

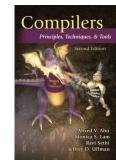
EGRE 591 Compiler Construction

Dr. Dan Resler
 Dept. of Electrical & Computer Engineering
 Virginia Commonwealth University

August 2019



Optional Reference Book



Compilers: Principles, Techniques, and Tools (2nd Edition)

by Aho/Lam/Sethi/Ullman (Addison Wesley: 2006)
 ISBN-13: 978-0321-48681-3

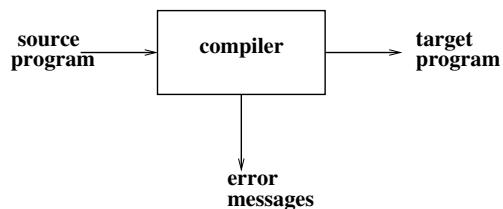
<http://tinyurl.com/y8ldtnr>



└ Introduction
 └ Definitions & History

Definition

A **compiler** is a program that reads a program written in one language and translates it into a logically equivalent program in another language.



└ Introduction
 └ Definitions & History

- Compilers usually translate human-readable computer programs into machine-readable programs
- Therefore they allow programs to be *portable* across machines
- The term compiler was coined in the early 50s by Grace Murray Hopper
- Prior to the late 1950s, many computer scientists considered compilers impossible to write
- First compiler: FORTRAN (took 18 staff-years to implement)



└ Introduction

└ Definitions & History

What Compilers Do

- Compilers can be distinguished in at least two ways:
 - ▶ by the kind of machine code they generate
 - ★ pure machine code
 - ★ augmented machine code
 - ★ virtual machine code
 - ▶ by the format of the target code they generate
 - ★ assembly language
 - ★ relocatable binary
 - ★ absolute binary



└ Introduction

└ Definitions & History

Syntax & Semantics

Definition

The definition of the **syntax** of a programming language specifies its *form* or structure; the specification of a programming language's **semantics** defines the *meaning* of the program

- Syntax is almost-universally defined using **context-free grammars** (CFGs)
- CFGs cannot define all aspects of a program, however
- Semantics of programming languages are commonly defined into two classes: **static** & **runtime** semantics



└ Introduction

└ Definitions & History

Interpreters

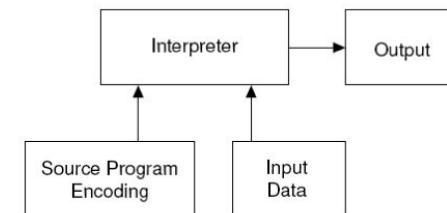


Figure 1.3: An interpreter.



└ Introduction

└ Definitions & History

Static Semantics

Definition

The **static semantics** of a language provide a set of rules that specify which syntactically legal programs are actually valid

- For example, in Java the statement
`a = b + c;`
is valid only if variables `a`, `b`, and `c` are defined and are of the appropriate type
- Static semantics can be specified formally or informally
- Checked at compile-time



Runtime Semantics

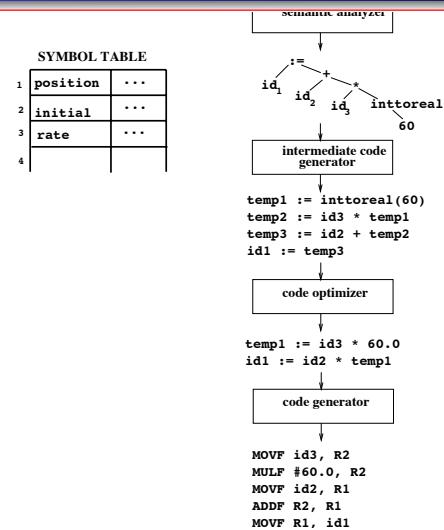
Definition

The **runtime semantics** of a language define what a valid program actually computes (or does) when executed

- Can be specified formally or informally
- Defining formal definitions of runtime semantics is outside the scope of this course



An Example



Organization of a Compiler

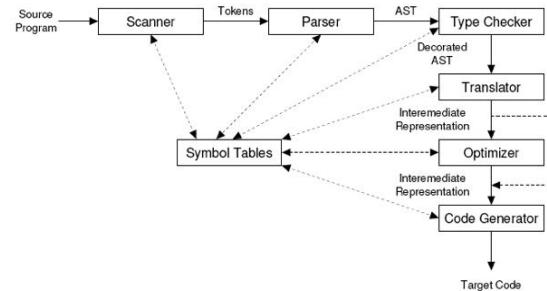
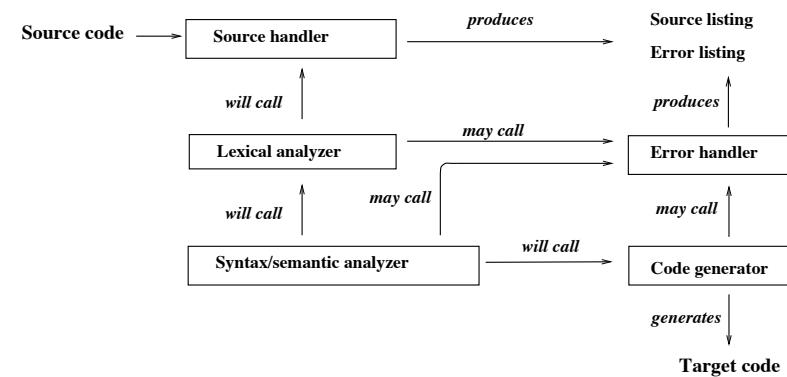


Figure 1.4: A syntax-directed compiler. AST denotes the Abstract Syntax Tree.



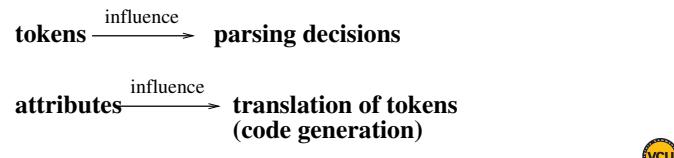
The Interleaved Approach



Definition

The character sequence forming a token is referred to as the **lexeme** for the token, e.g. for the identifier `counter` the token class would be ID and the lexeme "counter"

- Usually one token per keyword in a language (also called *reserved words*)
- Token classes (e.g. ID, REL_OP, BIN_OP) usually have *semantic values* (also called *attributes*) attached to them



- Few program errors are discernible by the scanner
- Example:
`fi (n == 13) System.out.print(n);`
- Can recognize:
 - illegal characters (e.g. % is never allowed in standard Pascal)
 - ill-formed tokens (e.g. in Pascal, .6 is an illegal floating point literal)
- Lexical errors often result in a warning message with the offending character(s) simply ignored



Other tasks of scanners:

- removing white space and comments
- correlating error messages from compiler with source
- producing a listing
- calling or invoking preprocessors

- Scanners are often generated automatically (by a program called a **scanner generator**) from a formal specification of the tokens



- Tokens can be described using natural language, e.g.

Java identifiers ... must start with a letter – including any currency symbol or connecting punctuation (such as _) – followed by letters, digits, or both.

- However, using natural language to describe tokens can be tedious and error-prone. In addition, it is not always obvious how to implement a machine to recognize tokens specified in this manner.
- Therefore scanners are specified using the formal language of **regular expressions** and implemented using **deterministic finite automata**, with mathematics to connect the two

Definition

A set of strings defined by a regular expression is called a **regular set**



Definition

- Regular expressions are defined by a **formal language**

Definition

A **formal language** is a set of strings defined by a **metalanguage**; a metalanguage is a language used to describe other languages



Formal Language Definitions

Definition

- A **grammar** is essentially a set of rules for choosing the subsets of A^* in which we are interested; it has four components:

N – a set of **non-terminal symbols**
T – a set of **terminal symbols**
S – a special **goal** or **start** symbol
P – a set of **productions**

- The set *N* is often spoken of as denoting the **syntactic classes** of the grammar. The union of sets *N* and *T* denotes the **vocabulary** *V* of the grammar.



Formal Language Definitions

Definition

- A **symbol** is an atomic entity, represented by a character, or sometimes by a reserved word or keyword
- An **alphabet** or **vocabulary** *A* is a non-empty finite set of symbols.
- A **word** or **string** 'over' an alphabet *A* is a sequence $\sigma = a_1 a_2 \dots a_n$ of symbols from the alphabet.
- The **null string** or **empty word** is a string of length zero usually denoted by ϵ (epsilon) (Fischer et al. use λ (lambda)).
- The set of all strings over an alphabet *A* is denoted A^* .
- The **language** *L* over alphabet *A* is a subset of A^* . Note that this involves no concept of meaning.



Operations on Languages

Operation	Notation	Definition
union of <i>L</i> and <i>M</i>	$L \cup M$	$L \cup M = \{s \mid s \in L \vee s \in M\}$
concatenation of <i>L</i> and <i>M</i>	LM	$LM = \{st \mid s \in L \wedge t \in M\}$
exponentiation of <i>L</i>	L^i	$L^i = \begin{cases} \epsilon, & \text{if } i = 0 \\ L^{i-1}L, & \text{if } i > 0 \end{cases}$
Kleene closure of <i>L</i>	L^*	$L^* = \bigcup_{i=0}^{\infty} L^i$
positive closure of <i>L</i>	L^+	$L^+ = \bigcup_{i=1}^{\infty} L^i$



Operations on Languages

- Example:
 $L = \{A, B, \dots, Z, a, b, \dots, z\}$ $D = \{0, 1, \dots, 9\}$
 - $L \cup D$ is the set of letters and digits
 - LD is the set of strings consisting of a letter followed by a digit
 - L^4 is the set of all four-letter strings over L
 - L^* is the set of all strings of letters over L , including ϵ
 - $L(L \cup D)^*$ is the set of all strings of letters and digits beginning with a letter
 - D^+ is the set of all strings of one or more digits
- Since a symbol is a string of length one, L & D can also be thought of as finite languages
- Operations on languages can be thought of as creating new languages



Overview

- Therefore a regular expression is built up out of simpler regular expressions according to a set of rules
- A regular expression is a set of strings L over an alphabet such that L can be defined by a regular expression
- Regular expression r is said to denote language $L(r)$
- Additional metacharacters often used in regular expressions:
 - | (vertical line) means 'or'
 - (dot) sometimes used for concatenation
 - ? (question mark) denotes 0 or 1 occurrence
 - * (star) denotes 0 or more occurrences
 - () (parenthesis) used to group terms
 - [] (brackets) denote character classes
 - not(X) denotes all characters except X
 - "?" quotes used to delimit metacharacter literals



Overview

Definition

A **regular expression** is a formal expression that is

- a single character in the alphabet Σ
- the empty string ϵ
- the empty set $\{\}$
- any string derived by a finite number of applications of the union, concatenation, or closure operations on any sets of strings in Σ^*



Examples

a	{a} (the language containing just the string a)
$a b$	{a, b} (a language containing 2 strings, a and b)
$a" "b$	{a b} (a language containing just the string a b)
$(a b) \cdot a$	{aa, ba}
$(a b)a$	{aa, ba}
$a?$	$a \epsilon$
$(a \cdot b) \epsilon$	{"", ab}, equivalent to $(a \cdot b) ?$
$[a, b, c]$	$a b c$
$[A-Za-z]$	$A B \dots Z a b \dots z$
$((a b)a)^*$	{"", aa, ba, aaaa, baaa, aaba, baba, aaaaaa, ... }
$(0 1)^* \cdot 0$	zero & all binary numbers that are multiples of 2
$a(\text{not}(a))a$	all 3 character strings beginning & ending with a except aaa
$b^*(abb^*)^*a$	strings of a's and b's with no consecutive a's



Examples from Programming Languages

- the keyword `while`:
`while`
- identifiers that begin with a letter followed by 0 or more letters or digits:
`[A-Za-z][A-za-z0-9]*`
- a literal floating point number that must begin with a digit:
`[-]?[0-9]+.[0-9]+`
- a Java comment (where EOL denotes end-of-line):
`//(not(EOL))*EOL`
- a comment delimited by @@ marks which allows single @'s within the comment body:
`@@(@not(@))|not(@)*@@`



Specifying Regular Expressions for Scanner Generators

```

if                                {return IF;}
[a-z][a-z0-9]*                  {return ID;}
[0-9]+                            {return NUM;}
([0-9]+."[0-9]*)|([0-9]*."[0-9]+) {return REAL;}
"--"[a-z]*\n)|(" "|\n|\t)+    { /* do nothing */ }
                                {error();}

```

FIGURE 2.2.

Regular expressions for some tokens.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel



Limitations of Regular Expressions

- Many languages cannot be described by regular expressions; for example: $L = \{a^n b^n \mid n \geq 1\}$
- Regular expressions can only be used to denote a fixed number of repetitions or an unspecified number—two arbitrary number of repetitions cannot be compared to see if they are equal



Overview

Definition

A **recognizer** for a language is a program that takes as input a string x and answers “yes” if x is a sentence of the language, “no” otherwise

Definition

- Finite automata** are the (theoretical or real) machines that recognize regular expressions
- A finite automata consists of:
 - a finite set of **states**, with one designated as the **start** state and one or more **final** states
 - labeled **edges** to connect states



Overview

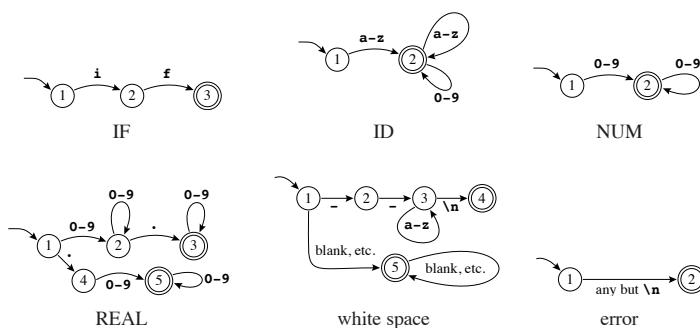


FIGURE 2.3. Finite automata for lexical tokens. The states are indicated by circles; final states are indicated by double circles. The start state has an arrow coming in from nowhere. An edge labeled with several characters is shorthand for many parallel edges.
From *Modern Compiler Implementation in Java*, Cambridge University Press, ©1998 Andrew W. Appel



Token Recognition

Last Final State	Current Input	Accept Action
0 1	if --not-a-com	
2 2	if --not-a-com	
3 3	if --not-a-com	
3 0	if --not-a-com	return IF
0 1	if --not-a-com	
12 12	if --not-a-com	
12 0	if --not-a-com	found white space; resume
0 1	if --not-a-com	
9 9	if --not-a-com	
9 10	if --not-a-com	
9 0	if --not-a-com	error, illegal token '-'; resume
0 1	if --not-a-com	
9 9	if --not-a-com	
9 0	if --not-a-com	error, illegal token '-'; resume

FIGURE 2.5. The automaton of Figure 2.4 recognizes several tokens. The symbol | indicates the input position at each successive call to the lexical analyzer, the symbol ⊥ indicates the current position of the automaton, and ⊤ indicates the most recent position in which the recognizer was in a final state.
From *Modern Compiler Implementation in Java*, Cambridge University Press, ©1998 Andrew W. Appel



Combining Finite Automata

```

if [a-z][a-z0-9]* {return IF;}
[0-9]+ {return ID;}
([0-9]+.“[0-9]*)|([0-9]+.“[0-9]+) {return NUM;}
(“--|[a-z]*\n)|( “|\n|\t)+ /* do nothing */
{error();}
.
.
.

```

FIGURE 2.2. Regular expressions for some tokens.
From *Modern Compiler Implementation in Java*, Cambridge University Press, ©1998 Andrew W. Appel

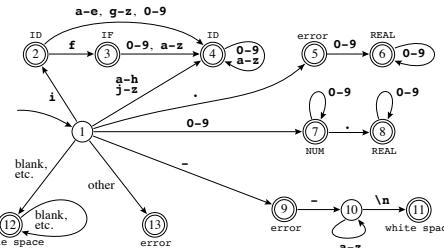
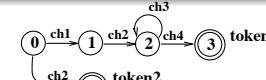


FIGURE 2.4. Combined finite automaton.
From *Modern Compiler Implementation in Java*, Cambridge University Press, ©1998 Andrew W. Appel



Implementing a Finite State Automata by Hand



```

state = 0; lexeme = ""; currChar = getChar();
while not finished do

```

```

case state of
  0: if currChar = ch1 then
      concatenate currChar to lexeme
      state = 1; currChar = getChar()
    elseif currChar = ch2 then
      concatenate currChar to lexeme
      return token2 // state 4
    else error
  1: if currChar = ch2 then
      concatenate currChar to lexeme
      state = 2; currChar = getChar()
    else error
  2: if currChar = ch3 then
      concatenate currChar to lexeme
      currChar = getChar()
    elseif currChar = ch4 then
      concatenate currChar to lexeme
      return token1 // state 3
    else error
endcase
endwhile

```

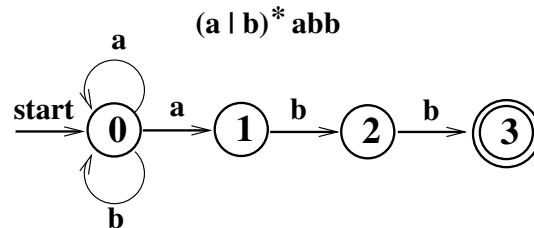


Types of Finite Automata

- Two types of finite automata:
 - non-deterministic (NFA)
 - deterministic (DFA)
- With NFAs, more than one transition may be possible on the same input symbol
- Both are capable of recognizing regular sets
- DFA faster, bigger than NFA
- Essentially only DFAs are used to implement scanners



Non-deterministic Finite Automata



	INPUT SYMBOL	
STATE	a	b
0	{0, 1}	{0}
1	-	{2}
2	-	{3}



Non-deterministic Finite Automata

Definition

- A **non-deterministic finite automaton** is a mathematical model that consists of
 - a set of states S
 - a set of input symbols Σ
 - a transition function *move* that maps state-symbol pairs to sets of states
 - a state s_0 designated as the *start* (or *initial*) state
 - a set of states F designated as *accepting* (or *final*) states



Non-deterministic Finite Automata

- An NFA *accepts* an input string x if and only if there is some path in the transition graph from the start to some accepting state such that the edges along the path spell out x
- A path can be represented by a sequence of state transitions called *moves*; the following diagram shows the moves made in accepting string $aabb$ using the NFA on the previous slide:



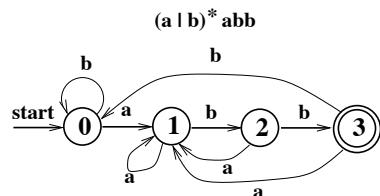
- NFAs allow ϵ -transitions
- The *language defined* by an NFA is the set of input strings it accepts



Deterministic Finite Automata

Definition

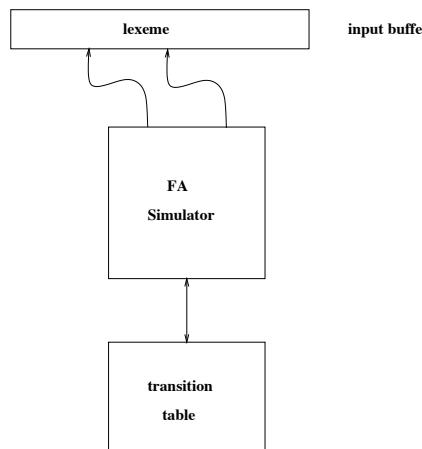
- A **deterministic finite automaton** is a special case of a non-deterministic finite automaton in which
 - ① no state has an ϵ -transition (a transition on input ϵ)
 - ② for each state s and input symbol a , there is at most one edge labeled a leaving s



