*

□

7.4 Peephole optimization

Some machine code optimizations are very easy to perform by scanning the code (using a sliding window of instructions, hence the name) for certain patterns that can be replaced by improved patterns.

Often, it is the case that successful replacements induce opportunities for further replacements; therefore scanning is usually repeated until no further matches are found.

Below we list some example patterns and replacements.

²One could also have a higher cost factor for instructions in an inner loop than in an encompassing loop.

| Pattern | Replacement | Description |
|---|--|-----------------------------|
| move x,y move y,x | move x,y | redundant load, store |
| jump(con) z jump y z: (seq1) y: (seq2) | jump(notcon) y z: (seq1) y: (seq2) | jumps over jumps |
| jump(con) z .. z: jump y | jump(con) y .. z: jump y | jumps to jumps |
| mul #1,x add #0,x | (delete) (delete) | algebraic simplification |
| mul #2,x | add x,x | machine specifics |
| mul #2,x | shiftright x,1 | machine specifics |
| add #1,x | inc x | machine specifics |

Appendix A

A Short Introduction to x86 Assembler Programming under Linux

A.1 Architecture

Memory is byte addressable. There are several general purpose registers: `%eax`, `%ebx`, `%ecx`, `%edx`, `%edi` and `%esi`.

The following registers have a special purpose:

- The “base pointer” `%ebp` is used to point to a stack frame for easy accessing to local variables and parameters, see Section A.4.
- The “stack pointer” `%esp` always points to the top of the stack which grows downwards. E.g. if the contents of `%esp` is 1000, then the top two elements of the stack are at addresses 1000 and 999.
- The instruction counter `%eip` points to the next instruction (in memory) to be executed.
- The register `%eflags` contains the condition codes which are set by several instructions and can be tested/used by others. Table A.1 shows some commonly used flags.

| code | name | description |
|------|----------|---|
| O | overflow | 1 if destination cannot hold result |
| S | sign | sign of last result |
| Z | zero | 1 if last result was 0 |
| P | parity | 1 if low byte of result has even number of 1 bits |
| C | carry | 1 if overflow bit |

Table A.1: Some x86 Flags

A.2 Instructions

A.2.1 Operands

The operands for instructions may be of the following types:

- *Immediate* (I), i.e. constant integer data. For the gnu assembler (*as*), such data is prefixed with '\$', e.g. \$100 (decimal) or \$0x80 (hexadecimal). Such data are encoded using 1, 2 or 4 bytes.
- *Register* (R), referring to a (4 byte) register, e.g. `%eax`.
- *Memory* (M), referring to 4 consecutive bytes of memory.

As an example, the `movl` moves a word (4 bytes) from the first to the second operand. The possibilities are illustrated in the following code fragment.

```

1 movl $100, %eax /* IR store 100 in register eax */
2 movl $100, (%eax) /* IM store 100 in memory pointed to by eax */
3 movl %eax, %ebx /* RR store content of eax in ebx */
4 movl %ebx, (%eax) /* RM store content of ebx in memory pointed to by eax */
5 movl (%eax), %ebx /* MR store content of memory pointed to by eax in ebx */

```

Note that single-instruction MM (memory to memory) transfers are impossible. In fact, no instruction can have two M operands.

A.2.2 Addressing Modes

The following *addressing modes* are available for referring to a memory location:

- *Normal* mode, e.g. `(%eax)` refers to the location with address the contents of the `%eax` register. Thus, using $*r$ to denote the contents of register r ,

$$(r) \rightarrow [*r]$$

- In *displacement* mode, an offset is specified together with a register. E.g. -4 (`%esp`) refers to the (32bit) word just “under” the top of the stack while 8 (`%esp`) refers to the location at $(\%esp) + 8$. Thus:

$$d(r) \rightarrow [*r + d]$$

- *Indexed* addressing mode involves a *base* (any) and an *index* (any except `%esp`) register as well as a constant *displacement* (1, 2 or 4) and a *scale* which must be 1, 2, 4 or 8:

$$d(r_b, r_i, s) \rightarrow [*r_b + s \times *r_i + d]$$

E.g. `movl 8(%ebx, %edi, 4), %eax` would move the word located at address $(\%ebx) + 4 \times (\%edi) + 8$ to register `%eax`.

One of the registers or the scale (default 1) or displacement (default 0) may be missing. E.g. *as* accepts all of the following:

```
1      movl $100, (%eax, %ebx)
2      movl $100, 4(%eax, %ebx)
3      movl $100, (%eax, %ebx, 4)
4      movl $100, (, %ebx, 4)
```

The tables in the following sections are adapted from [Bar03]. Also, o_1 and o_2 are sometimes used to refer to the first and second operand of the instruction.

A.2.3 Moving Data

| Data Transfer Instructions | | | |
|----------------------------|----------|-----------|--|
| name | operands | flags set | description |
| <code>movl</code> | IMR,IMR | OSZC | move word from o_1 to o_2 |
| <code>movb</code> | IMR,IMR | OSZC | move byte from o_1 to o_2 |
| <code>leal</code> | M,IMR | OSZC | store address of M in o_2 |
| <code>popl</code> | MR | OSZC | pop word from stack into operand |
| <code>pushl</code> | IMR | OSZC | push word on stack |
| <code>pushfl</code> | | | push <code>eflags</code> register on stack |
| <code>xchgl</code> | MR, MR | OSZC | exchange the words in o_1 and o_2 |

Table A.2: Data Transfer Instructions

| Integer Arithmetic Instructions | | | |
|---------------------------------|----------|-----------|--|
| name | operands | flags set | description |
| adcl | IMR,MR | OSZPC | add with carry $o_2 += (o_1 + C)$ |
| addl | IMR,MR | OSZPC | add $o_2 += o_1$ |
| cdq | | OSZPC | convert %eax to double word in %edx : %eax , see <code>idivl</code> |
| cmpl | IMR,MR | OSZPC | set flags according to $o_2 - o_1$ |
| decl | MR | OSZPC | decrement (-1) word |
| decbl | MR | OSZPC | decrement (-1) byte |
| divl | MR | OSZPC | unsigned divide %edx : %eax by o_1 , quotient in %eax , remainder in %edx |
| idivl | MR | OSZPC | signed division, see <code>divl</code> |
| imull | IMR, R | OSZPC | signed multiplication: $o_2 \times= o_1$ |
| imull | IMR | OSZPC | signed multiplication: %edx : %eax = %eax \times o_1 |
| incl | MR | OSZPC | increment $(+1)$ word |
| incbl | MR | OSZPC | increment $(+1)$ byte |
| negl | MR | OSZPC | negate (2's complement inversion) operand |
| sbbbl | IMR, MR | OSZPC | $o_2 -= o_1$, with borrowing (see <code>adcl</code>) |
| subl | IMR, MR | OSZPC | $o_2 -= o_1$, no borrowing |

Table A.3: Integer Arithmetic Instructions

A.2.4 Integer Arithmetic

See Table A.3 on page 142.

A.2.5 Logical Operations

See Table A.5 on page 145.

A.2.6 Control Flow Instructions

See Table A.6 on page 146.

A.3 Assembler Directives

The `gnu` assembler `as` supports a number of directives:

- A *section* switch directive is of the form

```
.section name
```

where *name* is usually *.text*, for code, *.data*, for data embedded in the program, or *.bss* for uninitialized data. It indicates which section the following directives and instructions belong to.

- An *equ* directive is of the form

```
.equ label, value
```

and assigns the *value* to *label*.

- Data are defined using an *.ascii*, *.byte* or *.long* directive, as in the following example

```
1 astring: /* label */  
2 .ascii "Hello world\0"  
3 bytes:  
4 .byte 0, 1, 2, 3, 4, 5  
5 words:  
6 .long 1024, 2048, 4096
```

- Data in the *.bss* section can be reserved using the *.lcomm* directive which has the form

```
.lcomm label, size_in_bytes
```

- The *.globl* directive exports a label so that it is accessible from other object files. The *.type* directive tells the linker that a label is a function, as in the following example.

```
1 .globl ackerman  
2 .type ackerman, @function  
3 ackerman:  
4 ...
```

The *.include* directive inserts another file at that point in the text. It has the form

```
.include "pathname"
```

- Finally, the *.rept* and *.endr* can be used to repeat everything between the directives a number of times, as in the following example. which reserves 100 data bytes that are initialized to 0.

```

1 .rept 100
2 .byte 0
3 .endr

```

A.4 Calling a function

The C calling convention is illustrated in the following example. Note how *%ebp* is used as a frame pointer.

```

1 .globl print
2 .type myfun, @function
3 # myfun(int a, int b) { int c = 0; int d = 2; ... }
4 myfun:
5     # stack: b a [return-addr]
6     pushl %ebp # save base(frame) pointer on stack
7     # stack: b a [return-addr] [saved-bp]
8     movl %esp, %ebp # base pointer is stack pointer
9     .equ a, 8
10    .equ b, 4
11    # a(%ebp) => a; b(%ebp) => b
12    pushl $0 # initialize c
13    pushl $2 # initialize d
14    # stack: b a [return-addr] [saved-bp] 0 2
15    .equ c, -4
16    .equ d, -8
17    # c(%ebp) => c; d(%ebp) => d
18
19    /* body of function, the return value is stored in %eax */
20
21    movl %ebp, %esp # restore stack pointer
22    # stack: b a [returnaddress] [savedebp]
23    popl %ebp # restore saved base (frame) pointer
24    /* stack: b a [return-addr]
25    ret

```

A.5 System calls

To obtain a service from Linux, it suffices to specify the service in register *eax*, set up the input parameters in various registers (see Table A.4 below) and then

generate an interrupt using the instruction **int int** $\$0 \times 80$. The return value and/or error code will be stored in register **eax**.

