

This document contains solutions to the exercises in Chapters 1 through 8 of *Engineering a Compiler, 2nd Edition*.
The solution for exercise 7.13(a) is incomplete.

Engineering a Compiler, 2e

Solutions to Exercises, V1.0
Chapters 1 through 8

Keith D. Cooper
Linda Torczon

Rice University
Houston, Texas

Limited Copies Distributed
Reproduction requires explicit permission
Copyright 2011, Morgan-Kaufmann Publishers and the Authors
All rights reserved

CHAPTER 1

The questions for Chapter 1 are all “soft” questions, in the sense that they have no single correct answer. Instead, they are intended to provoke thought and discussion. Our answers are more detailed than we would expect from a beginning student.

1. Consider a simple web browser that takes as input a textual string in HTML format and displays the specified graphics on the screen. Is the display process one of compilation or interpretation?

Answer: A browser can be viewed as a compiler that translates a web page from a textual language to a graphical language. It can also be viewed as an interpreter because it parses the input every time that it is displayed. Because the web page, viewed as a program, takes no input, the distinction between compilation and interpretation is blurred.

2. In designing a compiler, you will face many tradeoffs. What are the five qualities that you, as a user, consider most important in a compiler that you purchase? Does that list change when you are the compiler writer? What does your list tell you about a compiler that you would implement?

Answer: The answers are subjective, so no single set of answers is correct. The authors’ subjective answers are:

From the buyer’s perspective:

- (a) Compiled code is correct.
- (b) Compiled code adheres to system-wide calling conventions—that is, the code it produces operates correctly with code produced by other compilers and, if possible, other languages.
- (c) Compiler produces efficient code for my applications.
- (d) Compiler supports one or more of the standard interfaces for debug information, such as the DWARF standard.
- (e) Compiler produces clear, well-localized error messages.

From the compiler writer’s perspective:

- (a) Compiled code is correct.
- (b) Compiler produces efficient code for a broad variety of input programs.
- (c) Debugging individual passes in the compiler is relatively easy.
- (d) The compiler’s intermediate representation is human-readable.
- (e) The compiler’s design makes it easy to add and replace passes.

Insights lists: Optimization matters, as does compile time.

3. Compilers are used in many different circumstances. What differences might you expect in compilers designed for the following applications?

Many answers are possible.

- (a) A *just-in-time* compiler used to translate user interface code downloaded over a network

Answer: The time constraints on a just-in-time compiler make it unlikely that the compiler can afford to make many passes over the code. This may limit the amount of optimization that it can do.

- (b) A compiler that targets the embedded processor used in a cellular telephone

Answer: A compiler for an embedded processor may focus on issues other than speed. In particular, generating code that is small and code that uses less power may be more important. However, it is not clear that speed and low power are, necessarily, competing issues.

- (c) A compiler used in an introductory programming course at a high school

Answer: A compiler for an introductory programming course should give detailed feedback on syntax errors. It may also perform analysis to report on semantic errors that are not enforced by the language standard.

- (d) A compiler used to build wind-tunnel simulations that run on a massively parallel processor (where all processors are identical)

Answer: A compiler for a massively parallel machine will concentrate on performance issues, especially finding sections of code that can execute in parallel.

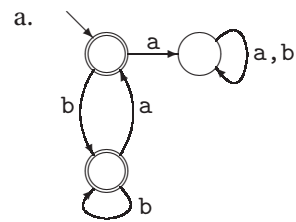
- (e) A compiler that targets numerically intensive programs to a large number of diverse machines

Answer: A multi-platform compiler for numerically-intensive codes will focus on retargetability in its code generator and on machine-independent optimizations—particularly those that help improve array-element references in nested loops.

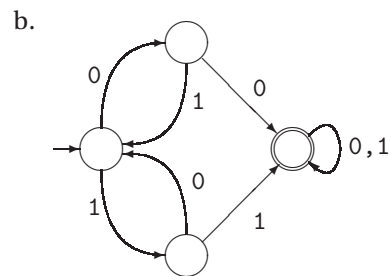
CHAPTER 2

Section 2.2

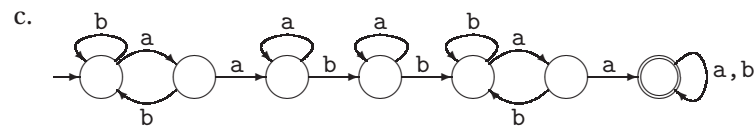
1. Describe informally the languages accepted by the following FAs:



Answer: The set of strings in which every 'a' is preceded by a 'b'



Answer: The set of strings that are not a concatenation of the strings '01' or '10'

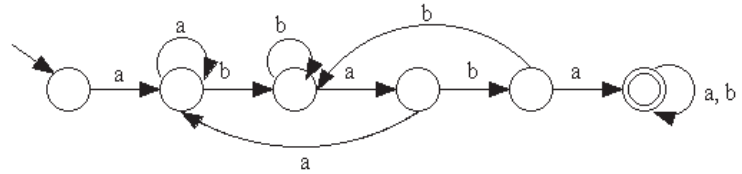


Answer: The set of strings that contain 'aab' as a substring and then 'baa' as a substring (not necessarily concatenated one with the other).

2. Construct an FA accepting each of the following languages:

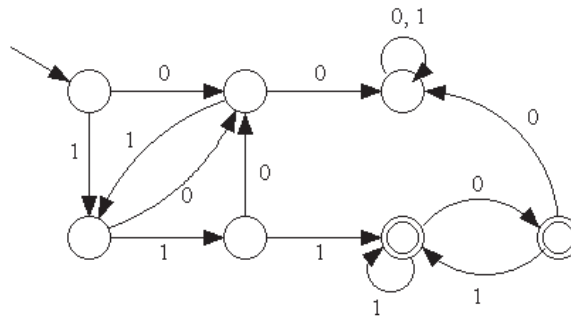
- a. $\{w \in \{a, b\}^* \mid w \text{ starts with 'a' and contains 'baba' as a substring}\}$

Answer:



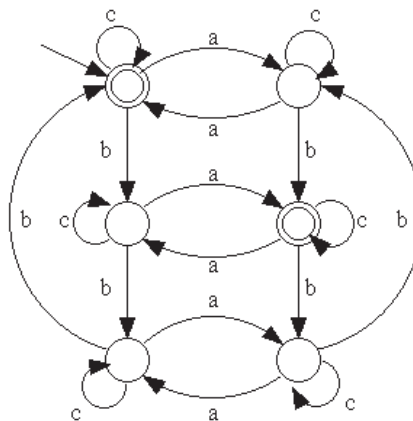
- b. $\{w \in \{0, 1\}^* \mid w \text{ contains '111' as a substring and does not contain '00' as a substring}\}$

Answer:



- c. $\{w \in \{a, b, c\}^* \mid \text{in } w \text{ the number of 'a's modulo 2 is equal to the number of 'b's modulo 3}\}$

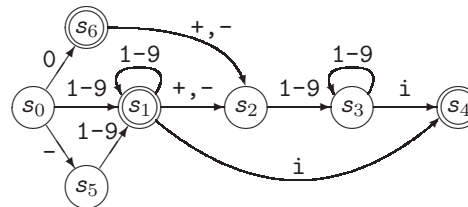
Answer:



3. Create FAs to recognize

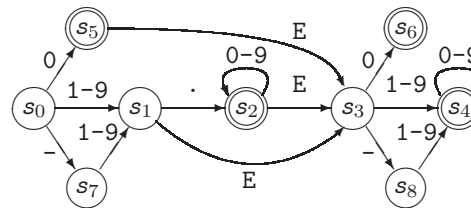
- a. words that represent complex numbers

Answer:



- b. words that represent decimal numbers written in scientific notation.

Answer:



Section 2.3

4. Different programming languages use different ways to represent integers. Construct a regular expression for each one of the following:

- a. Nonnegative integers in C represented in bases 10 and 16

Answer: $[0-9][0-9]^* | 0x0 | 0x([1-9][a-f][A-F])([0-9][a-f][A-F])^*$

- b. Nonnegative integers in VHDL that may include underscores

(An underscore cannot occur as the first or last character.)

Answer: $[0-9](_?[0-9])^*$

VHDL allows leading zeroes.

- c. Currency, in dollars, represented as a positive decimal number rounded to the nearest one-hundredth. Such numbers begin with the character \$, have commas separating each group of three digits to the left of the decimal point, and end with two digits to the right of the decimal point, for example, \$8,937.43 and \$7,777,777.77

Answer:

$\$([1-9](\epsilon|[0-9])(\epsilon|[0-9])(\epsilon|[0-9][0-9][0-9])^*)|0)\underline{[0-9][0-9]}$

(Underlines added to emphasize punctuation marks.)

Hint: not all the specifications describe regular languages.

5. Write a regular expression for each of the following languages.
- a. Given an alphabet $\Sigma = \{0, 1\}$, L is the set of all strings of alternating pairs of 0s and pairs of 1s.

Answer: $(1100(1100)^*(11|\epsilon))|(0011(0011)^*(00|\epsilon))$

- b. Given an alphabet $\Sigma = \{0, 1\}$, L is the set of all strings of 0s and 1s that contain an even number of 0s or an even number of 1s.

Answer: $(1^*01^*01^*)^+|(0^*10^*10^*)^+$

- c. Given the lowercase English alphabet, L is the set of all strings in which the letters appear in ascending lexicographical order.

Answer: $a^*b^*c^*d^*e^*f^*g^*h^*i^*j^*k^*l^*m^*n^*o^*p^*q^*r^*s^*t^*u^*v^*w^*x^*y^*z^*$

- d. Given an alphabet $\Sigma = \{a, b, c, d\}$, L is the set of strings $xyzwy$, where x and w are strings of one or more characters in Σ , y is any single character in Σ , and z is the character z , taken from outside the alphabet.

(Each string $xyzwy$ contains two words xy and wy built from letters in Σ . The words end in the same letter, y . They are separated by z .)

Answer: $(a|b|c|d)(a|b|c|d)^*a z (a|b|c|d)(a|b|c|d)^*a |$
 $(a|b|c|d)(a|b|c|d)^*b z (a|b|c|d)(a|b|c|d)^*b |$
 $(a|b|c|d)(a|b|c|d)^*c z (a|b|c|d)(a|b|c|d)^*c |$
 $(a|b|c|d)(a|b|c|d)^*d z (a|b|c|d)(a|b|c|d)^*d$

- e. Given an alphabet $\Sigma = \{+, -, \times, \div, (,), \text{id}\}$, L is the set of algebraic expressions using addition, subtraction, multiplication, division, and parentheses over id 's.

Answer: This language cannot be described in a regular expression because it requires matching parentheses.

6. Write a regular expression to describe each of the following programming language constructs:

- a. Any sequence of tabs and blanks (sometimes called *white space*)

Answer: $(b|t)^*$

- b. Comments in the programming language C

Answer: $/ * [(\Sigma - *)^* (*^* (\Sigma - \{ /, * \}))^*]^* *^* */$

where Σ represents the alphabet and $(\Sigma - x)$ represents the alphabet excluding the character 'x'.

- c. String constants (without escape characters)

Answer: $\|(\Sigma - \lambda I)^k\|$

- d. Floating-point numbers

Answer: The regular expression can become quite complex, depending on the specific notation used for exponents. See, for example, the FA in Section 2.2, question 3.b.

For simple decimal numbers, without scientific notation, the following RE would work:

$$(0 \mid [1 - 9][0 - 9]^*) \cdot [0 - 9]^*$$

It allows trailing zeros to indicate precision.

Section 2.4

7. Consider the three regular expressions:

$$(ab \mid ac)^*$$

$$(0 \mid 1)^* 1100 \ 1^*$$

$$(01 \mid 10 \mid 00)^* 11$$

We present the solution for each RE in its entirety.

Answer: For the RE $(ab \mid ac)^*$:

Part (a): Build an NFA for the RE using Thompson's construction.

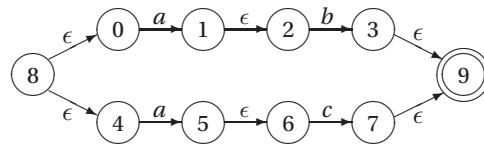
First, build DFAs for a , b , a , and c .



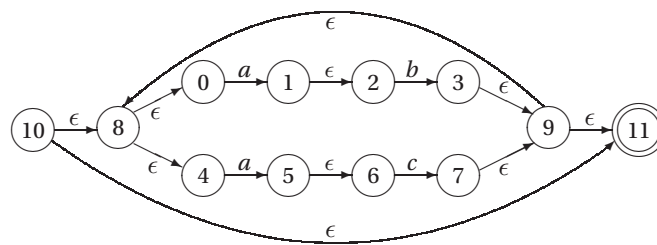
Next, combine them into NFAs for ab and ac



Combine the NFAs for ab and ac to form an NFA for $ab \mid ac$.



Finally, add two states and the ϵ -transitions for the closure to form an NFA for $(ab|ac)^*$.



Note: the NFAs can be built in any order that is consistent with the precedence among RE operators.

Part (b): Convert the NFA to a DFA using the subset construction.

Initialization:

Note: the “core” states, those generated by δ , are set in bold. The states added by ϵ -closure are set in normal face.

The NFA’s start state is 10.

$$s_0 \leftarrow \epsilon\text{-closure}(10) = \{\mathbf{10}, 8, 0, 4, 11\}$$

Round 1:

$$s_1 \leftarrow \epsilon\text{-closure}(\delta(s_0, a)) = \{\mathbf{1}, \mathbf{5}, 2, 6\}$$

There are no transitions out of s_0 on either b or c .

Round 2:

There is no transition out of s_1 on a .

$$s_2 \leftarrow \epsilon\text{-closure}(\delta(s_1, b)) = \{\mathbf{3}, 9, 8, 0, 4, 11\}$$

$$s_3 \leftarrow \epsilon\text{-closure}(\delta(s_1, c)) = \{\mathbf{7}, 9, 8, 0, 4, 11\}$$

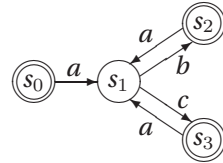
Round 3:

$$\epsilon\text{-closure}(\delta(s_2, a)) = s_1$$

$$\epsilon\text{-closure}(\delta(s_3, a)) = s_1$$

There are no transitions out of either s_2 or s_3 on b or c .

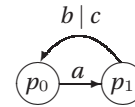
The worklist is empty. The resulting DFA has four states:



Part (c): Minimize the DFA:

Current Partition	Transition	Action
$\{ p_0 = \{s_0, s_2, s_3\}, p_1 = \{s_1\} \}$	—	—
	$\delta(p_0, a) = p_1$	none
	$\delta(p_0, b) = \emptyset$	none
	$\delta(p_0, c) = \emptyset$	none
	$\delta(p_1, a) = \emptyset$	none
	$\delta(p_1, b) = p_0$	none
	$\delta(p_1, c) = p_0$	none

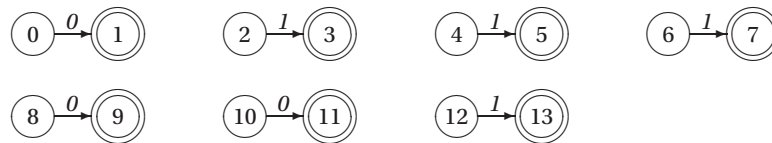
None of the potential transitions splits a set, so the minimal DFA has two states, as shown to the right.



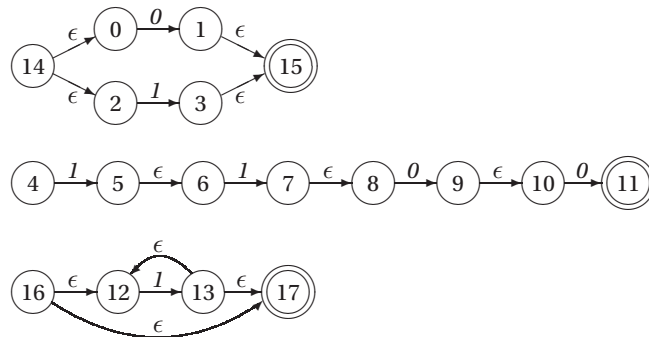
Answer: For the RE $(0|1)^* 1100 1^*$:

Part (a): Build an NFA for the RE using Thompson's construction.

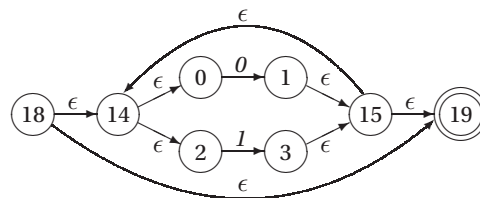
First build distinct DFAs for each of the digits:



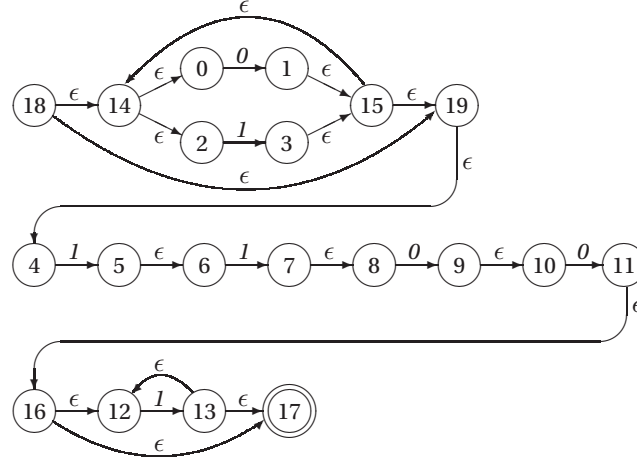
Next, combine these DFAs into NFAs for $0|1$, 1100 , and 1^* :



Next, build the NFA for $(0|1)^*$:



Now, connect the NFAs to form a single NFA for $(0|1)^* 1100 1^*$:



Part (b): Convert the NFA to a DFA using the subset construction.

Initialization:

The NFA's start state is 18.

$$s_0 \leftarrow \epsilon\text{-closure}(18) = \{\mathbf{18}, 14, 0, 2, 19, 4\}$$

Round 1:

$$s_1 \leftarrow \epsilon\text{-closure}(\delta(s_0, 0)) = \{\mathbf{1}, 15, 14, 0, 2, 19, 4\}$$

$$s_2 \leftarrow \epsilon\text{-closure}(\delta(s_0, 1)) = \{\mathbf{3}, \mathbf{5}, 15, 14, 0, 2, 19, 4, 6\}$$

Round 2:

$$\epsilon\text{-closure}(\delta(s_1, 0)) = \{\mathbf{1}, \dots\} = s_1$$

$$\epsilon\text{-closure}(\delta(s_1, 1)) = \{\mathbf{3}, \mathbf{5}, \dots\} = s_2$$

$$\epsilon\text{-closure}(\delta(s_2, 0)) = \{\mathbf{1}, \dots\} = s_1$$

$$s_3 \leftarrow \epsilon\text{-closure}(\delta(s_2, 1)) = \{\mathbf{3}, \mathbf{5}, \mathbf{7}, 15, 14, 0, 2, 19, 4, 6, 8\}$$

Round 3:

$$s_4 \leftarrow \epsilon\text{-closure}(\delta(s_3, 0)) = \{\mathbf{1}, \mathbf{9}, 15, 14, 0, 2, 19, 4, 10\}$$

$$\epsilon\text{-closure}(\delta(s_3, 1)) = \{\mathbf{3}, \mathbf{5}, \mathbf{7}, \dots\} = s_3$$

Round 4:

$$s_5 \leftarrow \epsilon\text{-closure}(\delta(s_4, 0)) = \{\mathbf{1}, \mathbf{11}, 15, 14, 0, 2, 19, 4, 16, 12, 17\}$$

$$\epsilon\text{-closure}(\delta(s_4, 1)) = \{\mathbf{3}, \mathbf{5}, \dots\} = s_2$$

Note: the "core" states, those generated by δ , are set in bold. The states added by ϵ -closure are set in normal face.

Round 5:

$$\epsilon\text{-closure}(\delta(s_5, 0)) = \{\mathbf{1}, \dots\} = s_1$$

$$s_6 \leftarrow \epsilon\text{-closure}(\delta(s_5, 1)) = \{\mathbf{3}, \mathbf{5}, \mathbf{13}, 15, 14, 0, 2, 19, 4, 6, 12, 17\}$$

Round 6:

$$\epsilon\text{-closure}(\delta(s_6, 0)) = \{\mathbf{1}, \dots\} = s_1$$

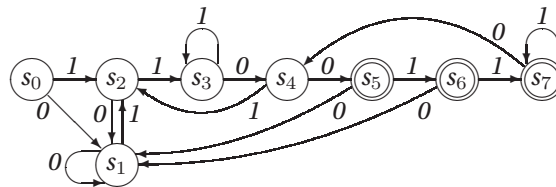
$$s_7 \leftarrow \epsilon\text{-closure}(\delta(s_6, 1)) = \{\mathbf{3}, \mathbf{5}, \mathbf{7}, \mathbf{13}, 15, 14, 0, 2, 19, 4, 6, 8, 12, 17\}$$

Round 7:

$$\epsilon\text{-closure}(\delta(s_7, 0)) = \{\mathbf{1}, \mathbf{9}, \dots\} = s_4$$

$$\epsilon\text{-closure}(\delta(s_7, 1)) = \{\mathbf{3}, \mathbf{5}, \mathbf{7}, \mathbf{13}, \dots\} = s_7$$

At this point, the worklist is empty. The resulting DFA has eight states.

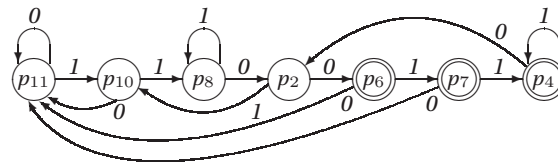


Part (c): Minimize the DFA.

Single element sets cannot split. We will omit the transitions for them.

Current Partition	Transition	Action
$\{ p_f = \{s_5, s_6, s_7\}, p_1 = \{s_0, s_1, s_2, s_3, s_4\} \}$	$\delta(p_0, 0) = p_1$ $\delta(p_0, 1) = p_0$ $\delta(p_1, 0) = \{p_0, p_1\}$	<i>none</i> <i>none</i> split p_1
$\{ p_0, p_2 = \{s_4\}, p_3 = \{s_0, s_1, s_2, s_3\} \}$	$\delta(p_0, 0) = \{p_2, p_3\}$	split p_0
$\{ p_4 = \{s_7\}, p_5 = \{s_5, s_6\}, p_2, p_3 \}$	$\delta(p_5, 0) = p_3$ $\delta(p_5, 1) = \{p_5, p_4\}$	<i>none</i> split p_5
$\{ p_4, p_6 = \{s_5\}, p_7 = \{s_6\}, p_2, p_3 \}$	$\delta(p_3, 0) = \{p_3, p_2\}$	split p_3
$\{ p_4, p_6, p_7, p_2, p_8 = \{s_3\}, p_9 = \{s_0, s_1, s_2\} \}$	$\delta(p_9, 0) = p_9$ $\delta(p_9, 1) = \{p_9, p_8\}$	<i>none</i> split p_9
$\{ p_4, p_6, p_7, p_2, p_8, p_{10} = \{s_2\}, p_{11} = \{s_0, s_1\} \}$	$\delta(p_{11}, 0) = p_{11}$ $\delta(p_{11}, 1) = p_{10}$	<i>none</i> <i>none</i>

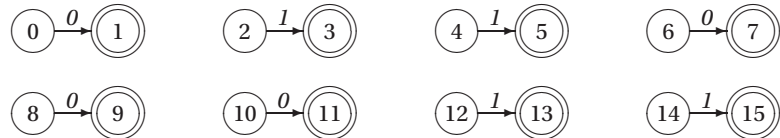
The minimization algorithm found seven distinct states. It combined states s_0 and s_1 in the original DFA. The minimal DFA is



Answer: For the RE $(01 \mid 10 \mid 00)^* 11$:

Part (a): Build an NFA for the RE using Thompson's construction.

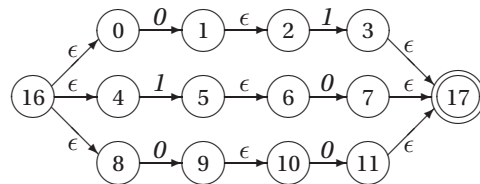
First, build distinct DFAs for each of the digits:



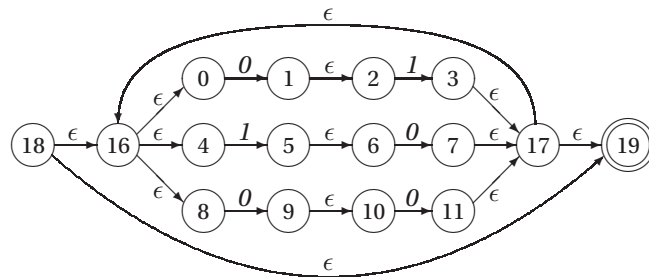
Next, combine the individual DFAs to form NFAs for the subterms 00, 10, 00, and 11:



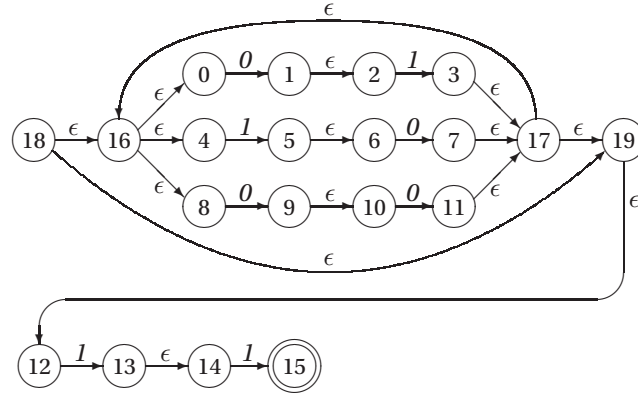
As the third step, combine the NFAs for 01, 10, and 00 to form an NFA for $(01 \mid 10 \mid 00)$:



As the fourth step, form the NFA for $(01 \mid 10 \mid 00)^*$:



And, finally, concatenate the NFA for 11 to the end of the NFA for $(01 \mid 10 \mid 00)^*$:



Part (b): Convert the NFA to a DFA using the subset construction.