STATE	INPUT SYMBOL	
	a	b
0	{1}	{0}
1	{1}	{2}
2	{1}	{3}
3	{1}	{0}



Using the Finite Automata Model



DFA Simulator

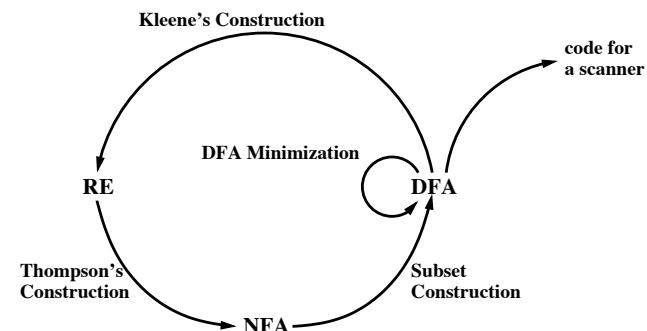
- Definition: **edge**(s, c) – the set of states reachable by following a *single* edge with label c from state s
- DFA simulation algorithm (given a start state of s_1 , the input string c_1, \dots, c_k , and the set of accepting states F):


```

d ← {s1}
for i ← 1 to k do
  d ← edge(d, ci)
endfor
if d ⊆ F then
  return "yes"
else
  return "no"
endif
```



Automatic Generation of Scanners from Regular Expressions¹



¹Diagram is from *Engineering a Compiler* by Cooper & Torczon, Morgan Kaufmann, 2004, pg. 45



Automatic Generation of Scanners from Regular Expressions



Converting a Regular Expression to an NFA

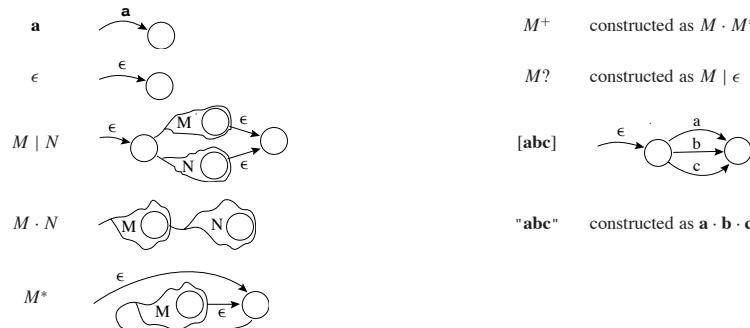


FIGURE 2.6. Translation of regular expressions to NFAs.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel

Converting a Regular Expression to an NFA

- Useful because it is easy to convert a regular expression to an NFA
 - Algorithm referred to as *Thompson's construction*
 - Each individual regular expression is translated into an NFA using a standard template, then templates are combined

Converting a Regular Expression to an NFA

```

if                                {return IF;}
[a-z][a-z0-9]*                  {return ID;}
[0-9]+                            {return NUM;}
([0-9]+"."[0-9]*)|([0-9]*"."[0-9]+) {return REAL;}
"--"[a-z]*\n)|(" \|n|\t)+      { /* do nothing */ }
                                  {error();}
```

FIGURE 2.2. Regular expressions for some tokens.
 From *Modern Compiler Implementation in Java*,
 Cambridge University Press, ©1998 Andrew W. Appel

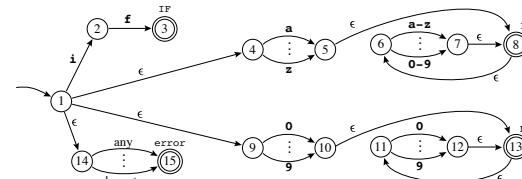
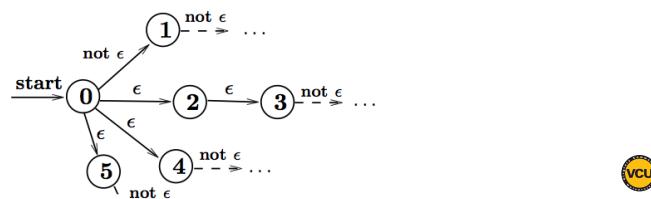


FIGURE 2.7. Four regular expressions translated to an NFA.
From Modern Compiler Implementation in Java,
 Cambridge University Press, ©1998 Andrew W. Appel

Converting an NFA to a DFA

- In the transition table for an NFA, each entry is a set of states; in the transition table for a DFA, each entry is a single state
- The general idea behind NFA-to-DFA translation is that each DFA state corresponds to a set of NFA states
- Process called **subset construction**
- Definition: **ϵ -closure(S)** – the set of NFA sets reachable from NFA state S (or set of states S) on ϵ -transitions alone



D. Resler © 2019-2024

EGRE591 Compiler Construction 49/243



Converting an NFA to a DFA

- Remember that the DFA generated from our NFA-to-DFA algorithm will (potentially) have single states representing multiple NFA states:

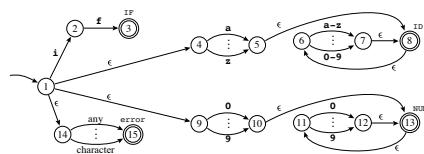


FIGURE 2.7.
Four regular expressions translated to an NFA.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel

FIGURE 2.8.
NFA converted to DFA.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel



D. Resler © 2019-2024

EGRE591 Compiler Construction 51/243

Converting an NFA to a DFA

- Formally, ϵ -closure(S) can be defined as the smallest set T such that

$$T = S \cup \left(\bigcup_{s \in T} \text{edge}(s, \epsilon) \right)$$

- Algorithm for calculating T :

 $T \leftarrow S$ **repeat** $T' \leftarrow T$ $T \leftarrow T' \cup (\bigcup_{s \in T'} \text{edge}(s, \epsilon))$ **until** $T = T'$ 

D. Resler © 2019-2024

EGRE591 Compiler Construction 50/243

Converting an NFA to a DFA

- Definition: **DFAedge(d, c)** – given the generated DFA state d will return the new DFA state reached (which represents a set of NFA states) on input c

- More formally:

$$\text{DFAedge}(d, c) = \epsilon\text{-closure} \left(\bigcup_{s \in d} \text{edge}(s, c) \right)$$



D. Resler © 2019-2024

EGRE591 Compiler Construction 52/243

Converting an NFA to a DFA

- NFA simulation algorithm (given an NFA start state of s_1 , the input string c_1, \dots, c_k , and the set F of accepting states):

```

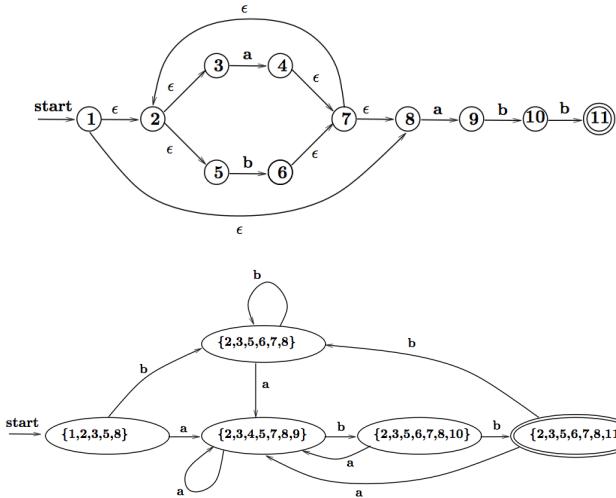
 $d \leftarrow \epsilon\text{-closure}(\{s_1\})$ 
for  $i \leftarrow 1$  to  $k$  do
     $d \leftarrow \text{DFAedge}(d, c_i)$ 
endfor
if ( $\exists e \mid (e \in d \wedge e \in F)$ ) then
    return "yes"
else
    return "no"
endif

```

- Manipulating sets in real time each time a character is inputted is prohibitively expensive; however, it is possible to do all of the sets-of-states calculations in advance



Converting an NFA to a DFA



Converting an NFA to a DFA

- NFA → DFA construction algorithm (assume that $\text{DFAedge}(d_i, c) = d_j$ if there is an edge from d_i to d_j labeled with c , the NFA start state is 1, and that Σ is the alphabet):

```

states[1]  $\leftarrow \{\}$ ; states[1]  $\leftarrow \epsilon\text{-closure}(s_1)$ 
p  $\leftarrow 1$ ; j  $\leftarrow 1$ 
while  $j \leq p$  do
    foreach  $c \in \Sigma$  do
        e  $\leftarrow \text{DFAedge}(\text{states}[j], c)$ 
        if  $e = \text{states}[i]$  for some  $i \leq p$ 
            then trans[j,c]  $\leftarrow i$ 
        else p  $\leftarrow p + 1$ 
            states[p]  $\leftarrow e$ 
            trans[j,c]  $\leftarrow p$ 
        endif
    endfor
    j  $\leftarrow j + 1$ 
endwhile

```



Syntax Analysis

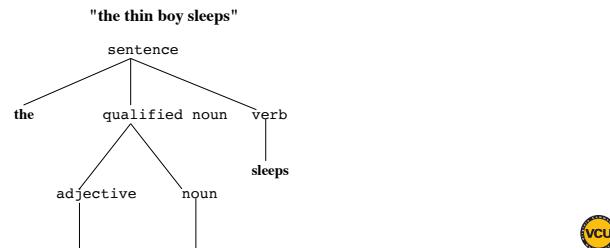
Definition

- The syntax of programming languages can be described by **grammars**
 - A grammar is a set of rules for *choosing* subsets of a language
 - Or it can be thought of a set of rules for *generating* subsets of a language
- The process of determining whether a sentence is legal for a given language is called **parsing**



Syntax Analysis

sentence	$\rightarrow \text{the qualified_noun verb} \mid \text{pronoun verb}$
qualified_noun	$\rightarrow \text{adjective noun}$
noun	$\rightarrow \text{man} \mid \text{girl} \mid \text{boy} \mid \text{lecturer}$
pronoun	$\rightarrow \text{he} \mid \text{she}$
verb	$\rightarrow \text{drinks} \mid \text{sleeps} \mid \text{mystifies}$
adjective	$\rightarrow \text{tall} \mid \text{thin} \mid \text{thirsty}$



Backus-Naur Form (BNF)

Metasymbol	Meaning
\rightarrow or $::=$	"is defined as"
$ $	"or"
< > or special notation	delimiters/indicators for non-terminals
non-delimited text or special notation	terminals
ϵ	empty right-hand side

- For example:

$A \rightarrow B .$	$\langle A \rangle ::= \langle B \rangle .$
$B \rightarrow x \mid (C) \mid [B] \mid \epsilon$	$\langle B \rangle ::= x \mid (\langle C \rangle) \mid [\langle B \rangle] \mid \epsilon$
$C \rightarrow BD$	$\langle C \rangle ::= \langle B \rangle \langle D \rangle$
$D \rightarrow +BD \mid \epsilon$	$\langle D \rangle ::= + \langle B \rangle \langle D \rangle \mid \epsilon$



Components of a Grammar

- Remember that we previously defined a grammar as:

N – a set of **non-terminal** symbols
 T – a set of **terminal** symbols
 S – a special **goal** or **start** symbol
 P – a set of **productions**

- For the grammar on the previous slide:

$N = \{ \text{sentence,qualified_noun,noun,pronoun,verb,adjective} \}$
 $T = \{ \text{the,man,girl,boy,lecturer, ..., tall,thin,thirsty} \}$
 $S = \text{sentence}$
 $P = \{ \text{'sentence} \rightarrow \text{the qualified_noun verb}', \text{'sentence} \rightarrow \text{pronoun verb}', \text{'qualified_noun} \rightarrow \text{adjective noun}', \text{'noun} \rightarrow \text{man}', \text{'noun} \rightarrow \text{girl}', \dots, \text{'adjective} \rightarrow \text{thin}', \text{'adjective} \rightarrow \text{thirsty}' \}$



BNF vs. EBNF

- BNF:

$$\begin{aligned} A &\rightarrow B . \\ B &\rightarrow x \mid (C) \mid [B] \mid \epsilon \\ C &\rightarrow BD \\ D &\rightarrow +BD \mid \epsilon \end{aligned}$$

- EBNF (Extended BNF):

$$\begin{aligned} A &\rightarrow B . \\ B &\rightarrow [x \mid (C) \mid [B]] \\ C &\rightarrow BD \\ D &\rightarrow \{ +BD \} \end{aligned}$$

- In addition EBNF often delimits tokens with double quotes, e.g. "x"

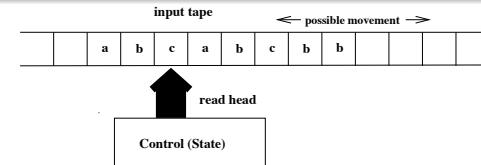


Chomsky's Heirarchy

Class	Grammar	Recognizer
0	Unrestricted $\alpha \rightarrow \beta$ (where α contains a nonterminal)	Turing Machine
1	Context-sensitive $\alpha B \gamma \rightarrow \alpha \beta \gamma$	Linear-bounded Automaton
2	Context-free $A \rightarrow \alpha$	Push-down Automaton
3	Regular $A \rightarrow a \mid aB \mid \epsilon$	Finite-state Automaton



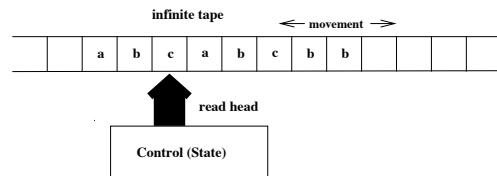
Grammars and Recognizers



- A basic automaton consists of:
 - ▶ a control mechanism with a finite number of states
 - ▶ a possibly infinite *tape* which may be read and advanced (plus possibly moved in either direction and written to)
 - ▶ A finite *alphabet*
 - ▶ A *start state* and a series of *transitions* to other states upon processing of the tape
 - ▶ Initial *input* represented by characters on the tape
 - ▶ Ideally it will have transitions to *halt* either upon it recognizing a legal string or upon it being *blocked* by illegal input
 - ▶ Our recognizers will be variations of this machine



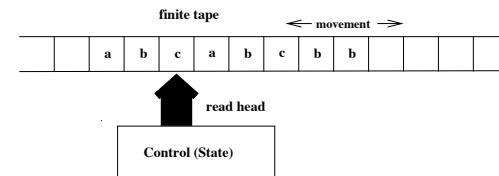
Unrestricted Languages



- Unrestricted automata model Turing machines
- Infinitely long read/write tape moveable in both directions with a start but no end
- Can compute all programs



Context-Sensitive Languages

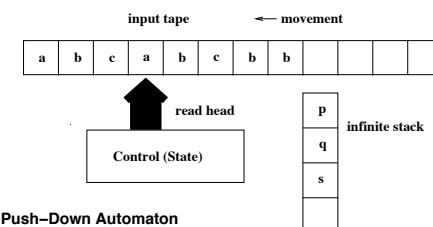


- Linear bounded automata recognize context-sensitive languages
- Finite read/write tape moveable in both directions
- Essentially a Turing machine running on a finite machine

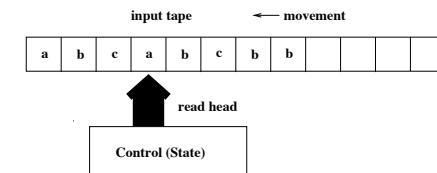


Context-Free Languages

- A language that can be generated by a context-free grammar is said to be a *context-free language*
- Recognizer for context-free languages: *push-down automata*



Regular Languages



Finite-State Automaton



Derivations

- Sentences can be **derived** from a grammar through a series of rewrite rules that specify the method for expanding non-terminals

1 $S \rightarrow S ; S$	4 $E \rightarrow id$
2 $S \rightarrow id := E$	5 $E \rightarrow num$
3 $S \rightarrow print (L)$	6 $E \rightarrow E + E$
	7 $E \rightarrow (S , E)$
	8 $L \rightarrow E$
	9 $L \rightarrow L , E$

GRAMMAR 3.1. A syntax for straight-line programs.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel

\underline{S}	$\underline{S} ; \underline{S}$
$\underline{S} ; \underline{id} := \underline{E}$	$\underline{id} := \underline{E} ; \underline{id} := \underline{E}$
$\underline{id} := \underline{num} ; \underline{id} := \underline{E}$	$\underline{id} := \underline{num} ; \underline{id} := \underline{E} + \underline{E}$
$\underline{id} := \underline{num} ; \underline{id} := \underline{E} + \underline{E}$	$\underline{id} := \underline{num} ; \underline{id} := \underline{E} + (\underline{S} , \underline{E})$
$\underline{id} := \underline{num} ; \underline{id} := \underline{E} + (\underline{S} , \underline{E})$	$\underline{id} := \underline{num} ; \underline{id} := \underline{id} + (\underline{id} := \underline{E} , \underline{E})$
$\underline{id} := \underline{num} ; \underline{id} := \underline{id} + (\underline{id} := \underline{E} , \underline{E})$	$\underline{id} := \underline{num} ; \underline{id} := \underline{id} + (\underline{id} := \underline{E} + \underline{E} , \underline{id})$
$\underline{id} := \underline{num} ; \underline{id} := \underline{id} + (\underline{id} := \underline{E} + \underline{E} , \underline{id})$	$\underline{id} := \underline{num} ; \underline{id} := \underline{id} + (\underline{id} := \underline{num} + \underline{E} , \underline{id})$
$\underline{id} := \underline{num} ; \underline{id} := \underline{id} + (\underline{id} := \underline{num} + \underline{num} , \underline{id})$	

DERIVATION 3.2.
From *Modern Compiler Implementation in Java*.



Sentential Forms

Definition

A **sentential form** is the goal or start symbol or any string that can be derived from it

- A string is said to be a *sentential form of a grammar* if it can be derived from it



Left-most vs. Right-most Derivations

$$\begin{array}{lll} 1. S \rightarrow S ; S & 4. E \rightarrow id & 8. L \rightarrow E \\ 2. S \rightarrow id := E & 5. E \rightarrow num & 9. L \rightarrow L , E \\ 3. S \rightarrow print (L) & 6. E \rightarrow E + E & \\ & 7. E \rightarrow (S , E) & \end{array}$$

GRAMMAR 3.1. A syntax for straight-line programs.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel

- For the sentence 'id := (id := num, id); print(id)':

► **left-most:**

$S \Rightarrow S ; S \Rightarrow id := E ; S \Rightarrow id := (S , E) ; S \Rightarrow$
 $id := (id := E ; E) ; S \Rightarrow id := (id := num ; E) ; S \Rightarrow$
 $id := (id := num ; id) ; S \Rightarrow id := (id := num ; id) ; print(L) \Rightarrow$
 $id := (id := num ; id) ; print(E) \Rightarrow id := (id := num ; id) ; print(id)$

► **right-most:**

$S \Rightarrow S ; S \Rightarrow S ; print(L) \Rightarrow S ; print(E) \Rightarrow S ; print(id) \Rightarrow$
 $id := E ; print(id) \Rightarrow id := (S , E) ; print(id) \Rightarrow$
 $id := (S , id) ; print(id) \Rightarrow id := (id := E , id) ; print(id) \Rightarrow$
 $id := (id := num , id) ; print(id)$



Ambiguous Grammars

- Every parse tree has associated with it a unique derivation
- However, not every sentence necessarily has only one possible parse tree
- A grammar that can derive more than one parse tree for a given sentence is said to be *ambiguous*

$$\begin{array}{l} E \rightarrow id \\ E \rightarrow num \\ E \rightarrow E * E \\ E \rightarrow E / E \\ E \rightarrow E + E \\ E \rightarrow E - E \\ E \rightarrow (E) \end{array}$$

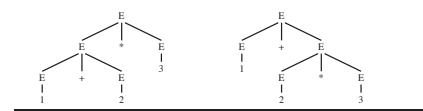


FIGURE 3.7. Two parse trees for the sentence $1+2*3$ in Grammar 3.5.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel



Parse Trees

- Essentially a graph of (or representing) a derivation
- Two different derivations can have the same parse tree

$\frac{S}{S : \underline{S}}$
 $\frac{\underline{S}}{\underline{id} := E}$
 $\frac{id := \underline{E}}{id := num ; id := \underline{E}}$
 $\frac{id := num ; id := \underline{E} + \underline{E}}{id := num ; id := \underline{E} + (\underline{S} , E)}$
 $\frac{id := num ; id := \underline{S} + \underline{E}}{id := num ; id := id + (\underline{id} := \underline{E} , E)}$
 $\frac{id := num ; id := id + (id := \underline{E} + E , \underline{E})}{id := num ; id := id + (id := \underline{E} + E , id)}$
 $\frac{id := num ; id := id + (id := num + \underline{E} , \underline{id})}{id := num ; id := id + (id := num + num , id)}$
 $\frac{id := num ; id := id + (id := num + num , id)}{id := num ; id := id + (id := num + num , id)}$

DERIVATION 3.2.