More information can be found in the man pages corresponding to the “function” column in Table A.4 which lists a few I/O related system calls.

| %eax | name | %ebx | %ecx | %edx | function |
|------|-------|------|--------|--------|----------------------------|
| 1 | exit | code | | | exit (code) |
| 3 | read | fd | buf | count | read (fd, buf, count) |
| 4 | write | fd | buf | count | write (fd, buf, count) |
| 5 | open | path | flags | mode | open (path, flags, mode) |
| 6 | close | fd | | | close (fd) |
| 19 | lseek | fd | offset | whence | lseek (fd, offset, whence) |

Table A.4: Linux System Calls

| Logical Operations | | | |
|--------------------|----------|-----------|---|
| name | operands | flags set | description |
| andl | IMR,MR | OSZPC | $o_2 \&= o_1$ $C = O = 0$ |
| orl | IMR,MR | OSZPC | $o_2 = o_1$ $C = O = 0$ |
| notl | MR | | $o_1 = o_1$ |
| rcll | I%cl, RM | OC | rotate $o_2 + C$ by o_1 bits left |
| rcrl | I%cl, RM | OC | rotate $o_2 + C$ by o_1 bits right |
| roll | I%cl, RM | OC | rotate o_2 by o_1 bits left |
| rorl | I%cl, RM | OC | rotate o_2 by o_1 bits right |
| sall | I%cl, RM | C | arithmetic shift o_1 bits left; sign bit to C, 0 added as lsb |
| sarl | I%cl, RM | C | arithmetic shift o_1 bits right; lsb bit to C, S to msb |
| shll | I%cl, RM | C | shift o_1 bits left; msb bit to C |
| shrl | I%cl, RM | C | shift o_1 bits right; lsb bit to C |
| testl | IRM, RM | OSZPC | set flags from $o_1 \& o_2$ |
| xorl | IRM, RM | OSZPC | $o_2 = o_2 \mathbf{xor} o_1$, $O = C = 0$ |

Table A.5: Logical Operations

| Control Flow Instructions | | | |
|---------------------------|----------|-----------|--|
| name | operands | flags set | description |
| call | address | OSZC | push %eip, jump to address |
| call | *R | OSZC | call function a (R), e.g. call *%eax |
| int | I | OSZC | cause interrupt number <i>I</i> , <i>I</i> = 0x080 for Linux kernel call |
| jmp | address | OSZC | unconditional jump |
| jmp | *R | OSZC | goto (R), e.g. jmp *%eax |
| ret | | OSZC | return: <code>popl %eip; jmp *%eip</code> |
| jecxz | address | OSZC | jump if %ecx is 0 |
| j[n][spoc] | address | OSZC | jump if [SPOC] is (not) set |
| j[n][gl][e] | address | OSZC | signed comparison: jump if (not) > (g), < (l), ≥ (ge) or ≤ (le) |
| j[n]z | address | OSZC | jump if (not) 0 |
| j[n]e | address | OSZC | jump if (not) equal |
| j[n][ab][e] | address | OSZC | unsigned comparison jump if (not) > (a), < (b), ≥ (ae) or ≤ (be) |

Table A.6: Control Flow Instructions

A.6 Example

Example 21 The listing below shows a function that prints its single parameter (an integer) on stdout, followed by '\n'.

```

1 # print an integer
2 .section .text
3 .globl print_int
4 .type print_int, @function
5 # print_it(int value) -- print value followed by \n to stdout.
6 print_int:
7     pushl %ebp
8     movl %esp, %ebp
9     .equ value, 8 # parameter offset
10    # initialize local variables:
11    .equ base, -4
12    pushl $10 # base = 10
13    .equ bufsize, -8
14    pushl $1 # bufsize = 1 ('\n')
15    .equ negative, -12
16    pushl $0 # negative = 0
17    # stack: .. value return-addr saved-ebp base bufsize

```

```

18     pushl $10 # push newline to be printed last
19     movl value(%ebp), %eax
20     jge .L1 # if value >= 0
21     # value < 0: remember
22     movl $1, negative(%ebp)
23     negl %eax # value = -value
24 .L1:
25     cdq # setup %edx:%eax for division
26     # aex = value/base, edx = value % base
27     divl base(%ebp)
28     # push remainder digit
29     pushl %edx
30     addl $48, (%esp)
31     incl bufsize(%ebp) # ++bufsize
32     cmpl $0, %eax
33     jne .L1 # loop if value != 0
34     # put sign if negative
35     cmpl $0, negative(%ebp)
36     je .L2
37     pushl $45 # '-'
38     incl bufsize(%ebp)
39 .L2:
40     # write(2): eax = 4, ebx = fd, ecx = buffer start, edx = buffer size
41     movl $4, %eax # code for write syscall
42     movl $1, %ebx # fd stdout = 1
43     movl %esp, %ecx # buffer start = top of stack
44     movl $4, %edx # bufsize * 4 bytes
45     imul bufsize(%ebp), %edx
46     int $0x80 # syscall
47     movl %ebp, %esp
48     popl %ebp # restore saved frame pointer
49     ret

```

Example 22 The function from Example 21 is used in the following program.

```

1  # print an integer
2  .section .text
3  .globl _start
4  _start:
5      call main
6      jmp exit
7  .include "print_int.s"
8  .globl main
9  .type main, @function
10 main:
11     pushl %ebp # save frame pointer
12     movl %esp, %ebp # set frame pointer
13     pushl $-100 # number to print
14     call print_int

```

```
15         movl %ebp, %esp
16         popl %ebp
17         ret
18 .type exit, @function
19 exit:
20         movl $0, %ebx
21         movl $1, %eax
22         int $0x80
```

It can be assembled and linked using the following command which translates *tst.s* into an object file *tst.o* which is then linked to an executable file *tst*

```
as tst.s -o tst.o && ld tst.o -o tst
```

Appendix B

Mc: the Micro-x86 Compiler

B.1 Lexical analyzer

```
1  %{
2  /* $Id: lex.l,v 1.1 2008/07/09 13:06:42 dvermeir Exp $
3  *
4  * Lexical analyzer for the toy language ``Micro``
5  */
6  #include <string.h> /* for strcmp, strdup & friends */
7  #include <stdlib.h> /* for atoi() */
8
9  #include "micro.tab.h" /* token type definitions from .y file */
10 #include "symbol.h" /* symbol table management */
11
12 extern int lineno; /* defined in micro.y */
13
14 void
15 lex_init() {
16     /* Initialize data structures etc. for scanner */
17     symbol_insert("declare",DECLARE); /*Insert keywords in symbol table */
18     symbol_insert("read",READ); /*Insert keywords in symbol table */
19     symbol_insert("write",WRITE);
20 }
21
22 /*
23 * The macro below will be called automatically when the generated scanner
24 * initializes itself.
25 */
26 #define YY_USER_INIT lex_init();
27
28 %}
29
```

```

30 alpha    [A-Za-z]
31 digit    [0-9]
32 alphanum [A-Za-z0-9]
33
34 %%
35 [ \t]    break;  /* ignore white space */
36 "\n"     ++lineno;
37
38 {alpha}{alphanum}*      {
39                           yyval.idx = symbol_find(yytext);
40
41                           if (yyval.idx<0) { /* new symbol: insert it */
42                               yyval.idx =symbol_insert(yytext, NAME);
43                               return NAME;
44                           }
45                           else
46                               return symbol_type(yyval.idx);
47                           }
48
49 {digit}+                {
50                           yyval.value = atoi(yytext);
51                           return NUMBER;
52                           }
53
54 "("                    return LPAREN;
55 ")"                    return RPAREN;
56 "{"                    return LBRACE;
57 "}"                    return RBRACE;
58 "="                    return ASSIGN;
59 ";"                    return SEMICOLON;
60 "+"                    return PLUS;
61 "-"                    return MINUS;
62
63 .                      {
64                           fprintf(stderr,
65                               "Illegal character \'%s\' on line #%d\n",
66                               yytext, lineno);
67                           exit(1);
68                           }
69
70 %%
71
72 int
73 yywrap() {
74     return 1; /* tell scanner no further files need to be processed */
75 }