Initialization:

The initial state is 18.

$$s_0 \leftarrow \epsilon\text{-closure}(18) = \{18, 16, 0, 4, 8, 19, 12\}$$

Round 1:

$$s_1 \leftarrow \epsilon\text{-closure}(\delta(s_0, 0)) = \{1, 9, 2, 10\}$$

$$s_2 \leftarrow \epsilon\text{-closure}(\delta(s_0, 1)) = \{5, 13, 6, 14\}$$

Round 2:

$$s_3 \leftarrow \epsilon\text{-closure}(\delta(s_1, 0)) = \{11, 17, 16, 0, 4, 8, 19, 12\}$$

$$s_4 \leftarrow \epsilon\text{-closure}(\delta(s_1, 1)) = \{3, 17, 16, 0, 4, 8, 19, 12\}$$

$$s_5 \leftarrow \epsilon\text{-closure}(\delta(s_2, 0)) = \{7, 17, 16, 0, 4, 8, 19, 12\}$$

$$s_6 \leftarrow \epsilon\text{-closure}(\delta(s_2, 1)) = \{15\}$$

Round 3:

$$\epsilon\text{-closure}(\delta(s_3, 0)) = \{1, 9, \dots\} = s_1$$

$$\epsilon\text{-closure}(\delta(s_3, 1)) = \{5, 13, \dots\} = s_2$$

$$\epsilon\text{-closure}(\delta(s_4, 0)) = \{1, 9, \dots\} = s_1$$

$$\epsilon\text{-closure}(\delta(s_4, 1)) = \{5, 13, \dots\} = s_2$$

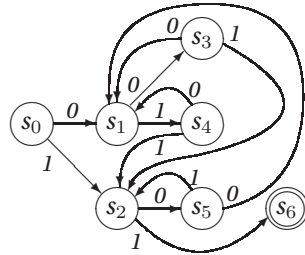
$$\epsilon\text{-closure}(\delta(s_5, 0)) = \{1, 9, \dots\} = s_1$$

$$\epsilon\text{-closure}(\delta(s_5, 1)) = \{5, 13, \dots\} = s_2$$

$$\epsilon\text{-closure}(\delta(s_6, 0)) = \emptyset$$

$$\epsilon\text{-closure}(\delta(s_6, 1)) = \emptyset$$

At this point, the worklist is empty. The resulting DFA has seven states.

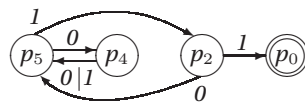


Part (c): Minimize the DFA.

Single element sets cannot split. We will omit the transitions for them.

Current Partition	Transition	Action
$\{ p_0 = \{s_6\}, p_1 = \{s_0, s_1, s_2, s_3, s_4, s_5\} \}$	$\delta(p_1, 0) = p_1$ $\delta(p_1, 1) = \{p_0, p_1\}$	<i>none</i> split p_1
$\{ p_0, p_2 = \{s_2\}, p_3 = \{s_0, s_1, s_3, s_4, s_5\} \}$	$\delta(p_3, 0) = p_3$ $\delta(p_3, 1) = \{p_2, p_3\}$	<i>none</i> split p_3
$\{ p_0, p_2, p_4 = \{s_1\} p_5 = \{s_0, s_3, s_4, s_5\} \}$	$\delta(p_5, 0) = p_4$ $\delta(p_5, 1) = p_2$	<i>none</i> <i>none</i>

The resulting DFA has four states.



8. One way of proving that two REs are equivalent is to construct their minimized DFAs and then to compare them. If they differ only by state names, then the REs are equivalent. Use this technique to check the following pairs of REs and state whether or not they are equivalent.

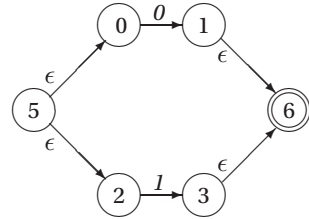
- a. $(0 \mid 1)^*$ and $(0^* \mid 10^*)^*$
 b. $(ba)^+ (a^* b^* \mid a^*)$ and $(ba)^* ba^+ (b^* \mid \epsilon)$

Answer: For $(0 \mid 1)^*$ and $(0^* \mid 10^*)^*$

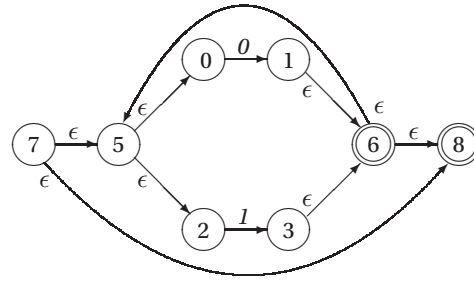
First, build a DFA for $(0 \mid 1)^*$. We start with NFAs for 0 and 1



Next, we connect them to form an NFA for $0 \mid 1$.



Next, we add the states and edges for the closure, to form an NFA for $(0 \mid 1)^*$.



Now, we convert the NFA to a DFA with the subset construction, using the same notation as in the solution to Problem 2.8.

Initialization:

$$s_0 \leftarrow \epsilon\text{-closure}(7) = \{7, 5, 0, 2, 8\}$$

Round 1:

$$s_1 \leftarrow \epsilon\text{-closure}(\delta(s_0, 0)) = \{1, 6, 5, 0, 2, 8\}$$

$$s_2 \leftarrow \epsilon\text{-closure}(\delta(s_0, 1a)) = \{3, 6, 5, 0, 2, 8\}$$

Round 2:

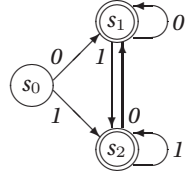
$$\epsilon\text{-closure}(\delta(s_1, 0)) = \{1, \dots\} = s_1$$

$$\epsilon\text{-closure}(\delta(s_1, 1)) = \{3, \dots\} = s_2$$

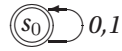
$$\epsilon\text{-closure}(\delta(s_2, 0)) = \{1, \dots\} = s_1$$

$$\epsilon\text{-closure}(\delta(s_2, 1)) = \{3, \dots\} = s_2$$

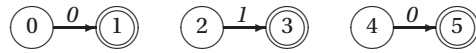
At this point the worklist is empty. The DFA has three states.



DFA minimization produces a single-state DFA, which a good compiler writer might have drawn from inspection.



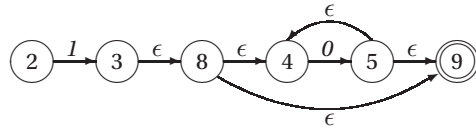
As the second step, we build the NFA for $(0^* \mid 10^*)^*$. We start with NFAs for 0, 1, and 0



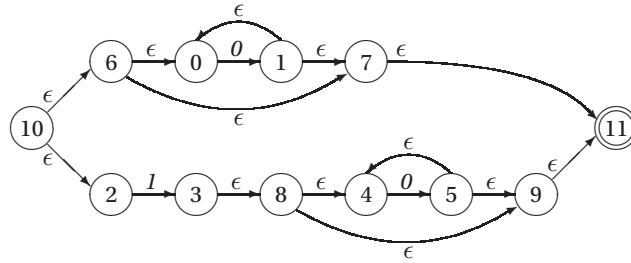
Next, we build the NFAs for 0^* and 0^* .



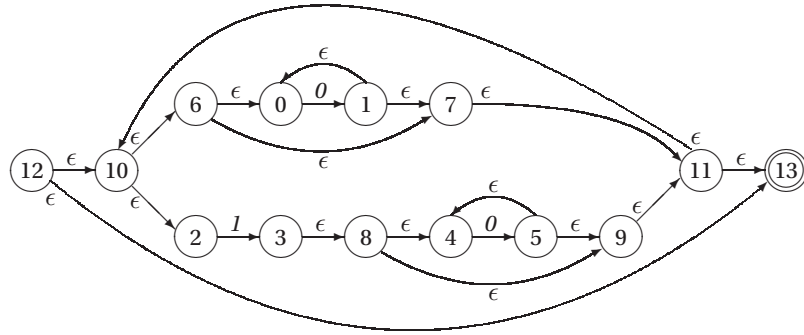
Now, we concatenate the NFAs for 1 and 0^* to form an NFA for 10^* .



Then, we form the NFA for $0^* \mid 10^*$.



And, finally, add the states to the NFA that implement the outermost closure.



This NFA is rather large, with lots of ϵ transitions. First, we apply the subset construction to find the equivalent DFA. Then, we will minimize the DFA.

Initialization:

$$s_0 = \epsilon\text{-closure}(12) = \{12, 10, 6, 0, 7, 11, 13, 2\}$$

Round 1:

$$s_1 = \epsilon\text{-closure}(\delta(s_0, 0)) = \{1, 7, 11, 10, 6, 0, 2, 13\}$$

$$s_2 = \epsilon\text{-closure}(\delta(s_0, 1)) = \{3, 8, 4, 9, 11, 10, 6, 0, 2, 13\}$$

Round 2:

$$\epsilon\text{-closure}(\delta(s_1, 0)) = \{1, \dots\} = s_1$$

$$\epsilon\text{-closure}(\delta(s_1, 1)) = \{3, \dots\} = s_2$$

$$s_3 = \epsilon\text{-closure}(\delta(s_2, 0)) = \{1, 5, 7, 11, 10, 6, 0, 2, 9, 13\}$$

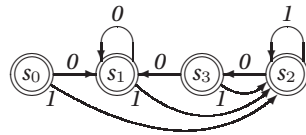
$$\epsilon\text{-closure}(\delta(s_2, 1)) = \{3, \dots\} = s_2$$

Round 3:

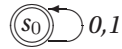
$$\epsilon\text{-closure}(\delta(s_3, 0)) = \{1, \dots\} = s_1$$

$$\epsilon\text{-closure}(\delta(s_3, I)) = \{3, \dots\} = s_2$$

At this point, the worklist is empty. The resulting DFA has four states.

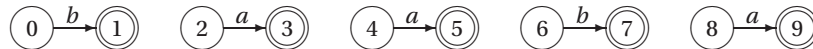


Minimization discovers that states s_0 , s_1 , s_2 , and s_3 form an equivalence class: each state is a final state and each has a transition back into the class on 0 and on 1. Thus, the minimized DFA has one state:

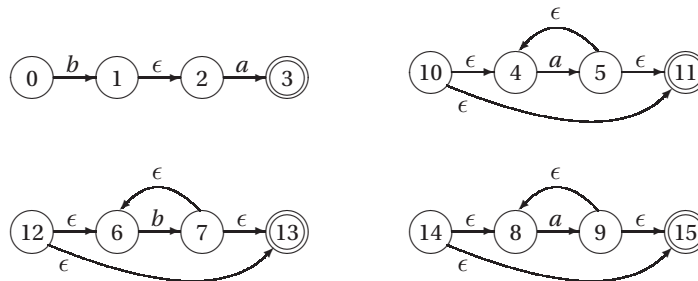


Answer: For $(ba)^+ (a^* b^* \mid a^*)$ and $(ba)^* ba^+ (b^* \mid \epsilon)$

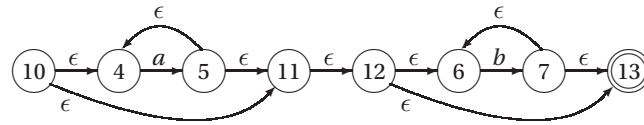
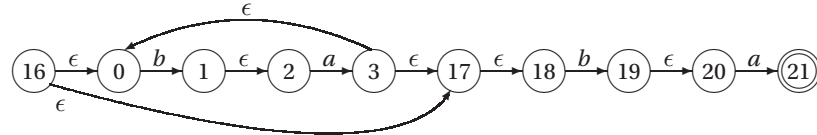
First, we build the NFA for $(ba)^+ (a^* b^* \mid a^*)$. We start with individual NFAs for each of the letters.



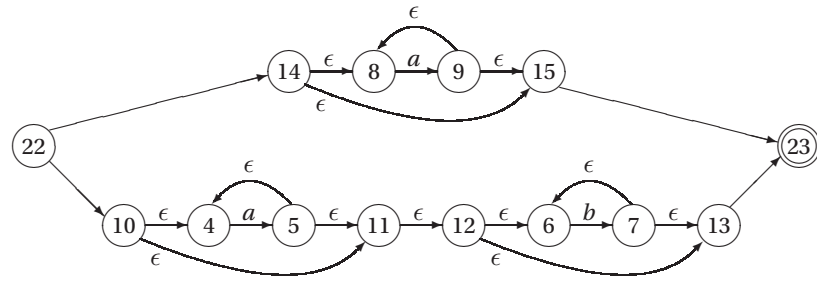
Next, we build NFAs for ba , a^* , b^* , and a^* .



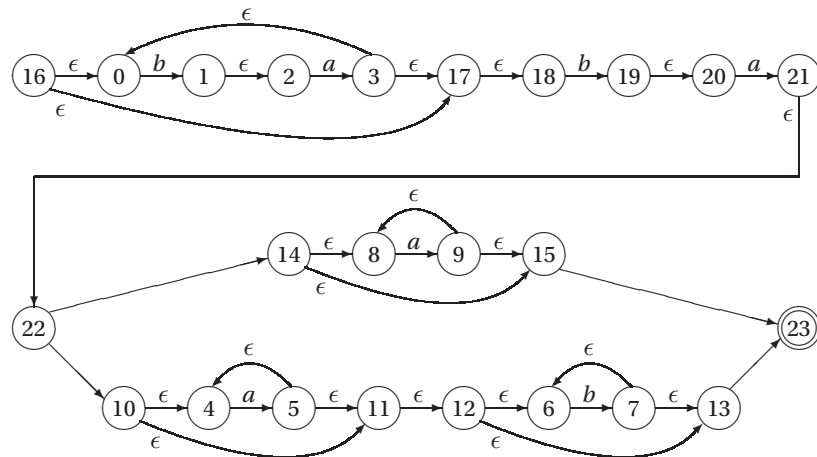
Thompson's construction does not give us a form for $(ba)^+$. We can either model it as $(ba)^*ba$, or we can build a more concise form by adding a single ϵ transition from state 0 to state 3. Since these solutions are intended to be explanatory, we will choose the former approach. Thus, we next produce NFAs for $(ba)^*ba$ and for a^*b^*



Next is the NFA for $(a^*b^* | a^*)$.



Finally, we tie together the NFAs for $(ba)^*ba$ and $(a^*b^* | a^*)$ by adding an ϵ transition from state 21 to state 22, to produce the final NFA.



The next step is to use the subset construction to convert the NFA to a DFA.

Initialization

The initial state is 16.

$$s_0 \leftarrow \epsilon\text{-closure}(16) = \{\mathbf{16}, 0, 17, 18\}$$

Round 1:

$$\epsilon\text{-closure}(\delta(s_0, a)) = \emptyset$$

$$s_1 \leftarrow \epsilon\text{-closure}(\delta(s_0, b)) = \{\mathbf{1}, \mathbf{19}, 2, 20\}$$

Round 2:

$$s_2 \leftarrow \epsilon\text{-closure}(\delta(s_1, a)) = \{\mathbf{3}, \mathbf{21}, 0, 17, 18, 22, 14, 8, 15, 23, 10, 4, 11, 12, 6, 13\}$$

$$\epsilon\text{-closure}(\delta(s_1, b)) = \emptyset$$

Round 3:

$$s_3 \leftarrow \epsilon\text{-closure}(\delta(s_2, a)) = \{\mathbf{9}, \mathbf{5}, 8, 15, 23, 4, 11, 12, 6, 13\}$$

$$s_4 \leftarrow \epsilon\text{-closure}(\delta(s_2, b)) = \{\mathbf{1}, \mathbf{19}, \mathbf{7}, 2, 20, 6, 13, 23\}$$

Round 4:

$$\epsilon\text{-closure}(\delta(s_3, a)) = \{\mathbf{9}, \mathbf{5}, 8, 15, 23, 4, 11, 12, 6, 13\} = s_3$$

$$s_5 \leftarrow \epsilon\text{-closure}(\delta(s_3, b)) = \{\mathbf{7}, 6, 13, 23\}$$

$$\begin{aligned} \epsilon\text{-closure}(\delta(s_4, a)) &= \{\mathbf{3}, \mathbf{21}, 0, 17, 18, 22, 14, 8, 15, 23, 10, 4, 11, 12, 6, 13\} \\ &= s_2 \end{aligned}$$

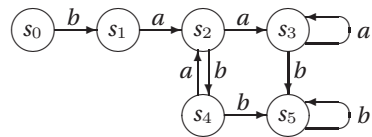
$$\epsilon\text{-closure}(\delta(s_4, b)) = \{\mathbf{7}, 6, 13, 23\} = s_5$$

Round 5:

$$\epsilon\text{-closure}(\delta(s_5, a)) = \emptyset$$

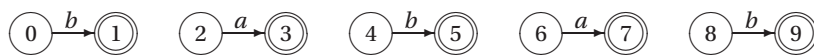
$$\epsilon\text{-closure}(\delta(s_5, b)) = \{\mathbf{7}, 6, 13, 23\} = s_5$$

The resulting DFA is:

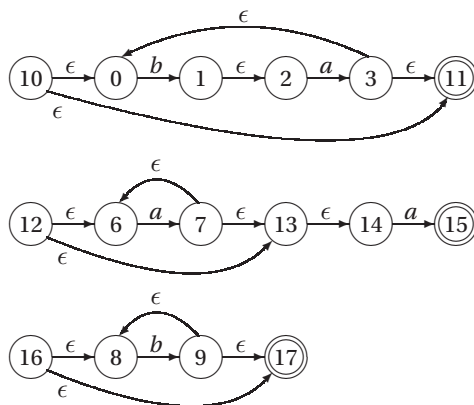


We will defer minimization, another intricate step, until we see the DFA for the other RE.

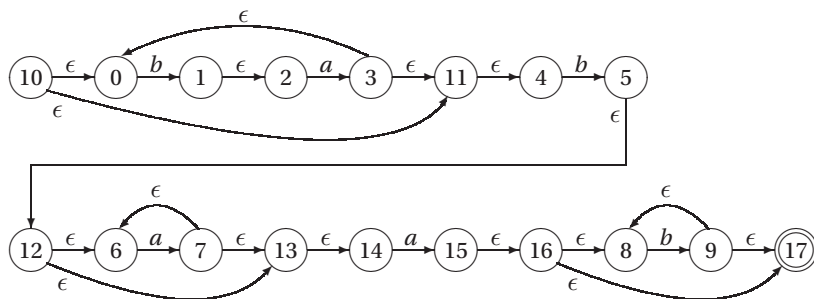
Next, we build the NFA for $(ba)^*ba^+(b^*|\epsilon)$. We will make one simplification, noting that the term $(b^*|\epsilon)$ is equivalent to b^* . We begin by building individual NFAs for each letter



Next, we build NFAs for the terms $(ba)^*$, a^+ , and b^* . (We skipped building ba . It occurs as shown for the first RE in this question. We will treat a^+ as a^*a .)



Finally, we tie these three NFAs together, along with the NFA for b (states 4 and 5 earlier) to form the final NFA



Next, we convert the NFA to a DFA with the subset construction.

Initialization

The initial state is 10.

$$s_0 \leftarrow \epsilon\text{-closure}(10) = \{\mathbf{10}, 0, 11, 4\}$$

Round 1

$$\epsilon\text{-closure}(\delta(s_0, a)) = \emptyset$$

$$s_1 \leftarrow \epsilon\text{-closure}(\delta(s_0, a)) = \{\mathbf{1}, \mathbf{5}, 2, 12, 6, 13, 14\}$$

Round 2

$$s_2 \leftarrow \epsilon\text{-closure}(\delta(s_1, a)) = \{\mathbf{3}, \mathbf{7}, \mathbf{15}, 0, 11, 4, 6, 13, 14, 16, 8, 17\}$$

$$\epsilon\text{-closure}(\delta(s_1, b)) = \emptyset$$

Round 3

$$s_3 \leftarrow \epsilon\text{-closure}(\delta(s_2, a)) = \{\mathbf{7}, \mathbf{15}, 6, 13, 14, 16, 8, 17\}$$

$$s_4 \leftarrow \epsilon\text{-closure}(\delta(s_2, b)) = \{\mathbf{1}, \mathbf{5}, \mathbf{9}, 2, 12, 6, 13, 14, 8, 17\}$$

Round 4

$$\epsilon\text{-closure}(\delta(s_3, a)) = \{\mathbf{7}, \mathbf{15}, 6, 13, 14, 16, 8, 17\} = s_3$$

$$s_5 \leftarrow \epsilon\text{-closure}(\delta(s_3, b)) = \{\mathbf{9}, 8, 17\}$$

$$\epsilon\text{-closure}(\delta(s_4, a)) = \{\mathbf{3}, \mathbf{7}, \mathbf{15}, 0, 11, 4, 6, 13, 14, 16, 8, 17\} = s_2$$

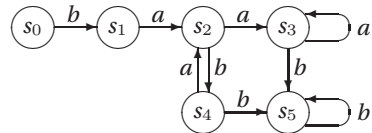
$$\epsilon\text{-closure}(\delta(s_4, b)) = \{\mathbf{9}, 8, 17\} = s_5$$

Round 5

$$\epsilon\text{-closure}(\delta(s_5, a)) = \emptyset$$

$$\epsilon\text{-closure}(\delta(s_5, b)) = \{\mathbf{9}, 8, 17\} = s_5$$

The resulting DFA is:



Notice that this DFA is identical to the DFA for the other RE. Thus, their minimized forms will be identical, within a renaming of states and we have shown that the two REs describe the same language.

9. In some cases, two states connected by an ϵ -move can be combined.

- a. Under what set of conditions can two states connected by an ϵ -move be combined?

Answer: Two states $s_i \xrightarrow{\epsilon} s_j$ can be combined if the ϵ -move is the only transition that enters s_j and the transitions leaving s_i and s_j , other than $s_i \xrightarrow{\epsilon} s_j$ are compatible. That is, if $s_i \xrightarrow{c} s_k$, then either s_k has no transition on c or $s_j \xrightarrow{c} s_k$.

These conditions hold for the ϵ transitions in the initial NFAs formed in Thompson's construction.

- b. Give an algorithm for eliminating ϵ -moves.

Answer: The following algorithm removes ϵ transitions of the kind found in NFAs created by Thompson's construction.

```

for each edge  $e \in E$ 
  if  $e = \langle q_i, q_j, \epsilon \rangle$  then
    add  $e$  to WorkList
while (WorkList  $\neq \emptyset$ )
  remove  $e = \langle q_i, q_j, \alpha \rangle$  from WorkList
  if  $\alpha = \epsilon$  then
    for each  $\langle q_j, q_k, \beta \rangle$  in  $E$ 
      add  $\langle q_i, q_k, \beta \rangle$  to  $E$ 
      if  $\beta = \epsilon$  then
        add  $\langle q_i, q_k, \beta \rangle$  to WorkList
    delete  $\langle q_i, q_j, \epsilon \rangle$  from  $E$ 
    if  $q_j \in F$  then
      add  $q_i$  to  $F$ 
for each state  $q_i \in N$ 
  if  $q_i$  has no entering edge then
    delete  $q_i$  from  $N$ 

```

If $s_i \xrightarrow{\epsilon} s_j$, the algorithm copies transitions that leave s_j to s_i to create an equivalent path through the transition graph that avoids the ϵ -move. After all such transitions have been copied, it deletes the ϵ -move. A final pass over the transition graph eliminates any unreachable states – states that were reachable only with an ϵ -move will be removed in this pass.

- c. How does your algorithm relate to the ϵ -closure function used to implement the subset construction?

Answer: Using the ϵ -closure function, the algorithm could directly discover the set of outbound transitions needed at each node. In some sense, the algorithm that we gave for part (b) computes ϵ -closure(s_i) on the fly, as it copies the edges.

An algorithm that used ϵ -closure directly would still need to create all the new edges and make a final pass to delete unreachable states.

10. Show that the set of regular languages is closed under intersection.

Answer: We can prove that regular languages are closed under union and under complement.

(Closure under union, sketch of proof: Take the DFAs for L_1 and L_2 , add a new start state with ϵ -moves to the start states of the DFAs for L_1 and L_2 . This new NFA accepts a string if it is accepted by the DFA for L_1 or if it is accepted by the DFA for L_2 , so it accepts $L_1 \cup L_2$.)

(Closure under complement, sketch of proof: Take the DFA for L_1 . Make all non-final states in the original DFA be final states of a new DFA and all the final states of the original DFA be non-final states in the new DFA. Now, the new DFA accepts any string except those that the original DFA would accept, so it accepts $\overline{L_1}$.)

Now, we know from DeMorgan's laws that $L_1 \cap L_2$ is equivalent to $\overline{\overline{L_1} \cup \overline{L_2}}$. Thus, we can compute $L_1 \cap L_2$ by finding the regular languages that are the complements of L_1 and L_2 , finding the language that is the union of those new languages, and taking its complement. At each step, the languages are regular because of closure under complement and union, so the resulting language is regular. This procedure guarantees a regular language from $L_1 \cap L_2$, so the set of regular languages is closed under intersection.

11. The DFA minimization algorithm given in Figure 2.9 is formulated to enumerate all the elements of P and all of the characters in Σ on each iteration of the while loop.
- a. Recast the algorithm so that it uses a worklist to hold the sets that must still be examined.