From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel

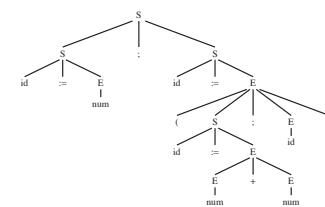


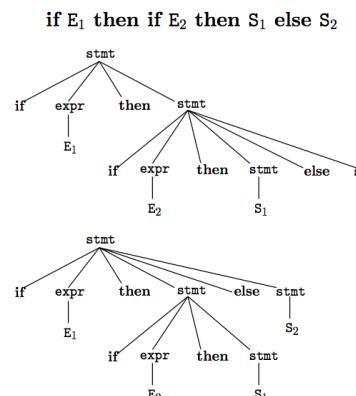
FIGURE 3.3. Parse tree.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel



Eliminating Ambiguity

- Sometimes a grammar can be re-written to eliminate ambiguity
- Example: the 'dangling else'

$stmt \rightarrow if\ expr\ then\ stmt$
 $| if\ expr\ then\ stmt\ else\ stmt$
 $| \dots$

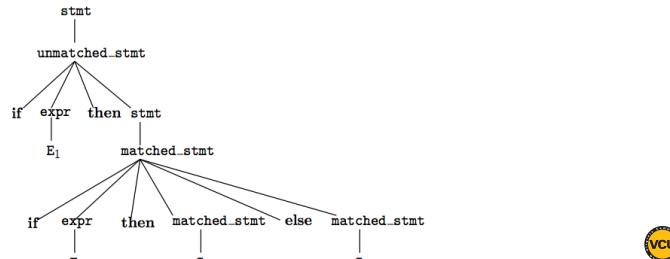


Eliminating Ambiguity

```

stmt      → matched_stmt
          | unmatched_stmt
matched_stmt → if expr then matched_stmt else matched_stmt
          | ...
unmatched_stmt → if expr then stmt
                 | if expr then matched_stmt else unmatched_stmt
                   if E1 then if E2 then S1 else S2

```



Top-down vs. Bottom-up Parsing

- **top-down:** start from the goal symbol and try to generate the sentence by substituting the *expanding* non-terminals at each step
- **bottom-up:** start from the input sentence and make *reductions* to non-terminals until you end up with the goal symbol
- top-down parsing is associated with left-most derivations, bottom-up with right-most



Syntax Error Recovery

- Simplest: parsing halts on first error
- Often can do better by restoring parser to some state where parsing can continue
- Strategies:
 - Panic-mode recovery – discard tokens until some synchronizing token is found
 - Phrase-level recovery – parser corrects token stream by changing or inserting tokens
 - Error productions – parser writer anticipates certain errors and inserts error productions into grammar
 - Global correction – analyze and modify entire program with purpose of altering it to be syntactically correct
- Can use combinations of strategies

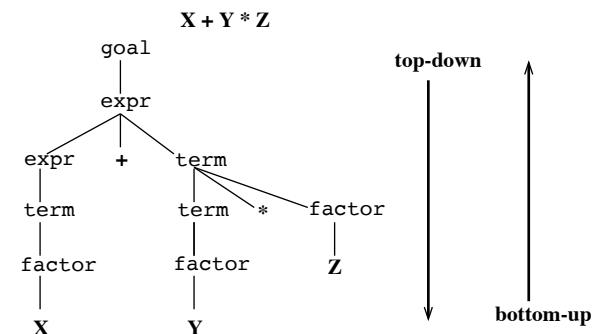


Top-down vs. Bottom-up Parsing

```

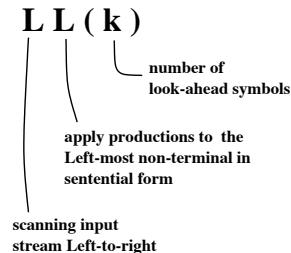
goal   → expr
expr  → term | expr + term
term   → factor | term * factor
factor → X | Y | Z

```



LL(k) Parsers

- LL(k) parsers parse LL(k) languages
- (usually deterministic) top-down parsers



- Most common type: LL(1)
- Informally, a language is LL(k) if there is enough information to choose the correct expansion of a non-terminal at any point in the parse given the current k look-ahead symbols
- Not all LL languages can be parsed with LL(1) parsers



Predictive Parsers

Definition

A **predictive parser** is a deterministic LL(k) parser

- With predictive parsing, the inability to proceed with the parsing of a sentence indicates an error (i.e. the sentence is not a legal sentence for the language)
- Most commonly used predictive parsers: **recursive descent parsers**



LL(k) Parsers

$A \rightarrow xB xC$	<u>Sentential Form</u>	<u>Input String</u>	
$B \rightarrow xB y$	$S = A$	$xxxz$	← try $A \rightarrow xB$
$C \rightarrow xC z$	xB	$xxxz$	
	xB	xxz	
	xB	xz	
	xB	z	
	y	z	?

- Solution: allow parser to backtrack
- Very expensive
- Need deterministic parsers, one where we can be sure which production to apply at any given moment in the parse



Overview

- **Recursive descent parsers** execute a set of recursive procedures to parse the input
 - ▶ a parsing procedure P is associated with each nonterminal A
 - ▶ procedure P is responsible for accomplishing one step of a derivation by choosing and applying one of A 's productions (i.e. it recognizes the language subphrase associated with A)
 - ▶ the parser chooses the appropriate production for A by inspecting the next k tokens in the input stream
 - ▶ the **predict set** for production $A \rightarrow \alpha$ is the set of tokens that trigger application of that production
 - ▶ the predict set for $A \rightarrow \alpha$ is determined primarily by α , the right-hand side (RHS) of the production



└ Top-Down Parsing

└ Recursive Descent Parsing

Overview

$$\begin{aligned} A &\rightarrow B \cdot \\ B &\rightarrow x \mid (C) \mid [B] \mid \epsilon \\ C &\rightarrow BD \\ D &\rightarrow +BD \mid \epsilon \end{aligned}$$

- A recursive descent parser for the above language would require 4 procedures
- Note that if we are recognizing subphrase B , if the next token is:
 - x expand using production $B \rightarrow x$
 - (expand using production $B \rightarrow (C)$
 - [expand using production $B \rightarrow [B]$
 - .)] + expand using production $B \rightarrow \epsilon$
- Therefore the predict set for production B : { x, (., ., .),]+, }



└ Top-Down Parsing

└ Recursive Descent Parsing

Overview

$$\begin{array}{lll} S \rightarrow E\$ & T \rightarrow T * F & F \rightarrow id \\ E \rightarrow E + T & T \rightarrow T / F & F \rightarrow num \\ E \rightarrow E - T & T \rightarrow F & F \rightarrow (E) \end{array}$$

GRAMMAR 3.10.

From Modern Compiler Implementation in Java,
Cambridge University Press, ©1998 Andrew W. Appel

```
void S() { E(); eat(EOF); }
void E() { switch(tok) {
    case ?: E(); eat(PLUS); T(); break;
    case ?: E(); eat(MINUS); T(); break;
    case ?: T(); break;
    default: error();
}
void T() { switch(tok) {
    case ?: T(); eat(TIMES); F(); break;
    case ?: T(); eat(DIV); F(); break;
    case ?: F(); break;
    default: error();
}}
```

- Recursive-descent parsing works only on grammars where the *first terminal symbol* of each subexpression provides enough information to choose which production to use



└ Top-Down Parsing

└ Recursive Descent Parsing

Overview

$$\begin{array}{ll} S \rightarrow \text{if } E \text{ then } S \text{ else } S & L \rightarrow \text{end} \\ S \rightarrow \text{begin } S \text{ } L & L \rightarrow ; \text{ } S \text{ } L \\ S \rightarrow \text{print } E & E \rightarrow \text{num} = \text{num} \end{array}$$

GRAMMAR 3.11.

From Modern Compiler Implementation in Java,
Cambridge University Press, ©1998 Andrew W. Appel

```
final int IF=1, THEN=2, ELSE=3, BEGIN=4, END=5, PRINT=6,
SEMI=7, NUM=8, EQ=9;

int tok = getToken();

void advance() {tok = getToken();}
void eat(int t) {if (tok==t) advance(); else error();}

void S() {switch(tok) {
    case IF: eat(IF); E(); eat(THEN); S(); break;
    case ELSE: eat(ELSE); S(); break;
    case BEGIN: eat(BEGIN); S(); L(); break;
    case PRINT: eat(PRINT); E(); break;
    default: error();
}}
void L() {switch(tok) {
    case END: eat(END); break;
    case SEMI: eat(SEMI); S(); L(); break;
    default: error();
}}
void E() {eat(NUM); eat(EQ); eat(NUM);}
```



└ Top-Down Parsing

└ Recursive Descent Parsing

Overview

- Recursive-descent parsing is so powerful and easy to implement that we'll consider the restrictions one would have to put on a grammar to allow one to use the recursive-descent algorithm
- Problems arise when we have alternative productions for the same non-terminal where each production has the same initial symbol, e.g.

$$\begin{aligned} A &\rightarrow xB \mid xC \\ B &\rightarrow xB \mid y \\ C &\rightarrow xC \mid z \end{aligned}$$

- It is possible to evaluate a grammar and determine if it can be parsed using recursive-descent parsing; to do so one must examine the FIRST and FOLLOW sets for grammar



LL(1) Grammars for Predictive Parsing

Rule 1

Given a production

$$A \rightarrow \xi_1 \mid \xi_2 \mid \dots \xi_n$$

the set of initial terminal symbols of all sentences that can be generated from each of the ξ_k 's must be disjoint, that is

$$\text{FIRST}(\xi_j) \cap \text{FIRST}(\xi_k) = \emptyset \text{ for all } j \neq k$$

where the set $\text{FIRST}(\xi_k)$ is the set of all terminal symbols that can appear in the first position of sentences derived from ξ_k . That is

$$a \in \text{FIRST}(\xi_k) \text{ if } \xi_k \rightarrow a\xi$$



LL(1) Grammars for Predictive Parsing

- when a grammar rule A can derive the empty set (i.e.
 $A \rightarrow \epsilon$),

$$\text{predict}(A) = \text{FIRST}(A) \cup \text{FOLLOW}(A)$$

- however, when a grammar rule A does not derive the empty set, $\text{predict}(A) = \text{FIRST}(A)$



LL(1) Grammars for Predictive Parsing

Rule 2

For every non-terminal A which can generate the null string, the set $\text{FIRST}(A)$ of the initial symbols which can be derived by using any alternative productions for A must be disjoint from the set $\text{FOLLOW}(A)$ of symbols that may follow any sequence generated from A , that is

$$\text{FIRST}(A) \cap \text{FOLLOW}(A) = \emptyset$$

where the set $\text{FOLLOW}(A)$ is computed by considering every production P_k of the form

$$P_k \rightarrow \xi_k A \zeta_k$$

and forming the sets $\text{FIRST}(\zeta_k)$, when

$$\text{FOLLOW}(A) = \text{FIRST}(\zeta_1) \cup \text{FIRST}(\zeta_2) \cup \dots \cup \text{FIRST}(\zeta_n)$$

with the addition that if at least one ζ_k is capable of generating the null string, then the set of $\text{FOLLOW}(P_k)$ has to be included in the set $\text{FOLLOW}(A)$ as well.



LL(1) Grammars for Predictive Parsing

$$S \rightarrow A C_1 \$$$

$$C \rightarrow c \mid \epsilon$$

$$A \rightarrow a B C d \mid B Q$$

$$B \rightarrow b B \mid \epsilon$$

$$Q \rightarrow q \mid \epsilon$$

FIRST & FOLLOW sets?

- $\text{FIRST}(C) = \{ c \}; \text{FIRST}(B) = \{ b \}; \text{FIRST}(Q) = \{ q \}$

- $\text{FOLLOW}(A) = \text{FIRST}(C) \cup \text{FOLLOW}(C_1) = \{ c, \$ \}$

- $\text{FIRST}(A) = \{ a \} \cup \text{FIRST}(B Q) = \{ a, b, q \}$

- $\text{FIRST}(S) = \text{FIRST}(A C_1 \$) = \{ a, b, c, q, \$ \}$

- $\text{FOLLOW}(C) = \{ d, \$ \}$

- $\text{FOLLOW}(B) = \text{FIRST}(C) \cup \text{FOLLOW}(C) \cup \text{FIRST}(Q) \cup \text{FOLLOW}(A) = \{ c, d, q, \$ \}$

- $\text{FOLLOW}(Q) = \text{FOLLOW}(A) = \{ c, \$ \}$

Is it LL(1)?



└ Top-Down Parsing

└ Recursive Descent Parsing

LL(1) Grammars for Predictive Parsing

```

1 S → A C $ 
2 C → c 
3 | λ 
4 A → a B C d 
5 | b Q 
6 B → b B 
7 | λ 
8 Q → q 
9 | λ

```

Figure 5.2: A CFGs.

Rule Number	A	$X_1 \dots X_m$	First($X_1 \dots X_m$)	Derives Empty?	Follow(A)	Answer
1	S	A C \$	a,b,q,c,\$	No		a,b,q,c,\$
2	C	c	c	No	c	c
3	λ			Yes	d,\$	d,\$
4	A	a B C d	a	No	a	a
5		B Q	b,q	Yes	c,\$	b,q,c,\$
6	B	b B	b	No	b	b
7	λ			Yes	q,c,d,\$	q,c,d,\$
8	Q	q	q	No	q	q
9	λ			Yes	c,\$	c,\$

Figure 5.3: Predict calculation for the grammar of Figure 5.2.

└ Top-Down Parsing

└ Recursive Descent Parsing

LL(1) Grammars for Predictive Parsing

- Eliminating left recursion:

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

is transformed into

$$A \rightarrow \beta A' \quad A' \rightarrow \alpha A' \mid \epsilon$$

- Note the right recursion is OK

- Generally:

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

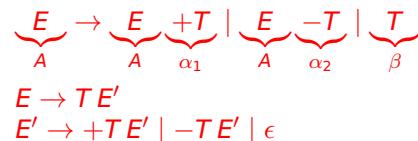
where no β_i begins with an A.

This becomes:

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A' \quad A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$$

$S \rightarrow E \$$	$T \rightarrow T * F$	$F \rightarrow id$
$E \rightarrow E + T$	$T \rightarrow T / F$	$F \rightarrow num$
$E \rightarrow E - T$	$T \rightarrow F$	$F \rightarrow (E)$
$E \rightarrow T$		

GRAMMAR 3.10. From Modern Compiler Implementation in Java, Cambridge University Press, ©1998 Andrew W. Appel



└ Top-Down Parsing

└ Recursive Descent Parsing

LL(1) Grammars for Predictive Parsing

- Problem: left recursion

$$S \rightarrow E \$$$

$$\begin{aligned} E \rightarrow & E + T \\ E \rightarrow & E - T \\ E \rightarrow & T \end{aligned}$$

$$\begin{aligned} T \rightarrow & T * F \\ T \rightarrow & T / F \\ T \rightarrow & F \end{aligned}$$

$$\begin{aligned} F \rightarrow & id \\ F \rightarrow & num \\ F \rightarrow & (E) \end{aligned}$$

GRAMMAR 3.10.

From Modern Compiler Implementation in Java, Cambridge University Press, ©1998 Andrew W. Appel



└ Top-Down Parsing

└ Recursive Descent Parsing

LL(1) Grammars for Predictive Parsing

└ Top-Down Parsing

└ Recursive Descent Parsing

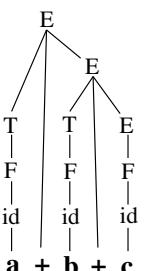
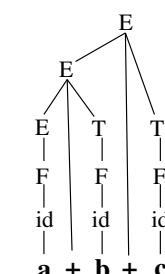
LL(1) Grammars for Predictive Parsing

- left recursive:

$$\begin{aligned} E \rightarrow & T \mid E + T \mid E - T \\ T \rightarrow & F \mid T * F \mid T / F \\ F \rightarrow & id \mid num \mid (E) \end{aligned}$$

- right recursive:

$$\begin{aligned} E \rightarrow & T \mid T + E \mid T - E \\ T \rightarrow & F \mid F * T \mid F / T \\ F \rightarrow & id \mid num \mid (E) \end{aligned}$$



└ Top-Down Parsing

└ Recursive Descent Parsing

LL(1) Grammars for Predictive Parsing

- The problem – the following grammar is not predictive

LL(1):

$$\begin{aligned} S &\rightarrow \text{if } E \text{ then } S \text{ else } S \\ S &\rightarrow \text{if } E \text{ then } S \end{aligned}$$

- Solution – **left factor** the grammar:

$$\begin{aligned} S &\rightarrow \text{if } E \text{ then } S \ X \\ X &\rightarrow \text{else } S \mid \epsilon \end{aligned}$$



└ Top-Down Parsing

└ Recursive Descent Parsing

Writing Recursive-Descent Parsers

- Shell of parsing procedure:

```
procedure S;
begin
    (* assert SYM ∈ FIRST(S) *)
    T(string)
    (* assert SYM ∈ FOLLOW(S) *)
end
```

where the transformation $T(\text{string})$ is governed by the following rules:

- if the production expands to a single terminal:

```
T(terminal) ⇒
  if SYM = terminal
    then get next symbol SYM
  else report error
e.g. Given  $S \rightarrow a$ ,
      if SYM = token_a
      then SYM = getNextToken()
      else reportError();
```



└ Top-Down Parsing

└ Recursive Descent Parsing

Writing Recursive-Descent Parsers²

- Need to write a procedure for each production in the LL grammar; this procedure usually is named the same as the non-terminal
- Procedure analyzes a sequence of symbols defined by the RHS of the production and verifies that they are correct for the given grammar, e.g. given a non-terminal S and a look-ahead buffer of SYM , procedure S will:
 - assume that S has been called with $\text{SYM} \in \text{FIRST}(S)$
 - will then parse a set of symbols that can be derived from production S
 - will relinquish parsing assuming that $\text{SYM} \in \text{FOLLOW}(S)$

² much of this section is from *Programming Language Translation* by Pat Terry, Addison-Wesley, 1986