```


B.2 Symbol table management

```

1  /* $Id: symbol.c,v 1.1 2008/07/09 13:06:42 dvermeir Exp $
2  *
3  * Symbol table management for toy 'Micro' language compiler.
4  * This is a trivial example: the only information kept
5  * is the name, the token type: READ, WRITE or NAME and, for NAMES
6  * whether they have been declared in the JVM code.
7  */
8  #include <stdio.h> /* for (f)printf(), std{out,int} */
9  #include <stdlib.h> /* for exit */
10 #include <string.h> /* for strcmp, strdup & friends */
11 #include "micro.tab.h" /* token type definitions from .y file */
12
13 #include "symbol.h"
14
15 typedef struct {
16     char *name;
17     int type; /* READ, WRITE, or NAME */
18     int declared; /* NAME only: 1 iff already declared in JVM code, 0 else */
19 } ENTRY;
20
21 #define MAX_SYMBOLS 100
22 static ENTRY symbol_table[MAX_SYMBOLS]; /* not visible from outside */
23 static int n_symbol_table = 0; /* number of entries in symbol table */
24
25 int
26 symbol_find(char* name) {
27     /* Find index of symbol table entry, -1 if not found */
28     int i;
29
30     for (i=0; (i<n_symbol_table); ++i)
31         if (strcmp(symbol_table[i].name, name)==0)
32             return i;
33     return -1;
34 }
35
36 int
37 symbol_insert(char* name,int type) {
38     /* Insert new symbol with a given type into symbol table,
39      * returns index new value */
40     if (n_symbol_table>=MAX_SYMBOLS) {
41         fprintf(stderr, "Symbol table overflow (%d) entries\n", n_symbol_table);
42         exit(1);
43     }
44     symbol_table[n_symbol_table].name = strdup(name);
45     symbol_table[n_symbol_table].type = type;
46     symbol_table[n_symbol_table].declared = 0;

```

```

47     return n_symbol_table++;
48 }
49
50 int
51 symbol_type(int i) {
52     /* Return type of symbol at position i in table. */
53     /* ASSERT ((0<=i)&&(i<n_symbol_table)) */
54     return symbol_table[i].type;
55 }
56
57 void
58 symbol_declare(int i) {
59     /* Mark a symbol in the table as declared */
60     /* ASSERT ((0<=i)&&(i<n_symbol_table)&&(symbol_table[i].type==NAME)) */
61     symbol_table[i].declared = 1;
62 }
63
64 int
65 symbol_declared(int i) {
66     /* Return declared property of symbol */
67     /* ASSERT ((0<=i)&&(i<n_symbol_table)&&(symbol_table[i].type==NAME)) */
68     return symbol_table[i].declared;
69 }
70
71 char*
72 symbol_name(int i) {
73     /* Return name of symbol */
74     /* ASSERT ((0<=i)&&(i<n_symbol_table)) */
75     return symbol_table[i].name;
76 }

```

B.3 Parser

```

1  %{
2  /* $Id: micro.y,v 1.9 2008/07/11 18:56:43 dvermeir Exp $
3  *
4  * Parser specification for Micro
5  */
6  #include <stdio.h> /* for (f)printf() */
7  #include <stdlib.h> /* for exit() */
8
9  #include "symbol.h"
10
11 int  lineno = 1; /* number of current source line */
12 extern int yylex(); /* lexical analyzer generated from lex.l */
13 extern char *yytext; /* last token, defined in lex.l */
14
15 void

```

```

16 yyerror(char *s) {
17     fprintf(stderr, "Syntax error on line %d: %s\n", lineno, s);
18     fprintf(stderr, "Last token was \"%s\"\n", yytext);
19     exit(1);
20 }
21
22 #define PROLOGUE "\
23 .section .text\n\
24 .globl _start\n\
25 \n\
26 _start:\n\
27     call main\n\
28     jmp exit\n\
29 .include \"../x86asm/print_int.s\"\n\
30 .globl main\n\
31 .type main, @function\n\
32 main:\n\
33     pushl %ebp /* save base(frame) pointer on stack */\n\
34     movl %esp, %ebp /* base pointer is stack pointer */\n\
35 "
36
37 #define EPILOGUE "\
38     movl %ebp, %esp\n\
39     popl %ebp /* restore old frame pointer */\n\
40     ret\n\
41 .type exit, @function\n\
42 exit:\n\
43     movl $0, %ebx\n\
44     movl $1, %eax\n\
45     int $0x80\n\
46 "
47
48 %}
49
50 %union {
51     int idx;
52     int value;
53 }
54
55 %token NAME
56 %token NUMBER
57 %token LPAREN
58 %token RPAREN
59 %token LBRACE
60 %token RBRACE
61 %token ASSIGN
62 %token SEMICOLON
63 %token PLUS
64 %token MINUS

```

```

65 %token DECLARE
66 %token WRITE
67 %token READ
68
69 %type <idx> NAME var
70 %type <value> NUMBER
71
72 %%
73 program          : LBRACE
74                   { puts(".section .data\n"); } declaration_list
75                   { puts(PROLOGUE); } statement_list
76                   RBRACE { puts(EPILOGUE); }
77                   ;
78 declaration_list : declaration SEMICOLON declaration_list
79                 | /* empty */
80                 ;
81
82 statement_list  : statement SEMICOLON statement_list
83                 | /* empty */
84                 ;
85
86 statement       : assignment
87                 | read_statement
88                 | write_statement
89                 ;
90
91 declaration     : DECLARE NAME
92                 {
93                     if (symbol_declared($2)) {
94                         fprintf(stderr,
95                             "Variable \"%s\" already declared (line %d)\n",
96                             symbol_name($2), lineno);
97                         exit(1);
98                     }
99                     else {
100                         printf(".lcomm %s, 4\n", symbol_name($2));
101                         symbol_declare($2);
102                     }
103                 }
104                 ;
105
106 assignment      : var ASSIGN expression
107                 {
108                     /* we assume that the expresion value is (%esp) */
109                     printf("\tpopl %s\n", symbol_name($1));
110                 }
111                 ;
112
113 read_statement  : READ var { /* not implemented */}

```

```

114             ;
115
116 write_statement : WRITE expression { puts("\tcall print_int\n"); }
117             ;
118
119 expression      : term
120                 | term PLUS term { puts("\tpopl %eax\n\taddl %eax, (%esp)\n"); }
121                 | term MINUS term { puts("\tpopl %eax\n\tsubl %eax, (%esp)\n"); }
122                 ;
123
124 term            : NUMBER { printf("\tpushl %d\n", $1); }
125                 | var { printf("\tpushl %s\n", symbol_name($1)); }
126                 | LPAREN expression RPAREN
127                 ;
128
129 var             : NAME
130                 {
131                     if (!symbol_declared($1)) {
132                         fprintf(stderr,
133                             "Variable \"%s\" not declared (line %d)\n",
134                             symbol_name($1), lineno);
135                         exit(1);
136                     }
137                     $$ = $1;
138                 }
139                 ;
140 %%
141
142 int
143 main(int argc, char *argv[]) {
144     return yyparse();
145 }

```

B.4 Driver script

```

1  #!/bin/sh
2  #
3  #      $Id: microc.sh,v 1.3 2008/07/10 16:44:14 dvermeir Exp $
4  #
5  #      Usage:  microc basename.mi
6  #
7  #      e.g. "microc tst.mi"  to compile tst.mi, resulting in
8  #
9  #      tst.s    assembler source code
10 #      tst.o    object file
11 #      tst      executable file
12 #
13 # determine basename

```

```

14 base=`basename $1 .mi`
15 # this checks whether $1 has .mi suffix
16 [ ${base} = $1 ] && { echo "Usage: microc basename.mi"; exit 1; }
17 # make sure source file exists
18 [ -f "$1" ] || { echo "Cannot open \"$1\""; exit 1; }
19 # compile to assembly code
20 ./micro <$1 >${base}.s || { echo "Errors in compilation of $1.mi"; exit 1; }
21 # assemble to object file: the --gdwarf2 option generates info for gdb
22 as --gdwarf2 ${base}.s -o ${base}.o || { echo "Errors assembling $1.s"; exit 1; }
23 # link
24 ld ${base}.o -o ${base} || { echo "Errors linking $1.o"; exit 1; }