Answer:

```

Worklist  $\leftarrow \{ \text{FinalStates}, \text{States} - \text{FinalStates} \}$ 
Partition  $\leftarrow \{ \text{FinalStates}, \text{States} - \text{FinalStates} \}$ 

while (Worklist  $\neq \emptyset$ ) do
  remove a set  $w$  from Worklist
  temp  $\leftarrow \text{Split}(w)$ 
  if temp  $\neq w$  then
    remove  $w$  from Partition // Split() refined  $w$ 
    Partition  $\leftarrow \text{Partition} \cup \text{temp}$  // Update Partition
    Worklist  $\leftarrow \text{Worklist} \cup \text{temp}$  // Update Worklist

```

- b. Recast the *Split* function so that it partitions the set around all of the characters in Σ .

Answer:

```

Split( $S$ ) {
  Result  $\leftarrow \{ S \}$ 
  for each  $c \in \Sigma$  do {
    for each  $s \in \text{Result}$  do {
       $T \leftarrow \emptyset$ 
      if  $c$  splits  $s$  into  $s_1$  and  $s_2$  then
        remove  $s$  from Result
         $T \leftarrow T \cup \{ s_1, s_2 \}$ 
    }
    Result  $\leftarrow \text{Result} \cup T$ 
  }
  return Result
}

```

- c. How does the expected case complexity of your modified algorithm compare to the expected case complexity of the original algorithm?

Answer: The total number of splits that the algorithm can perform is bounded by the number of states. Each state is in exactly one partition, so the number of splits cannot exceed the number of states.

However, the algorithm in Figure 2.19 can examine a set many more times than the modified algorithm. The original algorithm tries to split every set in every partition on each iteration of the outer while loop, so it does something like $O(S^2)$ calls to *Split*, where S is the number of states in the minimal DFA.

Ignore, for the moment, the modifications to *Split* in part (b) of the question. The worklist version of the algorithm, in part (a), calls *Split* once on each set that appears in *Worklist*. How many sets appear in *Worklist*? We know that every state in the minimal DFA appears in *Worklist*. Those states are either one of the two initial sets, or are split from some other set. If you think of the splitting as a binary tree that has S leaves, you can see that *Worklist* has, potentially, $O(S^2)$ distinct entries over time.

Are the two algorithms the same? No. The modified algorithm considers $\frac{1}{2}S^2$ sets where the original algorithm considered S^2 sets. In practice, that savings is worthwhile.

What about the modifications to *Split*? While the modifications would reduce the number of calls to *Split*, they would not reduce the total number of sets that the algorithm touches.

Section 2.5

- 12. Construct a DFA for each of the following C language constructs, and then build the corresponding table for a table-driven implementation for each of them:
 - a. Integer constants

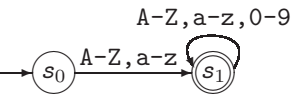
Answer:



δ	0-9	other
s_0	s_0	s_e
s_e	s_e	s_e

b. Identifiers

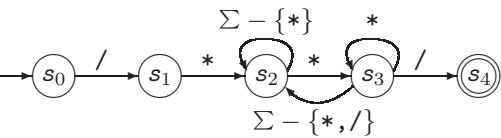
Answer:



δ	$A-Z, a-z$	$0-9$	<i>other</i>
s_0	s_1	s_e	s_e
s_1	s_1	s_1	s_e
s_e	s_e	s_e	s_e

c. Comments

Answer:



δ	$/$	$*$	<i>other</i>
s_0	s_1	s_e	s_e
s_1	s_e	s_2	s_e
s_2	s_2	s_3	s_2
s_3	s_4	s_3	s_2
s_4	s_e	s_e	s_e
s_e	s_e	s_e	s_e

Note: the transitions for s_4 reflect the fact that the DFA only accepts comments. It does not accept any text following the comment.

13. For each of the DFAs from the previous exercise, build the corresponding direct-coded scanner.

Answer:

```

goto s0
s0: char ← NextCharacter
    if ('0' ≤ char ≤ '9')
        then goto s0
        else if char = eof
            then report acceptance
            else goto se
se: print error message
    report failure

```

(a) Integer Constants

```

goto s0
s0: char ← NextCharacter
    if (char = '_' ) or
        ('a' ≤ char ≤ 'z') or
        ('A' ≤ char ≤ 'Z')
        then goto s1
        else goto se
s1: char ← NextCharacter
    if (char = '_' ) or
        ('a' ≤ char ≤ 'z') or
        ('A' ≤ char ≤ 'Z') or
        ('0' ≤ char ≤ '9')
        then goto s1
        else if char = eof
            then report acceptance
            else goto se
se: print error message
    report failure

```

(b) Identifiers

```

goto s0
s0: char ← NextCharacter
    if (char = '/')
        then goto s1
        else goto se
s1: char ← NextCharacter
    if (char = '*')
        then goto s2
        else goto se
s2: char ← NextCharacter
    if (char = '*')
        then goto s3
        else goto s2

```

```

s3: char ← NextCharacter
    if (char = '*')
        then goto s3
        else if (char = '/')
            then goto s4
            else goto s2
s4: report acceptance
se: print error message
    report failure

```

(c) Comments in C

14. This chapter described two styles of DFA implementations: a table-driven implementation that uses a case or switch statement, and a direct-coded scanner that uses branches and jumps. A third alternative is to use mutually recursive functions to implement a scanner. Discuss the advantages and disadvantages of such an implementation.

Answer: The primary advantage of implementing a DFA as a set of mutually recursive is that the implementation is conceptually simple.

The primary disadvantage is efficiency. Each procedure call must, at least, include a transfer of control that is equivalent to the one that occurs in the direct-coded scanner. In most implementations, the call will also allocate an activation record (see Chapter 6) and establish bindings for parameters. All of these make the cost per character higher than in the direct-coded scanner. (Space costs for activation records will be $O(n)$ where n is the number of characters. The extra overhead from parameter binding and such will be constant per call, with one call for each character.)

15. To reduce the size of the transition table, the scanner generator can use a character classification scheme. Generating the classifier table, however, seems expensive. The obvious algorithm would require $O(|\Sigma|^2 \cdot |states|)$ time. Derive an asymptotically faster algorithm for finding identical columns in the transition table.

Answer: Many solutions are possible to this problem. One simple solution is to perform a radix sort, state by state, on the columns, from s_0 through s_e . Characters whose columns end up in the same bucket belong in the same character class. The complexity of this approach should be $O(|\Sigma| \cdot |states|)$, a factor of $|\Sigma|$ faster than the naive algorithm. Many other approaches will work.

16. Figure 2.15 shows a scheme that avoids quadratic roll back behavior in a scanner built by simulating a DFA. Unfortunately, that scheme requires that the scanner know in advance the length of the input stream and that it maintain a bit-matrix, *Failed*, of size $|states| \times |input|$. Devise a scheme that avoids the need to know the size of the input stream in advance. Can you use the same scheme to reduce the size of the *Failed* table in cases where the worst case input does not occur?

Answer: In fact, this question is worded incorrectly. We do not know of a scheme that avoids, in general, the need to know the length of the input in the worst case. We can reduce the space needed by the worst case syntax by only storing every k^{th} state; that reduces the needed space by a factor of k while increasing the cost of scanning—the scanner may repeat k states before detecting that it is repeating a failed path in the DFA.

See the paper “‘Maximal-munch’ tokenization in linear time” by Thomas Reps, in ACM Transactions on Programming Languages and Systems (TOPLAS), 20(2), March 1998, pages 259–273.

CHAPTER 3

Section 3.2

1. Write a context-free grammar for the syntax of regular expressions.

Answer:

$$\begin{array}{lcl}
 R & \rightarrow & R R \\
 & | & R^* \\
 & | & (R) \\
 & | & R \mid R \\
 & | & \text{char}
 \end{array}$$

where char represents a single character

2. Write a context-free grammar for the Backus-Naur Form (BNF) notation for context-free grammars.

Answer:

$$\begin{array}{lcl}
 BNF & \rightarrow & LHSList \\
 LHSList & \rightarrow & Production \ LHSList \\
 & | & Production \\
 Production & \rightarrow & NonTerm \ Derives \ RHSList \\
 RHSList & \rightarrow & RHS \ Also \ RHSList \\
 & | & RHS \\
 RHS & \rightarrow & NonTerm \ RHS \\
 & | & Term \ RHS \\
 & | & NonTerm \\
 & | & Term
 \end{array}$$

where NonTerm is a nonterminal symbol, Term is a terminal symbol, Derives represents ' \rightarrow ', and Also represents '|'.

3. When asked about the definition of an *unambiguous context-free grammar* on an exam, two students gave different answers. The first defined it as “a grammar where each sentence has a unique syntax tree by leftmost derivation.” The second defined it as “a grammar where each sentence has a unique syntax tree by any derivation.” Which one is correct?

Answer: The two definitions are equivalent, as can be proved by induction. The basic idea is to show that any syntax tree generated by any derivation can be generated by a left-most derivation.

Section 3.3

4. The following grammar is not suitable for a top-down predictive parser. Identify the problem and correct it by rewriting the grammar. Show that your new grammar satisfies the LL(1) condition.

$$\begin{array}{lll}
 L \rightarrow R a & R \rightarrow a b a & Q \rightarrow b b c \\
 | Q b a & | c a b a & | b c \\
 & | R b c &
 \end{array}$$

Answer: The grammar, as written, contains a left-recursion on the non-terminal symbol R . We must rewrite to remove that left recursion. The productions

$$\begin{array}{lll}
 R \rightarrow a b a & & \\
 | c a b a & \text{can be rewritten as} & R \rightarrow a b a S \quad S \rightarrow b c S \\
 | R b c & & | c a b a S \quad | \epsilon
 \end{array}$$

At that point, the remaining problem lies in the productions for Q , which both begin with b . Left factoring the productions will cure that problem:

$$\begin{array}{lll}
 Q \rightarrow b b c & & Q \rightarrow b T \\
 | b c & \text{can be rewritten as} & T \rightarrow b c \\
 & & | c
 \end{array}$$

To show that the resulting grammar is LL(1), we need to compute, for each non-terminal symbol, the FIRST set of each alternative right-hand side. The grammar is, after the transformations:

$$\begin{array}{llllll}
 L \rightarrow R a & R \rightarrow a b a S & S \rightarrow b c S & Q \rightarrow b T & T \rightarrow b c \\
 | Q b a & | c a b a S & | \epsilon & & | c
 \end{array}$$

To show that this grammar is LL(1), we must show that the FIRST sets of the alternative right hand sides for each non-terminal symbol are disjoint:

$$\begin{array}{lll}
 L: & \text{FIRST}(Ra) = \{a, c\} & \text{FIRST}(Qba) = \{b\} & \{a, c\} \cap \{b\} = \emptyset \\
 R: & \text{FIRST}(abaS) = \{a\} & \text{FIRST}(cabaS) = \{c\} & \{a\} \cap \{c\} = \emptyset \\
 S: & \text{FIRST}(bcS) = \{b\} & \text{FIRST}(\epsilon) = \{\epsilon\} & \{b\} \cap \{\epsilon\} = \emptyset \\
 T: & \text{FIRST}(bc) = \{b\} & \text{FIRST}(c) = \{c\} & \{b\} \cap \{c\} = \emptyset
 \end{array}$$

5. Consider the following grammar:

$$\begin{array}{ll} A \rightarrow B a & C \rightarrow c B \\ B \rightarrow dab & | A c \\ & | C b \end{array}$$

Does this grammar satisfy the LL(1) condition? Justify your answer. If it does not, rewrite it as an LL(1) grammar for the same language.

Answer: The grammar has an indirect left recursion from the productions $A \rightarrow Ba$, $B \rightarrow Cb$, and $C \rightarrow Ac$.

The algorithm for removal of indirect left recursion, in Figure 3.6, cures this problem. Conceptually, it operates in two steps. First, it converts the indirect left recursion to a direct left recursion:

$$\begin{array}{lll} A \rightarrow daba & B \rightarrow dab & C \rightarrow cB \\ | cBba & | Cb & | Ac \\ | Acba & & \end{array}$$

Second, it eliminates the direct left recursion on A by rewriting the productions for A as follows:

$$\begin{array}{ll} A \rightarrow dabaD & D \rightarrow cbaD \\ | cBbaD & | \epsilon \end{array}$$

The algorithm, as written, interleaves these steps. Combining them produces the following grammar.

$$\begin{array}{llll} A \rightarrow dabaD & B \rightarrow dab & C \rightarrow cB & D \rightarrow cbaD \\ | cBbaD & | Cb & | Ac & | \epsilon \end{array}$$

The alternative right-hand sides for A and D have disjoint FIRST sets. Both B and C , however, have problems in this regard.

$$B: \quad \text{FIRST}(dab) = \{d\} \quad \text{FIRST}(Cb) = \{c, d\} \quad \{d\} \cap \{c, d\} = \{d\}$$

$$C: \quad \text{FIRST}(cB) = \{c\} \quad \text{FIRST}(Ac) = \{d, c\} \quad \{c\} \cap \{d, c\} = \{c\}$$

Left factoring can cure this problem. To see this point, substitute the two right-hand sides of A into the right-hand side for C , then substitute the right-hand sides of C into the right-hand side for B , to yield:

$$\begin{array}{l}
 B \rightarrow \text{dab} \\
 \quad | \text{c} B \\
 \quad | \text{daba} D c \\
 \quad | \text{c} B \text{ba} D c
 \end{array}$$

This grammar fragment has two right-hand sides that begin with d and two that begin with c. In each case, we can introduce a new non-terminal symbol for the common prefix, to yield:

$$\begin{array}{lll}
 B \rightarrow \text{c} B E & E \rightarrow \text{ba} D c & F \rightarrow \text{a} D c \\
 \quad | \text{dab} F & \quad | \epsilon & \quad | \epsilon
 \end{array}$$

Combining these new fragments with the remaining fragments of the grammar produces the following final grammar.

$$\begin{array}{lllll}
 A \rightarrow \text{daba} D & B \rightarrow \text{c} B E & D \rightarrow \text{cba} D & E \rightarrow \text{ba} D c & F \rightarrow \text{a} D c \\
 \quad | \text{cBba} D & \quad | \text{dab} F & \quad | \epsilon & \quad | \epsilon & \quad | \epsilon
 \end{array}$$

Inspecting the alternative right-hand sides for each non-terminal symbol shows that they have distinct FIRST and FIRST⁺ sets.

6. Grammars that can be parsed top-down, in a linear scan from left to right, with a k word lookahead are called LL(k) grammars. In the text, the LL(1) condition is described in terms of FIRST sets. How would you define the FIRST sets necessary to describe an LL(k) condition?

Answer: In an LL(1) grammar, the parser needs to look just one word ahead, so the FIRST sets have elements that are terminal symbols of the grammar. For an LL(k) grammar, the corresponding FIRST sets, called FIRST _{k} sets would need to contain elements representing combinations of up to k terminal symbols. Thus, FIRST _{k} (α) would be a subset of the set of all combinations of terminals that form words of length $\leq k$.

7. Suppose an elevator is controlled by two commands: \uparrow to move the elevator up one floor and \downarrow to move the elevator down one floor. Assume that the building is arbitrarily tall and that the elevator starts at floor x .

Write an LL(1) grammar that generates arbitrary command sequences that (1) never cause the elevator to go below floor x and (2) always return the elevator to floor x at the end of the sequence. For example, $\uparrow\uparrow\downarrow$ and $\uparrow\downarrow\uparrow\downarrow$ are

valid command sequences, but $\uparrow\downarrow\downarrow\uparrow$ and $\uparrow\downarrow\downarrow$ are not. For convenience, you may consider a null sequence as valid. Prove that your grammar is LL(1).

Answer: This problem is analogous to the parenthesis matching problem. Thus, one solution is

$$\begin{array}{lcl} E & \rightarrow & \uparrow E \downarrow E \\ & & | \quad \epsilon \end{array}$$

The grammar is trivially LL(1)

Section 3.4

8. Top-down and bottom-up parsers build syntax trees in different orders. Write a pair of programs, `TopDown` and `BottomUp`, that take a syntax tree and print out the nodes in order of construction. `TopDown` should display the order for a top-down parser, while `BottomUp` should show the order for a bottom-up parser.

Answer: A top-down parser constructs its parse tree in a root-first, left-child first order. Thus a preorder traversal, with left-to-right handling of children, should produce the order of construction for a top-down parser. A bottom-up parser constructs its parse tree in a root-last, left-child first order. Thus, a postorder traversal, with left-to-right handling of children, should produce the order of construction of a bottom-up parser.

9. The *ClockNoise* language (*CN*) is represented by the following grammar:

$$\begin{array}{lcl} Goal & \rightarrow & ClockNoise \\ ClockNoise & \rightarrow & ClockNoise \text{ tick } \text{tock} \\ & & | \quad \text{tick } \text{tock} \end{array}$$

a. What are the LR(1) items of *CN*?

Answer:

- $[Goal \rightarrow \bullet ClockNoise, EOF]$
- $[Goal \rightarrow \bullet ClockNoise, tick]$
- $[Goal \rightarrow \bullet ClockNoise, tock]$
- $[Goal \rightarrow ClockNoise \bullet, EOF]$
- $[Goal \rightarrow ClockNoise \bullet, tick]$
- $[Goal \rightarrow ClockNoise \bullet, tock]$
- $[ClockNoise \rightarrow \bullet ClockNoise tick tock, EOF]$
- $[ClockNoise \rightarrow \bullet ClockNoise tick tock, tick]$
- $[ClockNoise \rightarrow \bullet ClockNoise tick tock, tock]$
- $[ClockNoise \rightarrow ClockNoise \bullet tick tock, EOF]$
- $[ClockNoise \rightarrow ClockNoise \bullet tick tock, tick]$
- $[ClockNoise \rightarrow ClockNoise \bullet tick tock, tock]$
- $[ClockNoise \rightarrow ClockNoise tick \bullet tock, EOF]$
- $[ClockNoise \rightarrow ClockNoise tick \bullet tock, tick]$
- $[ClockNoise \rightarrow ClockNoise tick \bullet tock, tock]$
- $[ClockNoise \rightarrow ClockNoise tick tock \bullet, EOF]$
- $[ClockNoise \rightarrow ClockNoise tick tock \bullet, tick]$
- $[ClockNoise \rightarrow ClockNoise tick tock \bullet, tock]$
- $[ClockNoise \rightarrow \bullet tick tock, EOF]$
- $[ClockNoise \rightarrow \bullet tick tock, tick]$
- $[ClockNoise \rightarrow \bullet tick tock, tock]$
- $[ClockNoise \rightarrow tick \bullet tock, EOF]$
- $[ClockNoise \rightarrow tick \bullet tock, tick]$
- $[ClockNoise \rightarrow tick \bullet tock, tock]$
- $[ClockNoise \rightarrow tick tock \bullet, EOF]$
- $[ClockNoise \rightarrow tick tock \bullet, tick]$
- $[ClockNoise \rightarrow tick tock \bullet, tock]$

b. What are the FIRST sets of CN ?

<i>Answer:</i>	Symbol	FIRST
	$Goal$	tick
	$ClockNoise$	tick
	tick	tick
	tock	tock
	EOF	EOF

c. Construct the Canonical Collection of Sets of LR(1) Items for CN .

Answer:

$CC_0: \{ [Goal \rightarrow \bullet ClockNoise, EOF],$
 $[ClockNoise \rightarrow \bullet ClockNoise tick tock, \{EOF, tick\}],$
 $[ClockNoise \rightarrow \bullet tick tock, \{EOF, tick\}] \}$

$CC_1: \{ [Goal \rightarrow ClockNoise \bullet, EOF],$
 $[ClockNoise \rightarrow ClockNoise \bullet tick tock, \{EOF, tick\}] \}$

$CC_2: \{ [ClockNoise \rightarrow tick \bullet tock, \{EOF, tick\}] \}$

$CC_3: \{ [ClockNoise \rightarrow ClockNoise tick \bullet tock, \{EOF, tick\}] \}$

$CC_4: \{ [ClockNoise \rightarrow tick tock \bullet, \{EOF, tick\}] \}$

$CC_5: \{ [ClockNoise \rightarrow ClockNoise tick tock \bullet, \{EOF, tick\}] \}$

d. Derive the Action and Goto tables.

<i>Answer:</i>	ACTION			GOTO	
	EOF	tick	tock		$ClockNoise$
CC_0		s 2		CC_0	1
CC_1	acc	s 3		CC_1	0
CC_2			s 4	CC_2	0
CC_3			s 5	CC_3	0
CC_4	r 3	r 3		CC_4	0
CC_5	r 2	r 2		CC_5	0

10. Consider the following grammar:

$$\begin{aligned}
 \textit{Start} &\rightarrow S \\
 S &\rightarrow A a \\
 A &\rightarrow B C \\
 &\quad | B C f \\
 B &\rightarrow b \\
 C &\rightarrow c
 \end{aligned}$$

- a. Construct the canonical collection of sets of LR(1) items for this grammar.

Answer:

$$CC_0: \{ [Start \rightarrow \bullet S, EOF], [S \rightarrow \bullet A a, EOF], [A \rightarrow \bullet B C, a], [A \rightarrow \bullet B C f, a], [B \rightarrow \bullet b, c] \}$$

$$CC_1: \{ [Start \rightarrow S \bullet, EOF] \}$$

$$CC_2: \{ [S \rightarrow A \bullet a, EOF] \}$$

$$CC_3: \{ [A \rightarrow B \bullet C, a], [A \rightarrow B \bullet C f, a], [C \rightarrow \bullet c, \{a, f\}] \}$$

$$CC_4: \{ [B \rightarrow b \bullet, c] \}$$

$$CC_5: \{ [S \rightarrow A a \bullet, EOF] \}$$

$$CC_6: \{ [A \rightarrow B C \bullet, a], [A \rightarrow B C \bullet f, a] \}$$

$$CC_7: \{ [C \rightarrow c \bullet, \{a, f\}] \}$$

$$CC_8: \{ [A \rightarrow B C f \bullet, a] \}$$

b. Derive the Action and Goto tables.

Answer:

ACTION						GOTO				
	EOF	a	f	b	c		S	A	B	C
CC ₀				s 4		CC ₀	1	2	3	0
CC ₁	acc					CC ₁	0	0	0	0
CC ₂		s 5				CC ₂	0	0	0	0
CC ₃					s 7	CC ₃	0	0	0	6
CC ₄					r 5	CC ₄	0	0	0	0
CC ₅	r 2					CC ₅	0	0	0	0
CC ₆		r 3	s 8			CC ₆	0	0	0	0
CC ₇		r 6	r 6			CC ₇	0	0	0	0
CC ₈		r 4				CC ₈	0	0	0	0

c. Is the grammar LR(1)?

Answer: Yes. the table construction succeeds.

11. Consider a robot arm that accepts two commands: ∇ puts an apple in the bag and \triangle takes an apple out of the bag. Assume the robot arm starts with an empty bag.

A valid command sequence for the robot arm should have no prefix that contains more \triangle commands than ∇ commands. As examples, $\nabla\nabla\triangle\triangle$ and $\nabla\triangle\nabla$ are valid command sequences, but $\nabla\triangle\triangle\nabla$ and $\nabla\triangle\nabla\triangle\triangle$ are not.

- a. Write an LR(1) grammar that represents all the value command sequences for the robot arm.

Answer: A valid command sequence is a list of self-contained sets of commands, where each set consists of a ∇ , an optional set, and an optional \triangle . We will assume that a command sequence must have at least one command, a ∇ . We will assume that the grammar needs a unique start symbol. From that description, the grammar is straightforward.

$$\begin{array}{lll}
 \text{Start} \rightarrow \text{List} & \text{List} \rightarrow \text{Set List} & \text{Set} \rightarrow \nabla \text{Tail} \\
 & | \nabla \text{List} & \text{Tail} \rightarrow \text{Set } \triangle \\
 & | \epsilon & | \triangle
 \end{array}$$

b. Prove that the grammar is LR(1).

Answer: The easiest way to tell if a grammar is LR(1) is to build the canonical LR(1) tables for it. That can be done by hand or using a tool. Personally, I would accept a yacc or bison input file, along with the y.output file that showed the outcome. Students, however, tend to construct the table by hand.

First, we construct the FIRST Sets.