└ Top-Down Parsing

└ Recursive Descent Parsing

Writing Recursive-Descent Parsers

- if the production expands to a single non-terminal:
 $T(\text{single_non-terminal } < A >) \Rightarrow A$
 e.g. Given $S \rightarrow A$, just call the procedure A
- if the production expands into a number of alternative forms:

```
T( $\alpha_1 \mid \alpha_2 \mid \dots \alpha_n$ ) ⇒
  case SYM of
    FIRST( $\alpha_1$ ):  $T(\alpha_1)$ ;
    FIRST( $\alpha_2$ ):  $T(\alpha_2)$ ;
    ...
    FIRST( $\alpha_n$ ):  $T(\alpha_n)$ ;
    FOLLOW(S): (* do nothing *)
  end;
```

```
e.g. Given  $S \rightarrow aB \mid bC \mid \epsilon$ ,
      case SYM of
        token_a: begin SYM = getNextToken(); B; end;
        token_b: begin SYM = getNextToken(); C; end;
      end;
```



└ Top-Down Parsing

└ Recursive Descent Parsing

Writing Recursive-Descent Parsers

- ④ if the production allows for possible repetition of the RHS:

 $T(\{\alpha\}) \Rightarrow \text{while } \text{SYM} \in \text{FIRST}(\alpha) \text{ do } T(\alpha)$
e.g. Given $S \rightarrow \{aB\}$,

```
while SYM = token_a do
begin
    SYM = getNextToken(); B;
end;
```

- ⑤ if the production generates a sequence of terminals and non-terminals:

 $T(\alpha_1 \alpha_2 \dots \alpha_n) \Rightarrow \text{begin } T(\alpha_1); T(\alpha_2); \dots T(\alpha_n); \text{end}$
e.g. Given $S \rightarrow a \{+ A\} BC$,

```
if SYM = token_a then
begin
    SYM = getNextToken();
    while SYM = token_plus do
        begin SYM = getNextToken(); A; end;
    B; C;
end;
else reportError();
```



└ Top-Down Parsing

└ Recursive Descent Parsing

Example: A Simple Recursive-Descent Parser

```
private static void getSourceString() throws IOException {
    BufferedReader stdin =
        new BufferedReader( new InputStreamReader(System.in) );
    System.out.print("input: ");
    buffer = stdin.readLine()+"$"; // $ represents end-of-stream
}

private static int currentSym = 0;
private static char getNextToken() {
    char ch = buffer.charAt(currentSym++);
    if (debug)
        System.out.println("[SCANNER] "+ch);
    return ch;
}

private static void load_sym() { sym = getNextToken(); }

private static void accept(char terminal) throws Exception {
    if (sym == terminal)
        load_sym();
    else
        throw new Exception("syntax error");
}

private static void entering(String production) {
    if (debug)
        System.out.println("[PARSER] entering " + production);
}
```



└ Top-Down Parsing

└ Recursive Descent Parsing

Example: A Simple Recursive-Descent Parser

- Grammar:

 $A \rightarrow B .$ $B \rightarrow x | (C) | [B] | \epsilon$ $C \rightarrow B D$ $D \rightarrow \{ + B \}$

```
// Author: D. Resler 4/2004
// modified from "Programming Language Translation"
// by Pat Terry (Addison- Wesley, 1986)
```

```
import java.io.*;

public class RecursiveDescent {

    private static char sym;
    private static String buffer;
    private static boolean debug = true;

    public static void main(String [] args) throws Exception {
        getSourceString();
        sym = getNextToken(); A();
        System.out.println("Successful parse");
    }
}
```



└ Top-Down Parsing

└ Recursive Descent Parsing

Example: A Simple Recursive-Descent Parser

```
private static void exiting(String production) {
    if (debug)
        System.out.println("[PARSER] exiting " + production);
}

private static void A() throws Exception {
    entering("A"); B(); accept('.'); exiting("A");
}

private static void B() throws Exception {
    entering("B");
    switch (sym) {
        case 'x': load_sym(); break;
        case '(': load_sym(); C(); accept(')'); break;
        case '[': load_sym(); B(); accept(']'); break;
        case ')': case '+': case '.':
            /* followers of B - do nothing */
            break;
        default: throw new Exception("syntax error");
    }
    exiting("B");
}

private static void C() throws Exception {
    entering("C");
    B(); D();
    exiting("C");
}
```

$A \rightarrow B .$
$B \rightarrow x (C) [B] \epsilon$
$C \rightarrow B D$
$D \rightarrow \{ + B \}$

```
private static void D() throws Exception {
    entering("D");
    while (sym == '+') {
        load_sym(); B();
    }
    exiting("D");
}
```



Example: A Simple Recursive-Descent Parser

```

input: (x+x).
[SCANNER] (
[PARSER] entering A
[PARSER] entering B
[SCANNER] x ----- B
[PARSER] entering C
[PARSER] entering B
[SCANNER] + ----- B
[PARSER] exiting B
[PARSER] entering D
[SCANNER] x ----- B
[PARSER] entering B
[SCANNER] ) ----- B
[PARSER] exiting B
[PARSER] exiting D
[PARSER] exiting C
[SCANNER] . ----- B
[PARSER] exiting B
[SCANNER] $ ----- B
[PARSER] exiting A
Successful parse
  
```



Error Recovery in Recursive-Descent Parsers

- To avoid skipping too many symbols, the parser can pass a set FOLLOWERS to procedure A that contains all of the symbols you can skip to upon an error
- FOLLOWERS is constructed to include:
 - the set FOLLOW(A)
 - symbols that have been passed as FOLLOWERS to the procedure(s) that called A
 - so-called *beacon* symbols that should never be skipped



Error Recovery in Recursive-Descent Parsers

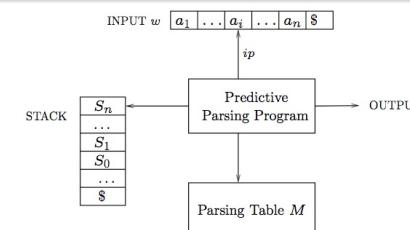
- Remember that when a procedure A that recognizes language phrase $A \rightarrow \alpha$ is called, the lookahead buffer should contain a symbol from $\text{FIRST}(\alpha)$
- Also, when you leave procedure A, the lookahead buffer should contain a symbol from $\text{FOLLOW}(A)$
- Possible strategy:


```

procedure A;
begin
  if not lookahead in FIRST(A)
    then report error and SKIPT0(FIRST(A) ∪ FOLLOW(A))
  if lookahead in FIRST(A)
    then parse phrase A
    if not lookahead in FOLLOW(A)
      then report an error and SKIPT0(FOLLOW(A))
end;
      
```
- This often skips too many symbols, however (why?)



Table-driven Predictive Parsing



(a = current input symbol, ip = instruction pointer)

```

set ip to the first symbol in w; push start non-terminal onto the stack
set X to the top of stack symbol
while (X ≠ $) /* stack is not empty */
  if (X is a) pop the stack and advance ip
  else if (X is a terminal) error()
  else if (M[X, a] is an error entry) error()
  else if (M[X, a] = X → Y1Y2...Yk) {
    output the production X → Y1Y2...Yk
    pop the stack
    push Yk, Yk-1, ..., Y1 onto the stack, with Y1 on top
  }
  set X to the top of stack symbol
}
  
```



Table-driven Predictive Parsing

- Grammar:

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

```

set ip to the first symbol in w;
push start non-terminal onto the stack
set X to the top of stack symbol
while ( X ≠ $ ) /* stack is not empty */
  if ( X is a ) pop the stack and advance ip
  else if ( X is a terminal ) error()
  else if ( M[X, a] is an error entry ) error()
  else if ( M[X, a] = X → Y1Y2...Yk )
    output the production X → Y1Y2...Yk
    pop the stack
    push Yk,Yk-1,...,Y1 onto the stack, with Y1 on top
  }
  set X to the top of stack symbol
}

```

- Parse table:

NON-TERMINAL	INPUT SYMBOL				
	id	+	*	()
E	E → TE'			E → TE'	
E'		E' → +TE'			E' → ε
T	T → FT'			T → FT'	
T'		T' → ε	T' → *FT'		T' → ε
F	F → id			F → (E)	



Error Recovery in Table-driven Predictive Parsers³

- A heuristic is used to determine a set of synchronizing tokens for productions
- New parse table for our example (where sync = set of sync tokens):

NON-TERMINAL	INPUT SYMBOL				
	id	+	*	()
E	E → TE'			E → TE'	sync
E'		E' → +TE'			E' → ε
T	T → FT'	sync		T → FT'	sync
T'		T' → ε	T' → *FT'		T' → ε
F	F → id	sync	sync	F → (E)	sync

- When the parser looks up table entry M[A, a]:
 - if it is blank: input symbol a is skipped
 - if it is sync: nonterminals on top of stack are popped until token on top of stack is in sync set; then if the token on the top of the stack does not match a, pop the token from the stack as before
- Another popular method: call error handling routine whenever error entry reached in table



Table-driven Predictive Parsing

- Parse of 'id + id * id':

MATCHED	STACK	INPUT	ACTION
	E\$	id + id * id\$	output E → TE'
	TE'\$	id + id * id\$	output T → FT'
	FT'E'\$	id + id * id\$	output F → id
id	idT'E'\$	+ id * id\$	match id
id	E'\$	+ id * id\$	output T' → ε
id	+ TE'\$	+ id * id\$	output E' → + TE'
id +	TE'\$	id * id\$	match +
id +	FT'E'\$	id * id\$	output T → FT'
id +	idT'E'\$	id * id\$	output F → id
id + id	T'E'\$	* id\$	match id
id + id	* FT'E'\$	id\$	output T' → * FT'
id + id *	FT'E'\$	id\$	match *
id + id *	idT'E'\$	id\$	output F → id
id + id * id	T'E'\$	\$	match id
id + id * id	E'\$	\$	output T' → ε
id + id * id	\$	\$	output E' → ε



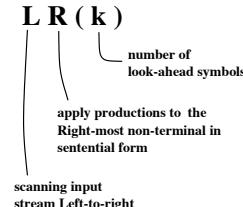
Properties of LL(1) Parsers

- a correct, left-most parse is constructed
- all grammars in the LL(1) class are unambiguous
- all table-driven LL(1) parsers operate in linear time and space with respect to the length of the parsed input



LR Parsing

- The fundamental concern of LL parsers is which production to choose in expanding a given nonterminal
- Weakness of LL(k) parsers: must predict which production to use having seen only k tokens of the right-hand-side (RHS)
- More powerful technique: LR(k) bottom-up parsing



- Look-ahead k assumed to be 1 if not explicitly stated



LR Parsing

- Also called **shift-reduce** parsing
- The parser has a **stack** and an **input**. The first k tokens of the input are the **look-ahead**. Based on the contents of the stack and the look-ahead, the parser performs one of 2 possible actions:

Shift push the first input token onto the stack

Reduce choose a grammar rule $X \rightarrow \alpha\beta\gamma$; pop γ , β , α from the stack; push X onto the stack

- Initially: stack is empty, parser is at beginning of input
- Parser accepts string when end-of-file marker (\$) is shifted and there is no more input



LR Parsing

- LR parsers:
 - ▶ begin with the parse tree's leaves and move towards its root
 - ▶ trace a rightmost derivation (in reverse)
 - ▶ uses grammar rules to replace a rule's right-hand side with its left hand side
- LR parsers can be constructed to recognize virtually all programming language constructs
- The class of grammars that can be parsed using LR parsing is a proper superset of the class of grammars that can be parsed with LL predictive parsing
- Principle drawback: difficult to construct LR parsers by hand
- Therefore usually done using a LR parser generator (a 'compiler-compiler')



LR Parsing

Stack	Input	Action
1 \$	a := 7 ; b := c + (d := 5 + 6 , d) \$	shift
1 id4	a := 7 ; b := c + (d := 5 + 6 , d) \$	shift
1 id4 :=6	a := 7 ; b := c + (d := 5 + 6 , d) \$	shift
1 id4 :=6 num10	a := 7 ; b := c + (d := 5 + 6 , d) \$	reduce E → num
1 id4 :=6 E11	a := 7 ; b := c + (d := 5 + 6 , d) \$	reduce S → id := E
1 S2	a := 7 ; b := c + (d := 5 + 6 , d) \$	shift
1 S2 ;	a := 7 ; b := c + (d := 5 + 6 , d) \$	shift
1 S2 ; id4	a := 7 ; b := c + (d := 5 + 6 , d) \$	shift
1 S2 ; id4 :=6	a := 7 ; b := c + (d := 5 + 6 , d) \$	shift
1 S2 ; id4 :=6 id20	a := 7 ; b := c + (d := 5 + 6 , d) \$	shift
1 S2 ; id4 :=6 E11	+ (d := 5 + 6 , d) \$	shift
1 S2 ; id4 :=6 E11 +16	+ (d := 5 + 6 , d) \$	shift
1 S2 ; id4 :=6 E11 +16 (s	d := 5 + 6 , d) \$	shift
1 S2 ; id4 :=6 E11 +16 (s id4	= 5 + 6 , d) \$	shift
1 S2 ; id4 :=6 E11 +16 (s id4 :=6	5 + 6 , d) \$	shift
1 S2 ; id4 :=6 E11 +16 (s id4 :=6 num10	+ 6 , d) \$	reduce E → num
1 S2 ; id4 :=6 E11 +16 (s id4 :=6 E11	+ 6 , d) \$	shift
1 S2 ; id4 :=6 E11 +16 (s id4 :=6 E11 +16	6 , d) \$	shift
1 S2 ; id4 :=6 E11 +16 (s id4 :=6 E11 +16 num10	, d) \$	reduce E → num
1 S2 ; id4 :=6 E11 +16 (s id4 :=6 E11 +16 E17	, d) \$	reduce E → E + E
1 S2 ; id4 :=6 E11 +16 (s id4 :=6 E11	, d) \$	reduce S → id := E
1 S2 ; id4 :=6 E11 +16 (s S12	, d) \$	shift
1 S2 ; id4 :=6 E11 +16 (s S12 +18 id4	, d) \$	reduce E → id
1 S2 ; id4 :=6 E11 +16 (s S12 +18 E21	, d) \$	shift
1 S2 ; id4 :=6 E11 +16 (s S12 +18 E21)22	, d) \$	reduce E → (S ; E)
1 S2 ; id4 :=6 E11 +16 E17	, d) \$	reduce E → E + E
1 S2 ; id4 :=6 E11	, d) \$	reduce S → id := E
1 S2 ; S3	, d) \$	reduce S → S ; S
1 S2	,	accept

FIGURE 3.18. Shift-reduce parse of a sentence. Numeric subscripts in the Stack are DFA state numbers; see Table 3.19.
From *Modern Compiler Implementation in Java*, Cambridge University Press, ©1998 Andrew W. Appel



└ Bottom-Up Parsing

└ Overview

LR Parsing Engine

- How does the LR parser know when to shift and when to reduce? **By using a DFA!**
- DFA applied to stack, not input
- Edges of the DFA are labeled by the symbols (terminals and nonterminals) that can appear on the stack



└ Bottom-Up Parsing

└ Overview

LR Parsing Algorithm

Step 1 Look up action in table using top stack state and input symbol

Step 2

If action is:	Do:
Shift(n)	Push n onto stack, advance input one token
Reduce(k)	1. Pop x symbols off stack, where $x = \text{number of symbols on the RHS of rule } k$ 2. Let X be the LHS symbol of rule k ; use state currently on top of stack and X to index into table to get 'goto n' 3. Push n onto stack
Accept	Stop parsing, report success
Error	Stop parsing, report failure

Step 3 Goto **Step 1**



└ Bottom-Up Parsing

└ Overview

LR Parsing Engine

	id	num	print	:	*	+	:=	()	S	S	E	L
1	s4		s7		s3					a			
2											g5		
3	s4		s7										
4													
5													
6	s20		s10							r1	r1		
7													
8													
9	s4		s7										
10													
11													
12													
13													
14													
15													
16	s20		s10										
17													
18													
19	s20		s10										
20	s20		s10							r4	r4		
21													
22													
23													

TABLE 3.19. LR parsing table for Grammar 3.1.
From *Modern Compiler Implementation in Java*, Cambridge University Press, ©1998 Andrew W. Appel

i $S \rightarrow S ; S$
 s $E \rightarrow \text{id}$
 z $E \rightarrow \text{id} := E$
 o $E \rightarrow E + E$
 r $L \rightarrow L , E$

GRAMMAR 3.1. A syntax for straight-line programs.
From *Modern Compiler Implementation in Java*, Cambridge University Press, ©1998 Andrew W. Appel

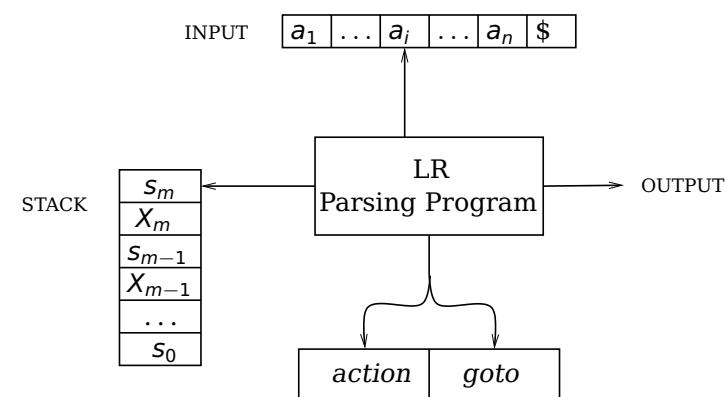
sn Shift into state n
 gn Goto state n
 rk Reduce by rule k
 a Accept
 Error (a blank entry in table)



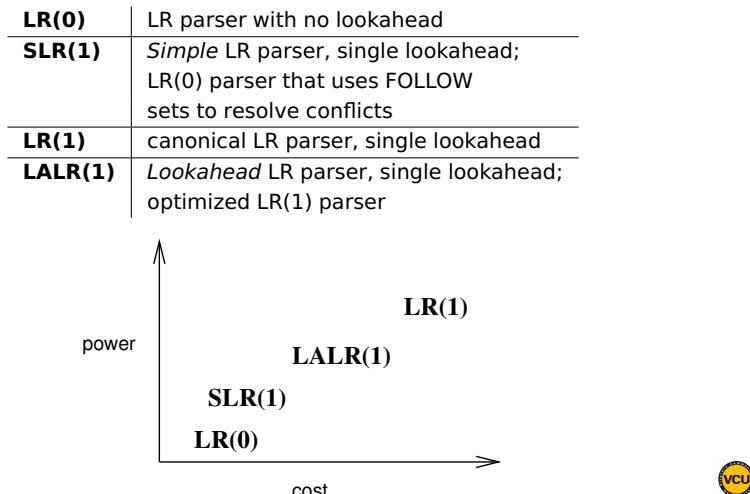
└ Bottom-Up Parsing

└ Overview

Model of a LR Parser



The LR Parsing Family



LR Parse Tables – The Big Picture

0: $s \rightarrow e$ 4: $t \rightarrow f$
 1: $e \rightarrow e + t$ 5: $f \rightarrow (e)$
 2: $e \rightarrow t$ 6: $f \rightarrow \text{NUM}$
 3: $t \rightarrow t * f$