```

B.5 Makefile

```

1 # $Id: Makefile,v 1.8 2008/07/10 16:44:14 dvermeir Exp $
2 #
3 CFLAGS=          -Wall -g
4 CC=              gcc
5 #
6 SOURCE=          microc.sh micro.y lex.l symbol.c symbol.h Makefile *.mi
7 #
8 all:             micro microc demo
9 micro:           micro.tab.o lex.yy.o symbol.o
10                 gcc $(CFLAGS) -o $@ $^
11
12 lex.yy.c:        lex.l micro.tab.h
13                 flex lex.l
14 #
15 #      Bison options:
16 #
17 #      -v      Generate micro.output showing states of LALR parser
18 #      -d      Generate micro.tab.h containing token type definitions
19 #
20 micro.tab.h\
21 micro.tab.c:     micro.y
22                 bison -v -d $^
23 ##
24 demo:           microc micro demo.mi
25                 ./microc demo.mi
26 #
27 CFILES= $(filter %.c, $(SOURCE)) micro.tab.c lex.yy.c
28 HFILES= $(filter %.h, $(SOURCE)) micro.tab.h
29 include make.depend
30 make.depend:    $(CFILES) $(HFILES)
31                 gcc -M $(CFILES) >$@
32
33 clean:
34                 rm -f lex.yy.c micro.tab.[hc] *.o microc micro *.jasm *.class micro

```

```
35 #
36 tar:
37         tar cvf micro.tar $(SOURCE)
```

B.6 Example

B.6.1 Source program

```
1 {
2 declare xyz;
3 xyz = (33+3)-35;
4 write xyz;
5 }
```

B.6.2 Assembly language program

```
1 .section .data
2
3 .lcomm xyz, 4
4 .section .text
5 .globl _start
6
7 _start:
8     call main
9     jmp exit
10 .include "../x86asm/print_int.s"
11 .globl main
12 .type main, @function
13 main:
14     pushl %ebp /* save base(frame) pointer on stack */
15     movl %esp, %ebp /* base pointer is stack pointer */
16
17     pushl $33
18     pushl $3
19     popl %eax
20     addl %eax, (%esp)
21
22     pushl $35
23     popl %eax
24     subl %eax, (%esp)
25
26     popl xyz
27     pushl xyz
28     call print_int
29
30     movl %ebp, %esp
31     popl %ebp /* restore old frame pointer */
```

```
32         ret
33 .type exit, @function
34 exit:
35         movl $0, %ebx
36         movl $1, %eax
37         int $0x80
```


Appendix C

Minic parser and type checker

C.1 Lexical analyzer

```
1  %{
2  /*      lex.l (1.6)      17:46:22      97/12/10
3  *
4  *      Lexical analyzer for the toy language ``minic``
5  */
6  #include      <string.h>      /* for strcmp, strdup & friends */
7  #include      <stdio.h>      /* for strcmp, strdup & friends */
8  #include      <stdlib.h>      /* for atoi() */
9
10 #include      "symtab.h"      /* symbol table management */
11 #include      "types.h"      /* symbol table management */
12 #include      "minic.tab.h"    /* token type definitions from .y file */
13 #include      "names.h"      /* string pool management */
14
15 extern int     lineno;        /* defined in minic.y */
16 extern SYM_TAB *scope;       /* defined in minic.y */
17
18 void
19 lex_init()
20 /*      Initialize data structures etc. for scanner */
21 {
22     scope    = symtab_open(0);    /* open topmost scope */
23 }
24
25 /*
26 *      The macro below will be called automatically when the generated scanner
27 *      initializes itself.
28 */
29 #define YY_USER_INIT    lex_init();
```

```

30
31 %}
32
33 alpha          [A-Za-z]
34 digit          [0-9]
35 alphanum       [A-Za-z0-9]
36
37 %%
38 [ \t]          break;           /* ignore white space */
39 "\n"           ++lineno;
40
41
42 int            return INT;       /* Keywords come before NAMES */
43 if             return IF;
44 else           return ELSE;
45 return         return RETURN;
46 float          return FLOAT;
47 struct         return STRUCT;
48
49 {alpha}{alphanum}* {
50     yyval.name      = names_find_or_add(yytext);
51     return NAME;
52 }
53
54 {digit}+         {
55     yyval.value     = atoi(yytext);
56     return NUMBER;
57 }
58
59 "("             return LPAR;
60 ")"             return RPAR;
61 "{"             return LBRACE;
62 "}"             return RBRACE;
63 "["             return LBRACK;
64 "]"             return RBRACK;
65 "=="            return EQUAL;
66 "="             return ASSIGN;
67 ";"             return SEMICOLON;
68 ","             return COMMA;
69 "."             return DOT;
70 "+"             return PLUS;
71 "-"             return MINUS;
72 "*"             return TIMES;
73 "/"             return DIVIDE;
74
75 .               {
76     fprintf(stderr,
77         "Illegal character \'%s\' on line #%d\n",
78         yytext,lineno);

```

```

79             exit(1);
80         }
81
82     %%
83
84     int
85     yywrap()
86     {
87     return 1; /* tell scanner no further files need to be processed */
88     }

```

C.2 String pool management

```

1  #ifndef NAMES_H
2  #define NAMES_H
3  /*      names.h(1.1)      09:25:08      97/12/08
4  *
5  *      String pool management.
6  */
7  char*  names_insert(char*);
8  char*  names_find(char*);
9  char*  names_find_or_add(char*);
10
11 #endif

1  /*      names.c(1.1)      09:25:08      97/12/08
2  *
3  *      String pool management.
4  */
5  #include      <stdio.h>      /* for fprintf and friends */
6  #include      <stdlib.h>     /* for malloc and friends */
7  #include      <assert.h>     /* for assert(invariant) */
8  #include      <string.h>     /* for strcmp(), strcpy() */
9
10 #include      "names.h"
11
12 #define INCREMENT_SIZE      1000
13
14 static char*  buf          = 0;
15 static int    buf_size     = 0; /* total size of buffer */
16 static size_t buf_avail    = 0; /* free bytes in buf */
17 static char*  next_name    = 0; /* address available for next string */
18
19 #define INVARIANT      assert(buf+buf_size==next_name+buf_avail); \
20                        assert(buf_size>=0); \
21                        assert(buf_avail>=0); \
22                        assert(next_name>=buf);
23 static void

```

```
24 names_grow(size_t size)
25 {
26     size_t next_offset = next_name - buf;
27
28     { INVARIANT }
29
30     buf = realloc(buf, buf_size + size); /* like malloc() if buf==0 */
31     if (!buf) {
32         fprintf(stderr, "Cannot expand name space (%d bytes)", buf_size + size);
33         exit(1);
34     }
35     buf_avail += size;
36     buf_size += size;
37     next_name = buf + next_offset;
38
39     { INVARIANT }
40 }
41
42 char*
43 names_insert(char* s)
44 {
45     char* ps;
46     size_t l = strlen(s) + 1; /* include trailing '\0' */
47
48     { INVARIANT }
49
50     while (l > buf_avail)
51         names_grow(INCREMENT_SIZE);
52
53     ps = strcpy(next_name, s);
54
55     buf_avail -= l;
56     next_name += l;
57
58     { INVARIANT }
59
60     return ps;
61 }
62
63 char*
64 names_find(char *s)
65 {
66     char *pc;
67
68     for (pc = buf; (pc != next_name); pc += strlen(pc) + 1)
69         if (!strcmp(pc, s))
70             return pc;
71     return 0;
72 }
```

```

73
74 char*
75 names_find_or_add(char *s)
76 {
77     char    *pc = names_find(s);
78
79     if (!pc)
80         pc = names_insert(s);
81     return pc;
82 }

```

C.3 Symbol table management

```

1  #ifndef SYMBOL_H
2  #define SYMBOL_H
3  /*      symtab.h(1.5)      10:07:41      97/12/10
4  *
5  *      Symbol table management
6  */
7
8  typedef struct syminfo {
9      char                *name;
10     struct type_info    *type;
11     } SYM_INFO;
12
13
14 typedef struct symcell {
15     SYM_INFO            *info;
16     struct symcell    *next;
17     } SYM_LIST;
18
19 typedef struct symtab {
20     struct symtab    *parent;
21     SYM_INFO        *function; /* enclosing this scope */
22     SYM_LIST        *list;
23     } SYM_TAB;
24
25 SYM_TAB*      symtab_open(SYM_TAB* enclosing_scope);
26 SYM_INFO*    symtab_find(SYM_TAB*, char*);
27 SYM_INFO*    symtab_insert(SYM_TAB*, char*, struct type_info*);
28 int          symtab_list_equal(SYM_LIST*, SYM_LIST*);
29 SYM_LIST*    symtab_list_insert(SYM_LIST*, SYM_INFO*);
30 SYM_INFO*    symtab_list_find(SYM_LIST*, char*);
31 void         symtab_list_release(SYM_LIST*);
32 SYM_INFO*    symtab_info_new(char*, struct type_info*);
33
34 void         symtab_print(FILE*, SYM_TAB*);
35 void         symtab_info_print(FILE*, SYM_INFO*);