Symbol	FIRST
<i>Start</i>	$\nabla \epsilon$
<i>List</i>	$\nabla \epsilon$
<i>Set</i>	∇
<i>Tail</i>	$\nabla \triangle$
∇	∇
\triangle	\triangle
EOF	EOF

Next, we construct the Canonical Collection of Sets of LR(1) Items:

CC₀: [*Start* → • *List*, EOF] [*List* → • *Set List*, EOF] [*List* → • ∇ *List*, EOF]
 [*List* → • ϵ , EOF] [*Set* → • ∇ *Tail*, EOF] [*Set* → • ∇ *Tail*, ∇]

CC₁: [*Start* → *List* •, EOF]

CC₂: [*List* → • *Set List*, EOF] [*List* → *Set* • *List*, EOF]
 [*List* → • ∇ *List*, EOF] [*List* → • ϵ , EOF] [*Set* → • ∇ *Tail*, EOF]
 [*Set* → • ∇ *Tail*, ∇]

CC₃: [*List* → • *Set List*, EOF] [*List* → • ∇ *List*, EOF]
 [*List* → ∇ • *List*, EOF] [*List* → • ϵ , EOF] [*Set* → • ∇ *Tail*, EOF]
 [*Set* → • ∇ *Tail*, ∇] [*Set* → • ∇ *Tail*, \triangle] [*Set* → ∇ • *Tail*, EOF]
 [*Set* → ∇ • *Tail*, ∇] [*Tail* → • *Set* \triangle , EOF] [*Tail* → • *Set* \triangle , ∇]
 [*Tail* → • \triangle , EOF] [*Tail* → • \triangle , ∇]

CC₄: [*List* → ϵ •, EOF]

CC₅: [*List* → *Set List* •, EOF]

CC₆: [*List* → ∇ *List* •, EOF]

CC₇: [*List* → • *Set List*, EOF] [*List* → *Set* • *List*, EOF]
 [*List* → • ∇ *List*, EOF] [*List* → • ϵ , EOF] [*Set* → • ∇ *Tail*, EOF]
 [*Set* → • ∇ *Tail*, ∇] [*Tail* → *Set* • \triangle , EOF] [*Tail* → *Set* • \triangle , ∇]

$CC_8: [List \rightarrow \bullet Set\ List, EOF] [List \rightarrow \bullet \nabla List, EOF]$
 $[List \rightarrow \nabla \bullet List, EOF] [List \rightarrow \bullet \epsilon, EOF] [Set \rightarrow \bullet \nabla Tail, EOF]$
 $[Set \rightarrow \bullet \nabla Tail, \nabla] [Set \rightarrow \bullet \nabla Tail, \triangle] [Set \rightarrow \nabla \bullet Tail, EOF]$
 $[Set \rightarrow \nabla \bullet Tail, \nabla] [Set \rightarrow \nabla \bullet Tail, \triangle] [Tail \rightarrow \bullet Set\ \triangle, EOF]$
 $[Tail \rightarrow \bullet Set\ \triangle, \nabla] [Tail \rightarrow \bullet Set\ \triangle, \triangle] [Tail \rightarrow \bullet \triangle, EOF]$
 $[Tail \rightarrow \bullet \triangle, \nabla] [Tail \rightarrow \bullet \triangle, \triangle]$

$CC_9: [Set \rightarrow \nabla Tail\ \bullet, EOF] [Set \rightarrow \nabla Tail\ \bullet, \nabla]$

$CC_{10}: [Tail \rightarrow \triangle \bullet, EOF] [Tail \rightarrow \triangle \bullet, \nabla]$

$CC_{11}: [Tail \rightarrow Set\ \triangle \bullet, EOF] [Tail \rightarrow Set\ \triangle \bullet, \nabla]$

$CC_{12}: [List \rightarrow \bullet Set\ List, EOF] [List \rightarrow Set\ \bullet List, EOF]$
 $[List \rightarrow \bullet \nabla List, EOF] [List \rightarrow \bullet \epsilon, EOF] [Set \rightarrow \bullet \nabla Tail, EOF]$
 $[Set \rightarrow \bullet \nabla Tail, \nabla] [Tail \rightarrow Set\ \bullet \triangle, EOF] [Tail \rightarrow Set\ \bullet \triangle, \nabla]$
 $[Tail \rightarrow Set\ \bullet \triangle, \triangle]$

$CC_{13}: [Set \rightarrow \nabla Tail\ \bullet, EOF] [Set \rightarrow \nabla Tail\ \bullet, \nabla] [Set \rightarrow \nabla Tail\ \bullet, \triangle]$

$CC_{14}: [Tail \rightarrow \triangle \bullet, EOF] [Tail \rightarrow \triangle \bullet, \nabla] [Tail \rightarrow \triangle \bullet, \triangle]$

$CC_{15}: [Tail \rightarrow Set\ \triangle \bullet, EOF] [Tail \rightarrow Set\ \triangle \bullet, \nabla] [Tail \rightarrow Set\ \triangle \bullet, \triangle]$

And, finally, we build the ACTION and GOTO tables:

	EOF	∇	ϵ	\triangle		List	Set	Tail
cc_0		s 3	s 4		cc_0	1	2	0
cc_1	acc				cc_1	0	0	0
cc_2		s 3	s 4		cc_2	5	2	0
cc_3		s 8	s 4	s 10	cc_3	6	7	9
cc_4	r 4				cc_4	0	0	0
cc_5	r 2				cc_5	0	0	0
cc_6	r 3				cc_6	0	0	0
cc_7		s 3	s 4	s 11	cc_7	5	2	0
cc_8		s 8	s 4	s 14	cc_8	6	12	13
cc_9	r 5	r 5			cc_9	0	0	0
cc_{10}	r 7	r 7			cc_{10}	0	0	0
cc_{11}	r 6	r 6			cc_{11}	0	0	0
cc_{12}		s 3	s 4	s 15	cc_{12}	5	2	0
cc_{13}	r 5	r 5		r 5	cc_{13}	0	0	0
cc_{14}	r 7	r 7		r 7	cc_{14}	0	0	0
cc_{15}	r 6	r 6		r 6	cc_{15}	0	0	0

The construction succeeds, so the grammar is LR(1).

12. The following grammar has no known LL(1) equivalent:

0	<i>Start</i>	\rightarrow	<i>A</i>
1			<i>B</i>
2	<i>A</i>	\rightarrow	(<i>A</i>)
3			<u>a</u>
4	<i>B</i>	\rightarrow	(<i>B</i> ≥
5			<u>b</u>

Show that the grammar is LR(1).

Answer: The easiest way to show that a grammar is LR(1) is to build the canonical LR(1) tables for it. That can be done by hand or using a tool. Personally, I would accept a yacc or bison input file, along with the y.output file that showed the outcome. Students, however, tend to construct the table by hand.

First, we construct FIRST sets.

Symbol	FIRST
EOF	EOF
<i>Start</i>	(a b
<i>A</i>	(a
<i>B</i>	(b
((
))
a	a
>	>
b	b

Next, we construct the Canonical Collection of Sets of LR(1) Items:

CC_0 : $[Start \rightarrow \bullet B, EOF]$ $[B \rightarrow \bullet (B >, EOF]$ $[B \rightarrow \bullet b, EOF]$

CC_1 : $[Start \rightarrow B \bullet, EOF]$

CC_2 : $[B \rightarrow \bullet (B >, >]$ $[B \rightarrow (\bullet B >, EOF]$ $[B \rightarrow \bullet b, >]$

CC_3 : $[B \rightarrow b \bullet, EOF]$

CC_4 : $[B \rightarrow (B \bullet >, EOF]$

CC_5 : $[B \rightarrow \bullet (B >, >]$ $[B \rightarrow (\bullet B >, >]$ $[B \rightarrow \bullet b, >]$

CC_6 : $[B \rightarrow b \bullet, >]$

CC_7 : $[B \rightarrow (B > \bullet, EOF]$

CC_8 : $[B \rightarrow (B \bullet >, >]$

CC_9 : $[B \rightarrow (B > \bullet, >]$

And, finally, we build the ACTION and GOTO tables:

ACTION						GOTO			
	EOF	()	a	>	b		A	B
cc_0		s 2				s 3	cc_0	0	1
cc_1	acc						cc_1	0	0
cc_2		s 5				s 6	cc_2	0	4
cc_3	r 6						cc_3	0	0
cc_4					s 7		cc_4	0	0
cc_5		s 5				s 6	cc_5	0	8
cc_6					r 6		cc_6	0	0
cc_7	r 5						cc_7	0	0
cc_8					s 9		cc_8	0	0
cc_9					5		cc_9	0	0

The fact that the table construction succeeds (i.e., it does not define any table location twice) implies that the grammar is LR(1).

Section 3.6

13. Write a grammar for expressions that can include binary operators (+ and \times), unary minus (-), autoincrement (++), and autodecrement (--) with their customary precedence. Assume that repeated unary minuses are not allowed, but that repeated autoincrement and autodecrement operators are allowed.

Answer:

$Expr \rightarrow Expr + Term$	$Factor \rightarrow (Expr)$	$Auto \rightarrow ++ Auto$
$Expr - Term$	$- (Expr)$	$-- Auto$
$Term$	$- num$	$Auto ++$
	num	$Auto --$
$Term \rightarrow Term \times Factor$	$- id$	id
$Term \div Factor$	$- Auto$	
$Factor$	$Auto$	

Section 3.7

14. Consider the task of building a parser for the programming language Scheme. Contrast the effort required for a top-down recursive-descent parser with that needed for a table-driven LR(1) parser. (Assume that you already have an LR(1) table generator.)

Answer: While the answer is subjective, a good case can be made for either technique. The syntax of Scheme is sufficiently simple that a recursive descent parser will be simple, compact, and fast. By the same argument, writing an LR(1) grammar for Scheme is easy.

Most of the difficulty lies in deriving a grammar that is suitable for the selected parsing technique. For the top-down parser, the grammar must avoid left recursion. That simple fact may shift the balance slightly in favor of the LR(1) solution.

15. The text describes a manual technique for eliminating useless productions in a grammar.

- a. Can you modify the LR(1) table-construction algorithm so that it automatically eliminates the overhead from useless productions?

Answer: The author of the table constructor can easily add code that recognizes any LR(1) item of the form $[A \rightarrow \bullet B]$, for a non-terminal symbol B , and replaces it with a series of LR(1) items that expand B into its various alternative right-hand sides. This transformation has the effect of eliminating productions of the form $A \rightarrow B$.

- b. Even though a production is syntactically useless, it may serve a practical purpose. For example, the compiler writer might associate a syntax-directed action (see Chapter 4) with the useless production. How should your modified table-construction algorithm handle an action associated with a useless production?

Answer: The transformation needs to check for the presence of a syntax-directed action before it eliminates a production. If a syntactically-useless production has an action that occurs when the parser reduces by it, the production serves a purpose and should not be eliminated.

CHAPTER 4

Section 4.2

1. In Scheme, the `+` operator is overloaded. Given that Scheme is dynamically typed, describe a method to type check an operation of the form `(+ a b)` where `a` and `b` may be of any type that is valid for the `+` operator.

Answer: The implementation can associate a tag with the run-time representation of each value. The implementation of `+` can verify the types of its operands by checking their tag fields. The results of those checks allow `+` to select the correct sequence of actions for any set of operands.

2. Some languages, such as APL or PHP, neither require variable declarations nor enforce consistency between assignments to the same variable. (A program can assign the integer 10 to `x` and later assign the string value “book” to `x` in the same scope.) This style of programming is sometimes called *type juggling*.

Suppose that you have an existing implementation of a language that has no declarations but requires type-consistent uses. How could you modify it to allow type juggling?

Answer: The implementation must consider each assignment as defining both a new value and, potentially, a new type. When a variable is used, its type is the type last assigned to the variable.

With control-flow, the programmer can write code that assigns distinct types to a single variable along different control-flow paths. Such use patterns may require run-time checking. The implementation should attach a tag (a run-time location that holds a value denoting the item’s data type) to each variable. The implementation can set the tag’s value appropriately at each assignment and explicitly check the tag’s value at each use. By tracking the type information at runtime, the implementation can ensure that each value is handled correctly.

Section 4.3

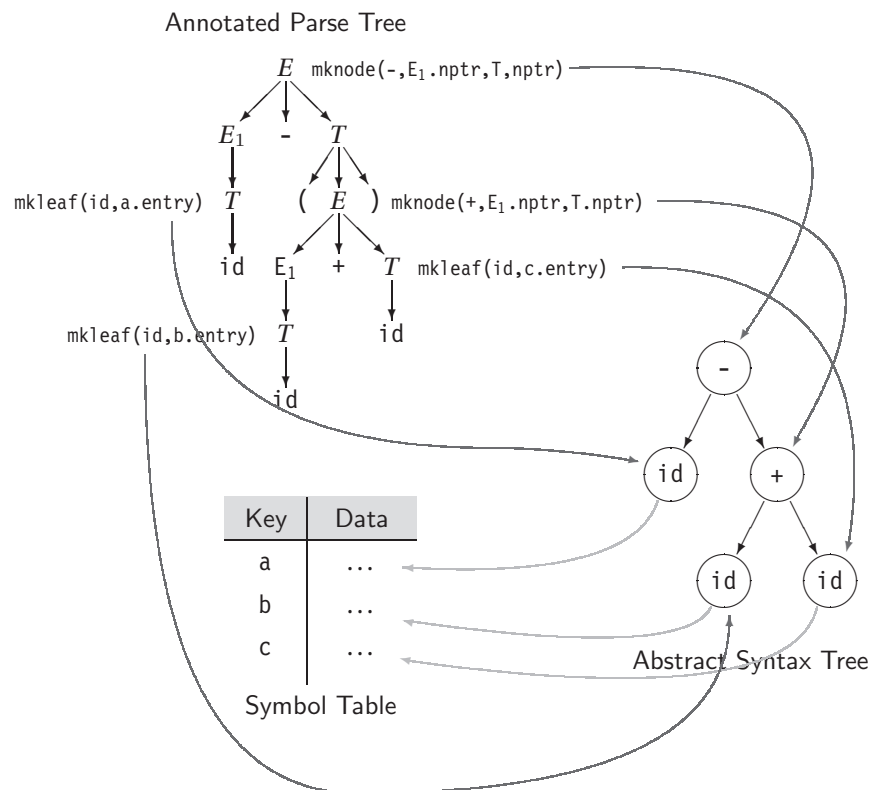
3. Based on the following evaluation rules, draw an annotated parse tree that shows how the syntax tree for $a - (b + c)$ is constructed.

Production

Evaluation Rules

 $E_0 \rightarrow E_1 + T \quad \{ E_0.nptr \leftarrow mknode(+, E_1.nptr, T.nptr) \}$
 $E_0 \rightarrow E_1 - T \quad \{ E_0.nptr \leftarrow mknode(-, E_1.nptr, T.nptr) \}$
 $E_0 \rightarrow T \quad \{ E_0.nptr \leftarrow T.nptr \}$
 $T \rightarrow (E) \quad \{ T.nptr \leftarrow E.nptr \}$
 $T \rightarrow id \quad \{ T.nptr \leftarrow mkleaf(id, id.entry) \}$

Answer:



4. Use the attribute-grammar paradigm to write an interpreter for the classic expression grammar. Assume that each name has a value attribute and a lexeme attribute. Assume that all attributes are already defined and that all values will always have the same type.

Answer: Assume the following grammar:

$Expr \rightarrow Expr + Term$	$Term \rightarrow Term \times Factor$	$Factor \rightarrow (Expr)$
$ Expr - Term$	$ Term \div Factor$	$ number$
$ Term$	$ Factor$	$ name$

Now, a set of attribution rules:

Production	Attribution Rule
1 $Expr \rightarrow Expr + Term$	$Expr_0.value \leftarrow Expr_1.value + Term.value$
2 $Expr \rightarrow Expr - Term$	$Expr_0.value \leftarrow Expr_1.value - Term.value$
3 $Expr \rightarrow Term$	$Expr.value \leftarrow Term.value$
4 $Term \rightarrow Term \times Factor$	$Term_0.value \leftarrow Term_1.value \times Factor$
5 $Term \rightarrow Term \div Factor$	$Term_0.value \leftarrow Term_1.value \div Factor$
6 $Term \rightarrow Factor$	$Term.value \leftarrow Factor.value$
7 $Factor \rightarrow (Expr)$	$Factor.value \leftarrow Expr.value$
8 $Factor \rightarrow number$	$Factor.value \leftarrow ValueOf(number.lexeme)$
9 $Factor \rightarrow name$	$Factor.value \leftarrow name.value$

Here, *ValueOf* is a simple function that computes a value from a lexeme—the text representation of the number.

5. Write a grammar to describe all binary numbers that are multiples of four. Add attribution rules to the grammar that will annotate the start symbol of a syntax tree with an attribute value that contains the decimal value of the binary number.

Answer: To simplify the grammar, we will assume that the scanner can recognize the final double zero in a binary number that is a multiple of four as a single word. (If we try to recognize it as two consecutive zeros in the grammar, we hit shift reduce problems that seriously complicate the grammar.) To recognize *DZero* requires two characters of lookahead, but it has a small constant additional cost. With that caveat, we can use the following LR(1) grammar:

$$\begin{aligned}
 \text{Number} &\rightarrow 0 \\
 &\quad | \quad 1 \text{ List} \\
 \text{List} &\rightarrow 1 \text{ List} \\
 &\quad | \quad 0 \text{ List} \\
 &\quad | \quad \text{DZero}
 \end{aligned}$$

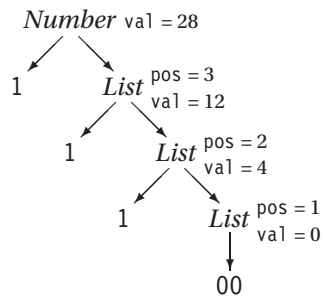
Next, we add attribution rules to the grammar:

Production	Attribution Rule
1 $\text{Number} \rightarrow 0$	$\text{Number.val} \leftarrow 0$
2 $\text{Number} \rightarrow 1 \text{ List}$	$\text{Number.val} \leftarrow 2^{(\text{List.pos}+1)} + \text{List.val}$
3 $\text{List}_0 \rightarrow 1 \text{ List}_1$	$\text{List}_0.\text{pos} \leftarrow \text{List}_1.\text{pos} + 1$ $\text{List}_0.\text{val} \leftarrow 2^{\text{List}_0.\text{pos}} + \text{List}_1.\text{val}$
4 $\text{List}_0 \rightarrow 0 \text{ List}_1$	$\text{List}_0.\text{pos} \leftarrow \text{List}_1.\text{pos} + 1$ $\text{List}_0.\text{val} \leftarrow \text{List}_1.\text{val}$
5 $\text{List} \rightarrow \text{DZero}$	$\text{List.val} \leftarrow 0$ $\text{List.pos} \leftarrow 1$

6. Using the grammar defined in the previous exercise, build the syntax tree for the binary number 11100.

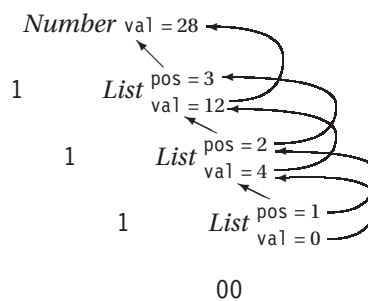
(a) Show all the attributes in the tree with their corresponding values.

Answer:



(b) Draw the attribute dependence graph for the syntax tree and classify all attributes as being either synthesized or inherited.

Answer:



All attributes are synthesized.

Section 4.4

7. A Pascal program can declare two integer variables a and b with the syntax

var a, b: int

This declaration might be described with the following grammar:

```

VarDecl  →  var IDList : TypeID
IDList   →  IDList, ID
          |  ID
  
```

where *IDList* derives a comma-separated list of variable names and *TypeID* derives a valid Pascal type. You may find it necessary to rewrite the grammar.

- (a) Write an attribute grammar that assigns the correct data type to each declared variable.

Answer:

$VarDecl \rightarrow \text{var } IDList : TypeID$	$TypeID.type \leftarrow \text{type of } TypeID;$ $IDList.type \leftarrow TypeID.type$
$IDList_0 \rightarrow IDList_1, ID$	$ID.type \leftarrow IDList_0.type;$ $IDList_1.type \leftarrow IDList_0.type$
$ ID$	$ID.type \leftarrow IDList_0.type;$

- (b) Write an ad hoc syntax-directed translation scheme that assigns the correct data type to each declared variable.

$VarDecl \rightarrow \text{var } IDList : TypeID$	for each ID_k in \$1 $ID_k \leftarrow \text{type of } TypeID$
Answer: $IDList_0 \rightarrow IDList_1, ID$	$$$ \leftarrow \text{cons}(\$1, \text{makenode}(ID))$
$ ID$	$$$ \leftarrow \text{makenode}(ID)$

The iterator used in the action for the first production is somewhat abstract. This can be done with an iterative list traversal, an abstract function such as Scheme's `map`, or a hand-coded recursive function

- (c) Can either scheme operate in a single pass over the syntax tree?

Answer: With an appropriate evaluation order (e.g., a right-to-left pre-order walk), an attribute grammar evaluator can process the rules from part (a) in a single pass.

The ad-hoc syntax-directed translation scheme must parse the list of *IDs* first, before it sees the declared type. The scheme in part (b) of this solution builds a list of those elements, discovers the type, and walks the list to set their types. Thus, it takes two passes over the syntax tree.

8. Sometimes, the compiler writer can move an issue across the boundary between context-free and context-sensitive analysis. Consider, for example, the classic ambiguity that arises between function invocation and array references in FORTRAN 77 (and other languages). These constructs might be added to the classic expression grammar using the productions:

$$\begin{array}{ll} \textit{Factor} & \rightarrow \text{ name (ExprList) } \\ \textit{ExprList} & \rightarrow \text{ ExprList , Expr } \\ & \quad | \text{ Expr } \end{array}$$

Here, the only difference between a function invocation and an array reference lies in how the name is declared.

In previous chapters, we have discussed using cooperation between the scanner and the parser to disambiguate these constructs. Can the problem be solved during context-sensitive analysis? Which solution is preferable?

Answer: To solve the problem during context sensitive analysis, the compiler can use a grammar that contains only a single derivation for the two constructs (function invocations and array references). The parser can build a representation that preserves all the information. During context-sensitive analysis, when the type information is naturally available, the compiler can rewrite the intermediate representation as a function invocation or an array reference, as appropriate.

Both approaches (scanner-parser cooperation and rewriting during context-sensitive analysis) work. The scanner-parser scheme isolates the problem in the syntax-oriented portions of the front end. It may add a small overhead to processing every name. The context-sensitive analysis scheme simplifies the scanner and parser, defers the issue until complete type information is available, and pushes the complication later into compilation.

Neither approach is a clear win. The best answer is to change the syntax of the source language and avoid the ambiguity. Unfortunately, that is not an option for FORTRAN. We are at least 40 years too late!

9. Sometimes, a language specification uses context-sensitive mechanisms to check properties that can be tested in a context-free way. Consider the grammar fragment in Figure 4.16 on page 208. It allows an arbitrary number of *StorageClass* specifiers when, in fact, the standard restricts a declaration to a single *StorageClass* specifier.

- (a) Rewrite the grammar to enforce the restriction grammatically.

Answer:

$$\begin{array}{lcl} \text{SpecifierList} & \rightarrow & \text{StorageClass TypeList} \\ & | & \text{TypeList} \\ \text{TypeList} & \rightarrow & \text{TypeSpecifier TypeList} \\ & | & \text{TypeSpecifier} \end{array}$$

- (b) Similarly, the language allows only a limited set of combinations of *TypeSpecifier*. *long* is allowed with either *int* or *float*; *short* is allowed only with *int*. Either *signed* or *unsigned* can appear with any form of *int*. *signed* may also appear on *char*. Can these restrictions be written into the grammar?

Answer: Yes. They can be added to the grammar at the cost of additional productions and states in the parser.

- (c) Propose an explanation for why the authors structured the grammar as they did.

Answer: Reducing the set of terminals (by combining all *StorageClass* specifiers into a single terminal and all *TypeSpecifiers* into a single terminal) shrinks the tables and produces a more compact parser, at the cost of checking the restrictions during context-sensitive analysis.

- (d) Do your revisions to the grammar change the overall speed of the parser? In building a parser for C, would you use the grammar like the one in Figure 4.16, or would you prefer your revised grammar? Justify your answer.

Answer: The version shown in Figure 4.16 also avoids one reduction on each use of a *StorageClass* or *TypeSpecifier*, so it makes the parser marginally faster. As to which version of the grammar is preferred, that is a subjective matter.

Encoding the restriction in the grammar avoids the need to write code that checks the restrictions. Doing the check during context-sensitive analysis allows the compiler writer to produce a more useful diagnostic message when an error is discovered.

Hint: The scanner returned a single token type for any of the *StorageClass* values and another token type for any of the *TypeSpecifiers*.

Section 4.5

10. Object-oriented languages allow operator and function overloading. In these languages, the function name is not always a unique identifier, since you can have multiple related definitions, as in

```
void Show(int);
void Show(char *);
void Show(float);
```

For lookup purposes, the compiler must construct a distinct identifier for each function. Sometimes, such overloaded functions will have different return types, as well. How would you create distinct identifiers for such functions?

Answer: Assume that the compiler has a character, such as \$, that is not allowed in an identifier. The compiler can create distinct function names by concatenating the function name, the return type, and the parameter types, separating each with \$. Applying this scheme to the declarations for Show yields three distinct identifiers:

```
Show$void$int
Show$void$char*
Show$void$float
```

11. Inheritance can create problems for the implementation of object-oriented languages. When object type *A* is a parent of object type *B*, a program can assign a “pointer to *B*” to a “pointer to *A*,” with syntax such as $a \leftarrow b$. This should not cause problems since everything that *A* can do, *B* can also do. However, one cannot assign a “pointer to *A*” to a “pointer to *B*,” since object class *B* can implement methods that object class *A* does not.

Design a mechanism that can use ad hoc syntax-directed translation to determine whether or not a pointer assignment of this kind is allowed.

Answer: If the declarations of *A* and *B* must be present at compile time, the compiler can build a tree to represent the inheritance hierarchy. Given an inheritance tree and an assignment “ $x \leftarrow y$ ” the compiler can determine if the type of *x* is a parent of the type of *y* by walking up the inheritance tree from the node for type of *y*.

Section 5.2

1. A parse tree contains much more information than an abstract syntax tree.

- a. In what circumstances might you need information that is found in the parse tree but not the abstract syntax tree?

Answer: In any application where the compiler needs to generate output that the programmer will read, the parse tree may contain information that helps produce results that more closely resemble the original input program.

- b. What is the relationship between the size of the input program and its parse tree? Its abstract syntax tree?

Answer: The parse tree has, to a first approximation, size proportional to the grammatical derivation of the input program (recall Chapter 3). That is, it has (roughly) a node for each word in the input program, along with a node for every reduction in the parse. Of course, some multi-word constructs, such as a C `for` loop, may be represented with fewer nodes than they have words.

The abstract syntax tree has, to a first approximation, a node for each word in the input program. It may have more nodes if it expands source-level abstractions, such as an array reference (see Chapter 7), into an address calculation. It may have more nodes if it represents control-flow structures implicitly—for example, a single node for a loop construct such as the C `for` loop, rather than a node for each of the parentheses and semicolons.