	\$	NUM	+	*	()	s	e	t	f
0	-	s1	-	-	s2	-	-	g3	g4	g5
1	r6	-	r6	r6	-	r6	-	-	-	-
2	-	s1	-	-	s2	-	-	g6	g4	g5
3	acc	-	s7	-	-	-	-	-	-	-
4	r2	-	r2	s8	-	r2	-	-	-	-
5	r4	-	r4	r4	-	r4	-	-	-	-
6	-	-	s7	-	-	s9	-	-	-	-
7	-	s1	-	-	s2	-	-	-	-	g10
8	-	s1	-	-	s2	-	-	-	-	g11
9	r5	-	r5	r5	-	r5	-	-	-	-
10	r1	-	r1	s8	-	r1	-	-	-	-
11	r3	-	r3	-	r3	-	-	-	-	-

	Stack		Input	Next Action
	States	Symbols		
(1)	0		1 * (2 + 3) \$	shift NUM
(2)	0 1	NUM	* (2 + 3) \$	reduce by 6
(3)	0	f	* (2 + 3) \$	goto 5
(4)	0 5		* (2 + 3) \$	reduce 4
(5)	0	t	* (2 + 3) \$	goto 4
(6)	0 4		* (2 + 3) \$	shift *
(7)	0 4 8	t*	(2 + 3) \$	shift (
(8)	0 4 8 2	t*(2 + 3) \$	shift NUM
(9)	0 4 8 2 1	t*(NUM	+ 3) \$	reduce 6
(10)	0 4 8 2	t*(f	+ 3) \$	goto 5
(11)	0 4 8 2 5	t*(f	+ 3) \$	reduce 4
(12)	0 4 8 2	t*(t	+ 3) \$	goto 4

LR Parse Tables – The Big Picture

- Example grammar L_1 & parse table:

0: $s \rightarrow e$ 4: $t \rightarrow f$
 1: $e \rightarrow e + t$ 5: $f \rightarrow (e)$
 2: $e \rightarrow t$ 6: $f \rightarrow \text{NUM}$
 3: $t \rightarrow t * f$

	\$	NUM	+	*	()	s	e	t	f
0	-	s1	-	-	s2	-	-	g3	g4	g5
1	r6	-	r6	r6	-	r6	-	-	-	-
2	-	s1	-	-	s2	-	-	g6	g4	g5
3	acc	-	s7	-	-	-	-	-	-	-
4	r2	-	r2	s8	-	r2	-	-	-	-
5	r4	-	r4	r4	-	r4	-	-	-	-
6	-	-	s7	-	-	s9	-	-	-	-
7	-	s1	-	-	s2	-	-	-	-	g10
8	-	s1	-	-	s2	-	-	-	-	g11
9	r5	-	r5	r5	-	r5	-	-	-	-
10	r1	-	r1	s8	-	r1	-	-	-	-
11	r3	-	r3	-	r3	-	-	-	-	-

LR Parse Tables – The Big Picture

0: $s \rightarrow e$ 4: $t \rightarrow f$
 1: $e \rightarrow e + t$ 5: $f \rightarrow (e)$
 2: $e \rightarrow t$ 6: $f \rightarrow \text{NUM}$
 3: $t \rightarrow t * f$

	\$	NUM	+	*	()	s	e	t	f
0	-	s1	-	-	s2	-	-	g3	g4	g5
1	r6	-	r6	r6	-	r6	-	-	-	-
2	-	s1	-	-	s2	-	-	g6	g4	g5
3	acc	-	s7	-	-	-	-	-	-	-
4	r2	-	r2	s8	-	r2	-	-	-	-
5	r4	-	r4	r4	-	r4	-	-	-	-
6	-	-	s7	-	-	s9	-	-	-	-
7	-	s1	-	-	s2	-	-	-	-	g10
8	-	s1	-	-	s2	-	-	-	-	g11
9	r5	-	r5	r5	-	r5	-	-	-	-
10	r1	-	r1	s8	-	r1	-	-	-	-
11	r3	-	r3	-	r3	-	-	-	-	-

	Stack	Symbols	Input	Next Action
	States			
(12)	0 4 8 2	t*(t	+ 3) \$	goto 4
(13)	0 4 8 2 4	t*(t	+ 3) \$	reduce 2
(14)	0 4 8 2	t*(e	+ 3) \$	goto 6
(15)	0 4 8 2 6	t*(e	+ 3) \$	shift +
(16)	0 4 8 2 6 7	t*(e +	3) \$	shift NUM
(17)	0 4 8 2 6 7 1	t*(e + NUM) \$	reduce 6
(18)	0 4 8 2 6 7	t*(e + f) \$	goto 5
(19)	0 4 8 2 6 7 5	t*(e + f) \$	reduce 4
(20)	0 4 8 2 6 7	t*(e + t) \$	goto 10
(21)	0 4 8 2 6 7 10	t*(e + t) \$	reduce 1
(22)	0 4 8 2	t*(e) \$	goto 6
(23)	0 4 8 2 6	t*(e) \$	shift)

LR Parse Tables – The Big Picture

0: $s \rightarrow e$ 4: $t \rightarrow f$
 1: $e \rightarrow e + t$ 5: $f \rightarrow (e)$
 2: $e \rightarrow t$ 6: $f \rightarrow \text{NUM}$
 3: $t \rightarrow t * f$

	\$	NUM	+	*	()	s	e	t	f
0	-	s1	-	-	s2	-	-	g3	g4	g5
1	r6	-	r6	r6	-	r6	-	-	-	-
2	-	s1	-	-	s2	-	-	g6	g4	g5
3	acc	-	s7	-	-	-	-	-	-	-
4	r2	-	r2	s8	-	r2	-	-	-	-
5	r4	-	r4	r4	-	r4	-	-	-	-
6	-	-	s7	-	-	s9	-	-	-	-
7	-	s1	-	-	s2	-	-	-	g10	g5
8	-	s1	-	-	s2	-	-	-	-	g11
9	r5	-	r5	r5	-	r5	-	-	-	-
10	r1	-	r1	s8	-	r1	-	-	-	-
11	r3	-	r3	r3	-	r3	-	-	-	-

	Stack			Input	Next Action
	States	Symbols	Input		
(23)	0 4 8 2 6	t * (e) \$	shift)	
(24)	0 4 8 2 6 9	t * (e)	\$	reduce 5	
(25)	0 4 8	t * f	\$	goto 11	
(26)	0 4 8 11	t * f	\$	reduce 3	
(27)	0	t	\$	goto 4	
(28)	0	e	\$	reduce 2	
(29)	0	e	\$	goto 3	
(30)	0 3	e	\$	accept	



Creating LR Parse Tables – The Big Picture

- Relating the table to the graph:

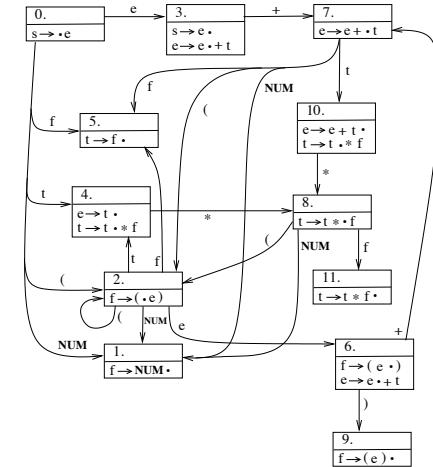
- Each box is a state with the state number being the row in the parse table
- Edges are labeled with either a terminal or nonterminal (never both)
- Edges labeled with a terminal relate to shift directives in the table
- Use the goto portion of the table to determine the next state when traversing edges labeled with a nonterminal

- Productions at the bottom of each state are called **LR(0) items**; LR(0) items are productions augmented by a dot
- Dots are used to indicate the present position of the parse, i.e. what is on the parse stack



Creating LR Parse Tables – The Big Picture

- LR State-Machine S_1 for language L_1 :



Creating LR Parse Tables – The Big Picture

- A dot all the way to the right of a production indicates the parser can reduce using that production
- Shifts cause the dot to move from the left of the symbol labeling the traversed edge to the right of that symbol
- Reductions first require moving backwards (against the arrows) one edge for every symbol on the right hand side of the production followed by moving to the indicated state



Creating LR Parse Tables

- How do we create an LR parse table?
- Main steps:
 - ① create LR(0) item sets
 - ② generate the table from items



Creating LR(0) Items

Step 1 Initialization

- Create state 0 with the start production with a **dot** at far left of the RHS
- Initial set of items in any state is called the **kernel** or **seed** items

0
s → ·e

0: s → e 4: t → f
1: e → e + t 5: f → (e)
2: e → t 6: f → NUM
3: t → t * f



LR(0) Parsing

- LR(0) grammars are those that can be parsed looking only at the stack
- Too weak of a parsing method to be of much use
- However, generating an LR(0) parse table is the usual starting point for many LR parsers
- Producing an LR(0) table involves first building LR(0) item sets



Creating LR(0) Items

Step 2 Close the kernel items

- Kernel items for next states are created with the **closure** operation
- Closure is performed on kernel items that have a nonterminal to the right of a dot
- Begin the closure process by adding all productions of the nonterminal that follows the dot

0
s → ·e
e → ·e + t
e → ·t

0: s → e 4: t → f
1: e → e + t 5: f → (e)
2: e → t 6: f → NUM
3: t → t * f



Creating LR(0) Items

- Closing state 0:
 - ▶ Since the dot is to the left of e , add the e productions

0
$s \rightarrow \cdot e$
$e \rightarrow \cdot e + t$
$e \rightarrow \cdot t$

0: $s \rightarrow e$	4: $t \rightarrow f$
1: $e \rightarrow e + t$	5: $f \rightarrow (e)$
2: $e \rightarrow t$	6: $f \rightarrow \text{NUM}$
3: $t \rightarrow t * f$	

- ▶ Note we have new items with a dot to the left of a nonterminal, so we need to add the t productions

0
$s \rightarrow \cdot e$
$e \rightarrow \cdot e + t$
$e \rightarrow \cdot t$
$t \rightarrow \cdot t * f$
$t \rightarrow \cdot f$



Creating LR(0) Items

Step 3 Partition the kernel items

- Group together all items that have the same symbol (called the **transition symbol**) to the right of the dot
- All items with the same transition symbol form a **partition**

0
$s \rightarrow \cdot e$
$e \rightarrow \cdot e + t$
$e \rightarrow \cdot t$
$t \rightarrow \cdot t * f$
$t \rightarrow \cdot f$



Creating LR(0) Items

- Closing state 0 *continued*:

- ▶ Repeat this process until no more items can be added
- ▶ Result:

0
$s \rightarrow \cdot e$
$e \rightarrow \cdot e + t$
$e \rightarrow \cdot t$
$t \rightarrow \cdot t * f$
$t \rightarrow \cdot f$

0: $s \rightarrow e$	4: $t \rightarrow f$
1: $e \rightarrow e + t$	5: $f \rightarrow (e)$
2: $e \rightarrow t$	6: $f \rightarrow \text{NUM}$
3: $t \rightarrow t * f$	



Creating LR(0) Items

Step 4 Form new states and next-state transitions as necessary

- Partitions form new states *except* those with dots at the far right (all ϵ productions are part of this latter partition)
- Modify productions in the potentially new state by moving the dot one symbol to the right
- If a state already exists that has the same set of kernel items as the modified partition, add an edge labeled with the transition symbol from the current state to the existing state

0
$s \rightarrow \cdot e$
$e \rightarrow \cdot e + t$
$e \rightarrow \cdot t$



Creating LR(0) Items

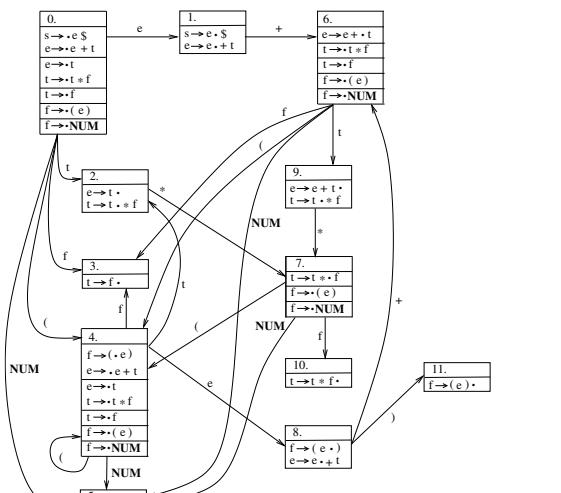
Step 4 continued

- The transition symbol is the symbol you moved the dot over
- If there is no existing state with the same kernel items as a given partition, create a new state with an edge from the current state labeled with the transition symbol

Step 5 Repeat steps 2–4 until no more new states can be added to the diagram



Creating LR(0) Items



Creating LR(0) Items

0. $s \rightarrow \cdot e \$$ $e \rightarrow \cdot e + t$ $e \rightarrow \cdot t$ $t \rightarrow \cdot t * f$ $t \rightarrow \cdot f$ $f \rightarrow \cdot (e)$ $f \rightarrow \cdot \text{NUM}$	2. $e \rightarrow t \cdot$ $t \rightarrow t \cdot * f$	6. $e \rightarrow e \cdot + t$ $t \rightarrow t \cdot f$ $t \rightarrow \cdot f$ $f \rightarrow \cdot (e)$ $f \rightarrow \cdot \text{NUM}$	9. $e \rightarrow e \cdot + t \cdot$ $t \rightarrow t \cdot * f$
	3. $t \rightarrow f \cdot$		10. $t \rightarrow t \cdot f \cdot$
	4. $f \rightarrow (e \cdot)$ $e \rightarrow e \cdot + t$		11. $f \rightarrow (e \cdot)$
1. $s \rightarrow e \cdot \$$ $e \rightarrow e \cdot + t$			
		8. $f \rightarrow (e \cdot)$ $e \rightarrow e \cdot + t$	
		5. $f \rightarrow \cdot \text{NUM} \cdot$	

- Item set numbering is irrelevant

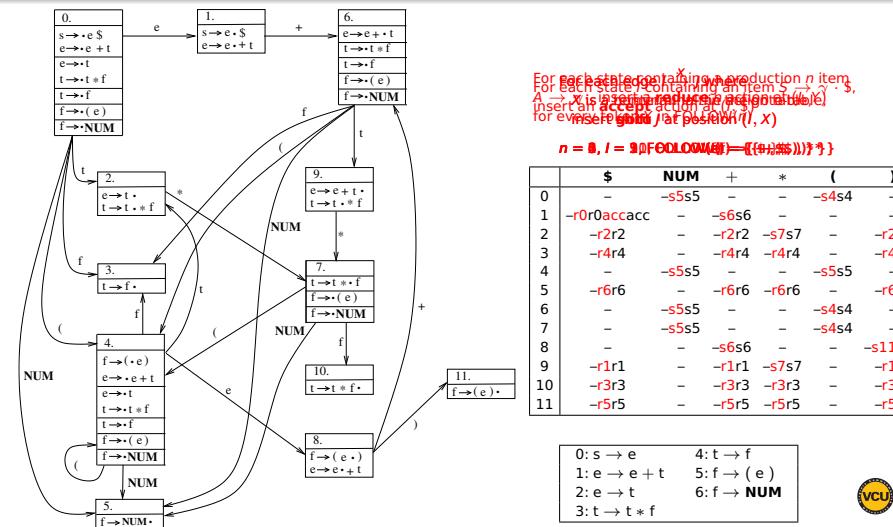


Creating LR Parse Tables from Item Sets

- For each edge $I \xrightarrow{X} J$ where
 - X is a terminal: in the action table, insert **shift** J at position (I, X)
 - X is a nonterminal: in the goto table, insert **goto** J at position (I, X)
- For each state containing a production n item $A \rightarrow \gamma \cdot$, insert a **reduce** n action at (I, Y) for every token Y in $\text{FOLLOW}(n)$
- For each state I containing an item $S \rightarrow \gamma \cdot \$$, insert an **accept** action at $(I, \$)$

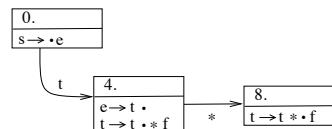


Creating LR Parse Tables from Item Sets



Problems with LR(0) State Machines

- Consider state 4 in S_1 , reproduced below:



- If you are currently in state 4, should you shift or reduce?
- Such situations are called **shift/reduce conflicts**
- Can also have **reduce/reduce conflicts**—two or more productions in the same state with dots all the way to the right
- States that contain these conflicts are called **inadequate** states
- A language is considered LR(0) if its state machine has no inadequate states
- Inadequate states produce multiple entries in an LR(0) parse table

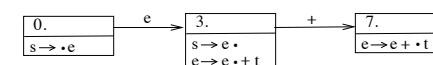


Creating LR Parse Tables from Item Sets

	\$	NUM	+	*	()	s	e	t	f
0	-	$\rightarrow s5$	-	-	$\rightarrow s4$	-	-	$g1$	$g2$	$g3$
1	acc	-	$\rightarrow s6$	-	-	-	-	-	-	-
2	$r2$	-	$\rightarrow r2$	$\rightarrow s7$	-	$\rightarrow r2$	-	-	-	-
3	$r4$	-	$\rightarrow r4$	$\rightarrow r4$	-	$\rightarrow r4$	-	-	-	-
4	-	$\rightarrow s5$	-	-	$\rightarrow s5$	-	-	$g8$	$g2$	$g3$
5	$r6$	-	$\rightarrow r6$	$\rightarrow r6$	-	$\rightarrow r6$	-	-	-	-
6	-	$\rightarrow s5$	-	-	$\rightarrow s4$	-	-	$g9$	$g3$	-
7	-	$\rightarrow s5$	-	-	$\rightarrow s4$	-	-	-	-	$g10$
8	-	-	$\rightarrow s6$	-	-	$\rightarrow s11$	-	-	-	-
9	$\rightarrow r1$	-	$\rightarrow r1$	$\rightarrow s7$	-	$\rightarrow r1$	-	-	-	-
10	$\rightarrow r3$	-	$\rightarrow r3$	$\rightarrow r3$	-	$\rightarrow r3$	-	-	-	-
11	$\rightarrow r5$	-	$\rightarrow r5$	$\rightarrow r5$	-	$\rightarrow r5$	-	-	-	-

SLR(1) Grammars

- Can solve most LR(0) shift/reduce errors by considering the next token in the input stream
- In an LR(0) parser we always reduce when a state contains a production in the form of $S \rightarrow \gamma \cdot$; in an SLR parser we reduce for states that contain $S \rightarrow \gamma \cdot$ and the lookahead symbol is in $FOLLOW(S)$
- Consider state 3 in S_1 , reproduced below:



- Reduce when lookahead is in $FOLLOW(S)$, shift otherwise
- Works if $+$ is not a member of $FOLLOW(S)$
- If all LR(0) conflicts can be resolved using this scheme, language is said to be SLR(1)



└ Bottom-Up Parsing

└ LR Parsing

LR(1) Grammars

- Many grammars are not SLR—a FOLLOW set may contain symbols that also label an outgoing edge
- Note that a nonterminal's FOLLOW set contains *all* symbols that can follow that nonterminal in *every* context; however, we only care about the follow symbols for the current state
- This set of relevant lookahead symbols is called the **lookahead set**; typically it is a subset of the complete FOLLOW set
- An **LR(1) item** is an LR(0) item augmented with a lookahead symbol from this lookahead set
- The lookahead symbol is actually part of an item—two items with the same production and dot position but with different lookahead symbols are different items



└ Bottom-Up Parsing

└ LR Parsing

LALR(1) Grammars

- Problem with LR(1) parsers: the state machine (and therefore the parse table) is typically twice the size of an equivalent LR(0) machine
- Many states in LR(1) machines differ only in their lookahead symbols; LALR parsers merge many of these states
- Though uncommon, these merges sometimes produce inadequate states
- LALR(1) parsers have the same number of states as LR(0)/SLR parsers with almost the power of LR(1) parsers
- Almost all modern programming languages are LALR(1)
- Most compiler generators produce LALR parsers



└ Bottom-Up Parsing

└ LR Parsing

Creating LR(1) Items

- LR(1) item $[s \rightarrow \alpha \cdot \beta, D]$ consists of a production, a dot, and a lookahead D
- Process same as creating LR(0) items except LR(1) items created during closure
- Initial item for L_1 : $[s \rightarrow \cdot e, \$]$
- Given the following LR(1) kernel item and production,

$$\begin{aligned} s &\rightarrow \alpha \cdot x\beta, C \\ x &\rightarrow \cdot \gamma \end{aligned}$$

create the new LR(1) item using

$$[x \rightarrow \cdot \gamma, \text{FIRST}(\beta, C)]$$

where α and β can be null, and

$$\text{FIRST}(\beta, C) = \begin{cases} \text{FIRST}(\beta) \cup C, & \text{if } \beta == \text{null} \\ \text{FIRST}(\beta), & \text{otherwise} \end{cases}$$

- Continue in the same manner as with LR(0) until no more LR(1) items can be created



└ Bottom-Up Parsing

└ LR Parsing

Hierarchy of Grammar Classes

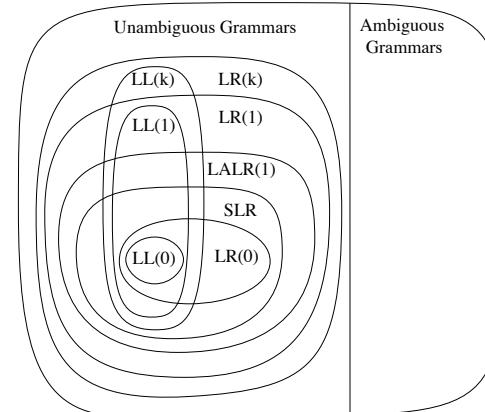


FIGURE 3.30.