```

```

36 void          symtab_list_print(FILE*,SYM_LIST*,char* separator);
37 #endif

1  /*          symtab.c(1.5)      10:07:40          97/12/10
2  *
3  *          Symbol table management
4  */
5  #include      <stdlib.h>          /* for malloc() and friends */
6  #include      <stdio.h>          /* for fprintf() and friends */
7
8  #include      "util.h"
9  #include      "symtab.h"
10 #include      "types.h"
11
12 SYM_TAB*
13 symtab_open(SYM_TAB * enclosing_scope)
14 {
15     SYM_TAB *st      = fmalloc(sizeof(SYM_TAB));
16
17     st->parent      = enclosing_scope;
18     st->list        = 0;
19     if (enclosing_scope)
20         st->function = enclosing_scope->function;
21     return st;
22 }
23
24 SYM_INFO*
25 symtab_find(SYM_TAB *st,char *name)
26 {
27     SYM_INFO      *i;
28
29     for ( ; (st); st = st->parent)
30         if ((i=symtab_list_find(st->list,name)))
31             return i;
32     return 0;
33 }
34
35 SYM_INFO*
36 symtab_insert(SYM_TAB *st,char *name,T_INFO* t)
37 {
38     SYM_INFO* i      = symtab_info_new(name,t);
39     st->list          = symtab_list_insert(st->list,i);
40
41     return i;
42 }
43
44 static int
45 symtab_info_equal(SYM_INFO* i1,SYM_INFO* i2)
46 {

```

```
47  /* rely on names and types being stored in a pool ``without duplicates`` */
48  return ((i1->name==i2->name)&&(i1->type==i2->type));
49  }
50
51  int
52  symtab_list_equal(SYM_LIST* l1, SYM_LIST* l2)
53  {
54      if (l1==l2)
55          return 1;
56
57      while (l1&&l2)
58          if (symtab_info_equal(l1->info, l2->info))
59              {
60                  l1 = l1->next;
61                  l2 = l2->next;
62              }
63          else
64              return 0;
65      if (l1) /* l2 == 0 */
66          return 0;
67      else
68          return (l1==l2);
69  }
70
71  SYM_LIST*
72  symtab_list_insert(SYM_LIST* l, SYM_INFO* i)
73  {
74      SYM_LIST* nl      = fmalloc(sizeof(SYM_LIST));
75      nl->info           = i;
76      nl->next           = l;
77      return nl;
78  }
79
80  SYM_INFO*
81  symtab_list_find(SYM_LIST* l, char* name)
82  {
83      for (; (l); l = l->next)
84          if (l->info->name==name) /* this works if all names in string pool */
85              return l->info;
86      return 0;
87  }
88
89  void
90  symtab_list_release(SYM_LIST* l)
91  {
92      if (l)
93          {
94              symtab_list_release(l->next);
95              free(l);
```

```

96         }
97     }
98
99     SYM_INFO*
100     symtab_info_new(char* name, T_INFO* t)
101     {
102         SYM_INFO* i      = fmalloc(sizeof(SYM_INFO));
103         i->name = name;
104         i->type = t;
105         return i;
106     }
107
108     void
109     symtab_info_print(FILE* f, SYM_INFO* info)
110     {
111         types_print(f, info->type);
112         fprintf(f, " %s", info->name);
113     }
114
115     void
116     symtab_list_print(FILE* f, SYM_LIST* l, char* separator)
117     {
118         while (l)
119         {
120             symtab_info_print(f, l->info);
121             if (l->next)
122                 fprintf(f, "%s", separator);
123             l = l->next;
124         }
125     }
126
127     void
128     symtab_print(FILE *f, SYM_TAB *tab)
129     {
130         if (!tab)
131             fprintf(f, "<null symtab>");
132         else
133             while (tab)
134             {
135                 symtab_list_print(f, tab->list, "\n");
136                 if (tab->parent)
137                     fprintf(f, "-----");
138                 tab = tab->parent;
139             }
140     }

```

C.4 Types library


```

1  #ifndef TYPES_H
2  #define TYPES_H
3  /*
4   *      types.h(1.4)      17:11:16      97/12/08
5   *
6   *      Type pool management
7   */
8  #include      "symtab.h"
9
10 typedef enum { int_t, float_t, fun_t, record_t, array_t } T_CONS; /* type construct
11
12 typedef struct {
13     SYM_LIST      *fields;
14     } T_RECORD;
15
16 typedef struct {
17     struct type_info      *target;
18     struct types_list      *source;
19     } T_FUN;
20
21 typedef struct array_type {
22     struct type_info      *base;
23     } T_ARRAY;
24
25
26 typedef struct type_info {
27     T_CONS      cons;
28     union {
29         T_FUN      fun;
30         T_ARRAY      array;
31         T_RECORD      record;
32     } info;
33     } T_INFO;
34
35 typedef struct types_list {
36     T_INFO      *type;
37     struct types_list      *next;
38     } T_LIST;
39
40 T_INFO* types_simple(T_CONS c);
41 T_INFO* types_fun(T_INFO* target, T_LIST *source);
42 T_INFO* types_record(SYM_LIST *fields);
43 T_INFO* types_array(T_INFO *base);
44 T_LIST* types_list_insert(T_LIST*, T_INFO*);
45 int      types_list_equal(T_LIST*, T_LIST*);
46 void      types_print(FILE*, T_INFO*);
47 void      types_list_print(FILE*, T_LIST*, char*);
48 void      types_list_release(T_LIST*);
49 void      types_print_all(FILE*);

```

```

50
51 #endif

1  /*      types.c(1.5)      17:11:14      97/12/08
2  *
3  *      Type pool management
4  */
5  #include      <stdlib.h>      /* for malloc() and friends */
6  #include      <stdio.h>      /* for fprintf() and friends */
7  #include      <assert.h>      /* for assert(condition) */
8
9  #include      "util.h"
10 #include      "symtab.h"
11 #include      "types.h"
12
13 static T_LIST  *types  = 0;
14
15 static int
16 types_equal(T_INFO *t1, T_INFO *t2)
17 {
18     if (t1==t2)
19         return 1;
20
21     if (t1->cons==t2->cons)
22         switch (t1->cons)
23         {
24             case int_t:
25             case float_t:
26                 return 1;
27             case fun_t:
28                 if (!types_equal(t1->info.fun.target, t2->info.fun.target))
29                     return 0;
30                 else
31                     return types_list_equal(t1->info.fun.source, t2->info.fun.source);
32             case array_t:
33                 return types_equal(t1->info.array.base, t2->info.array.base);
34             case record_t:
35                 return symtab_list_equal(t1->info.record.fields, t2->info.record.fields);
36             default:
37                 assert(0);
38         }
39     return 0;
40 }
41
42 void
43 types_list_release(T_LIST* l)
44 {
45     if (l) /* free memory held by list nodes, not by types in the nodes */
46         {

```

```
47         types_list_release(l->next);
48         free(l);
49     }
50 }
51
52 int
53 types_list_equal(T_LIST* t1,T_LIST *t2)
54 {
55     if (t1==t2)
56         return 1;
57
58     while (t1&& t2)
59         if (types_equal(t1->type,t2->type))
60             {
61                 t1      = t1->next;
62                 t2      = t2->next;
63             }
64         else
65             return 0;
66
67     if (t1) /* t2 == 0 */
68         return 0;
69     else /* t1 == 0 */
70         return (t1==t2);
71 }
72
73 static T_INFO*
74 types_find(T_INFO *t)
75 {
76     T_LIST *tl;
77
78     for (tl=types;(tl);tl=tl->next)
79         if (types_equal(t,tl->type))
80             return tl->type;
81     return 0;
82 }
83
84 static T_INFO*
85 types_new(T_CONS c)
86 {
87     T_INFO *t      = fmalloc(sizeof(T_INFO));
88     T_LIST *tl     = fmalloc(sizeof(T_LIST));
89     t->cons        = c;
90     tl->type       = t;
91     tl->next       = types;
92     types         = tl;
93     return t;
94 }
95
```

```
96 T_INFO*
97 types_simple(T_CONS c)
98 {
99     T_INFO t,*pt;
100     t.cons = c;
101
102     if ((pt=types_find(&t)))
103         return pt;
104     else
105         return types_new(c);
106 }
107
108 T_INFO*
109 types_fun(T_INFO* target,T_LIST *source)
110 {
111     T_INFO t,*pt;
112
113     t.cons = fun_t;
114     t.info.fun.target = target;
115     t.info.fun.source = source;
116
117     if (!(pt=types_find(&t)))
118     {
119         pt = types_new(fun_t);
120         pt->info.fun.source = source;
121         pt->info.fun.target = target;
122     }
123     else
124         types_list_release(source);
125     return pt;
126 }
127
128 T_INFO*
129 types_record(SYM_LIST *fields)
130 {
131     T_INFO t,*pt;
132
133     t.cons = record_t;
134     t.info.record.fields = fields;
135
136     if (!(pt=types_find(&t)))
137     {
138         pt = types_new(record_t);
139         pt->info.record.fields = fields;
140     }
141     else
142         symtab_list_release(fields);
143     return pt;
144 }
```

```

145
146 T_INFO*
147 types_array(T_INFO *base)
148 {
149     T_INFO  t,*pt;
150     t.cons      = array_t;
151     t.info.array.base      = base;
152
153     if (!(pt=types_find(&t)))
154     {
155         pt = types_new(array_t);
156         pt->info.array.base      = base;
157     }
158     return pt;
159 }
160
161 T_LIST*
162 types_list_insert(T_LIST* l,T_INFO *t)
163 {
164     T_LIST  *nl      = fmalloc(sizeof(T_LIST));
165     nl->type      = t;
166     nl->next      = l;
167     return nl;
168 }
169
170 void
171 types_print(FILE* f,T_INFO *t)
172 {
173     if (!t)
174         fprintf(f,"<null_type>");
175     else
176         switch (t->cons) {
177             case int_t:
178                 fprintf(f,"int");
179                 break;
180             case float_t:
181                 fprintf(f,"float");
182                 break;
183             case fun_t:
184                 types_print(f,t->info.fun.target);
185                 fprintf(f,"(");
186                 types_list_print(f,t->info.fun.source,",");
187                 fprintf(f,")");
188                 break;
189             case record_t:
190                 fprintf(f,"struct {");
191                 symtab_list_print(f,t->info.record.fields,";");
192                 fprintf(f,"}");
193                 break;