- c. Propose an algorithm to recover a program's parse tree from its abstract syntax tree.

Answer: The simplest algorithm would be to regenerate the syntax from the AST and to parse it again.

2. Write an algorithm to convert an expression tree into a DAG.

Answer: To facilitate comparisons among subexpressions, we need a hash table that maps subexpressions (or nodes in the expression tree) into unique integers (which can be as simple as the index of the corresponding table entry.) The algorithm traverses the expression tree in a post-order walk. Assume, without loss of generality, that each node has no more than two children.

```

BuildDAG( root ) {
    initialize hash table entries to null
    DAG ← Walk( root )
    return DAG
}

Walk( node ) {
    if node is a leaf then
        if lexeme(node) is not in the hash table then
            create an entry for it in the table
            create a node x for it in the new DAG
            store x in the table entry
        else let x be the DAG node stored in the table entry
    else {
        if LeftChild(node) exists
            then Left ← Walk(LeftChild(node))
            else Left ← null
        if RightChild(node) exists
            then Right ← Walk(RightChild(node))
            else Right ← null
        if the key "op(node),Left,Right" is not in the hash table then
            create an entry for it in the table
            create a node x in the new DAG
            set LeftChild(x) to be Left
            set RightChild(x) to be Right
        else let x be the DAG node stored in the table entry
    }
    return x
}

```

Section 5.3

3. Show how the following code fragment

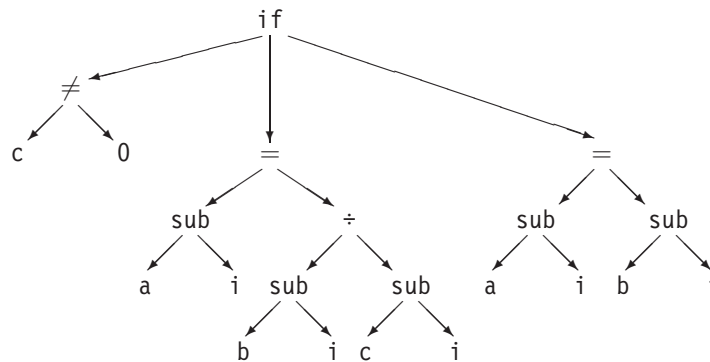
```

if (c[i]  $\neq$  0)
    then a[i]  $\leftarrow$  b[i]  $\div$  c[i];
    else a[i]  $\leftarrow$  b[i];

```

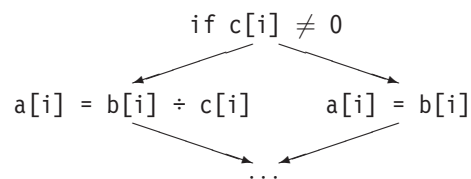
might be represented in an abstract syntax tree, in a control-flow graph, and in quadruples. Discuss the advantages of each representation. For what applications would one representation be preferable to the others?

Answer: *Abstract syntax tree:* The abstract syntax tree has a number of advantages. Because it is close to the actual syntax of the original code, it is easy to regenerate source code in a walk of the AST. The AST typically represents code at a near-source level of abstraction, which makes it a good representation for some kinds of analysis and transformation. For example, it is easier to move a string-copy operation as a compact subtree in an AST than as a long sequence of near-assembly level operations.



Among the disadvantages of an AST are the fact that it does not represent control-flow in any explicit way.

Control-flow graph: The control-flow graph is useful for reasoning about the flow of control in the original code. In the example, the CFG makes it obvious that one of the two assignments executes, where the AST does not.



Among the disadvantages of a CFG are the fact that it does not represent the computation in blocks in any explicit way.

Quadruples: A low-level representation such as quadruples has the advantage that it exposes *all* of the computation. The representation is easy to understand and easy to manipulate. Simple data structures can be used to represent quadruples.

```

r1      ←  i × 4
r2      ←  r1 + @c
r3      ←  Mem(r2)
r4      ←  0
cmp       r3 = r4
bneq     →  L1, L2

L1: r5    ←  i × 4
r6      ←  r5 + @b
r7      ←  Mem(r6)
r8      ←  r5 + @c
r9      ←  Mem(r8)
r10     ←  r7 ÷ r9
r11     ←  r5 + @a
Mem(r11) ←  r10
b         →  L3

L2: r12   ←  i × 4
r13     ←  r12 + @b
r14     ←  Mem(r13)
r15     ←  r12 + @a
Mem(r15) ←  r14
b         →  L3

L3: nop

```

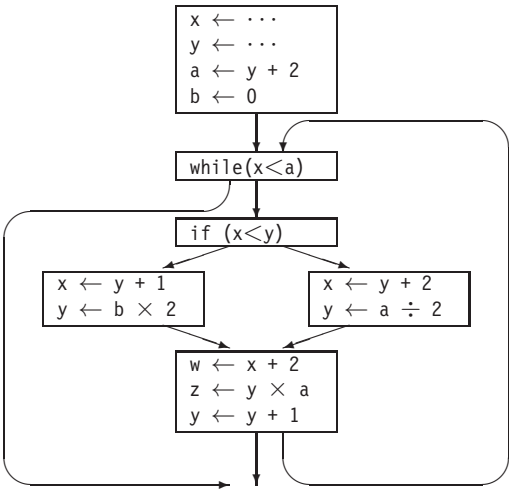
One disadvantage of a low-level linear representation such as quadruples is the fact that all detail is exposed. Transformations that try to manipulate groups of operations, such as moving an entire assignment statement or a string copy operation, must manipulate each operation. In the example, the block of operations starting at L₁ is a single source-level statement, as is the block at L₂.

```
...
x ← ...
y ← ...
a ← y + 2
b ← 0
while(x < a)
  if (y < x)
    x ← y + 1
    y ← b × 2
  else
    x ← y + 2
    y ← a ÷ 2;
w ← x + 2
z ← y × a
y ← y + 1
```

FIGURE 5.13 Code Fragment for Exercise 4

4. Examine the code fragment shown in Figure 5.13. Draw its CFG and show its SSA form as a linear code.

Answer:



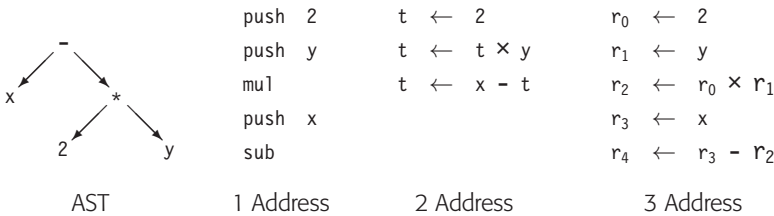
Control-flow Graph

```
x0 ← ...
y0 ← ...
a0 ← y0
b0 ← 0
if (x0 > a0) goto 12
11: x1 ← φ(x0, x4)
y1 ← φ(y0, y5)
if (x1 < y1)
  x2 ← y1 + 1
  y2 ← b0 × 2
else
  x3 ← y1 + 2
  y3 ← a0 / 2
x4 ← φ(x2, x3)
y4 ← φ(y2, y3)
w0 ← x4 + 2
z0 ← y4 × a0
y5 ← y4 + 1
if (x4 < a0) goto 11
12: ...
```

SSA Form

5. Show how the expression $x - 2 \times y$ might be translated into an abstract syntax tree, one-address code, two-address code, and three-address code.

Answer:



Note that the 1-address code assumes that sub computes $\text{stack}[\text{top}] - \text{stack}[\text{top}-1]$.

6. Given a linear list of ILOC operations, develop an algorithm that finds the basic blocks in the ILOC code. Extend your algorithm to build a control-flow graph to represent the connections between blocks.

Taken from § 9.2 of EaC, 1st Edition.

Answer: At its simplest, the CFG identifies the beginning and the end of each basic block and connects the resulting blocks with edges that describe the possible transfers of control among blocks. Initially, let's assume that the CFG-builder receives, as input, a simple, linear IR that represents a procedure in a classic, Algol-like language. At the end of this section, we briefly explore some of the complications that can arise in CFG construction.

To begin, the compiler must find the beginning and the end of each basic block. In a linear IR, the initial operation of a block is sometimes called a *leader*. An operation is a leader if it is the first operation in the procedure, or if it has a label that is, potentially, the target of some branch. The compiler can identify all the leaders in a single pass over the IR, shown on the left-hand side of the unnumbered Figure labelled "Building a Control-flow Graph". It iterates over the operations in the program, in order, finds the labelled statements, and records them as leaders.¹

The second pass finds every operation that ends a block. It assumes that every block ends with a branch or a jump, and that branches specify labels for both outcomes— a "branch taken" label and a "branch not taken" label. This simplifies the handling of blocks and allows the compiler's back end to

1. If the code contains spurious labels, this may unnecessarily split blocks. A more complex algorithm would only add branch and jump targets to the leader set. If the code contains any ambiguous jumps—that is, a jump to an address in a register—then it must include all labelled statements as leaders anyway.

<pre> next ← 1 leader[next++] ← 1 for i ← 1 to n if op_i has a label l_i then leader[next++] ← i create a CFG node for l_i </pre>	<pre> for i ← 1 to next - 1 j ← leader[i] + 1 while (j ≤ n and op_j ∉ leader) j ← j + 1 j ← j - 1 last[i] ← j if op_j is "cbr r_k→l₁,l₂" then add edge from j to node for l₁ add edge from j to node for l₂ else if op_j is "jump l→l₁" then add edge from j to node for l₁ else if op_j is "jump→r₁" then add edges from j to all other nodes </pre>
<i>Finding Leaders</i>	<i>Finding Last & Adding Edges</i>

Building a Control-flow Graph

choose which path will follow the “fall through” case of a branch. (For the moment, assume branches have no delay slots.)

To find the end of each block, the algorithm iterates through the blocks, in order of their appearance in the *leader* array. It walks forward through the IR until it finds the leader of the following block. The operation immediately before that leader ends the current block. The algorithm records that operation’s index in *last[i]*, so that the pair $\langle \text{leader}[i], \text{last}[i] \rangle$ describes block *i*. It adds the needed edges to the CFG. This pass takes time proportional to the size of the IR program.

As we have assumed throughout the book, the CFG should have a unique entry node n_0 and a unique exit node n_f . The underlying code should have this shape. If it does not, a simple postpass over the graph can create n_0 and n_f .

Complications in CFG Construction Features of the IR, the target architecture, and even the source language can complicate this process. For example, an ambiguous jump—ILOC’s *jump* operation—may force the compiler to add an edge from the jump to every labelled block in the procedure. This can add edges to the CFG that never occur at run time, degrade the quality of the data-flow information that the compiler derives, and reduce the effectiveness of optimization. In general, the compiler writer should avoid creating ambiguous jumps when the targets can be known, even if this requires

extra analysis during IR generation. The ILOC examples in this book avoid using the jump-to-register operation, `jump`, in favor of the immediate jump, `jumpI`.

The compiler writer can improve this situation, to some extent, by including features in the IR that record potential jump targets. ILOC includes the `tbl` pseudo-operation to let the compiler record the potential targets of an ambiguous jump. Anytime that the compiler generates a `jump`, it should follow the branch with a set of `tbl` operations that record the possible targets for the branch. When it builds the CFG, the compiler can use these hints to limit the number of spurious edges. Any `jump` followed by a set of `tbls` gets an edge to each block named in those `tbl` operations. Any `jump` that has no `tbls` generates an edge to each labelled block.

The target architecture can complicate the process. The algorithm given in the unnumbered figure assumes that all leaders, except the first, are labelled. If the target machine has fall-through branches, the algorithm must be extended to recognize unlabeled statements that receive control on a fall-through path. If the code has been scheduled and it models branch-delay slots, the problem becomes harder. A labelled statement that sits in a branch delay slot is a member of two distinct blocks. The compiler can cure this by replication—creating a new (unlabeled) copy of the operations in the delay slots.

Delay slots also complicate finding the end of a block. If a branch or jump can occur in a branch delay slot, the CFG builder must walk forward from the leader to find the block-ending branch—the first branch it encounters. Branches in the delay slot of a block-ending branch can, themselves, be pending on entry to the target block. They can split the target block and force creation of new blocks and new edges. The analysis required to create a CFG in these circumstances is much more complex.

Some languages, such as Pascal and Algol, allow jumps to labels outside the current procedure. This transfer can be modelled in the current procedure with a branch to a CFG node created to represent the target. The complication arises on the other end of the branch, where a transfer from an unknown source can reach a block. For this reason, non-local gotos are usually restricted to lexically nested procedures, where the compiler can see the CFGs for all related procedures at once. In that situation, it can insert the correct edge and model the effects correctly. Otherwise any block whose label can be the target of a non-local goto must be treated as if it has an unknown predecessor in the CFG.

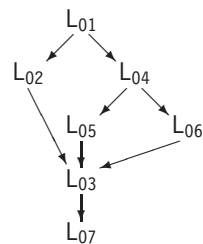
L01: add $r_a, r_b \Rightarrow r_1$	L05: add $r_9, r_b \Rightarrow r_{11}$
add $r_c, r_d \Rightarrow r_2$	add $r_a, r_b \Rightarrow r_{12}$
add $r_1, r_2 \Rightarrow r_3$	add $r_c, r_d \Rightarrow r_{13}$
add $r_a, r_b \Rightarrow r_4$	i2i $r_a \Rightarrow r_{13}$
cmp_LT $r_1, r_2 \Rightarrow r_5$	add $r_{13}, r_b \Rightarrow r_{14}$
cbr $r_5 \rightarrow L02, L04$	multI $r_{12}, 17 \Rightarrow r_{15}$
L02: add $r_a, r_b \Rightarrow r_6$	jumpI $\rightarrow L03$
multI $r_6, 17 \Rightarrow r_7$	L06: add $r_1, r_2 \Rightarrow r_{16}$
jumpI $\rightarrow L03$	i2i $r_2 \Rightarrow r_{17}$
L03: add $r_a, r_b \Rightarrow r_{22}$	i2i $r_1 \Rightarrow r_{18}$
multI $r_{22}, 17 \Rightarrow r_{23}$	add $r_{17}, r_{18} \Rightarrow r_{19}$
jumpI $\rightarrow L07$	add $r_{18}, r_{17} \Rightarrow r_{20}$
L04: add $r_c, r_d \Rightarrow r_8$	multI $r_1, 17 \Rightarrow r_{21}$
i2i $r_a \Rightarrow r_9$	jumpI $\rightarrow L03$
cmp_LT $r_9, r_d \Rightarrow r_{10}$	L07: nop
cbr $r_{10} \rightarrow L05, L06$	

FIGURE 5.14 Code Fragment for Exercise 7

Section 5.4

7. For the code shown in Figure 5.14 find the basic blocks and construct the CFG.

Answer: The basic blocks are easy to find. Each labelled statement is a leader. Each block ends at the first cbr or jumpI. The CFG is shown below.



```

static int max = 0;
void A(int b, int e)
{
    int a, c, d, p;
    a = B(b);
    if (b > 100) {
        c = a + b;
        d = c * 5 + e;
    }
    else
        c = a * b;
    *p = c;
    C(&p);
}

int B(int k)
{
    int x, y;
    x = pow(2, k);
    y = x * 5;
    return y;
}

void C(int *p)
{
    if (*p > max)
        max = *p;
}

```

FIGURE 5.15 Code for Exercise 8

8. Consider the three C procedures shown in Figure 5.15.

- a. Suppose a compiler uses a register-to-register memory model. Which variables in procedures A, B, and C would the compiler be forced to store in memory? Justify your answers.

Answer: In a register-to-register model, the compiler would need to store *p* in procedure A in memory, because A passes *p*'s address to C. The other variables are all used locally or passed as values to other procedures. Thus, the compiler can safely keep them in registers.

- b. Suppose a compiler uses a memory-to-memory model. Consider the execution of the two statements that are in the *if* clause of the *if-else* construct. If the compiler has two registers available at that point in the computation, how many loads and stores would the compiler need to issue in order to load values in registers and store them back to memory during execution of those two statements? What if the compiler has three registers available?

Answer: With two registers, the compiler might generate the following code:

<i>for</i> c = a + b	<i>for</i> d = c * 5 + e
load a \Rightarrow r ₁	loadI 5 \Rightarrow r ₁
load b \Rightarrow r ₂	load c \Rightarrow r ₂
add r ₁ , r ₂ \Rightarrow r ₃	mult r ₁ , r ₂ \Rightarrow r ₁
store r ₁ \Rightarrow c	load e \Rightarrow r ₂
	add r ₁ , r ₂ \Rightarrow r ₁
	store r ₁ \Rightarrow d

With a third register, the compiler could keep the value of *c* alive, avoiding the load before the second use (in the computation of *c* * 5). It cannot avoid the store to *c* because of the later use.

9. In FORTRAN, two variables can be forced to begin at the same storage location with an equivalence statement. For example, the following statement forces *a* and *b* to share storage:

equivalence (a,b)

Can the compiler keep a local variable in a register throughout the procedure if that variable appears in an equivalence statement? Justify your answer.

Answer: In general, the compiler cannot keep either *a* or *b* in a register if they appear together in an equivalence statement. More specifically, it cannot keep *a* in a register across an assignment to *b* (and vice-versa) because the assignment must change the value of both variables.

Typically, compilers handle this problem by keeping in memory any variable mentioned in an equivalence statement, as any other variable that might have an ambiguous name. (Examples include call-by-reference parameters and variables involved in pointer-based references.)

Section 5.5

10. Some part of the compiler must be responsible for entering each identifier into the symbol table.
- a. Should the scanner or the parser enter identifiers into the symbol table? Each has an opportunity to do so.

Answer: Either the scanner or the parser can create the initial symbol-table entry for an identifier. The scanner can create the entry but cannot enter any information about its type. Depending on the syntax of the language, the parser may be able to enter type information. (For example, in C's declaration syntax, when the parser encounters an identifier in a declaration, it has already seen the syntax that assigns it a base type.)

- b. Is there an interaction between this issue, declare-before-use rules, and disambiguation of subscripts from function calls in a language with the FORTRAN 77 ambiguity?

Answer: If the scanner has access to the symbol table, it can use type information from the symbol table to disambiguate certain kinds of ambiguity. For example, the scanner can return a different syntactic category for a function name than for an array name in FORTRAN, removing the ambiguity in `fee(a, b)`.

11. The compiler must store information in the IR version of the program that allows it to get back to the symbol table entry for each name. Among the options open to the compiler writer are pointers to the original character strings and subscripts into the symbol table. Of course, the clever implementor may discover other options.

What are the advantages and disadvantages of each of these representations for a name? How would you represent the name?

Answer: Pointers to the character string allow the various parts of the compiler to directly access the strings without further address calculations. This produces a lower overhead on access to names. Unless the compiler can guarantee that it only keeps one copy of each string, it cannot use the pointer for an equality test, so equality tests are expensive. Symbol table lookups would require a hash-computation on each lookup.

Symbol table indices provide direct access to the information in the symbol table, without the need for hash-computation on lookup. They also provide for an inexpensive equality test, which pointers may not provide. Access to the actual name of the symbol (its lexeme) requires more work under this scheme than it would in an implementation that stored the pointer.

To make this choice, the compiler should consider which form of access occurs more often: a symbol table access or an access to the lexeme. If space is not an overriding concern, the compiler writer might consider storing both the pointer and the index.

12. You are writing a compiler for a simple lexically-scoped language. Consider the example program shown in Figure 5.16.

- a. Draw the symbol table and its contents at line 11.

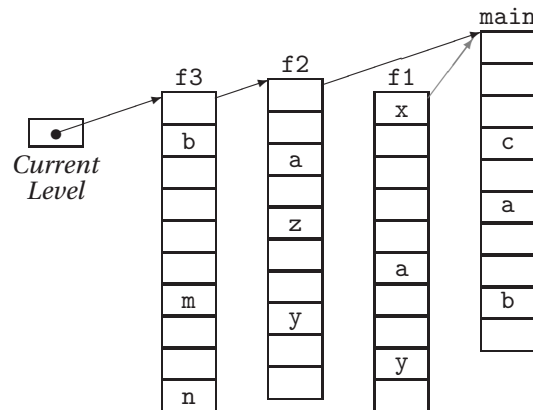
Answer: Many answers are possible for this question, depending on the scheme that the student proposes for implementing a lexically scoped table. The simplest scheme is the “sheaf of tables” implementation, which might produce the following state. (f 1 has been discarded.)

```

1  procedure main
2    integer a, b, c;
3    procedure f1(w,x);
4      integer a,x,y;
5      call f2(w,x);
6    end;
7    procedure f2(y,z)
8      integer a,y,z;
9      procedure f3(m,n)
10       integer b, m, n;
11       c = a * b * m * n;
12     end;
13     call f3(c,z);
14   end;
15   ...
16   call f1(a,b);
17 end;

```

FIGURE 5.16 Program for Exercise 12



- b. What actions are required for symbol table management when the parser enters a new procedure and when it exits a procedure?

Answer: In the "sheaf of tables" implementation, entering a new procedure (or scope) entails creating a new table and linking it into the chain of tables. (The compiler sets the new table's chain to the contents of *Current Level* and then sets *Current Level* to point at the new table.) On exit from a procedure (or scope), it removes the first table from the

chain of tables (setting *Current Level* to the chain leading from the *Current Level* table to the next scope). The de-linked table may be freed; alternately, the compiler may keep it around for use in the debugger or later stages of compilation.

In other table organizations, the detailed actions will be different. However, the basic actions are: on entry to a scope, mark the “pre-scope” state of the table, and on exit from a scope, restore it to the appropriate “pre-scope” state.

13. The most common implementation technique for a symbol table uses a hash table, where insertion and deletion are expected to have $O(1)$ cost.

- a. What is the worst-case cost for insertion and for deletion in a hash table?

Answer: The worst case cost for insertion depends on the implementation. In an open-hashing table, the worst case cost might be $O(1)$, if the insertion code does not perform a lookup first. If it does perform a lookup, the worst case for insertion becomes $O(n)$, for a table that contains n names.

In an open-addressing table, the worst-case cost will be $O(n)$ for a table that contains n names.

If deletion requires a lookup to find the record for the name being deleted, then the worst-case cost is $O(n)$ for a table that contains n names.

- b. Suggest an alternative implementation scheme that guarantees $O(1)$ insertion and deletion.

Answer: See, for example, the sidebar labelled *An Alternative to Hashing* on page 256 of Chapter 5, or the discussion in § 2.6.3 of closure-free regular expressions that starts on page 77.

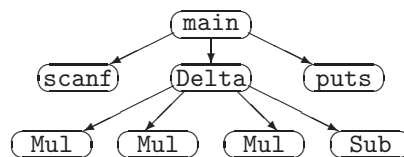
Section 6.2

1. Show the call tree and execution history for the following C program:

```
int Sub(int i, int j) {
    return i - j;
}
int Mul(int i, int j) {
    return i * j;
}
int Delta(int a, int b, int c) {
    return Sub(Mul(b,b), Mul(Mul(4,a),c));
}
void main() {
    int a, b, c, delta;
    scanf("%d %d %d", &a, &b, &c);
    delta = Delta(a, b, c);
    if (delta == 0)
        puts("Two equal roots");
    else if (delta > 0)
        puts("Two different roots");
    else
        puts("No root");
}
```

Answer:

Call Tree:



Note that main calls puts at most once, even though it has three distinct call sites that can invoke puts.

Execution History:

1. main calls scanf
2. scanf returns to main
3. main calls Delta
4. Delta calls Mul
5. Mul returns to Delta
6. Delta calls Mul
7. Mul returns to Delta
8. Delta calls Mul
9. Mul returns to Delta
10. Delta calls Sub
11. Sub returns to Delta
12. Delta returns to main
13. main calls puts (from 1 of 3 sites)
14. puts returns to main

2. Show the call tree and execution history for the following C program:

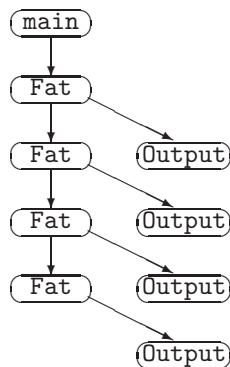
```
void Output(int n, int x) {
    printf("The value of %d! is %s.\n", n, x);
}

int Fat(int n) {
    int x;
    if (n > 1)
        x = n * Fat(n - 1);
    else
        x = 1;
    Output(n, x);
    return x;
}

void main() {
    Fat(4);
}
```

Answer:

Call Tree



Execution History:

1. main calls Fat (4)
2. Fat (4) calls Fat (3)
3. Fat (3) calls Fat (2)
4. Fat (2) calls Fat (1)
5. Fat (1) calls Output
6. Output returns to Fat (1)
7. Fat (1) returns to Fat (2)
8. Fat (2) calls Output
9. Output returns to Fat (2)
10. Fat (2) returns to Fat (3)
11. Fat (3) calls Output
12. Output returns to Fat (3)
13. Fat (3) returns to Fat (4)
14. Fat (4) calls Output
15. Output returns to Fat (4)
16. Fat (4) returns to main

Section 6.3

3. Consider the following Pascal program, in which only procedure calls and variable declarations are shown:

```
1  program Main(input, output);
2    var a, b, c : integer;
3    procedure P4; forward;
4    procedure P1;
5      procedure P2;
6        begin
7          end;
8      var b, d, f : integer;
9      procedure P3;
10       var a, b : integer;
11       begin
12         P2;
13       end;
14     begin
15       P2;
16       P4;
17       P3;
18     end;
19   var d, e : integer;
20   procedure P4;
21     var a, c, g : integer;
22     procedure P5;
23       var c, d : integer;
24       begin
25         P1;
26       end;
27   var d : integer;
28   begin
29     P1;
30     P5;
31   end;
32 begin
33   P1;
34   P4;
35 end.
```

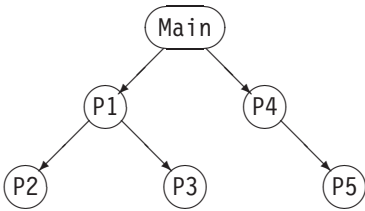
a. Construct a static coordinate table, similar to the one in Figure 6.3.