A hierarchy of grammar classes.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel



Error-recovery in LR Parsers

- Errors detected when lookup in parsing table produces error entry
- General error-recovery scheme:
 - scan down stack until a state s with a goto on a particular nonterminal A is found (might be more than one choice for A)
 - zero or more input symbols are then discarded until a symbol a is found that can legitimately follow A
 - parser then stacks the state $\text{GOTO}(s, A)$ and resumes normal parsing
- This method of recovery attempts to eliminate the phrase containing the syntax error
- As in table-drive predicative parsing, can also call special error handling routines when error entry encountered



An Example

- Consider again the context-sensitive language L :

$$L = \{a^n b^n c^n \mid n \geq 1\}$$

- No context-free grammar generates L !

- An attempt:

```
string → a_seq b_seq c_seq
a_seq → a | a_seq a
b_seq → b | b_seq b
c_seq → c | c_seq c
```

- This context-free grammar generates the language L' :

$$L' = a^+ b^+ c^+ = \{a^k b^m c^n \mid k \geq 1, m \geq 1, n \geq 1\}$$



Definitions

- An attribute grammar is a context-free grammar that has been extended to provide context-sensitive information by appending attributes to some of its nonterminals
- Each distinct symbol in the grammar has associated with it a finite, possibly empty, set of attributes
 - Each attribute has a domain of possible values
 - An attribute may be assigned values from its domain during parsing
 - Attributes can be evaluated in assignments or conditions
- Two Classes of Attributes:

Synthesized attributes: An attribute that gets its values from the attributes attached to children of its nonterminal

Inherited attributes: An attribute that gets its values from the attributes attached to the parent (or siblings) of its nonterminal



Using an Attribute Grammar

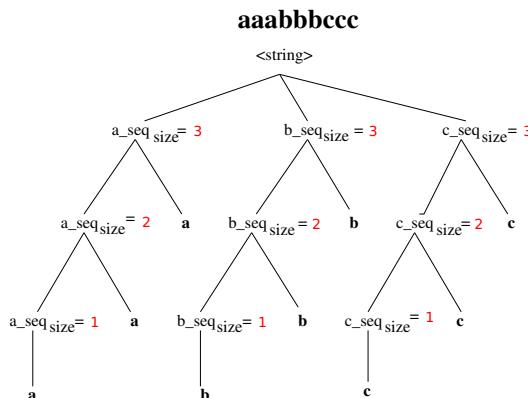
- Attach an integer synthesized attribute, size , to each of the nonterminals a_seq , b_seq , c_seq

- Impose a condition on the first production:

```
string → a_seqsize b_seqsize c_seqsize
  condition
    a_seqsize = b_seqsize = c_seqsize
  a_seqsize → a
    a_seqsize ← 1
    | a_seq2size a
      a_seq1size ← a_seq2size + 1
  b_seqsize → b
    b_seqsize ← 1
    | b_seq2size b
      b_seq1size ← a_seq2size + 1
  c_seqsize → c
    c_seqsize ← 1
    | c_seq2size c
      c_seq1size ← c_seq2size + 1
```



Decorating a Parse Tree



string → a_seq_size b_seq_size c_seq_size
condition: a_seq_size = b_seq_size = c_seq_size
 a_seq_size → a
 a_seq1_size ← 1
 | a_seq2_size a
 a_seq1_size ← a_seq2_size + 1
 b_seq_size → b
 b_seq1_size ← 1
 | b_seq2_size b
 b_seq1_size ← b_seq2_size + 1
 c_seq_size → c
 c_seq1_size ← 1
 | c_seq2_size c
 c_seq1_size ← c_seq2_size + 1



Using an Inherited Attribute

- Attach a synthesized attribute `size` to a_seq and and inherited attributes `inSize` to b_seq and c_seq

string → a_seq_size b_seq_inSize c_seq_inSize
 b_seq_inSize ← a_seq_size
 c_seq_inSize ← a_seq_size

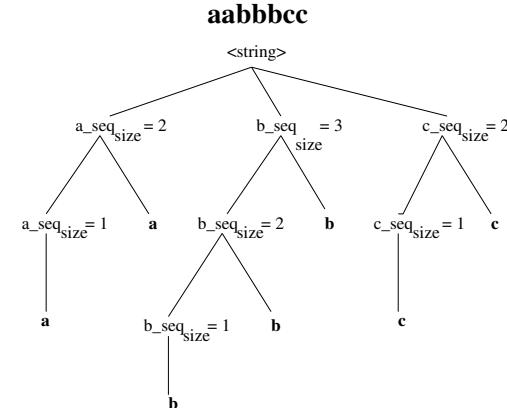
a_seq_size → a
 a_seq_size ← 1
 | a_seq2_size a
 a_seq_size ← a_seq2_size + 1

b_seq_inSize → b
condition: b_seq_inSize = 1
 | b_seq2_inSize b
 b_seq2_inSize ← b_seq_inSize - 1

c_seq_inSize → c
condition: c_seq_inSize = 1
 | c_seq2_inSize c
 c_seq2_inSize ← c_seq_inSize - 1



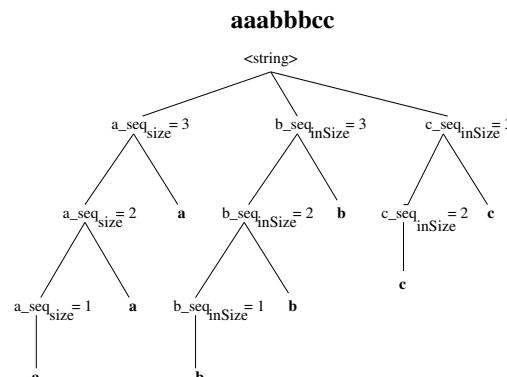
A problem?



string → a_seq_size b_seq_size c_seq_size
condition: a_seq_size = b_seq_size = c_seq_size
 a_seq_size → a
 a_seq1_size ← 1
 | a_seq2_size a
 a_seq1_size ← a_seq2_size + 1
 b_seq_size → b
 b_seq1_size ← 1
 | b_seq2_size b
 b_seq1_size ← b_seq2_size + 1
 c_seq_size → c
 c_seq1_size ← 1
 | c_seq2_size c
 c_seq1_size ← c_seq2_size + 1



Using Inherited & Synthesized Attributes



string → a_seq_size b_seq_inSize c_seq_inSize
 b_seq_inSize ← a_seq_size
 c_seq_inSize ← a_seq_size

a_seq_size → a
 a_seq_size ← 1
 | a_seq2_size a
 a_seq_size ← a_seq2_size + 1

b_seq_inSize → b
condition: b_seq_inSize = 1
 | b_seq2_inSize b
 b_seq2_inSize ← b_seq_inSize - 1

c_seq_inSize → c
condition: c_seq_inSize = 1
 | c_seq2_inSize c
 c_seq2_inSize ← c_seq_inSize - 1



Recognizer & Translator

- Compilers typically recognize legal programs for a given language *and* translate a program into another form
- Uses context-free grammar (CFG) to specify syntax for the recognizer
- Source program is also translated to target code as phrases are being recognized by parser
- Each grammar symbol in the CFG has a (possibly empty) set of **attributes** associated with it; these attributes usually are bound to values in the domain of the target language



Semantic Actions

- Remember that parsers generally have two main tasks – to recognize legal sentences for a given language *and* to transform those sentences into intermediate or target code
- Semantic actions are actions performed by the parser as it recognizes phrases in the source program
- In a recursive-descent parser, semantic action code is interspersed with the control flow of the parsing actions
- In the specification file for an LR parser generator, semantic actions are fragments of program code attached to grammar productions



Syntax-Directed Translation Definitions

- Each production in the CFG has a (possibly empty) set of **semantic rules** for computing attributes
- A **syntax-directed definition** (SDD) is a CFG plus a set of semantic rules
- An **attribute grammar** is a SDD in which the semantic rules cannot have side effects
- A SDT produces an **annotated parse tree**: a parse tree containing attributes and their values at each node
- Semantic rules that can be embedded anywhere in the right-hand side of a production are often referred to as **semantic actions**

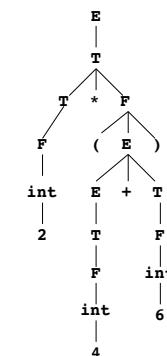


Example 1 – A Simple Calculator

- Goal: design a SDD that translates an arithmetic expression into its integer value, e.g. $2 * (4 + 6)$ would result in 20.

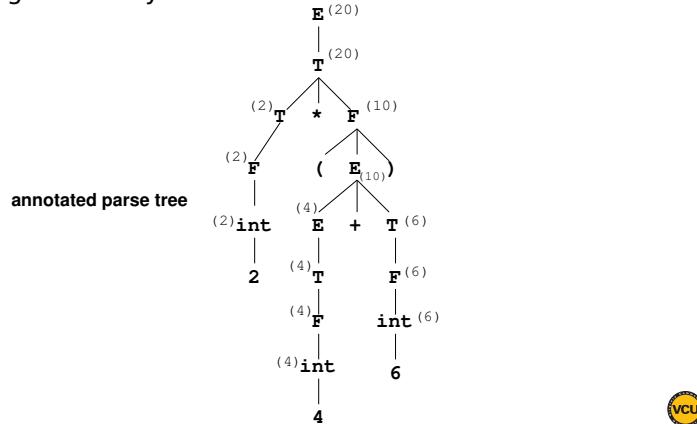
Step 1 Design a CFG for arithmetic expressions

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



Example 1 – A Simple Calculator

Step 2 Determine the attributes associated with each grammar symbol

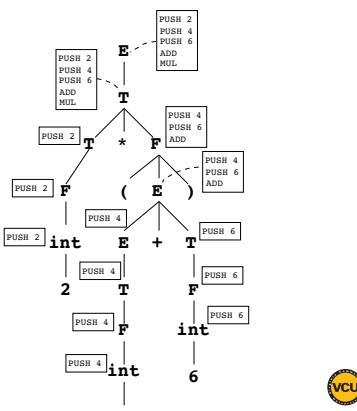


Example 2 – A Simple Translator

- Goal: Design a SDD that translates an arithmetic expression into a series of stack machine instructions, e.g.

$$2 * (4 + 6) \implies \begin{array}{l} \text{PUSH } 2 \\ \text{PUSH } 4 \\ \text{PUSH } 6 \\ \text{ADD} \\ \text{MUL} \end{array}$$

Step 1 Define the attributes



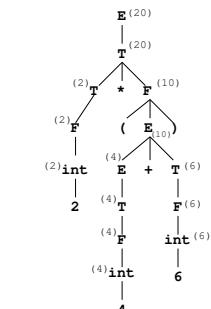
Example 1 – A Simple Calculator

Step 3 Write the semantic rules that will be used to determine the attributes.

CFG

$$\begin{aligned} E_1 &\rightarrow E_2 + T & E_1.\text{val} = E_2.\text{val} + T.\text{val} \\ E &\rightarrow T & E.\text{val} = T.\text{val} \\ T_1 &\rightarrow T_2 * F & T_1.\text{val} = T_2.\text{val} * F.\text{val} \\ T &\rightarrow F & T.\text{val} = F.\text{val} \\ F &\rightarrow \text{int} & F.\text{val} = \text{int}.\text{val} \\ F &\rightarrow (E) & F.\text{val} = E.\text{val} \end{aligned}$$

Semantic Rules



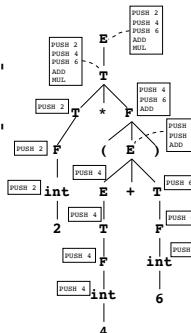
Example 2 – A Simple Translator

Step 2 Write the semantic rules that will be used to determine the attributes.

CFG

$$\begin{aligned} E_1 &\rightarrow E_2 + T & E_1.\text{code} = E_2.\text{code} + T.\text{code} + \text{" ADD"} \\ E &\rightarrow T & E.\text{code} = T.\text{code} \\ T_1 &\rightarrow T_2 * F & T_1.\text{code} = T_2.\text{code} + F.\text{code} + \text{" MUL"} \\ T &\rightarrow F & T.\text{code} = F.\text{code} \\ F &\rightarrow \text{int} & F.\text{code} = \text{"PUSH "} + \text{int}.\text{val} \\ F &\rightarrow (E) & F.\text{code} = E.\text{code} \end{aligned}$$

Semantic Rules



Question

What parsing method (i.e. LL or LR) is being used here?



Synthesized vs. Inherited Attributes

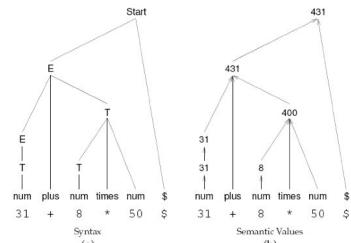


Figure 7.1: (a) Parse tree for the displayed expression;
(b) Synthesized attributes transmit values up the parse tree toward the root.

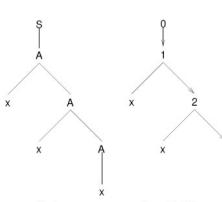


Figure 7.2: (a) Parse tree for the displayed input string; (b) Inherited attributes pass from parent to child.



Bottom-Up Syntax-Directed Translation

- 1 Start \rightarrow Digs_{ans} \$ call (ans)
- 2 Digs_{up} \rightarrow Digs_{below} d_{next}
up \leftarrow below $\times 10 +$ next
- 3 | d_{first}
up \leftarrow first

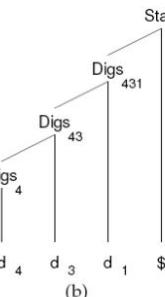


Figure 7.3: (a) Grammar with semantic actions; (b) Parse tree and propagated semantic values for the input 4 3 1 \$.



Bottom-Up Syntax-Directed Translation

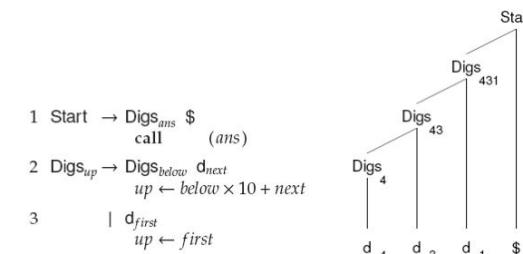
- Semantic actions are ‘triggered’ as reductions take place
- Abstractly, bottom-up parsers maintain two stacks:
 - ▶ the **syntactic (or parse) stack** containing terminal & nonterminal symbols
 - ▶ the **semantic stack** containing values associated with grammar symbols



Bottom-Up Syntax-Directed Translation

- Example for string 431:

Sentential Form	Reduction	Action
d ₄ d ₃ d ₁ \$	Digs _{up} \rightarrow d _{first}	up \leftarrow first
Digs ₄ d ₃ d ₁ \$	Digs _{up} \rightarrow Digs _{below} d _{next}	up \leftarrow below $\times 10 +$ next
Digs ₄₃ d ₁ \$	Digs _{up} \rightarrow Digs _{below} d _{next}	up \leftarrow below $\times 10 +$ next
Digs ₄₃₁ \$	Digs _{up} \rightarrow Digs _{below} d _{next}	up \leftarrow below $\times 10 +$ next
Start	Start \rightarrow Digs _{ans} \$	call PRINT(ans431)



Bottom-Up Syntax-Directed Translation

- Problem: what if we wish to trigger (execute) a semantic action *before* an actual reduction takes place?

1 Start → Num \$
 2 Num → o Digs
 3 | Digs
 4 Digs → Digs d
 5 | d

(a)

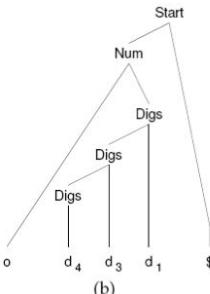


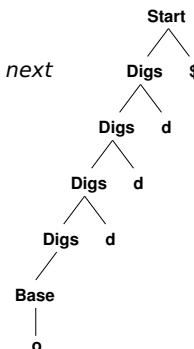
Figure 7.4: (a) Grammar and (b) parse tree for the input o 4 3 1 \$.



Bottom-Up Syntax-Directed Translation

- Better solution avoiding a global variable:

Start → Digs_{ans} \$
call PRINT(ans.val)
 Digs_{up} → Digs_{below} d_{next}
 $up.base \leftarrow below.base$
 $up.val \leftarrow below.val \times below.base + next$
 | Base_{base}
 $up.base \leftarrow base$
 $up.val \leftarrow 0$
 Base_n → 0
 $n \leftarrow 8$
 | ε
 $n \leftarrow 10$



Bottom-Up Syntax-Directed Translation

- Solution: use “dummy” ϵ -rules to force semantic actions

1 Start → Num_{ans} \$
call (ans)
 2 Num_{ans} → SignalOctal Digs_{octans}
 $ans \leftarrow octans$
 | SignalDecimal Digs_{decans}
 $ans \leftarrow decans$
 4 SignalOctal → o
 $base \leftarrow 8$
 5 SignalDecimal → λ
 $base \leftarrow 10$
 6 Digs_{up} → Digs_{below} d_{next}
 $up \leftarrow below \times base + next$
 7 | d_{first}
 $up \leftarrow first$

Figure 7.6: Use of λ -rules to force semantic action.