```

```

194         case array_t:
195             types_print(f,t->info.array.base);
196             fprintf(f,"*");
197             break;
198         default:
199             assert(0);
200     }
201     /* fprintf(f,"%x]",(unsigned int)t); */
202 }
203
204 void
205 types_list_print(FILE* f,T_LIST* tl,char* separator)
206 {
207     /* fprintf(f,"%x]",(unsigned int)tl); */
208     while (tl)
209     {
210         types_print(f,tl->type);
211         if (tl->next)
212             fprintf(f,"%s",separator);
213         tl = tl->next;
214     }
215 }
216
217 void
218 types_print_all(FILE* f)
219 {
220     types_list_print(f,types,"\n");
221 }

```

C.5 Type checking routines

```

1  #ifndef CHECK_H
2  #define CHECK_H
3  /*
4   *      check.h (1.3)      10:31:17      97/12/10
5   *
6   *      Semantic checks.
7   */
8  #include "types.h"
9  #include "symtab.h"
10
11 void      check_assignment(T_INFO*,T_INFO*);
12 T_INFO*   check_record_access(T_INFO* t,char* field);
13 T_INFO*   check_array_access(T_INFO* ta,T_INFO* ti);
14 T_INFO*   check_arith_op(int token,T_INFO* t1,T_INFO* t2);
15 T_INFO*   check_relop(int token,T_INFO* t1,T_INFO* t2);
16 T_INFO*   check_fun_call(SYM_TAB*,char*,T_LIST**);
17 SYM_INFO* check_symbol(SYM_TAB* scope,char* name);

```

```

18
19 #endif

1  /*
2  *      check.c(1.4)      17:46:20      97/12/10
3  */
4  #include      <stdio.h>          /* for fprintf() and friends */
5  #include      "check.h"
6
7  #include      "minic.tab.h"      /* for tokens */
8
9  extern int      lineno;          /* defined in minic.y */
10
11 static void
12 error(char *s1,char *s2,T_INFO* t1,char *s3,char *s4,T_INFO* t2)
13 {
14     fprintf(stderr,"type error on line %d: ",lineno);
15     if (s1) fprintf(stderr,"%s",s1);
16     if (s2) fprintf(stderr,"%s",s2);
17     if (t1) types_print(stderr,t1);
18     if (s3) fprintf(stderr,"%s",s3);
19     if (s4) fprintf(stderr,"%s",s4);
20     if (t2) types_print(stderr,t2);
21     fprintf(stderr,"\n");
22     exit(1);
23 }
24
25 void
26 check_assignment(T_INFO* tlexp,T_INFO* texp)
27 {
28     if (tlexp!=texp)
29         error("cannot assign ",0,texp," to ",0,tlexp);
30 }
31
32 T_INFO*
33 check_record_access(T_INFO* t,char *field)
34 {
35     SYM_INFO      *i;
36
37     if (t->cons!=record_t)
38         error("not a record: ",0,t," for field ",field,0);
39     if ((i=symtab_list_find(t->info.record.fields,field)))
40         return i->type;
41     error("record type ",0,t," has no field ",field,0);
42     return 0;
43 }
44
45 T_INFO*
46 check_array_access(T_INFO* ta,T_INFO* ti)

```

```

47 {
48 if (ta->cons!=array_t)
49     error("not an array type: ",0,ta,0,0,0);
50 if (ti->cons!=int_t)
51     error("index for ",0,ta," must be integer, not ",0,ti);
52 return ta->info.array.base;
53 }
54
55 T_INFO*
56 check_arith_op(int token,T_INFO* t1,T_INFO* t2)
57 {
58 if (t1!=t2)
59     error("type ",0,t1," does not match ",0,t2);
60 if ((t1->cons!=int_t)&&(t2->cons!=float_t))
61     error("type ",0,t1," is not numeric",0,0);
62 return t1;
63 }
64
65 T_INFO*
66 check_relop(int token,T_INFO* t1,T_INFO* t2)
67 {
68 if (t1!=t2)
69     error("type ",0,t1," does not match ",0,t2);
70 return types_simple(int_t);
71 }
72
73 SYM_INFO*
74 check_symbol(SYM_TAB* scope,char* name)
75 {
76 SYM_INFO* i = symtab_find(scope,name);
77
78 if (!i)
79     error("undeclared variable \"",name,0,"\"",0,0);
80 return i;
81 }
82
83 T_INFO*
84 check_fun_call(SYM_TAB* scope,char* name,T_LIST** args)
85 {
86 SYM_INFO*      i = symtab_find(scope,name);
87 T_INFO*        ft;
88
89 if (!i)
90     error("undeclared function \"",name,0,"\"",0,0);
91
92 ft = i->type;
93
94 if (ft->cons!=fun_t)
95     error(name," is not a function",0,0,0,0);

```



```

96
97 if (types_list_equal(ft->info.fun.source, (args?*args:0)))
98     {
99         /* release type_list from args, replace by equal list
100         * from function type
101         */
102         if (args) {
103             types_list_release(*args);
104             *args = ft->info.fun.source;
105         }
106         return ft->info.fun.target;
107     }
108 error("bad type of arguments for ",name,0,0,0,0);
109 return 0;
110 }

```

C.6 Parser with semantic actions

```

1  %{
2  /*      minic.y(1.9)      17:46:21      97/12/10
3  *
4  *      Parser demo of simple symbol table management and type checking.
5  */
6  #include      <stdio.h>      /* for (f)printf() */
7
8  #include      "syntab.h"
9  #include      "types.h"
10 #include      "check.h"
11
12 int          lineno = 1;      /* number of current source line */
13 extern int    yylex();        /* lexical analyzer generated from lex.l */
14 extern char   *yytext;        /* last token, defined in lex.l
15 */
16 SYM_TAB      *scope;          /* current symbol table, initialized in lex.l */
17 char          *base;          /* basename of command line argument */
18
19 void
20 yyerror(char *s)
21 {
22     fprintf(stderr, "Syntax error on line %d: %s\n", lineno, s);
23     fprintf(stderr, "Last token was \"%s\"\n", yytext);
24     exit(1);
25 }
26 %}
27
28 %union {
29     char*      name;

```

```

30      int          value;
31      T_LIST*      tlist;
32      T_INFO*      type;
33      SYM_INFO*    sym;
34      SYM_LIST*    slist;
35  }
36
37  %token  INT FLOAT NAME STRUCT IF ELSE RETURN NUMBER LPAR RPAR LBRACE RBRACE
38  %token  LBRACK RBRACK ASSIGN SEMICOLON COMMA DOT PLUS MINUS TIMES DIVIDE EQUAL
39
40  %type   <name>   NAME
41  %type   <value>  NUMBER
42  %type   <type>   type parameter exp lexp
43  %type   <tlist>  parameters more_parameters exps
44  %type   <sym>    field var
45  %type   <slist>  fields
46
47  /*      associativity and precedence: in order of increasing precedence */
48
49  %nonassoc      LOW  /* dummy token to suggest shift on ELSE */
50  %nonassoc      ELSE /* higher than LOW */
51
52  %nonassoc      EQUAL
53  %left          PLUS      MINUS
54  %left          TIMES     DIVIDE
55  %left          UMINUS    /* dummy token to use as precedence marker */
56  %left          DOT       LBRACK /* C compatible precedence rules */
57
58  %%
59  program          : declarations
60                  ;
61
62  declarations     : declaration declarations
63                  | /* empty */
64                  ;
65
66  declaration      : fun_declaration
67                  | var_declaration
68                  ;
69
70  fun_declaration  : type NAME { /* this is $3 */
71                      <sym>$ = symtab_insert(scope,$2,0);
72                      scope = symtab_open(scope); /* open new scope */
73                      scope->function = <sym>$; /* attach to this function */
74                      }
75                  LPAR parameters RPAR { /* this is $7 */
76                      <sym>3->type = types_fun($1,$5);
77                      }
78                  block { scope = scope->parent; }