Answer:

Scope	a	b	c	d	e	f	g
Main	$\langle 1,0 \rangle$	$\langle 1,4 \rangle$	$\langle 1,8 \rangle$	$\langle 1,12 \rangle$	$\langle 1,16 \rangle$	—	—
P1	$\langle 1,0 \rangle$	$\langle 2,0 \rangle$	$\langle 1,8 \rangle$	$\langle 2,4 \rangle$	—	$\langle 2,8 \rangle$	—
P2	$\langle 1,0 \rangle$	$\langle 1,4 \rangle$	$\langle 1,8 \rangle$	—	—	—	—
P3	$\langle 3,0 \rangle$	$\langle 3,4 \rangle$	$\langle 1,8 \rangle$	$\langle 2,4 \rangle$	—	$\langle 2,8 \rangle$	—
P4	$\langle 2,0 \rangle$	$\langle 1,4 \rangle$	$\langle 2,4 \rangle$	$\langle 1,12 \rangle$	$\langle 1,16 \rangle$	—	$\langle 2,8 \rangle$
P5	$\langle 2,0 \rangle$	$\langle 1,4 \rangle$	$\langle 3,0 \rangle$	$\langle 3,4 \rangle$	$\langle 1,16 \rangle$	—	$\langle 2,8 \rangle$

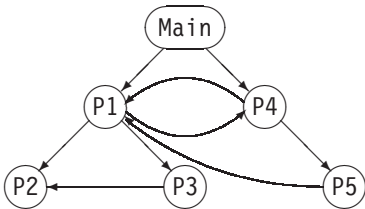
b. Construct a graph to show the nesting relationships in the program.

Answer:



c. Construct a graph to show the calling relationships in the program.

Answer:



4. Some programming languages allow the programmer to use functions in the initialization of local variables but not in the initialization of global variables.

- a. Is there an implementation rationale to explain this seeming quirk of the language definition?

Answer: Most implementations use static initialization (with an assembly-language directive) to initialize global variables. Since a function must execute to return its value, this mechanism excludes any initial value that is not known at compile time.

- b. What mechanisms would be needed to allow initialization of a global variable with the result of a function call?

Answer: If an implementation allowed a function call as the initializer of a global variable, it would need to defer initialization to runtime, when it could execute the function. It would also need to restrict the domain of the functions used in the initialization expression.

For example, the following program is not well formed because the implementation cannot resolve the initial value of `fee`.

```
int fee = foe();

int foe() {
    return fee + 5;
}
```

5. The compiler writer can optimize the allocation of ARs in several ways. For example, the compiler might:

- Allocate ARs for leaf procedures statically.
- Combine the ARs for procedures that are always called together. (When α is called, it always calls β .)
- Use an arena-style allocator in place of heap allocation of ARs.

For each scheme, consider the following questions:

- a. What fraction of the calls might benefit? In the best case? In the worst case?

Answer: Estimating these fractions is tenuous, but that is part of the point of the question. People make sweeping statements about the benefits of various call-site optimizations without much data as to how often the phenomenon in question occurs.

Leaf procedures: Estimating the fraction of calls that invoke leaf procedures is a tenuous business, at best. It can range from none (a program

can have no leaf procedures—think of recursive fibonacci) to all of the calls.

Always called together: If *p* always calls *q*, then every call to *p* will benefit from combining *p* and *q*'s ARs. However, the fraction of all calls that benefit will depend on how often *p* is called.

Arena allocation: This option benefits all calls to procedures implemented with heap-allocated ARs. If all procedures have heap-allocated ARs, then all calls benefit.

- b. What is the impact on runtime space utilization?

Answer: Leaf procedures: This scheme requires one activation record, sized large enough to hold the largest AR of a leaf procedure. (Only one leaf procedure can be active at a given point in execution.) That storage is permanently allocated, so it forms a constant overhead on the program's space requirement.

The extra space can be bounded by the size of the statically allocated AR.

Always called together: Assume the *p* and *q* have their ARs merged and *p* always calls *q*. This scheme requires no extra space, since the original program would allocate ARs for both *p* and *q*.

Arena allocation: This scheme might require significant amounts of extra space. The arena allocators have cheap allocation and bulk deallocation. The latter may be a poor match for AR allocation, since the dynamic behavior of the program will determine when ARs may be freed. With a contrary program, deallocation is never possible, so all ARs remain allocated from the point of their creation until the end of execution.

6. Draw the structures that the compiler would need to create to support an object of type Dumbo, defined as follows:

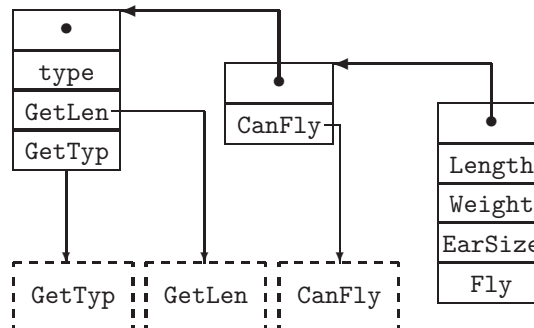
```
class Elephant {
    private int Length;
    private int Weight;
    static int type;

    public int GetLen();
    public int GetTyp();
}

class Dumbo extends Elephant {
    private int EarSize;
    private boolean Fly;

    public boolean CanFly();
}
```

Answer:



Dashed boxes represent executable code

7. In a programming language with an open class structure, the number of method invocations that need runtime name resolution, or dynamic dispatch, can be large. A method cache, as described in Section 6.3.4, can reduce the runtime cost of these lookups by short-circuiting them.

As an alternative to a global method cache, the implementation might maintain a single entry method cache at each call site—an inline method cache that records the address of the method most recently dispatched from that site, along with its class.

Develop pseudocode to use and maintain such an inline method cache. Explain the initialization of the inline method caches and any modifications to the general method lookup routine required to support inline method caches.

Answer:

The inline method cache consists of a pair $\langle \text{class tag}, \text{code} \rangle$. It is initialized to a known illegal value in the cache tag. Each method invocation is replaced with the following pseudo-code:

<pre> code ← lookup(obj.method) invoke code(obj, arguments) </pre>	<pre> if obj.class_tag = inline_cache.class_tag then inline_cache.code(obj, arguments) else { code ← lookup(obj.method) inline_cache.class_tag ← obj.class_tag inline_cache.code ← code inline_cache.code(obj, arguments) } </pre>
---	--

Original Code

With Inline Cache

As long as the general lookup routine returns a reference to the executable method in some compact form, such as a code pointer, this scheme should work. It would be good if the general lookup routine handled errors, such as an attempt to lookup a method that does not map to any known implementation. If, instead, that error check must be performed on the result returned by the lookup routine, it adds code to every method invocation, causing the code to be larger. (In either scheme, the code executes. If the code is replicated at each call, then the executable will have myriad copies of that check instead of one.)

```

1  procedure main;
2    var a : array[1...3] of int;
3    i : int;
4    procedure p2(e : int);
5      begin
6        e := e + 3;
7        a[i] := 5;
8        i := 2;
9        e := e + 4;
10     end;
11  begin
12    a := [1, 10, 77];
13    i := 1;
14    p2(a[i]);
15    for i := 1 to 3 do
16      print(a[i]);
17    end.

```

FIGURE 6.13 Program for Problem 8

Section 6.4

8. Consider the program written in Pascal-like pseudo code shown in Figure 6.13. Simulate its execution under call-by-value, call-by-reference, call-by-name, and call-by-value-result parameter binding rules. Show the results of the print statements in each case.

Answer:

Call-by-value should print 5, 10, and 77.

Call-by-reference should print 9, 10, and 77.

Call-by name should print 5, 14, and 77.

Call-by-value-result should print 8, 10, and 77.

9. The possibility that two distinct variables refer to the same object (memory area) is considered undesirable in programming languages. Consider the following Pascal procedure, with parameters passed by reference:

```

procedure mystery(var x, y : integer);
begin
  x := x + y;
  y := x - y;
  x := x - y;
end;

```

If no overflow or underflow occurs during the arithmetic operations:

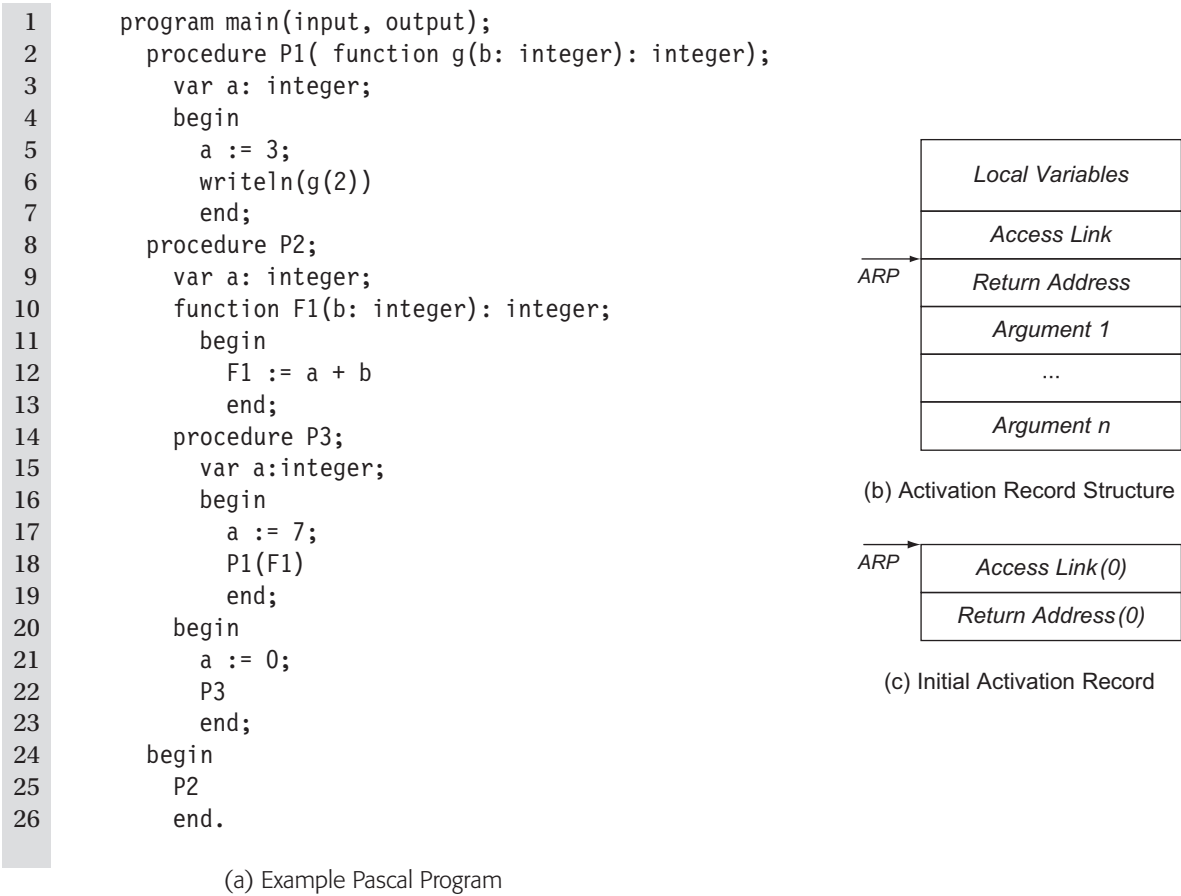
- a. What result does `mystery` produce when it is called with two distinct variables, `a` and `b`?

Answer: As long as no overflow or underflow occurs, `mystery` exchanges the values of `a` and `b`.

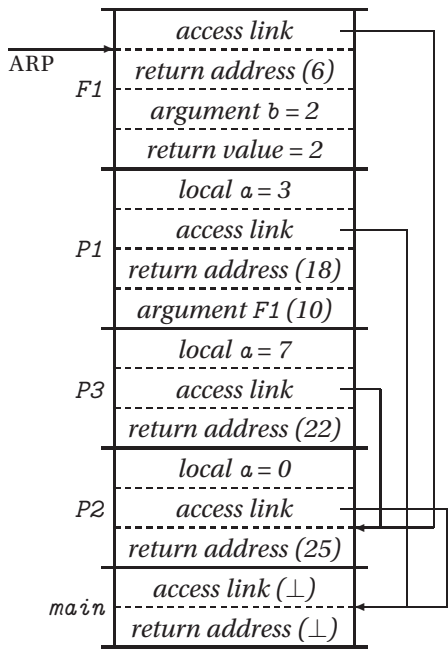
- b. What would be the expected result if `mystery` is invoked with a single variable `a` passed to both parameters? What is the actual result in this case?

Answer: Someone unfamiliar with the implementation of call-by-reference parameters might reasonably expect that a “swap” routine would leave its argument’s value unchanged when it is passed the same name as both actual parameters. For example, the standard three-variable sequence (`t ← x`, `x ← y`, and `y ← t`) has that behavior.

However, `mystery` sets its actual parameter to zero when it is passed the same argument in both parameter positions, assuming that no overflow or underflow occurs.



Answer:



11. Assume that the compiler is capable of analyzing the code to determine facts such as “*from this point on, variable v is not used again in this procedure*” or “*variable v has its next use in line 11 of this procedure,*” and that the compiler keeps all local variables in registers for the following three procedures:

```

procedure main
  integer a, b, c
  b = a + c;
  c = f1(a,b);
  call print(c);
end;
procedure f1(integer x, y)
  integer v;
  v = x * y;
  call print(v);
  call f2(v);
  return -x;
end;
procedure f2(integer q)
  integer k, r;
  ...
  k = q / r;
end;

```

- a. Variable x in procedure `f1` is live across two procedure calls. For the fastest execution of the compiled code, should the compiler keep it in a caller-saves or callee-saves register? Justify your answer.

Answer: Keeping x in a callee-saves register would let the two callees, `print` and `f2`, determine whether or not to save its value. Were the compiler to assign x to a caller-saves register in `f1`, the linkage sequence would save its value and restore it at each call site, even if the callee did not use the full set of registers. Thus, callee-saves is (probably) the best choice for x .

Of course, if the compiler knew that both `print` and `f2` saved their full complement of callee-saves registers, it might choose to keep x in a caller-saves register, save it before the call to `print` and restore it after the call to `f2`. This would produce fewer loads and stores, but requires knowledge about the behavior of both `print` and `f2` during compilation of `f1`.

- b. Consider variables a and c in procedure `main`. Should the compiler keep them in caller-saves or callee-saves registers, again assuming that the

compiler is trying to maximize the speed of the compiled code? Justify your answer.

Answer: Neither *a* nor *c* has a value that must be preserved across the call to *f1*. The value of *a* has no use after *main* evaluates it as an actual parameter.

The situation with *c* is slightly different. The use before the call to *f1* is disjoint from the definition and use after the call. However, no value passes from the first part to the second.

Since neither value must be preserved across the call, the compiler should assign them both to caller-saves registers. In this way, it can avoid the stores to save them and the loads to restore them.

12. Consider the following Pascal program. Assume that the ARs follow the same layout as in problem 10, with the same initial condition, *except* that the implementation uses a global display rather than access links.

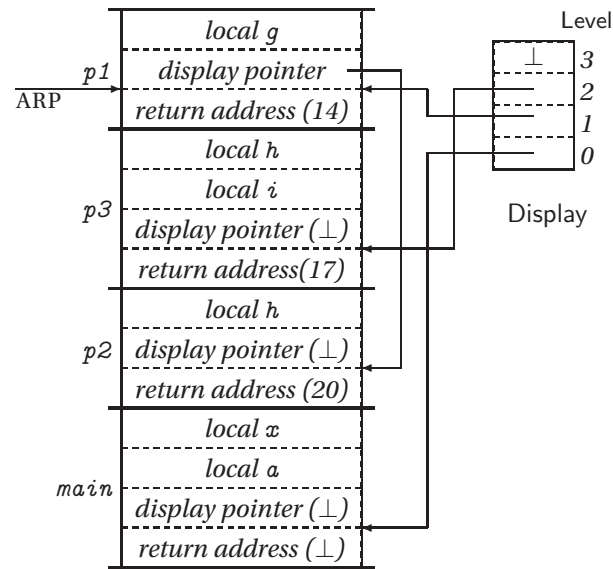
```

1  program main(input, output);
2      var x : integer;
3          a : float;
4      procedure p1();
5          var g:character;
6          begin
7              ...
8          end;
9      procedure p2();
10         var h:character;
11         procedure p3();
12             var h,i:integer;
13             begin
14                 p1();
15             end;
16         begin
17             p3();
18         end;
19     begin
20         p2();
21     end

```

Draw the set of ARs that are on the runtime stack when the program reaches line 7 in procedure *p1*.

Answer:



Section 7.2

1. Memory layout affects the addresses assigned to variables. Assume that character variables have no alignment restriction, short integer variables must be aligned to halfword (2 byte) boundaries, integer variables must be aligned to word (4 byte) boundaries, and long integer variables must be aligned to doubleword (8 byte) boundaries. Consider the following set of declarations:

```
char a;
long int b;
int c;
short int d;
long int e;
char f;
```

Draw a memory map for these variables:

- (a) Assuming that the compiler cannot reorder the variables

Answer:

a	1 byte
<i>empty</i>	7 bytes
b	8 bytes
c	4 bytes
d	2 bytes
<i>empty</i>	2 bytes
e	8 bytes
f	1 byte

- (b) Assuming the compiler can reorder the variables to save space

Answer: Many answers are possible. The general algorithm would produce the layout shown on the left; a greedy rearrangement might produce the layout shown on the right.

b	8 bytes
e	8 bytes
c	4 bytes
d	2 bytes
a	1 byte
f	1 byte

a	1 byte
f	1 byte
d	2 bytes
c	4 bytes
b	8 bytes
e	8 bytes

2. As demonstrated in the previous question, the compiler needs an algorithm to lay out memory locations within a data area. Assume that the algorithm receives as input a list of variables, their lengths, and their alignment restrictions, such as

$\langle a, 4, 4 \rangle, \langle b, 1, 3 \rangle, \langle c, 8, 8 \rangle, \langle d, 4, 4 \rangle, \langle e, 1, 4 \rangle, \langle f, 8, 16 \rangle, \langle g, 1, 1 \rangle.$

The algorithm should produce, as output, a list of variables and their offsets in the data area. The goal of the algorithm is to minimize unused, or wasted, space.

- (a) Write down an algorithm to lay out a data area with minimal wasted space.

Answer: The following pseudo-code assumes the presence of a global integer variable, *Offset*, and that it can destroy a copy of the list, stored in *List*. (It could, of course, mark list elements as processed rather than removing them from the list.)

```
Assign() {
    Offset ← 0
    NextSize ← FindBoundary( List )    // find largest size
    while (NextSize > 0)
        LayoutClass( List, NextSize )  // assign offsets
        NextSize ← FindBoundary( List ) // find next size
    }

FindBoundary( List ) {
    boundary ← 0
    for each element, e, in List
        if alignment(e) > boundary
            then boundary ← alignment(e)
    return boundary
}

LayoutClass( List, Size ) {
    for each element, e, in List
        if alignment(e) = Size then
            assign name(e) the storage at Offset
            print name(e), Offset
            Offset ← Offset + Size
            remove element from List
    }
```

- (b) Apply your algorithm to the example list above and two other lists that you design to demonstrate the problems that can arise in storage layout.

Answer: For the list in the problem, the algorithm produces the following output

f,0, c,16, a,24, d,28, e,32, b,36, g,39

- (c) What is the complexity of your algorithm?

Answer: The algorithm takes time proportional to the length of the list and the number of distinct alignment-boundaries present in the list. It makes two passes over the list for each distinct boundary size. In the worst case, each list element has a distinct boundary size and the algorithm will require $O(N^2)$ where N is the number of list elements. In practice, the number of boundary sizes will be limited to five or six.

3. For each of the following types of variable, state where in memory the compiler might allocate the space for such a variable. Possible answers include registers, activation records, static data areas (with different visibilities), and the runtime heap.

- a. A variable local to a procedure

Answer: The compiler can allocate it to a register or in the activation record.

- b. A global variable

Answer: The compiler can allocate it in a static global data area.

- c. A dynamically allocated global variable

Answer: The compiler can allocate it on the run-time heap.

- d. A formal parameter

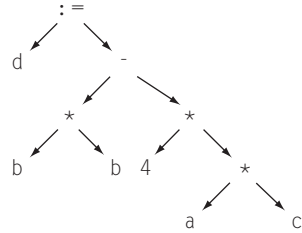
Answer: The compiler can allocate it to a register or in the activation record.

- e. A compiler-generated temporary variable

Answer: The compiler can allocate it to a register or in the activation record.

Section 7.3

4. Use the treewalk code-generation algorithm from Section 7.3 to generate naive code for the following expression tree. Assume an unlimited set of registers.



Answer: There are many possible answers. The one on the left follows the style of the book, while the one on the right is more concise.

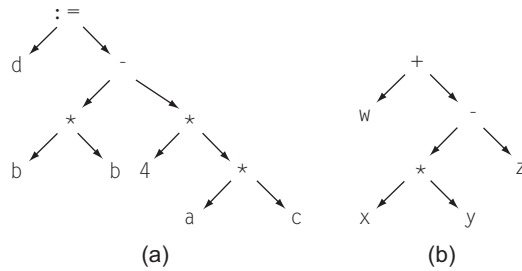
```

loadI  @b      ⇒ r1
loadAO  rarp,r1 ⇒ r2
mult    r2,r2 ⇒ r3
loadI   4      ⇒ r4
loadI  @a      ⇒ r5
loadAO  rarp,r5 ⇒ r6
loadI  @c      ⇒ r7
loadAO  rarp,r7 ⇒ r8
mult    r6,r8 ⇒ r9
mult    r4,r9 ⇒ r10
sub     r3,r10 ⇒ r11
loadI  @d      ⇒ r12
store  r11    ⇒ r12
  
```

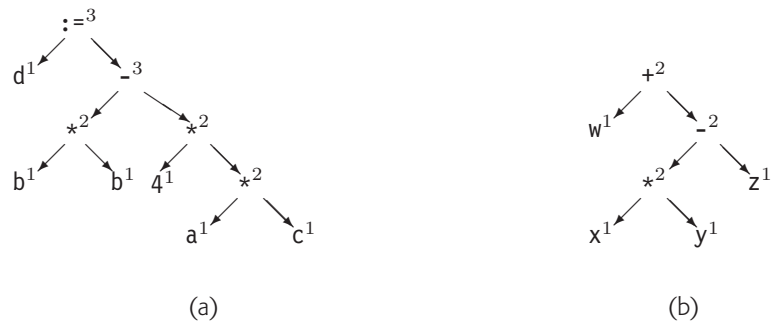
```

loadAI  rarp,@b ⇒ r1
mult    r1,r1 ⇒ r2
loadI   4      ⇒ r3
loadAI  rarp,@a ⇒ r4
loadAI  rarp,@c ⇒ r5
mult    r4,r5 ⇒ r6
mult    r3,r6 ⇒ r7
sub     r2,r7 ⇒ r8
loadI  @d      ⇒ r9
store  r8     ⇒ r9
  
```

5. Find the minimum number of registers required to evaluate the following trees using the ILOC instruction set. For each nonleaf node, indicate which of its children must be evaluated first in order to achieve this minimum number of registers.



Answer: Nodes are superscripted with the number of registers required to evaluate the subtree rooted at that node.



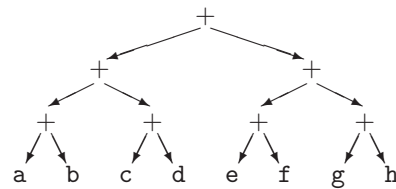
In each case, the code should evaluate the more demanding subtree—the one requiring more registers—first. In tree (a), notice that an immediate multiply would eliminate the register needed to hold 4, but would not decrease the overall demand for registers.

6. Build expression trees for the following two arithmetic expressions, using standard precedence and left-to-right evaluation. Compute the minimum number of registers required to evaluate each of them using the ILOC instruction set.

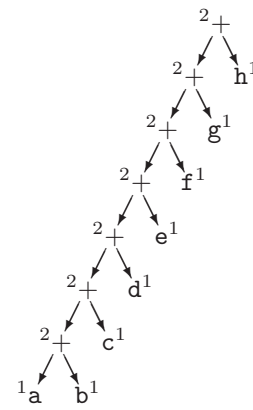
(a) $((a + b) + (c + d)) + ((e + f) + (g + h))$

(b) $a + b + c + d + e + f + g + h$

Answer:



(a)



(b)

Section 7.4

7. Generate predicated ILOC for the following code sequence. (No branches should appear in the solution.)

```
if (x < y)
  then z = x * 5;
  else z = y * 5;
w = z + 10;
```

Answer:

```
cmp_LT rx, ry ⇒ r1
not    r1      ⇒ r2
(r1)? loadI 5   ⇒ r3
(r1)? mult  rx, r3 ⇒ rz
(r2)? loadI 5   ⇒ r3
(r1)? mult  ry, r3 ⇒ rz
loadI 10   ⇒ r4
add      rz, r4 ⇒ rw
```

Note that the immediate load into r_3 is identical on both paths. The predicate could be dropped and one fewer operation generated. One way to catch this inefficiency is *code hoisting* (see § 10.3.2).