Bottom-Up Syntax-Directed Translation

- Bottom-up SDDs are naturally implemented using LR parsers
- A partial implementation of the grammar on the previous slide is given on the next page using the syntax of CUPS



Top-Down Syntax-Directed Translation

- Top-down parsers are usually written by hand with semantic actions embedded in the code
- Recall that in a recursive-descent parser there is one subroutine for each grammar production
- Semantic values can be naturally passed via subroutine parameters
- Since values can easily be passed in *and* out via parameters, it is common in recursive-descent top-down parsers to make use of both synthesized and inherited values



Top-Down Syntax-Directed Translation

```
import java.util.Scanner;

public class CMSC504SDTexampleTOPDOWN2 {

    // assumes octal digits are legal

    private static String source;
    private static int current = 0;
    private static final char EOF = '\0'; // arbitrary

    private static char buf;

    private static char getToken() {
        return (current < source.length() ? source.charAt(current++) : EOF);
    }

    public static void main(String [] args) {
        Scanner s = new Scanner(System.in);
        System.out.print("Input: "); source = s.nextLine();
        buf = getToken();
        start();
    }

    private static void start() {
        int b = base(); IntObj i = new IntObj();
        System.out.println(digs(b,i));
    }

    private static int base() {
        if (buf == '0') {
            buf = getToken(); return 8;
        }
        else return 10;
    }

    private static class IntObj {
        public int value;
        public IntObj(){}; // necessary 'cause class private
    }
}
```

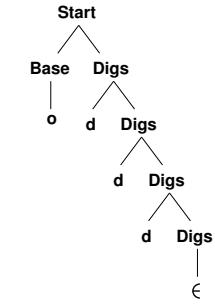
Start → Base Digs
 Digs → d Digs | ε
 Base → o | ε



Top-Down Syntax-Directed Translation

- Let's rewrite our example using both synthesized and inherited attributes (with synthesized (out) attributes in red, inherited (in) attributes in blue):

Start	→ Base _{base} Digs _{base,mult,ans}
	call PRINT(ans)
Digs _{base,mult,up}	→ d _{next} Digs _{base,mult,below}
	up ← next × mult + below
	mult ← mult × base
ε	
	up ← 0
	mult ← 1
Base _{base}	→ o
	base ← 8
ε	
	base ← 10



Top-Down Syntax-Directed Translation

```
private static int digs(int base, IntObj mult) {
    if (buf == EOF) { // Digs --> epsilon
        mult.value = 1;
        return 0;
    } else if (Character.isDigit(buf)) {
        int next = ((char) buf - (char) '0');
        buf = getToken();
        int below = digs(base,mult);
        int up = next * mult.value + below;
        mult.value = mult.value * base;
        return up;
    }
    else throw new RuntimeException("illegal character - "+buf);
}

private static int base() {
    if (buf == '0') {
        buf = getToken(); return 8;
    }
    else return 10;
}

private static class IntObj {
    public int value;
    public IntObj(){}; // necessary 'cause class private
}
```

Start → Base Digs
 Digs → d Digs | ε
 Base → o | ε

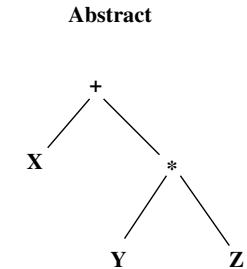
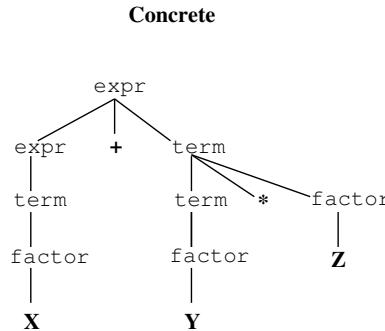


Overview

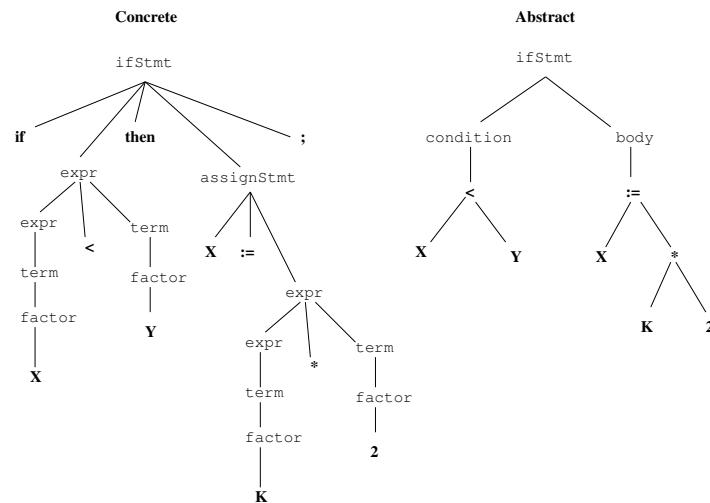
- Derivation trees from BNF grammars contain more information than is required to specify the semantics of a language
- The **abstract syntax** of a language describes the essential components of a language while ignoring inessential details
- The abstract syntax of a language can be used to create an **abstract syntax tree (AST)**
- Specifying the abstract syntax of a language means providing the possible templates for the various constructs of the language.
- Abstract syntax generally allows ambiguity since parsing has already been done
- Semantic specifications based on abstract syntax are much simpler



Abstract Syntax Trees



Abstract Syntax Trees



Concrete Syntax for Wren

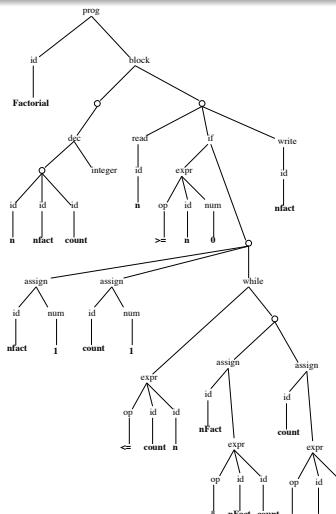
```

<program> → program ident is <block>
<block> → <dec seq> begin <cmd-seq> end
<dec seq> → ε | <dec><dec seq>
<dec> → var <var list> : <type>;
<type> → integer | boolean
<var list> → <var> | , <var list>
<cmd seq> → <cmd> | <cmd>; <cmd seq>
<cmd> → skip | ident := <expr>
| read ident | write <int expr>
| while <bool expr> do <cmd seq> end while
| if <bool expr> then <cmd seq> end if
| if <bool expr> then <cmd seq> else <cmd seq> end if
<expr> → <int expr> | <bool expr>
<int expr> → <term> | <int expr> weakop <term>
<term> → <element> | <term> strongop <element>
<element> → number | ident | ( <int expr> ) | - <element>
<bool expr> → <bool term> | <bool expr> or <bool term>
<bool term> → <bool element> | <bool term> and <bool element>
<bool element> → true | false | ident | <comparison>
| not ( <bool expr> ) | ( <bool expr> )
<comparison> → <int expr> relationop <int expr>
  
```

Abstract Syntax of Wren

Program \triangleq *prog*(Identifier, Block)
 Block \triangleq *block*(Declaration⁺, Command⁺)
 Declaration \triangleq *dec*(Identifier⁺, Type)
 Type \triangleq *integer* | *boolean*
 Command \triangleq *assign*(Identifier, Expression) | *skip*
 | *read*(Identifier) | *write*(Expression)
 | *while*(Expression, Command⁺)
 | *if*(Expression, Command⁺)
 | *ifelse*(Expression, Command⁺, Command⁺)
 Expression \triangleq Numeral | Identifier | *true* | *false*
 | *expr*(Operator, Expression, Expression)
 | *minus*(Expression)
 | *not*(Expression)
 Operator \triangleq + | - | * | / | *or* | *and*
 | <= | < | = | > | >= | <>

Concrete vs. Abstract Representation



Concrete vs. Abstract Representation

```

program Factorial is
var n,nFact,count:integer;
begin
  read n;
  if n >= 0 then
    nFact := 1; count := 1;
    while count <= n do
      nFact := nFact * count;
      count := count + 1;
    end while;
    write nFact;
  end if
end
  
```



Concrete vs. Abstract Representation

Creating the Abstract Syntax Tree

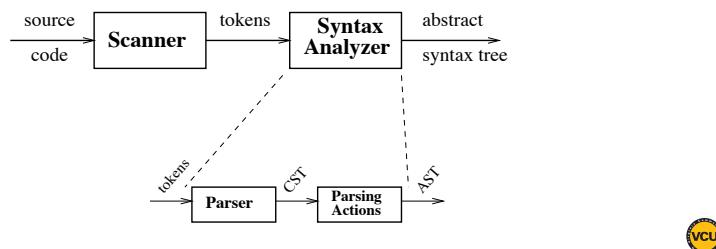
- The AST is typically constructed bottom-up; the AST data structure should therefore support tree construction from the leaves to the root
- Construction of the AST is usually done by semantic actions in the parser



Definition

Definition

A **canonical form** is the simplest or most symmetrical form to which all entities of a certain class can be reduced or transformed to without loss of generality



Overview

- Used to keep track of scope and binding information for names
- Searched everytime a name is encountered in code
- Required operations at compile-time: create S.T., remove S.T., add name, locate name
- Often must search multiple symbol tables for name
- Keywords usually handled separately
- Often implemented using a hash table or “searchable stack”



Overview

- Checks source program for semantic errors
- Gathers information for subsequent code-generation phase
- Main activity is usually type checking
- Relies heavily on information stored in symbol tables
- Semantic information is often specified using attribute grammars



Example

```

1 function f(a:int, b:int, c:int) =
2   ( print_int(a+c);
3     let var j := a+b
4       var a := "hello"
5     in print(a); print_int(j);
6   end;
7   print_int(b)
8 )

```

line 1:	$\sigma_0 = \sigma_{global} \uplus \{f \mapsto \text{func}\}$	\uplus ≡ overriding union
	$\sigma_1 = \sigma_0 \uplus \{a \mapsto \text{int}, b \mapsto \text{int}, c \mapsto \text{int}\}$	
line 2:	look up <i>a</i> , <i>b</i> in σ_1	
line 3,4:	$\sigma_2 = \sigma_1 \uplus \{j \mapsto \text{int}, a \mapsto \text{string}\}$	
line 5:	look up <i>a</i> , <i>j</i> in σ_2	
line 6:	discard σ_2	
line 7:	look up <i>b</i> in σ_1	
line 8:	discard σ_1	



Multiple Symbol Tables

```
structure M = struct
  structure E = struct
    val a = 5;
  end
  structure N = struct
    val b = 10;
    val a = E.a + b;
  end
  structure D = struct
    val d = E.a + N.a;
  end
end
```

(a) An example in ML

```
package M;
class E {
  static int a = 5;
}
class N {
  static int b = 10;
  static int a = E.a + b;
}
class D {
  static int d = E.a + N.a;
}
```

(b) An example in Java

$\sigma_0 = \text{base environment}$

$$\begin{aligned}\sigma_1 &= \{a \mapsto \text{int}\} \\ \sigma_2 &= \{E \mapsto \sigma_1\} \\ \sigma_3 &= \{b \mapsto \text{int}, a \mapsto \text{int}\} \\ \sigma_4 &= \{N \mapsto \sigma_3\} \\ \sigma_5 &= \{d \mapsto \text{int}\} \\ \sigma_6 &= \{D \mapsto \sigma_5\} \\ \sigma_7 &= \{M \mapsto (\sigma_2 + \sigma_4 + \sigma_6)\}\end{aligned}$$

FIGURE 5.1. Several active environments at once.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel



Overview

- Why generate intermediate code?
 - ▶ re-targeting code is facilitated
 - ▶ easier to optimize code (can be done independently of target machine instruction set)
- Will look at two kinds:
 - ▶ Graphical representation (syntax trees)
 - ▶ Three-address code (TAC) (quads, triples)



Using Hash Tables

```
class Bucket {String key; Object binding; Bucket next;
  Bucket(String k, Object b, Bucket n) {key=k; binding=b; next=n;}
}

class HashT {
  final int SIZE = 256;
  Bucket table[] = new Bucket[SIZE];

  int hash(String s) {
    int h=0;
    for(int i=0; i<s.length(); i++)
      h+=65539*s.charAt(i);
    return h;
  }

  void insert(String s, Binding b) {
    int index=hash(s)%SIZE
    table[index]=new Bucket(s,b,table[index]);
  }

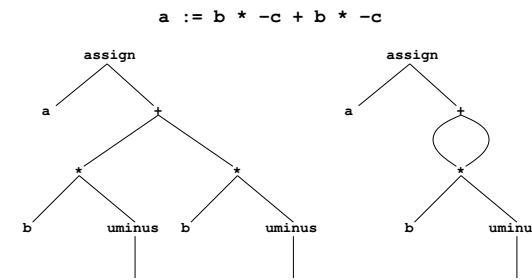
  Object lookup(String s) {
    int index=hash(s)%SIZE
    for (Binding b = table[index]; b!=null; b=b.next)
      if (s.equals(b.key)) return b.binding;
    return null;
  }

  void pop(String s) {
    int index=hash(s)%SIZE
    table[index]=table[index].next;
  }
}
```

PROGRAM 5.2. Hash table with external chaining.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel

Syntax Trees

- Natural way to represent the hierarchical structure of a program
- Can have more than one representation:



Syntax tree

DAG



Building Syntax Trees

- Can be generated by appropriate semantic actions:

Production	Semantic Rules
$S \rightarrow \mathbf{id} := E$	$S.\text{nptr} := \text{mknode}(\text{'assign'}, \text{mkleaf}(\mathbf{id}, \mathbf{id}.\text{place}), E.\text{nptr})$
$E \rightarrow E_1 + E_2$	$E.\text{nptr} := \text{mknode}(\text{'+'}, E_1.\text{nptr}, E_2.\text{nptr})$
$E \rightarrow E_1 * E_2$	$E.\text{nptr} := \text{mknode}(\text{*'}, E_1.\text{nptr}, E_2.\text{nptr})$
$E \rightarrow -E_1$	$E.\text{nptr} := \text{mknode}(\text{'uminus'}, E_1.\text{nptr})$
$E \rightarrow (E_1)$	$E.\text{nptr} := E_1.\text{nptr}$
$E \rightarrow \mathbf{id}$	$E.\text{nptr} := \text{mkleaf}(\mathbf{id}, \mathbf{id}.\text{place})$

- Can create DAGS if node creation routines try to return pointers to existing nodes rather than always creating new ones



Three-Address Code (TAC)

- General form: $x := y \text{ op } z$
- Can build complicated expressions using multiple TAC instructions
- For example, $x + y * z$ could be translated to:

$$\begin{aligned} t_1 &:= y * z \\ t_2 &:= x + t_1 \end{aligned}$$
- Using names for intermediate values allows TAC to be easily rearranged



Building Syntax Trees with a Compiler-Compiler

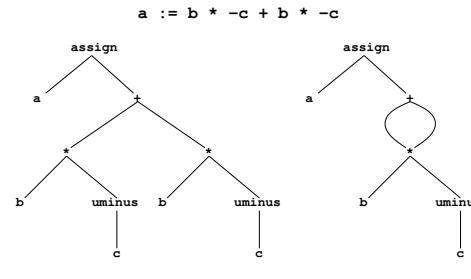
- CUP specification:

```
expr ::= expr:e1 PLUS expr:e2
       {: RESULT = new Abstract.BinaryArithmeticExpr(e1,sym.PLUS,e2); :}
| expr:e1 MINUS expr:e2
  {: RESULT = new Abstract.BinaryArithmeticExpr(e1,sym_MINUS,e2); :}
| expr:e1 TIMES expr:e2
  {: RESULT = new Abstract.BinaryArithmeticExpr(e1,sym.TIMES,e2); :}
| expr:e1 DIVIDE expr:e2
  {: RESULT = new Abstract.BinaryArithmeticExpr(e1,sym.DIVIDE,e2); :}
...
```



TAC Code & Syntax Trees

- TAC is a linear-representation of a syntax tree, e.g.



$t_1 := -c$	$t_1 := -c$
$t_2 := b * t_1$	$t_2 := b * t_1$
$t_3 := -c$	$t_3 := -c$
$t_4 := b * t_1$	$t_4 := b * t_1$
$t_5 := t_2 + t_4$	$t_5 := t_2 + t_4$
$a := t_5$	$a := t_5$



Examples of Three-Address Code

Binary Operators:	$a := b + c$
	$a := b - c$
	$a := b * c$
	$a := b / c$
Unary operators:	$a := -b$
Assignment:	$a := b$
Unconditional jump:	<code>goto a</code>
Conditional jump:	<code>ifz b goto a</code> <code>ifnz b goto a</code>
Function call & return:	<code>arg a</code> <code>a := call b</code> <code>return a</code>



Syntax Directed Translation into Three-Address Code continued

Production	Semantic Rules
$S \rightarrow \text{while } E \text{ do } S_1$	$S.begin := \text{newlabel};$ $S.after := \text{newlabel};$ $S.code := \text{gen}(S.begin ':') \parallel$ $E.code \parallel$ $\text{gen('if' } E.place '=' 0 \text{ 'goto' } S.after) \parallel$ $S_1.code \parallel$ $\text{gen('goto' } S.begin) \parallel$ $\text{gen}(S.after ':')$

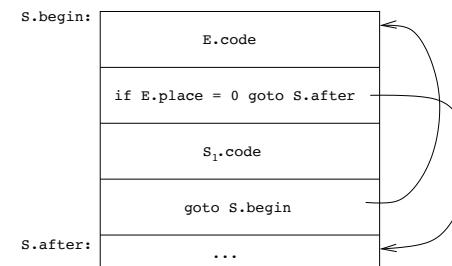


Syntax Directed Translation into Three-Address Code

Production	Semantic Rules
$S \rightarrow \text{id} := E$	$S.code := E.code \parallel \text{gen}(\text{id.place} ':=' E.place)$
$E \rightarrow E_1 + E_2$	$E.place := \text{newtemp};$ $E.code := E_1.code \parallel E_2.code \parallel$ $\text{gen}(E.place ':=' E_1.place '+' E_2.place)$
$E \rightarrow E_1 * E_2$	$E.place := \text{newtemp};$ $E.code := E_1.code \parallel E_2.code \parallel$ $\text{gen}(E.place ':=' E_1.place '*' E_2.place)$
$E \rightarrow -E_1$	$E.place := \text{newtemp};$ $E.code := E_1.code \parallel$ $\text{gen}(E.place ':=' 'uminus' E_1.place)$
$E \rightarrow (E_1)$	$E.place := E_1.place;$ $E.code := E_1.code$
$E \rightarrow \text{id}$	$E.place := \text{id.place};$ $E.code := "$



TAC for While Loop



```

S.begin := newlabel;
S.after := newlabel;
S.code := gen(S.begin ':') ||
E.code ||
gen('if' E.place '=' 0 'goto' S.after) ||
S1.code ||
gen('goto' S.begin) ||
gen(S.after ':')
  
```



Implementations of Three-Address Code

- Three-address code is an abstract form of intermediate code which can be implemented in many different ways
- We will look at **quads**, **triples**



Triples

- **triples**: record structure with three fields

$\begin{array}{c} \text{op1} \text{ arg1} \text{ arg2} \\ \underbrace{\hspace{1cm}}_{\text{S.T. entry or code position}} \end{array}$

- to avoid entering temps into S.T., can refer to temp value by position

$a := b * -c + b * -c$

	op	arg1	arg2
(0)	uminus	c	
(1)	*	b	(0)
(2)	uminus	c	
(3)	*	b	(2)
(4)	+	(1)	(3)
(5)	:=	(4)	a



Quads

- **quads**: record structure with four fields

$\begin{array}{c} \text{op1} \text{ arg1} \text{ arg2} \text{ result} \\ \underbrace{\hspace{1cm}}_{\text{S.T. entries}} \end{array}$

S.T. entries

- temporary variables are created as needed and placed in the S.T.