```

```

79          ;
80
81 parameters      : more_parameters      { $$ = $1; }
82                  |                      { $$ = 0; }
83                  ;
84
85 more_parameters : parameter COMMA more_parameters
86                  { $$ = types_list_insert($3,$1); }
87                  | parameter          { $$ = types_list_insert(0,$1); }
88                  ;
89
90 parameter       : type NAME {
91                   symtab_insert(scope,$2,$1); /* insert in symbol table */
92                   $$ = $1; /* remember type info */
93                   }
94                   ;
95
96 block           : LBRACE                { scope = symtab_open(scope); }
97                   var_declarations statements RBRACE
98                   { scope = scope->parent; /* close scope */ }
99                   ;
100
101 var_declarations: var_declaration var_declarations
102                  |
103                  ;
104
105 var_declaration : type NAME SEMICOLON   { symtab_insert(scope,$2,$1); }
106                  ;
107
108 type            : INT                   { $$ = types_simple(int_t); }
109                  | FLOAT                { $$ = types_simple(float_t); }
110                  | type TIMES           { $$ = types_array($1); }
111                  | STRUCT LBRACE fields RBRACE /* record */
112                  { $$ = types_record($3); }
113                  ;
114
115 fields          : field fields          { $$ = symtab_list_insert($2,$1); }
116                  |                      { $$ = 0; }
117                  ;
118
119 field           : type NAME SEMICOLON   { $$ = symtab_info_new($2,$1); }
120                  ;
121
122 statements      : statement SEMICOLON statements
123                  | /* empty */
124                  ;
125
126 statement       : IF LPAR exp RPAR statement %prec LOW
127                  | IF LPAR exp RPAR statement ELSE statement /* shift */

```

```

128         | lexp ASSIGN exp          { check_assignment($1,$3); }
129         | RETURN exp /* statements always in scope with function */
130         | { check_assignment(scope->function->type->info.fun.target,$1,$3); }
131         | block
132         ;
133
134 lexp      : var                    { $$ = $1->type; }
135         | lexp LBRACK exp RBRACK { $$ = check_array_access($1,$3); }
136         | lexp DOT NAME           { $$ = check_record_access($1,$3); }
137         ;
138
139 exp       : exp DOT NAME           { $$ = check_record_access($1,$3); }
140         | exp LBRACK exp RBRACK { $$ = check_array_access($1,$3); }
141         | exp PLUS exp             { $$ = check_arith_op(PLUS,$1,$3); }
142         | exp MINUS exp             { $$ = check_arith_op(MINUS,$1,$3); }
143         | exp TIMES exp            { $$ = check_arith_op(TIMES,$1,$3); }
144         | exp DIVIDE exp           { $$ = check_arith_op(DIVIDE,$1,$3); }
145         | exp EQUAL exp            { $$ = check_relop(EQUAL,$1,$3); }
146         | LPAR exp RPAR { $$ = $2; }
147         | MINUS exp %prec UMINUS /* this will force a reduce */
148         | { $$ = check_arith_op(UMINUS,$2,0); }
149         | var                      { $$ = $1->type; }
150         | NUMBER                   { $$ = types_simple(int_t); }
151         | NAME LPAR RPAR           { $$ = check_fun_call(scope,$1,0); }
152         | NAME LPAR exps RPAR      { $$ = check_fun_call(scope,$1,&$3); }
153         ;
154
155 exps      : exp                    { $$ = types_list_insert(0,$1); }
156         | exp COMMA exps           { $$ = types_list_insert($3,$1); }
157         ;
158
159 var       : NAME                   { $$ = check_symbol(scope,$1); }
160 %%
161
162 int
163 main(int argc,char *argv[])
164 {
165     if (argc!=2) {
166         fprintf(stderr,"Usage: %s base_file_name",argv[0]);
167         exit(1);
168     }
169     base = argv[1];
170     return yyparse();
171 }

```

C.7 Utilities

```
1 #ifndef UTIL_H
```

```

2 #define UTIL_H
3 /*
4 *      util.h(1.1)      11:24:32      97/12/08
5 *
6 *      General utility functions
7 */
8 #include      <stdlib.h>
9
10 void      *fmalloc(size_t);      /* malloc() version that aborts on failure */
11
12 #endif

1 /*
2 *      util.c(1.2)      10:13:01      97/12/10
3 *
4 *      Utility functions.
5 */
6
7 #include      <stdio.h>      /* for fprintf() and friends */
8
9 #include      "util.h"
10
11 /* aborts on failure */
12 void      *fmalloc(size_t s)
13 {
14 void      *ptr = malloc(s);
15 if (!ptr) {
16         fprintf(stderr, "Out of memory in fmalloc(%d)\n", s);
17         exit(1);
18     }
19 return ptr;
20 }

```

C.8 Driver script

```

1 #!/bin/sh
2 #
3 #      Usage:  mct basename
4 #
5 #      e.g. "mct tst"  to compile tst.c,
6 #
7 #      No output except errors on stderr.
8 #
9 #      First make sure source file exists
10 #
11 if [ ! -f "$1.c" ]
12 then
13     echo "Cannot open \"$1.c\""

```

```

14         exit 1
15 fi
16 #
17 minic $1 <$1.c

```

C.9 Makefile

```

1  #           %M%(%I%)           %U%           %E%
2  #
3  CFLAGS=      -Wall -g
4  CC=          gcc
5  #
6  HH_FILES=    util.h names.h symtab.h types.h check.h
7  CC_FILES=    util.c names.c symtab.c types.c check.c
8  SOURCE=      Makefile $(CC_FILES) $(HH_FILES) minic.y lex.l mct.sh
9  #
10 #
11 all:         minic mct
12 minic:       minic.tab.o lex.yy.o $(CC_FILES:%.c=%.o)
13             gcc $(CFLAGS) -o $@ $^
14
15 lex.yy.c:    lex.l minic.tab.h
16             flex lex.l
17 #
18 include      make_dependencies
19 #
20 #           Bison options:
21 #
22 #           -v           Generate minic.output showing states of LALR parser
23 #           -d           Generate minic.tab.h containing token type definitions
24 #
25 minic.tab.h\
26 minic.tab.c: minic.y
27             bison -v -d $^
28 #
29 clean:
30             rm -f lex.yy.c minic.tab.[hc] *.o mct minic *.jasm *.class minic.o
31 #
32 tar:
33             tar cvf minic.tar $(SOURCE)
34 source:
35             @echo $(SOURCE)
36 #
37 ##### creating dependencies
38 #
39 dep:         lex.yy.c minic.tab.c $(CC_FILES)
40             $(CC) -M $(CC_FILES) lex.yy.c minic.tab.c >make_dependencies
41 #

```

```
42 ##### sccs rules
43 #
44 delta:
45     @echo "Comment: \c"; read inp; \
46     for f in $(SOURCE); do\
47         [ -f SCCS/p.$$f ] && \
48         { echo "##$f:"; sccs delget -y\"$$inp\" $$f; } done;\
49     true
50 edit:
51     @for f in $(SOURCE); do\
52         [ ! -f SCCS/p.$$f ] && \
53         { echo "##$f:"; sccs edit $$f; } done; true
54 create:
55     @for f in $(SOURCE); do\
56         [ ! -f SCCS/s.$$f ] && \
57         { echo "##$f:"; sccs create $$f; } done; rm -f ,*; true
```

Index

- \Rightarrow_G^n , 36
- \Rightarrow_G^* , 36
- \Rightarrow_G , 36
- $|w|$, 24
- $pref_l(w)$, 24
- \vdash_M^* , 26, 29
- \vdash_M , 26, 29
- δ^* , 26
- ϵ -closure, 30
- ϵ -move, 28
- $\$1$, 76
- $\$\$$, 76
- $GEN_e(B)$, 107
- $KILL_{use}(B)$, 112
- $KILL_{live}(B)$, 110
- $GEN(B)$, 103
- $IN_e(B)$, 107
- $IN_{use}(B)$, 112
- $IN_{live}(B)$, 110
- $IN(B)$, 103
- $KILL_e(B)$, 107
- $KILL(B)$, 103
- $OUT_e(B)$, 107
- $OUT_{use}(B)$, 112
- $OUT_{live}(B)$, 110
- $OUT(B)$, 103
- $UD(a, u)$, 105
- $GEN_{use}(B)$, 112
- $GEN_{live}(B)$, 110
- absolute code, 130
- abstract syntax trees, 76
- action table, 53
- activation record, 132
- algorithm optimization, 93
- aliases, 126, 129
- aliasing, 102
- alphabet, 24
- assembler, 7
- assembler source, 130
- attribute, 65
- attribute grammar, 65
- attributes
 - inherited, 66
 - synthesized, 66
- available expression, 106
- awk, 8
- back edge, 117
- back end, 74
- backpatching, 83
- basic block, 94
- basic induction variable, 123
- bison
 - attributes, 68
- cfg
 - see context-free grammar, 35
- cfront, 7
- code simplification, 99
- code-generator generator, 78
- compiler, 6
- computation by a DFA, 26
- configuration, 26
- conflict
 - reduce-reduce, 57
 - shift-reduce, 57
- context-free grammar, 14, 35