8. As mentioned in Section 7.4, short-circuit code for the following expression in C avoids a potential division-by-zero error:

```
a != 0 && b / a > 0.5
```

If the source-language definition does not specify short-circuit evaluation for boolean-valued expressions, can the compiler generate short-circuit code as an optimization for such expressions? What problems might arise?

Answer: If short-circuit code is generated as an optimization and the source-language definition requires that the compiler preserve exceptions, the optimization cannot be applied to any expression that might generate an exception. In the expression shown above, the subexpression b / a can generate an exception, so the short-circuit code will (illegally) alter the behavior when a is zero.

Section 7.5

9. For a character array $A[10 \dots 12, 1 \dots 3]$ stored in row-major order, calculate the address of the reference $A[i, j]$, using at most four arithmetic operations in the generated code.

Answer: Using the “false zero” scheme, $@A_0$ becomes $@A - 31$, so the address of $A[i, j]$ is $@A_0 + i \times 3 + j$.

10. What is a dope vector? Give the contents of the dope vector for the character array in the previous question. Why does the compiler need a dope vector?

Answer: A dope vector is a structure that the compiler builds to pass information about an array’s size at a call site. The dope vector typically contains the array’s base address and the sizes of each dimension. It may, depending on language issues, need to contain the number of dimensions, too.

For the array $A[10 \dots 12, 1 \dots 3]$ in the previous exercise, the dope vector would contain $\langle @A_0, 3, 3 \rangle$. If the implementation needed the number of dimensions, then the dope vector would be $\langle @A_0, 2, 3, 3 \rangle$.

Dope vectors are necessary because a single procedure can be invoked from many different call sites with different array-valued actual parameters. Using a dope-vector lets the callee generate appropriate code without knowing the size of each dimension or, in some cases, the number of such dimensions.

11. When implementing a C compiler, it might be advisable to have the compiler perform range checking for array references. Assuming range checks are used and that all array references in a C program have successfully passed them, is it possible for the program to access storage outside the range of an array, for example, accessing $A[-1]$ for an array declared with lower bound zero and upper bound N ?

Answer: *This question is not a particularly useful question. It will be replaced in EaC3e.*

Range checking is useful for debugging and to detect unexpected or malicious behavior. Good techniques exist to decrease the runtime cost of range checks; see, for example, reference [257] in the bibliography.

In an implementation that includes range checking, a program can use pointer arithmetic to access storage outside the valid range of an array. For example, with an array declared as $A[0:N]$, the code could load the value from the storage that would be $A[-1]$ into i with a code sequence such as:

```
int i, *p;
p = &A[0];
p--;
i = *p;
```

Section 7.6

12. Consider the following character-copying loop from Section 7.6.2:

```

                                loadI @b    ⇒ r@b // get pointers
                                loadI @a    ⇒ r@a
                                loadI NULL ⇒ r1 // terminator
do {
                                L1: cload r@b ⇒ r2 // get next char
                                cstore r2 ⇒ r@a // store it
                                addI r@b, 1 ⇒ r@b // bump pointers
                                addI r@a, 1 ⇒ r@a
                                cmp_NE r1, r2 ⇒ r4
                                cbr r4 → L1, L2
                                L2: nop // next stmt
```

Modify the code so that it branches to an error handler at L_{sov} on any attempt to overrun the allocated length of a . Assume that the allocated length of a is stored as an unsigned four-byte integer at an offset of -8 from the start of a .


```

Answer:  loadI  @b      ⇒ r@b    // get pointers
         loadI  @a      ⇒ r@a
         loadAI r@a,-8   ⇒ rlen(a) // bound on a
         add    r@a,rlen(a) ⇒ r_ub
         loadI  NULL    ⇒ r1     // terminator
L1: cload r@b      ⇒ r2     // get next char
         cstore r2      ⇒ r@a    // store it
         addI   r@b,1    ⇒ r@b    // bump pointers
         addI   r@a,1    ⇒ r@a
         cmp_GT r@a,r_ub ⇒ r4     // check overflow
         cbr    r4       → Lsov,L2
         cmp_NE r1,r2    ⇒ r5
         cbr    r5       → L1,L3
L3: nop                // next stmt

```

This solution introduces a second compare and branch into the loop body. If branch overhead is a concern, and it should be, the following version might be worth generating.

```

         loadI  @b      ⇒ r@b    // get pointers
         loadI  @a      ⇒ r@a
         loadAI r@a,-8   ⇒ rlen(a) // bound on a
         add    r@a,rlen(a) ⇒ r_ub
         loadI  NULL    ⇒ r1     // terminator
L1: cload r@b      ⇒ r2     // get next char
         cstore r2      ⇒ r@a    // store it
         addI   r@b,1    ⇒ r@b    // bump pointers
         addI   r@a,1    ⇒ r@a
         cmp_LE r@a,r_ub ⇒ r4     // check overflow
         cmp_NE r1,r2    ⇒ r5     // check terminator
         and    r4,r5    ⇒ r6     // combine them
         cbr    r6       → L1,L2
L2: cbr    r4       → L3,Lsov // overflow exit?
L3: nop                // next stmt

```

This version has the same number of operations inside the loop, but replaces the `cbr` for overflow with a logical `and`. On exit from the loop, at `L2`, the code tests to determine if the exit occurred because the loop encountered a `NULL` or because the loop tried to overrun the allocated length of `a`. Finally, the first statement after the loop has moved to `L3`.

13. Arbitrary string assignments can generate misaligned cases.

- (a) Write the ILOC code that you would like your compiler to emit for an arbitrary PL/I-style character assignment, such as

```
fee(i:j) = fie(k:l);
```

where $j-i = l-k$. This statement copies the characters in `fie`, starting at location `k` and running through location `l` into the string `fee`, starting at location `i` and running through location `j`.

Include versions using character-oriented memory operations and versions using word-oriented memory operations. You may assume that `fee` and `fie` do not overlap in memory.

Answer: Assume that the values of `i`, `j`, `k`, and `l` are already in local variables with the obvious names, `ri`, `rj`, `rk`, and `rl`. We assume that $(j - i) = (k - l)$, so the code will not check that condition.

```

sub      rj, ri      ⇒ rcount
loadI    0           ⇒ rzero
loadI    0           ⇒ roffset
loadI    @fee        ⇒ r@fee
loadI    @fie        ⇒ r@fie

cmp      rcount, rzero ⇒ r1
cbr_LE   r1          → Lexit

L0: cloadAI rfee, roffset ⇒ rchar
cstoreAI rchar          ⇒ rfie, roffset
addI     roffset, 1     roffset
subI     rcount, 1     rcount
cmp      rcount, rzero ⇒ r1
cbr_GT   r1          → L0

Lexit: nop
```

Character-oriented operations

Improvements on this code are possible. For example, we could shift the loop-counter tests and branches from `rcount` to `roffset` (see, for example, the discussion of *linear function test replacement* in § 10.7.2, starting on page 589).

To implement the same functionality with word oriented operations, we need an initial loop to read characters up to a word boundary, then an loop that loads a word, adjusts its alignments, and writes a word.



The version that uses full-word load and store is only slightly more complex, if we assume that load and store can access arbitrarily-aligned locations in memory. With that assumption, the code simply performs full-word load and store operations until less than a word remains. At that point, it performs one, two, or three character-oriented operations. The truly complex case occurs with full-word load and store operations that only work on word-aligned addresses. We will provide ILOC code for that case in the next release of the solutions — the set that includes all thirteen chapters.

- (b) The programmer can create character strings that overlap. In PL/I, the programmer might write

```
fee(i:j) = fee(i+1:j+1);
```

or, even more diabolically,

```
fee(i+k:j+k) = fee(i:j);
```

How does this complicate the code that the compiler must generate for the character assignment?

Answer: First, the language must define what behavior is expected from assignments involving overlapping storage. We will assume that the language expects overlapped assignment to behave in the naïve way, that is, to read data from a location before writing new data to it.

(a). The easiest way, albeit the least efficient, would be to copy the source to an internally allocated buffer, and then to copy the data from the buffer to its destination location.

(b). The more complex version breaks into two separate cases: the destination string runs ahead of the source string, which is easy; and the destination string runs behind the source string (so that we must defer writing to it).

In the first case, we can use the character-copying loop and let it run. It should work. If the distance between the start of the source and destination is large enough, the word-oriented version will also work.

In the second case, we need to buffer the characters in the overlap region. The simplest approach is to use an internally-allocated, full-size buffer, as in (a).

- (c) Are there optimizations that the compiler could apply to the various character-copying loops that would improve runtime behavior? How would they help?

Answer: If the target machine supports a *load multiple/store multiple* instruction pair, then the loops might do better by clearing out a set of adjacent registers and using the multiword memory operations.

Because the loops contain a fair amount of code, the compiler might reduce code size by generating calls to a library routine, such as the Linux `memcpy` and `memmove` routines. Those library implementations should be optimized for the target machine.

Section 7.7

14. Consider the following type declarations in C:

```

struct S2 {          union U {          struct S1 {
    int i;             float r;           int a;
    int f;             struct S2;         double b;
};                     };                union U;
                                int d;
};

```

Build a structure-element table for S1. Include in it all the information that a compiler would need to generate references to elements of a variable of type S1, including the name, length, offset, and type of each element.

Answer:

<i>Name</i>	<i>Length</i>	<i>Offset</i>	<i>Type</i>
a	4	0	int
b	8	4	double
r	4	12	float
i	4	12	int
f	4	16	int
d	4	20	int

Structure Element Table for S1

This table assumes that the compiler introduces no padding between structure elements. In practice, the compiler might add padding between a and b to ensure doubleword alignment of b (required on some machines). Such padding would, of course, change the offsets of many of the names.

15. Consider the following declarations in C:

```
struct record {
    int StudentId;
    int CourseId;
    int Grade;
} grades[1000];
int g, i;
```

Show the code that a compiler would generate to store the value in variable g as the grade in the i^{th} element of $grades$, assuming the following:

(a) The array $grades$ is stored as an array of structures.

Answer:

```
multi ri,12      ⇒ r1    // record is 12 bytes long
addI  r1,8       ⇒ r2    // offset of Grade is 8
addI  r2,@grades ⇒ r3    // add base address
store rg        ⇒ r3
```

(b) The array $grades$ is stored as a structure of arrays.

Answer:

```
loadI @grades ⇒ r1    // get base address
addI  r1,8000 ⇒ r2    // start of grade array
multi ri,4    ⇒ r3    // offset of ith element
add   r2,r3  ⇒ r4    // add base address
store rg     ⇒ r4
```

Section 7.8

16. As a programmer, you are interested in the efficiency of the code that you produce. You recently implemented, by hand, a scanner. The scanner spends most of its time in a single `while` loop that contains a large case statement.
- (a) How would the different case statement implementation techniques affect the efficiency of your scanner?

Answer:

Linear Search: The order in which the case labels appear may greatly affect the efficiency of the scanner. If there are a large number of cases that appear with roughly equal frequency, linear search may produce worse results than the other implementation techniques.

Binary Search: This technique provides roughly equal access times for all the labels. If there are a large number of equi-probable labels, and they form a sparse set, then this technique is a good choice.

Jump Table: In this case, the compiler arranges for the code to directly compute an offset into a table of labels, sometimes called a *jump table*. If the set of labels forms a dense set, this technique can eliminate many of the table lookups required by the binary search method (without wasting space). The case lookup takes constant time, but requires a memory reference and a jump to register. For extremely small label sets, linear search may be faster. (In general, this technique is faster than the other methods, but it requires a dense set of labels.)

- (b) How would you change your source code to improve the runtime performance under each of the case statement implementation strategies?

Answer:

Linear Search: Since the compiler tests the labels in order of appearance, the programmer should rearrange them so that commonly occurring labels appear first, with rarely used labels appearing at the end.

Binary Search: Since the compiler must reorder the labels into lexicographic order, the order in which they appear in the code is irrelevant to the actual performance.

Jump Table: This technique has a constant cost per execution, so it is insensitive to the order in which the labels appear in the source code.

17. Convert the following C tail-recursive function to a loop:

```
List * last(List *l) {  
    if (l == NULL)  
        return NULL;  
    else if (l->next == NULL)  
        return l;  
    else  
        return last(l->next); }
```

Answer:

```
List * last(List *l) {  
    while(1) {  
        if (l == NULL)  
            break;  
        else if (l->next == NULL)  
            break;  
        else l = l->next;  
    }  
    return l;  
}
```


Section 7.9

18. Assume that x is an unambiguous, local, integer variable and that x is passed as a call-by-reference actual parameter in the procedure where it is declared. Because it is local and unambiguous, the compiler might try to keep it in a register throughout its lifetime. Because it is passed as a call-by-reference parameter, it must have a memory address at the point of the call.

- (a) Where should the compiler store x ?

Answer: Ideally, the compiler should place x in a caller-saves register, as long as there are enough such registers to accommodate x . In a caller-saves register, it will be spilled at the call.

- (b) How should the compiler handle x at the call site?

Answer: If x is in a caller-saves register, the compiler should emit the code to save x into the register save area of the current AR. When it evaluates x as an actual parameter, it can use the address in the register save area as x 's memory address. After the call, it should restore x from the register save area. (There is no need to create a second copy of x for the call site.)

- (c) How would your answers change if x was passed as a call-by-value parameter?

Answer: With a call-by-value parameter, the calling sequence needs a value for x rather than an address, so the sequence need not (necessarily) store x . Thus, the bias in favor of a caller-saves register is removed. The compiler should assign x to a register based on whatever other criteria it uses to select caller-saves versus callee-saves registers.

To evaluate x as an argument, it can copy the value from x into the appropriate register or AR location. There is no need to restore the value after the call.

19. The linkage convention is a contract between the compiler and any outside callers of the compiled code. It creates a known interface that can be used to invoke a procedure and obtain any results that it returns (while protecting the caller's runtime environment). Thus, the compiler should only violate the linkage convention when such a violation cannot be detected from outside the compiled code.

- (a) Under what circumstances can the compiler be certain that using a variant linkage is safe? Give examples from real programming languages.

Answer: If the procedure cannot be invoked from outside the current unit of compilation, then the compiler can safely use a non-standard linkage. Examples include `static` procedures in C and nested local procedures in Pascal. Since these procedures cannot be named outside the compilation unit, they cannot be invoked from outside.

- (b) In these circumstances, what might the compiler change about the calling sequence and the linkage convention?

Answer: The compiler can change almost anything in such a case. For example, it might pass all of the parameters to a procedure in registers. It might collect information about the number of registers used in each procedure in the compilation unit and use that knowledge to avoid unnecessary register save & restore code at some calls. It might merge the ARs of some of the procedures, saving a few operations involved in allocation and deallocation of the ARs.

Section 8.4

1. Apply the algorithm from Figure 8.4 to each of the following blocks:

$$\begin{aligned} t_1 &\leftarrow a + b \\ t_2 &\leftarrow t_1 + c \\ t_3 &\leftarrow t_2 + d \\ t_4 &\leftarrow b + a \\ t_5 &\leftarrow t_3 + e \\ t_6 &\leftarrow t_4 + f \\ t_7 &\leftarrow a + b \\ t_8 &\leftarrow t_4 - t_7 \\ t_9 &\leftarrow t_8 * t_6 \end{aligned}$$

Block b_0

$$\begin{aligned} t_1 &\leftarrow a \times b \\ t_2 &\leftarrow t_1 \times 2 \\ t_3 &\leftarrow t_2 \times c \\ t_4 &\leftarrow 7 + t_3 \\ t_5 &\leftarrow t_4 + d \\ t_6 &\leftarrow t_5 + 3 \\ t_7 &\leftarrow t_4 + e \\ t_8 &\leftarrow t_6 + f \\ t_9 &\leftarrow t_1 + 6 \end{aligned}$$

Block b_1

Answer: Value numbers are shown as superscripts. An asterisk before an operation indicates that it should be replaced with the operation listed to its right.

$$\begin{aligned} t_1^3 &\leftarrow a^1 + b^2 \\ t_2^5 &\leftarrow t_1^3 + c^4 \\ t_3^7 &\leftarrow t_2^5 + d^6 \\ * t_4^3 &\leftarrow b^2 + a^1 & t_4^3 &\leftarrow t_1^3 \\ t_5^9 &\leftarrow t_3^7 + e^8 \\ t_6^{11} &\leftarrow t_4^3 + f^{10} \\ * t_7^3 &\leftarrow a^1 + b^2 & t_7^3 &\leftarrow t_1^3 \\ * t_8^{12} &\leftarrow t_4^3 - t_7^3 & t_8^{12} &\leftarrow 0^{12} \\ * t_9^{13} &\leftarrow t_8^{12} * t_6^{11} & t_9^{13} &\leftarrow 0^{12} \end{aligned}$$

Block b_0

$$\begin{aligned}
t_1^3 &\leftarrow a^1 \times b^2 \\
t_2^5 &\leftarrow t_1^3 \times 2^4 \\
t_3^7 &\leftarrow t_2^5 \times c^6 \\
t_4^9 &\leftarrow 7^8 + t_3^7 \\
t_5^{11} &\leftarrow t_4^9 + d^{10} \\
t_6^{13} &\leftarrow t_5^{11} + 3^{12} \\
t_7^{15} &\leftarrow t_4^9 + e^{14} \\
t_8^{17} &\leftarrow t_6^{13} + f^{16} \\
t_9^{19} &\leftarrow t_1^3 + 6^{18}
\end{aligned}$$

Block b_1

Local value numbering finds no improvements in block b_1 .

2. Consider a basic block, such as b_0 or b_1 in question 1 above. It has n operations, numbered 0 to $n - 1$.
 - a. For a name x , $\text{USES}(x)$ contains the index in b of each operation that uses x as an operand. Write an algorithm to compute the USES set for every name mentioned in block b . If $x \in \text{LIVEOUT}(b)$, then add two dummy entries ($> n$) to $\text{USES}(x)$.

Answer: Assumptions: (1) We have a hash table, as in local value numbering. (2) The *lookup()* routine creates an entry for its argument if no entry exists. (3) *Uses* is a field in the hash table. (4) *LeftOp(i)* returns the first operand in operation i . *RightOp(i)* returns the second operand. *Target(i)* returns the name defined by operation i .

```

initialize the hash table
for  $i \leftarrow 0$  to  $n-1$  do
     $\text{index} \leftarrow \text{lookup}(\text{LeftOp}(i))$ 
    add  $i$  to  $\text{Uses}[\text{index}]$ 

     $\text{index} \leftarrow \text{lookup}(\text{RightOp}(i))$ 
    add  $i$  to  $\text{Uses}[\text{index}]$ 

for each  $i$  in the block's  $\text{LIVEOUT}$  set do
     $\text{index} \leftarrow \text{lookup}(i)$ 
    add  $n+1$  to  $\text{Uses}[\text{index}]$ 
    add  $n+2$  to  $\text{Uses}[\text{index}]$ 

```

b. Apply your algorithm to blocks b_0 and b_1 above.

Answer: To make the drawing understandable, assume that *lookup* allocates space in the table in a simple way. The first entry is at index 0, the second at 1, and so on. Since the question supplies no LIVEOUT sets, we will assume that LIVEOUT is empty.

The algorithm produces the following table for block b_0 :

Index	Name	Uses	Index	Name	Uses
0	a	1, 4, 7	8	e	5
1	b	1, 4, 7	9	t ₅	
2	t ₁	2,	10	f	6
3	c	2,	11	t ₆	9
4	d	3	12	t ₇	8
5	t ₂	3	13	t ₈	9
6	t ₃	5	14	t ₉	
7	t ₄	6, 8			

Applied to block b_1 , the algorithm produces the following table:

Index	Name	Uses	Index	Name	Uses
0	a	1	8	t ₅	6
1	b	2	9	t ₆	8
2	t ₁	2, 9	10	e	7
3	t ₂	3	11	t ₇	
4	c	3	12	f	8
5	t ₃	4	13	t ₈	
6	t ₄	5, 7	14	t ₉	
7	d	5			

Note that the algorithm is defined for *names*. For that reason, we did not create entries for the literal constants.

- c. For a reference to x in operation i of block b , $\text{DEF}(x,i)$ is the index in b where the value of x visible at operation i was defined. Write an algorithm to compute $\text{DEF}(x,i)$ for each reference x in b . If x is upward exposed at i , then $\text{DEF}(x,i)$ should be -1 .

Answer: We will assume a hash table, as in part (a) of the question, with an additional field, `LastDef`. When `lookup` creates an entry, it initializes `LastDef` to -1 .

initialize the hash table

for $i \leftarrow 0$ to $n-1$ do

<i>index</i> \leftarrow <i>lookup</i> (<i>LeftOp</i> (i))	$\text{DEF}(\text{LeftOp}(i), i)$
<i>mark LeftOp</i> (i) with <i>LastDef</i> (<i>index</i>)	$\leftarrow \text{LastDef}(\text{index})$
<i>index</i> \leftarrow <i>lookup</i> (<i>RightOp</i> (i))	$\text{DEF}(\text{RightOp}(i), i)$
<i>mark RightOp</i> (i) with <i>LastDef</i> (<i>index</i>)	$\leftarrow \text{LastDef}(\text{index})$
<i>index</i> \leftarrow <i>lookup</i> (<i>Target</i> (i))	
<i>LastDef</i> (<i>index</i>) $\leftarrow i$	

- d. Apply your algorithm to blocks b_0 and b_1 above.

Answer: The DEF information is harder to represent, because x in $\text{DEF}(x,i)$ refers to a specific reference—that is, a name and an operation. To simplify the issue, we will annotate the code with numerical superscripts. The superscript refers to the operation number in which the value is defined. Thus, c^2 indicates that this specific reference to c uses a value defined in operation 2. The superscript -1 indicates an upward exposed use.

1	$t_1 \leftarrow a^{-1} + b^{-1}$
2	$t_2 \leftarrow t_1^1 + c^{-1}$
3	$t_3 \leftarrow t_2^2 + d^{-1}$
4	$t_4 \leftarrow b^{-1} + a^{-1}$
5	$t_5 \leftarrow t_3^3 + e^{-1}$
6	$t_6 \leftarrow t_4^4 + f^{-1}$
7	$t_7 \leftarrow a^{-1} + b^{-1}$
8	$t_8 \leftarrow t_4^4 - t_7^7$
9	$t_9 \leftarrow t_8^8 * t_6^6$

Block b_0

1	$t_1 \leftarrow a^{-1} \times b^{-1}$
2	$t_2 \leftarrow t_1^1 \times 2$
3	$t_3 \leftarrow t_2^2 \times c^{-1}$
4	$t_4 \leftarrow 7 + t_3^3$
5	$t_5 \leftarrow t_4^4 + d^{-1}$
6	$t_6 \leftarrow t_5^5 + 3$
7	$t_7 \leftarrow t_4^4 + e^{-1}$
8	$t_8 \leftarrow t_6^6 + f^{-1}$
9	$t_9 \leftarrow t_1^1 + 6$

Block b_1

3. Apply the tree-height balancing algorithm from Figures 8.7 and 8.8 to the two blocks in problem 1. Use the information computed in problem 2b above. In addition, assume that $\text{LIVEOUT}(b_0)$ is $\{t_3, t_9\}$, that $\text{LIVEOUT}(b_1)$ is $\{t_7, t_8, t_9\}$, and that the names a through f are upward-exposed in the blocks.

Answer: For block b_0 , we first compute the USES set for each name, as shown in the table on the left below. Notice that t_5 is unused and that we have used an asterisk to indicate an arbitrary large integer—added to the USES set of any name in LIVEOUT.

t_5 is unused, or dead. Watch its progress through the algorithm. Similarly, t_4 and t_7 are identical and redundant with t_1 which is not a root. The algorithm performs tree-height balancing, not common subexpression elimination.

The next step is to find the roots of subexpressions that must be balanced and scheduled. The algorithm enqueues the $\langle \text{root}, \text{priority} \rangle$ pairs $\langle t_3, + \rangle$, $\langle t_4, + \rangle$, $\langle t_6, + \rangle$, $\langle t_7, + \rangle$, $\langle t_8, - \rangle$, and $\langle t_9, \times \rangle$. The roots and their trees are shown in the table to the right below.

#	Uses	Code	Root	Tree
1	2	$t_1 \leftarrow a + b$	t_9	$(\times t_8, t_6)$
2	3	$t_2 \leftarrow t_1 + c$	t_3	$(+ a, b, c, d)$
3	*, *, 5	$t_3 \leftarrow t_2 + d$	t_4	$(+ a, b)$
4	6, 8	$t_4 \leftarrow b + a$	t_6	$(+ t_4, f)$
5	—	$t_5 \leftarrow t_3 + e$	t_7	$(+ a, b)$
6	9	$t_6 \leftarrow t_4 + f$	t_8	$(- t_4, t_7)$
7	8	$t_7 \leftarrow a + b$		
8	9	$t_8 \leftarrow t_4 - t_7$		
9	*, *	$t_9 \leftarrow t_8 * t_6$		

Next, the algorithm pulls trees out of the queue, in precedence order. It flattens the tree and rebuilds it. If it encounters a root in this process, it recurs on the root to ensure that the result of that subtree is available in the current computation.

For block b_0 , the first subtree dequeued is t_9 . It recurs on t_8 (which recurs on t_4 and t_7) and t_6 (which recurs on t_4 and discovers that it has already been rebuilt). Finally, it dequeues t_3 and processes it. The algorithm produces the following code:

Generated Code	
1	$t_4 \leftarrow a + b$
2	$t_7 \leftarrow a + b$
3	$t_8 \leftarrow t_4 - t_7$
4	$t_6 \leftarrow t_4 + f$
5	$t_9 \leftarrow t_8 \times t_6$
6	$temp_0 \leftarrow a + b$
7	$temp_1 \leftarrow c + d$
8	$t_3 \leftarrow temp_0 + temp_1$

A good scheduler could issue operations (1, 2, 6, and 7) in parallel, (3, 4, and 8) in parallel, and finish the block with operation 5.

For block b_1 , we first compute the USES set for each name, as shown in the table to the left below. Again, we have used an asterisk to indicate an arbitrary large integer—added to the USES set of any name in LIVEOUT.

The next step is to find the roots of subexpressions that must be balanced and scheduled. The algorithm enqueues the $\langle root, priority \rangle$ pairs $\langle t_1, \times \rangle$, $\langle t_3, \times \rangle$, $\langle t_4, + \rangle$, $\langle t_7, + \rangle$, $\langle t_8, + \rangle$, and $\langle t_9, + \rangle$. The roots and their trees are shown in the table to the right below.

#	Uses	Code	Root	Tree
1	2, 9	$t_1 \leftarrow a \times b$	t_1	$(\times a, b)$
2	3	$t_2 \leftarrow t_1 \times 2$	t_3	$(\times t_1, 2, c)$
3	4	$t_3 \leftarrow t_2 \times c$	t_4	$(+ 7, t_3)$
4	5, 7	$t_4 \leftarrow 7 + t_3$	t_7	$(+ t_4, e)$
5	6	$t_5 \leftarrow t_4 + d$	t_8	$(+ t_4, d, 3, f)$
6	8	$t_6 \leftarrow t_5 + 3$	t_9	$(+ t_1, 6)$
7	*, *	$t_7 \leftarrow t_4 + e$		
8	*, *	$t_8 \leftarrow t_6 + f$		
9	*, *	$t_9 \leftarrow t_1 + 6$		

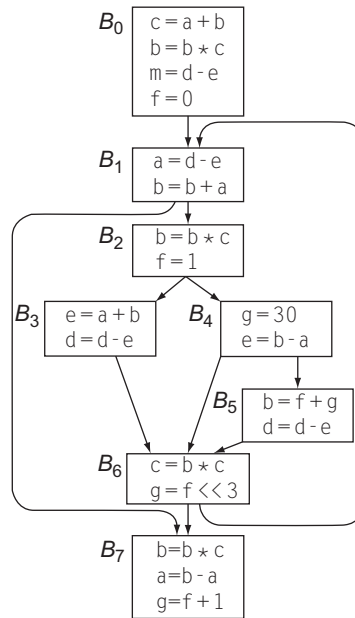
Next, the algorithm pulls trees out of the queue, in precedence order. It flattens the tree and rebuilds it. If it encounters a root in this process, it recurs on the root to make the result of that subtree available.

For block b_1 , the first subtree dequeued is t_1 . Next, the algorithm dequeues t_3 , which recurs on t_1 . It then dequeues the addition subtrees, t_4 , t_7 , t_8 , and t_9 . That process produces the following code:

Generated Code	
1	$t_1 \leftarrow a \times b$
2	$\text{temp}_0 \leftarrow 2 \times c$
3	$t_3 \leftarrow t_1 \times \text{temp}_0$
4	$t_4 \leftarrow 7 + t_3$
5	$t_7 \leftarrow t_4 + e$
6	$\text{temp}_1 \leftarrow t_4 + d$
7	$\text{temp}_2 \leftarrow 3 + f$
8	$t_8 \leftarrow \text{temp}_1 + \text{temp}_2$
9	$t_9 \leftarrow t_1 + 6$

Section 8.5

4. Consider the following control-flow graph:



a. Find the extended basic blocks and list their distinct paths.

Answer: The EBBs are $\{A\}$, $\{B, C, D, E, F\}$, $\{G\}$, and $\{H\}$.

b. Apply local value numbering to each block.

Answer:

B₀: $c_3 = a_1 + b_2$ *no redundancies in this block*
 $b_4 = b_2 * c_3$
 $m_7 = d_5 - e_6$
 $f_8 = 0_8$

B₁: $a_3 = d_1 - e_2$ *no redundancies in this block*
 $b_5 = b_4 + a_3$

B₂: $b_3 = b_1 * c_1$
 $f_4 = 1_4$

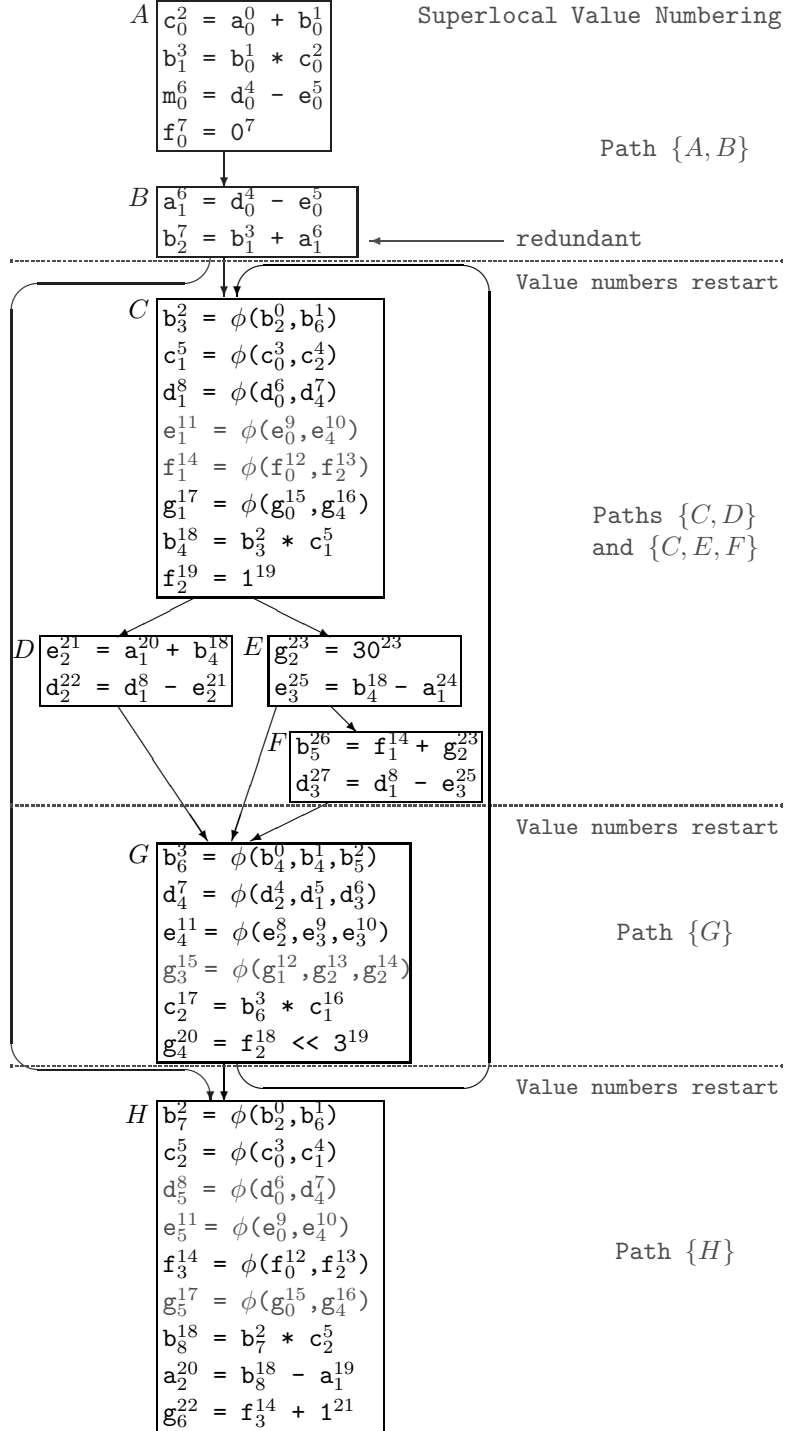
B ₃ : $e_3 = a_1 + b_2$ $d_5 = d_4 - e_3$	<i>no redundancies in this block</i>
B ₄ : $g_1 = 30_1$ $e_4 = b_2 - a_3$	<i>no redundancies in this block</i>
B ₅ : $b_3 = f_1 + g_2$ $d_6 = d_4 - e_5$	<i>no redundancies in this block</i>
B ₆ : $c_3 = b_1 * c_2$ $g_6 = f_4 < 3_5$	<i>no redundancies in this block</i>
B ₇ : $b_3 = b_1 * c_2$ $a_5 = b_3 - a_4$ $g_8 = f_6 + 1_7$	<i>no redundancies in this block</i>

- c. Apply superlocal value numbering to the EBBs and note any improvements that it finds beyond those found by local value numbering.

Answer: First, convert the code to SSA form. The answer (on the next page) uses minimal SSA form. Dead ϕ -functions are shown in gray. Some of the ϕ -functions that are live have their sole uses in other ϕ -functions, including dead ones.

Superlocal value numbering finds the extended basic blocks $\{A, B\}$, $\{C, D, E, F\}$, $\{G\}$, and $\{H\}$. All of these EBBs except $\{C, D, E, F\}$ have a single path, while $\{C, D, E, F\}$ has two paths ($\{C, D\}$ and $\{C, E, F\}$). Thus, the algorithm processed five paths: $\{A, B\}$, $\{C, D\}$, $\{C, E, F\}$, $\{G\}$, and $\{H\}$.

SSA names are indicated with subscripts. Value numbers appear as superscripts.



5. Consider the following simple five-point stencil computation:

```

do 20 i = 2, n-1, 1
  t1 = A(i,j-1)
  t2 = A(i,j)
  do 10 j = 2, m-1, 1
    t3 = A(i,j+1)
    A(i,j) = 0.2 × (t1 + t2 + t3 + A(i-1,j) + A(i+1,j))
    t1 = t2
    t2 = t3
  10 continue
20 continue

```

Each iteration of the loop executes two copy operations.

- a. Loop unrolling can eliminate the copy operations. What unroll factor is needed to eliminate all copy operations in this loop?

Answer: The inner loop should be unrolled by three, as shown in the following code:

```

do 20 i = 2, n-1, 1
  t1 = A(i,j-1)
  t2 = A(i,j)
  do 10 j = 2, m-1, 3
    t3 = A(i,j+1)
    t4 = A(i-1,j)
    t5 = A(i+1,j)
    A(i,j) = (t1+t2+t3+t4+t5)/5

    t1 = A(i,j+2)
    t4 = A(i-1,j+1)
    t5 = A(i+1,j+1)
    A(i,j) = (t1+t2+t3+t4+t5)/5

    t2 = A(i,j+3)
    t4 = A(i-1,j+2)
    t5 = A(i+1,j+2)
    A(i,j) = (t1+t2+t3+t4+t5)/5
  10 continue
20 continue

```

- b. In general, if a loop contains multiple cycles of copy operations, how can you compute the unroll factor needed to eliminate all of the copy operations?

Answer: Find the length of each cycle of copy operations. The desired unroll factor is the least common multiple of the complete set of unique cycle lengths.

Section 8.6

6. At some point p , $\text{LIVE}(p)$ is the set of names that are *live* at p . $\text{LIVEOUT}(b)$ is just the LIVE set at the end of block b .

- a. Develop an algorithm that takes as input a block b and its LIVEOUT set and produces as output the LIVE set for each operation in the block.

Answer: Assume that block b_c consists of a list of three-address operations, numbered from 0 to $n-1$. For the sake of illustration, we will compute a *LiveNow* set for each operation, plus one for a virtual n^{th} operation.

```

LiveNow( $n$ )  $\leftarrow$  LIVEOUT( $b$ )
for  $i = n-1$  to 0 do
    LiveNow( $i$ )  $\leftarrow$  LiveNow( $i+1$ )
    remove Target( $i$ ) from LiveNow( $i$ )
    add LeftOp( $i$ ) to LiveNow( $i$ )
    add RightOp( $i$ ) to LiveNow( $i$ )

```

In practice, most algorithms that use *LiveNow* compute it and use it incrementally, so that they do not need to store a set of each operation.

- b. Apply your algorithm to blocks b_0 and b_1 in problem 1, using $\text{LIVEOUT}(b_0) = \{t_3, t_9\}$ and $\text{LIVEOUT}(b_1) = \{t_7, t_8, t_9\}$.

Answer:

	Code	<i>LiveNow</i>
1	$t_1 \leftarrow a + b$	$\{a, b, c, d, e, f\}$
2	$t_2 \leftarrow t_1 + c$	$\{a, b, c, d, e, f\}$
3	$t_3 \leftarrow t_2 + d$	$\{a, b, d, e, f\}$
4	$t_4 \leftarrow b + a$	$\{t_3, a, b, e, f\}$
5	$t_5 \leftarrow t_3 + e$	$\{t_3, t_4, a, b, e, f\}$
6	$t_6 \leftarrow t_4 + f$	$\{t_3, t_4, a, b, f\}$
7	$t_7 \leftarrow a + b$	$\{t_3, t_4, t_6, a, b\}$
8	$t_8 \leftarrow t_4 - t_7$	$\{t_3, t_4, t_6, t_7\}$
9	$t_9 \leftarrow t_8 * t_6$	$\{t_3, t_6, t_8\}$
10	$\text{LIVEOUT}(b_0)$	$\{t_3, t_9\}$

Results for Block b_0

	Code	<i>LiveNow</i>
1	$t_1 \leftarrow a \times b$	$\{t_1, a, b, c, d, e, f\}$
2	$t_2 \leftarrow t_1 \times 2$	$\{t_1, c, d, e, f\}$
3	$t_3 \leftarrow t_2 \times c$	$\{t_1, t_2, c, d, e, f\}$
4	$t_4 \leftarrow 7 + t_3$	$\{t_1, t_3, d, e, f\}$
5	$t_5 \leftarrow t_4 + d$	$\{t_1, t_4, d, e, f\}$
6	$t_6 \leftarrow t_5 + 3$	$\{t_1, t_4, t_5, e, f\}$
7	$t_7 \leftarrow t_4 + e$	$\{t_1, t_4, t_6, e, f\}$
8	$t_8 \leftarrow t_6 + f$	$\{t_1, t_6, t_7, f\}$
9	$t_9 \leftarrow t_1 + 6$	$\{t_1, t_7, t_8\}$
10	$\text{LIVEOUT}(b_1)$	$\{t_7, t_8, t_9\}$

Results for Block b_1

7. Figure 8.16 shows an algorithm for constructing hot paths in the CFG.

- a. Devise an alternate hot-path construction that pays attention to ties among equal-weight edges.

Answer: Many solutions are possible. In our solution, we will assume a function, $MaxOut(\langle i, j \rangle)$ that, applied to an edge $\langle i, j \rangle$ e , returns the highest frequency count of an edge whose source is j . $MaxOut$ is straightforward to write.

One tie-breaking variant of the algorithm from Figure 8.16 is as follows:

```

 $E \leftarrow | \text{edges} |$ 
for each block  $b$ 
    make a degenerate chain,  $d$ , for  $b$ 
     $priority(d) \leftarrow E$ 
// Create a priority queue,  $EdgeQ$ , of edges ordered by frequency
 $EdgeQ \leftarrow \text{new priority queue}$ 
for each edge  $\langle x, y \rangle$ 
     $Enqueue(EdgeQ, \langle x, y \rangle, Frequency(\langle x, y \rangle))$ 
 $P \leftarrow 0$ 
// Create a priority queue,  $Q$ , of edges with same frequency,
// ordered by  $MaxOut$ 
 $\langle x, y \rangle \leftarrow Dequeue(EdgeQ)$ 
while  $EdgeQ$  is not empty do
     $Q \leftarrow \text{new priority queue}$ 
     $\langle a, b \rangle \leftarrow Dequeue(EdgeQ)$ 
    while  $Frequency(\langle a, b \rangle) = Frequency(\langle x, y \rangle)$  do
         $Enqueue(Q, \langle a, b \rangle, MaxOut(\langle a, b \rangle))$ 
        if  $EdgeQ$  is empty
            then break
     $x \leftarrow a$  // set up next iteration of outer loop
     $y \leftarrow b$ 
// process edges in  $Q$ , in order by  $MaxOut$  values
while  $Q$  is not empty do
     $\langle m, n \rangle \leftarrow Dequeue(Q)$ 
    if  $m$  is the tail of chain  $c_a$  and  $n$  is the head of chain  $c_b$  then
         $t \leftarrow priority(c_a)$ 
        append  $c_b$  onto  $c_a$ 
         $priority(c_a) \leftarrow \min(t, priority(c_b), P++)$ 

```


The version in Figure 8.16 takes equal-priority edges in the order that the queue presents them. This version constructs, for each distinct edge frequency count, a priority queue that contains all edges with that count. For each such edge, $\langle x, y \rangle$, it finds the largest frequency count of the edges that leave y and it uses that value to order the queue.

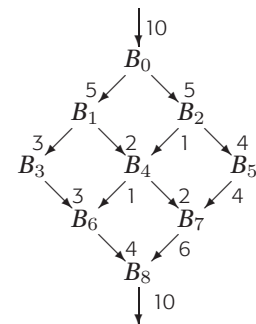
Thus, it will break ties between edges of equal frequency count by considering the frequencies of the edges successors in the CFG, and choosing the remaining edge with the “hottest” successor to process next.

- b. Construct two examples where your algorithm leads to a code layout that improves on the layout produced by the book’s algorithm. Use the code-layout algorithm from Figure 8.17 with the chains constructed by your algorithm and those built by the book’s algorithm.

Answer:

Of course, myriad examples are possible. The CFG shown to the right is one example. The first edge chosen by the algorithm will be $\langle B_7, B_8 \rangle$.

At that point, it must choose one of the edges leaving B_0 . Without tie breaking, the choice will be made on some arbitrary, but typically deterministic, criteria, such as the name of the destination node. If it chooses $\langle B_0, B_1 \rangle$ instead of $\langle B_0, B_2 \rangle$, then the hot path, $(B_0, B_2, B_5, B_7, B_8)$ will contain one or more taken branches. With the other choice, all those jumps will be fall-through cases.



Example CFG

Section 8.7

8. Consider the following code fragment. It shows a procedure fee and two call sites that invoke fee.

```
static int A[1000,1000], B[1000];
...
x = A[i,j] + y;
call fee(i,j,1000);
...
call fee(1,1,0);
...
fee(int row; int col; int ub) {
    int i, sum;
    sum = A[row,col];
    for (i=0; i<ub; i++) {
        sum = sum + B[i];
    }
}
```

- a. What optimization benefits would you expect from inlining fee at each of the call sites? Estimate the fraction of fee's code that would remain after inlining and subsequent optimization.

Answer: At the first call site, the compiler would realize some benefits from eliminating the call. For example, immediately before the call, the code references the value of $A[i, j]$ while the first executable statement inside fee refers to the same location. After inlining, the compiler should recognize that it can reuse the value loaded for the first statement in the second. Of course, some of the overhead of the procedure call and return will vanish, as well. The larger part of the work, however, occurs in the loop. The only impact that inlining will have on the loop is to make the upper bound explicit.

At the second call site, the compiler would realize that the loop is dead code. Since the only use of sum is in the loop, it might eliminate that assignment as well and eliminate all of the code in fee.

Of course, fee is, in some sense, dead code. It has no externally visible effects. Thus, neither call can affect the answers produced by the overall program, except by making the application run longer. However, the compiler cannot tell, at the calls, that fee has no effect. Unless the compiler performs fairly sophisticated interprocedural analysis, it cannot declare those calls dead.

Inlining, at either call site, would move all of the facts needed for dead code elimination into a single procedure. Thus, a compiler with a reasonable dead code eliminator, such as the *Dead* and *Clean* techniques described in Chapter 10, would eliminate all of the code that originated in fee after it inlined the calls.

- b. Based on your experience in part a, sketch a high-level algorithm for estimating the benefits of inlining a specific call site. Your technique should consider both the call site and the callee.

Answer: See answer for question nine, below. We asked the same basic question twice.

9. In Problem 8, features of the call site and its context determined the extent to which the optimizer could improve the inlined code. Sketch, at a high level, a procedure for estimating the improvements that might accrue from inlining a specific call site. (With such an estimator, the compiler could inline the call sites with the highest estimated profit, stopping when it reached some threshold on procedure size or total program size.)

Answer: The “profit” from inline substitution arises from several distinct sources.

- First, inlining a procedure eliminates the overhead of the actual calling sequence—the pre-call sequence, the prolog, the epilog, and the post-return sequence. Most of the work entailed in those sequences is eliminated.
- Second, if any of the actual parameters are constants, they can be folded into the body of the callee after inlining. Two kinds of effects occur: direct simplification of computations and simplification of control flow. If all operands of a computation are (now) known constants, the computation can be folded away. If some operands are known constants and others are not, algebraic simplifications may still be possible. (See, for example, Figure 8.3 on page 424.)

Control-flow constructs that depend on parameters with known constant values can sometimes be simplified. The loop in procedure fee of the example for Question 8.8 (above) does not execute when $ub < 1$.

We can calculate a rough estimate of the amount of computation that the compiler can eliminate by considering the savings in the procedure linkage, the savings in direct computations, and the simplifications of control flow. The estimate can be computed in an ad-hoc fashion. A slightly more formal methodology can be found in reference [31] in the EaC2e bibliography.

10. When the procedure placement algorithm, shown in Figure 8.21, considers an edge $\langle p, q \rangle$ it always places p before q .
- a. Formulate a modification of the algorithm that would consider placing the sink of an edge before its source.

Answer: In principle, the solution is simple. In the graph, each node represents one or more procedures. Initially, each procedure holds one node. When the algorithm coalesces together nodes that are connected by edges in the graph, the resulting node represents all the procedures in the two nodes that it coalesced.

We will call the first procedure in each node its head. We will keep a table of all procedures that contains, for each procedure, the head of its current node, the size of the procedure, and the procedure's distance from the head (from the first procedure in the node). With the entry for each procedure, we will keep a cumulative length for all the procedures in a node headed by that procedure.

These data structures are easy to initialize. The algorithm must update them when it merges two nodes.

Now, when the placement algorithm considers an edge $\langle x, y \rangle$, it should compare the distance between x and y with the node for x placed first against the distance with the node for x placed last. The calculation is easy, given that we know the distance from x to its head and the length of the code associated with that head, along with the same information for y . If

$$\text{dist}(x) + \text{length}(\text{head}(y)) - \text{dist}(y) \leq \text{dist}(y) + \text{length}(\text{head}(x)) - \text{dist}(x)$$

then x should follow y , otherwise y should follow x .

All that remains is to discuss how to maintain the data structures. When the algorithm merges two nodes, say i and j , and places i before j , it must update the information for every procedure contained in j . For each p in j , it must set:

$$\begin{aligned} \text{head}(p) &\leftarrow \text{head}(i) \\ \text{dist}(p) &\leftarrow \text{dist}(p) + \text{length}(\text{head}(i)) \end{aligned}$$

It must also update the length of $\text{head}(i)$. If we update the code from Figure 8.21 with these changes, it might look as follows:

```

build the call multi-graph  $G$  // Initialization work
 $Q \leftarrow$  a new priority queue

for each node  $x \in G$  // Set up new structures
     $list(x) \leftarrow \{x\}$ 
     $head(x) \leftarrow x$ 
     $dist(x) \leftarrow 0$ 
     $size(x) \leftarrow$  actual code length // size is for  $x$ 
     $length(x) \leftarrow size(x)$  // cumulative length of node  $x$ 

for each edge  $(x,y) \in G$  // Add weights to the edges
    if  $(x = y)$  // Self loop is irrelevant
        then delete  $(x,y)$  from  $G$ 
        else  $weight((x,y)) \leftarrow$  estimated execution frequency for  $(x,y)$ 

for each node  $x \in G$ 
     $list(x) \leftarrow \{x\}$  // Initialize placement lists
    if multiple edges exist from  $x$  to  $y$ 
        then combine them and their weights
    for each edge  $(x,z) \in G$  // Put each edge into  $Q$ 
        Enqueue( $Q, (x,z), weight((x,z))$ )

// Iterative reduction of the graph
while  $Q$  is not empty
     $(x,y) \leftarrow$  Dequeue( $Q$ ) // Take highest priority edge
    if  $dist(x) + length(head(y)) - dist(y) \leq$ 
        $dist(y) + length(head(x)) - dist(x)$ 
        then  $src \leftarrow x$  and  $dst \leftarrow y$ 
        else  $src \leftarrow y$  and  $dst \leftarrow x$ 

    for each edge  $(dst,z) \in G$  // Move source from  $dst$  to  $src$ 
        ReSource( $(dst,z), src$ )
    for each edge  $(z,dst) \in G$  // Move target from  $dst$  to  $src$ 
        ReTarget( $(z,dst), src$ )

    for each procedure  $p$  in  $list(dst)$ 
         $head(p) \leftarrow head(src)$ 
         $dist(p) \leftarrow dist(p) + length(head(src))$ 

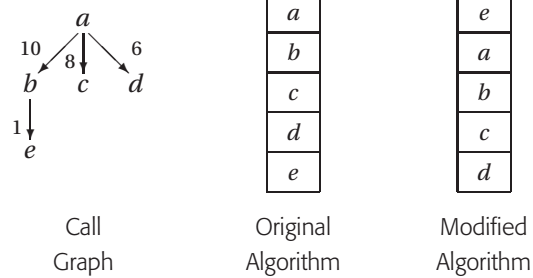
    append  $list(dst)$  to  $list(src)$  // Update the placement list
    delete  $dst$  and its edges from  $G$  // Clean up  $G$ 

```

More efficient schemes can be constructed, for example, using the disjoint-set union-find algorithm.

- b. Construct an example where this approach places two procedures closer together than the original algorithm. Assume that all procedures are of uniform size.

Answer: The number of correct solutions is, essentially, unbounded. Here is one simple graph where considering placement before and placement afterwards results in a shorter distance between two procedures connected by a call.



The original algorithm produces the order shown in the center panel. The modified algorithm produces the order shown in the right panel. If we assume that all procedures have unit size, then the modified algorithm puts *b* closer to *e*.