- derives, 36
- generated language, 36
- left recursive, 41
- left sentential form, 38
- leftmost derivation, 38
- LL(1), 43
- LL(k), 44
- LR(1), 59
- right sentential form, 50
- rightmost derivation, 50
- context-free language, 36
- cross-compiler, 7
- dag, 95
- dead variable, 110
- definition-use chain, 112
- dependency graph, 67
- derivation
 - leftmost, 38
 - rightmost, 50
- deterministic finite automaton, 26
- DFA, 26
 - computation, 26
 - computation step, 26
 - configuration, 26
- directed acyclic graph, 95
- dominate, 117
- empty string, 24
- exit node, 117
- expression
 - available, 106
- family of an induction variable, 123
- finite automaton
 - deterministic, 26
 - nondeterministic, 28
- first, 45
- flow graph, 94
- follow, 45
- formal language, 7, 24
- front end, 74
- goto table, 53
- gprof, 116
- grammar, 10
 - attribute, 65
- handle, 50
- header of a loop, 117
- immediate left recursion, 41
- induction variable, 123
 - basic, 123
 - family, 123
- inherently ambiguous languages, 38
- inherited attributes, 66
- initial basic block, 102
- inner loop, 117
- interpreter, 8
- item
 - LR(1), 54
- Java virtual machine, 8
- JVM, 8
- language
 - formal, 7, 24
 - little, 8
 - regular, 24
 - source, 6
 - target, 6
- languages
 - inherently ambiguous, 38
- leader, 94
- left recursion, 41
 - immediate, 41
- left sentential form, 38
- leftmost derivation, 38
- length of string, 24
- lexical analyze, 12
- lisp, 8
- little language, 8
- live variable, 110
- LL(1) grammar, 43

- LL(1) parse table construction, 45
- LL(k) grammars, 44
- loop
 - header, 117
 - inner, 117
 - invariant, 120
- LR(1) grammar, 59
- LR(1) item, 54
- lvalue, 89
- Micro, 9
- Minic, 72
- natural loop, 117
- NFA, 28
- nondeterministic finite automaton, 28
- nonterminal symbol, 35
- parse tree, 13, 37
- parser, 13
 - predictive, 44
- parsing
 - predictive, 38
 - top-down, 38
- path between points, 102
- phrase, 50
 - simple, 50
- point, 102
- predictive parser, 44
 - operation, 45
- predictive parsing, 38
- prefix, 24
- preheader, 121
- product, 24
- production, 35
- prolog, 8
- reduce operation, 51
- reduce-reduce conflict, 57
- register allocation, 17
- regular expressions, 24
- regular language, 24
- relocatable code, 130
- retargeting, 74
- right sentential form, 50
- rightmost derivation, 50
- scanner, 12
- scope, 69
- semantic action, 15, 65
- semantics
 - static, 71
- shift operation, 51
- shift-reduce conflict, 57
- simple phrase, 50
- simple postfix, 76
- Smalltalk, 8
- source language, 6
- SQL, 8
- start symbol, 35
- state, 26, 28
 - final, 26
 - initial, 26
- static semantics, 71
- strength reduction, 124
- successor, 102
- symbol table, 12
- syntax, 7
- syntax-directed translation, 15, 67
- synthesized attributes, 66
- target language, 6
- terminal symbol, 35
- three-address code, 15
 - instructions, 78
- token, 12
- top-down parsing, 38
- topological sort, 67
- transition function, 26
- type rules, 72
- type system, 71
- use-definition chain, 105
- used variable, 110

viable prefix, 50

x86, 10

yacc

 attributes, 68

yield, 14

List of Figures

| | | |
|-----|--|----|
| 1.1 | A source text in the C language | 6 |
| 1.2 | The syntax of the Micro language | 9 |
| 1.3 | A Micro program | 10 |
| 1.4 | X86 assembly code generated for the program in Figure 1.3 | 11 |
| 1.5 | Result of lexical analysis of program in Figure 1.3 | 13 |
| 1.6 | Parse tree of program in Figure 1.3 | 14 |
| 1.7 | three-address code corresponding to the program of Figure 1.3 . . | 16 |
| 1.8 | Optimized three-address code corresponding to the program of Figure 1.3 | 17 |
| 2.1 | A declaration for TOKEN and lex() | 18 |
| 2.2 | The transition diagram for <i>lex()</i> | 23 |
| 2.3 | A DFA for <i>NAME</i> | 27 |
| 2.4 | DFA implementation | 28 |
| 2.5 | M_1 | 29 |
| 2.6 | M_\emptyset , M_ϵ and M_a | 31 |
| 2.7 | $M_{r_1+r_2}$, $M_{r_1r_2}$ and $M_{r_1}^*$ | 32 |
| 2.8 | A generated scanner | 34 |
| 3.1 | Parse trees in the ambiguous context-free grammar from Example 5 | 37 |
| 3.2 | A simple top-down parse | 39 |
| 3.3 | Eliminating immediate left recursion | 42 |
| 3.4 | The prediction problem | 43 |
| 3.5 | A predictive parser | 44 |

| | | |
|------|---|-----|
| 3.6 | Predictive parser operation | 46 |
| 3.7 | Naive shift-reduce parsing procedure | 53 |
| 3.8 | A shift-reduce parser | 54 |
| 3.9 | Shift-reduce parser algorithm | 55 |
| 3.10 | A (partial) rightmost derivation | 56 |
| 3.11 | States of the viable prefix DFA of G_3 | 58 |
| 3.12 | Transition function of the viable prefix DFA of G_3 | 59 |
| 3.13 | LR(1) tables for G_3 | 60 |
| 4.1 | Computing synthesized attribute values | 66 |
| 4.2 | Computing inherited attribute values | 68 |
| 4.3 | A simple symbol table using a string pool and hashing | 70 |
| 4.4 | A stacked symbol table | 71 |
| 5.1 | Abstract syntax tree corresponding to the program in Figure 1.3 | 77 |
| 5.2 | Tree-rewriting rules | 78 |
| 5.3 | Three-address code instructions | 79 |
| 5.4 | Conditional control flow statements | 83 |
| 5.5 | Exit sets after “and” | 84 |
| 6.1 | Flow graph of example code | 95 |
| 6.2 | Dag of example code | 98 |
| 6.3 | Dag of code with array manipulation | 100 |
| 6.4 | Corrected dag of code with array manipulation | 101 |
| 6.5 | Dag using algebraic identities | 102 |
| 6.6 | Flow graph with cycle | 104 |
| 6.7 | Example flow graph for available expressions | 108 |
| 6.8 | Code before first subexpression elimination | 113 |
| 6.9 | Code after first subexpression elimination | 114 |
| 6.10 | Flow graph with loops | 117 |
| 6.11 | Flow graph of Example 20 | 121 |
| 6.12 | Flow graph of Example 20 after moving loop invariants | 122 |

| | | |
|------|---|-----|
| 6.13 | Flow graph of Example 20 after strength reduction | 125 |
| 6.14 | Flow graph of Figure 6.13 after copy propagation and constant folding | 126 |
| 6.15 | optimized flow graph for Example 20 | 127 |
| 7.1 | Memory layout in C | 131 |
| 7.2 | Activation records | 132 |
| 7.3 | Multiplication on IBM/370 | 135 |

List of Tables

| | | |
|-----|---|-----|
| 2.1 | Regular expressions describing Micro tokens | 25 |
| A.1 | Some x86 Flags | 140 |
| A.2 | Data Transfer Instructions | 141 |
| A.3 | Integer Arithmetic Instructions | 142 |
| A.4 | Linux System Calls | 145 |
| A.5 | Logical Operations | 145 |
| A.6 | Control Flow Instructions | 146 |

Bibliography

- [ALSU07] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, techniques and tools, second edition*. Pearson, Addison Wesley, 2007.
- [ASU86] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, techniques and tools*. Addison Wesley, 1986.
- [Bar03] Jonathan Bartlett. *Programming from the Ground Up*. <http://www.bartlettpublishing.com/> and <http://savannah.nongnu.org/projects/pgubook/>, 2003.
- [FL91] Charles N. Fischer and Richard J. LeBlanc. *Crafting a compiler with C*. Addison-Wesley, 1991.
- [HU69] John E. Hopcroft and Jeffrey D. Ullman. *Formal languages and their relation to automata*. Addison Wesley, 1969.
- [Sal73] A. Salomaa. *Formal languages*. ACM Monograph Series. Academic Press, 1973.