$a := b * -c + b * -c$

	op	arg1	arg2	result
(0)	uminus	c		t_1
(1)	*	b	t_1	t_2
(2)	uminus	c		t_3
(3)	*	b	t_3	t_4
(4)	+	t_2	t_4	t_5
(5)	:=	t_5		a



Local Declarations

$P \rightarrow D \quad \{ \text{offset} := 0 \}$

$D \rightarrow D ; D$

$D \rightarrow id : T \quad \{ \text{enter}(id.name, T.type, offset);$
 $\text{offset} := \text{offset} + T.width; \}$

$T \rightarrow \text{integer} \quad \{ T.type := \text{integer};$
 $T.width := 4; \}$

$T \rightarrow \text{real} \quad \{ T.type := \text{real};$
 $T.width := 8; \}$

$T \rightarrow \text{array} [num] \text{ of } T_1 \quad \{ T.type := \text{array}(num.val, T_1.type);$
 $T.width := num.val \times T_1.width \}$

$T \rightarrow \uparrow T_1 \quad \{ T.type := \text{pointer}(T_1.type);$
 $T.width := 4 \}$



└ Code Generation

└ Intermediate Languages

Declarations in Nested Procedures

```

P → M D      { addwidth(top(tblptr), top(offset)); }
                pop(tblptr); pop(offset); }

M → ε         { t := mkttable(nil);
                push(t,tblptr); push(0,offset); }

D → D1 ; D2

D → proc id ; N D1 ; S { t := top(tblptr);
                            addwidth(t,top(offset));
                            pop(tblptr); pop(offset);
                            enterproc(top(tblptr), id.name, t) }

D → id : T   { enter(top(tblptr),id.name, T.type, top(offset));
                top(offset) := top(offset) + T.width }

N → ε         { t := mkttable(top(tblptr));
                push(t,tblptr); push(0,offset); }

```

mkttable(previous) creates a new symbol table
enter(table.name.type,offset) creates a new entry for *name* in S.T. pointed to by *table*
addwidth(table,width) records cumulative width of all entries in *table*
enterproc(table.name,newtable) creates new entry for procedure *name* in S.T. pointed to by *table*



└ Code Generation

└ Intermediate Languages

Assignments & Arithmetic Expressions

```

S → id := E   { p := lookup(id.name);
                  if p ≠ nil then emit(p ':=' E.place)
                  else error }

E → E1 + E2 { E.place := newtemp;
                      emit(E.place ':=' E1.place +' E2.place) }

E → E1 * E2 { E.place := newtemp;
                      emit(E.place ':=' E1.place '*' E2.place) }

E → -E1     { E.place := newtemp;
                  emit(E.place ':=' 'uminus' E1.place) }

E → ( E1 )   { E.place := E1.place }

E → id        { p := lookup(id.name);
                  if p ≠ nil then E.place := p
                  else error }

```



└ Code Generation

└ Intermediate Languages

Records

```

T → record L D end { T.type := record(top(tblptr));
                        T.width := top(offset);
                        pop(tblptr); pop(offset) }

L → ε               { t := mkttable(nil);
                        push(t,tblptr); push(0,offset) }

```



└ Code Generation

└ Intermediate Languages

Boolean Expressions

- Grammar:

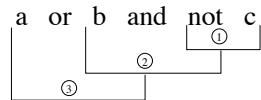
E → E **or** E | E **and** E | **not** E | (E) | id relop id |
true | **false**

- Two methods of translating:

- encode true and false numerically, then evaluate boolean expressions as arithmetic expressions
- flow of control: represent value of boolean expression by a position reached in a program



Numerical Representation



- Three-address code sequence for and/or/not:

$t_1 := \text{not } c$
 $t_2 := b \text{ and } t_1$
 $t_3 := a \text{ or } t_2$

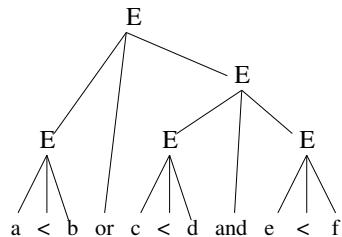
- Three-address code sequence for relational operators:

$a \text{ relop } b \equiv \text{if } a \text{ relop } b \text{ then } 1 \text{ else } 0$

```
100: if a relop b goto 103
101: t := 0
102: goto 104
103: t := 1
104:
```



Example Expression



```
100: if a < b goto 103 107: t2 := 1
101: t1 := 0 108: if e < f goto 111
102: goto 104 109: t3 := 0
103: t1 := 1 110: goto 112
104: if c < d goto 107 111: t3 := 1
105: t2 := 0 112: t4 := t2 and t3
106: goto 108 113: t5 := t1 or t4
```



Numerical Representation Translation Scheme

$E \rightarrow E_1 \text{ or } E_2$	{ E.place := newtemp; emit(E.place ':=> E ₁ .place 'or' E ₂ .place) }
$E \rightarrow E_1 \text{ and } E_2$	{ E.place := newtemp; emit(E.place ':=> E ₁ .place 'and' E ₂ .place) }
$E \rightarrow \text{not } E_1$	{ E.place := newtemp; emit(E.place ':=> 'not' E ₁ .place) }
$E \rightarrow (E_1)$	{ E.place := E ₁ .place }
$E \rightarrow id_1 \text{ relop } id_2$	{ E.place := newtemp; emit('if' id ₁ .place relop op id ₂ .place 'goto' nextstat + 3); emit(E.place ':=> '0'); emit('goto' nextstat + 2); emit(E.place ':=> '1') }
$E \rightarrow \text{true}$	{ E.place := newtemp; emit(E.place ':=> '1') }
$E \rightarrow \text{false}$	{ E.place := newtemp; emit(E.place ':=> '0') }



Flow of Control Method

- Translate expression into series of TAC instructions involving a sequence of unconditional and conditional jumps to one of two labels: E.true and E.false

- Example #1:



$a < b \equiv \text{if } a < b \text{ goto E.true}$
 goto E.false

- Example #2: $E \rightarrow E_1 \text{ or } E_2$

- if E_1 is true: $E_1.\text{true} \equiv E.\text{true}$
- if E_1 is false: E_2 must be evaluated, therefore $E_1.\text{false}$ is the label for the first statement in the code for E_2

- Example #3: $E \rightarrow \text{not } E_1$

No code is generated: $E_1.\text{true} \equiv E.\text{false}$
 $E_1.\text{false} \equiv E.\text{true}$



└ Code Generation

└ Intermediate Languages

Flow of Control Translation Scheme

```

E → E1 or E2 { E1.true := E.true;
                     E1.false := newlabel;
                     E2.true := E.true;
                     E2.false := E.false;
                     E.code := E1.code || gen(E1.false ':') || E2.code }

E → E1 and E2 { E1.true := newlabel;
                     E1.false := E.false;
                     E2.true := E.true;
                     E2.false := E.false;
                     E.code := E1.code || gen(E1.true ':') || E2.code }

E → not E1 { E1.true := E.false;
                   E1.false := E.true;
                   E.code := E1.code }

E → ( E1 ) { E1.true := E.true;
                  E1.false := E.false;
                  E.code := E1.code }

E → id1 relop id2 { E.code := gen('if' id1.place relop.op id2.place 'goto' E.true) ||
                           gen('goto' E.false) }

E → true { E.code := gen('goto' E.true) }

E → false { E.code := gen('goto' E.false) }

```



└ Code Generation

└ Intermediate Languages

Loop Statements

```

while-do:
BEGIN_LOOP:
    <code for test>
    if false, goto END_LOOP:
    <statements>
    goto BEGIN_LOOP:
END_LOOP: ...

```



└ Code Generation

└ Intermediate Languages

Branching Statements

if-then:

```

...
<code for test>
if false, goto END:
<statements>
END: ...

```

if-then-else:

```

...
<code for test>
if false, goto FALSE_PART:
<statements>
goto END:
FALSE_PART:
<statements>
END: ...

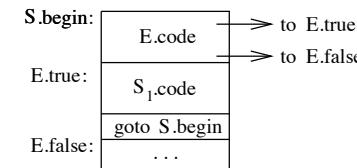
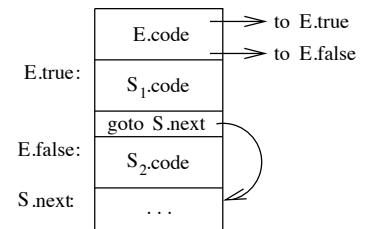
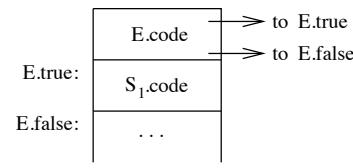
```



└ Code Generation

└ Intermediate Languages

Flow of Control Formats



Branching and Loop Statements Translation Scheme

```
S → if E then S1 { E.true := newlabel;
E.false := S.next;
S1.next := S.next;
S.code := E.code || gen(E.true ':') || S1.code }
```

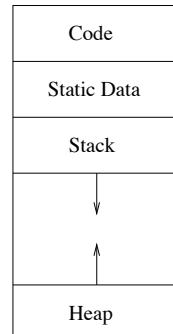
```
S → if E then S1 else S2 { E.true := newlabel;
E.false := newlabel;
S1.next := S.next;
S2.next := S.next;
S.code := E.code ||
gen(E.true ':') || S1.code ||
gen('goto' S.next) ||
gen(E.false ':') || S2.code }
```

```
S → while E do S1 { S.begin := newlabel;
E.true := newlabel;
E.false := S.next;
S1.next := S.begin;
S.code := gen(S.begin ':') || E.code ||
gen(E.true ':') || S1.code ||
gen('goto' S.begin) }
```



Run-Time Storage Organization

- Memory subdivided to hold:
 - generated target code
 - data objects
 - stack of subroutine activations
- Typical structure:



Storage Allocation Strategies

- Three primary strategies:
 - Static allocation
 - Stack allocation
 - Heap allocation
- Adopted strategies often combination of all three
- Will concentrate on stack allocation

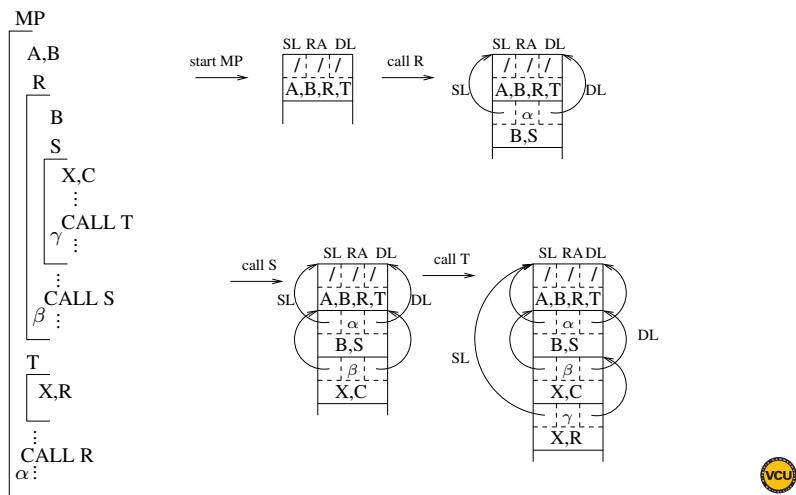


Activation Records

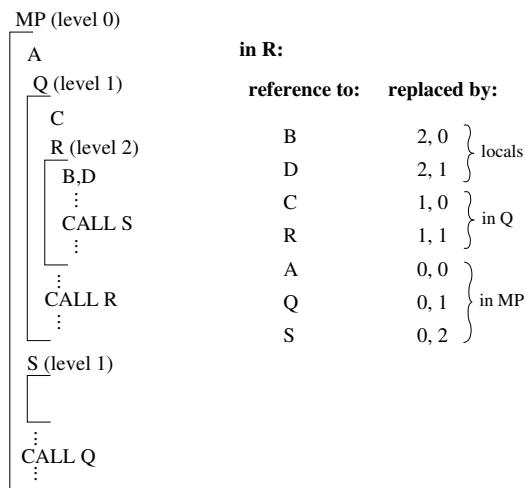
- Each subprogram call generates an *activation record* (AR) or *frame* on the run-time stack
- Activation records typically contain:
 - Temporary values
 - Local data
 - Saved machine status (registers)
 - Stack control links (dynamic links)
 - Variable access links (static links)
 - Subprogram parameters
 - Subprogram return values
 - Subprogram return address



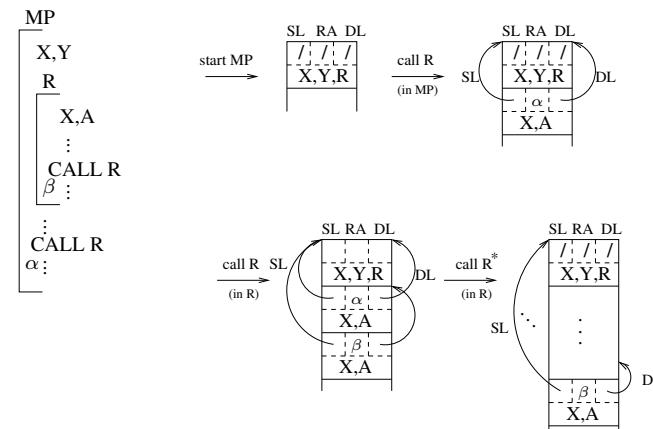
The Run-Time Stack



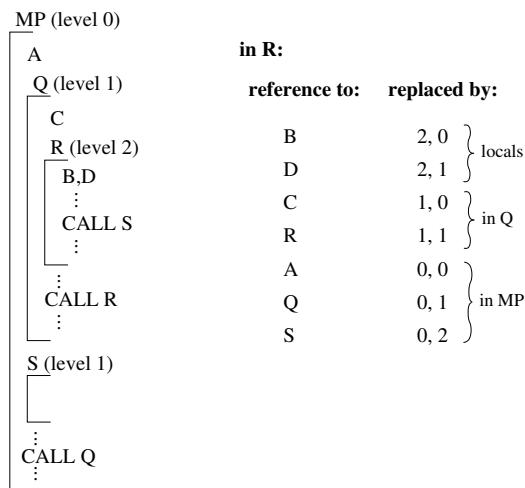
Addressing for Non-local Scopes



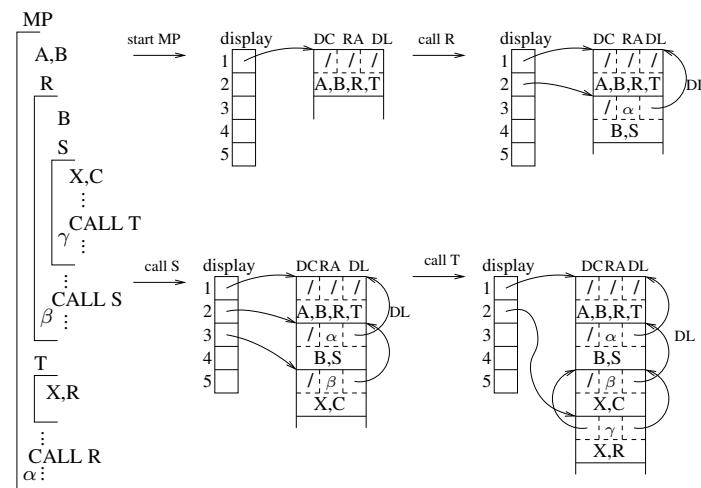
Recursion



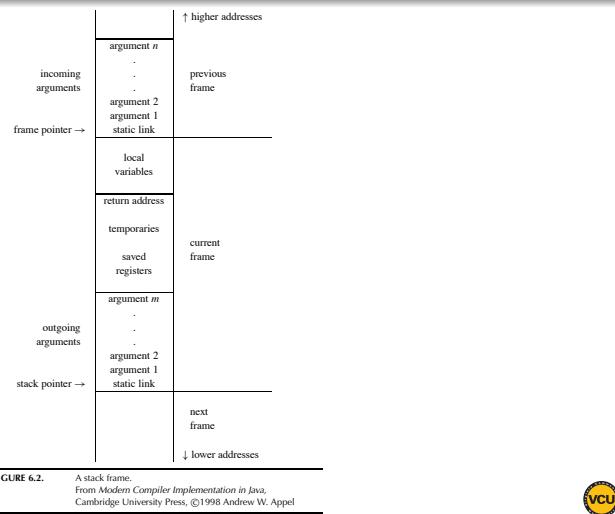
Addressing for Non-local Scopes



Using a Display



Parameter Passing



Overview

- Intermediate representations usually use an unbounded number of temporaries
- Target machine, however, has a bounded number of registers
- Two temporaries t_1 and t_2 can be stored in the same register if they are never “in use” at the same time
- Therefore to properly map an unbounded number of temporaries to a bounded set of registers, a compiler must be able to determine which variables are in use at any given time; excess temporaries can be stored in memory
- A variable is said to be **live** if it holds a value that may be needed in the future
- Liveness analysis begins with a **control-flow graph**



Optimizations

- Use registers to minimize memory traffic:
 - ▶ for parameters: studies show few routines have more than 6 parameters
 - ▶ for locals & temporaries
 - ▶ return address
 - ▶ displays
- Dynamic link can be eliminated if one has a set frame size



Liveness of Variables

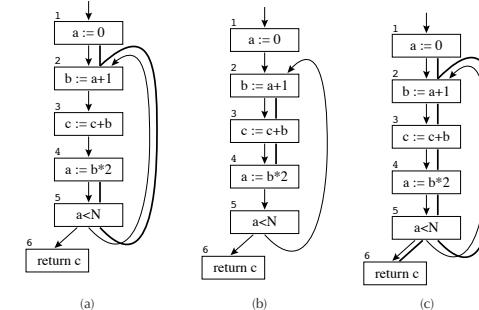


FIGURE 10.2. Liveness of variables a , b , c .
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel

- Determining the live range of a variable is an example of a **dataflow** problem



Dataflow Terminology & Definitions

- A flow graph has:
 - ▶ out edges that lead to **successor** nodes; the set **succ[n]** is the set of successors of node n
 - ▶ in edges that come from **predecessor** nodes; the set **pred[n]** is the set of predecessors of node n
- An assignment to a variable or temporary **defines** that variable; an occurrence of a variable in an expression **uses** the variable
- The sets **def[n]** and **use[n]** are the set of variables defined and used in node n
- A variable is **live** on an edge if there is a directed graph from that edge to a **use** of the variable that does not go through a **def** of that variable
- A variable is **live-in** at a node if it is live on any of the in-edges of the node; it is **live-out** at a node if it is live on any out-edges of the node



Computation of Liveness

```

for each  $n$ 
   $in[n] \leftarrow \{ \}; out[n] \leftarrow \{ \}$ 
repeat
  for each  $n$ 
     $in'[n] \leftarrow in[n]; out'[n] \leftarrow out[n]$ 
     $in[n] \leftarrow use[n] \cup (out[n] - def[n])$ 
     $out[n] \leftarrow \bigcup_{s \in succ[n]} in[s]$ 
  until  $in'[n] = in[n]$  and  $out'[n] = out[n]$  for all  $n$ 

```

ALGORITHM 10.4. Computation of liveness by iteration.
 From *Modern Compiler Implementation in Java*,
 Cambridge University Press, ©1998 Andrew W. Appel



Dataflow Equations

- ① If a variable is in **use[n]**, then it is **live-in** at node n
- ② If a variable is **live-in** at a node n , then it is **live-out** at all nodes in **pred[n]**
- ③ If a variable is **live-out** at node n and not in **def[n]**, then the variable is also **live-in** at n

$$in[n] = use[n] \cup (out[n] - def[n])$$

$$out[n] = \bigcup_{s \in succ[n]} in[s]$$

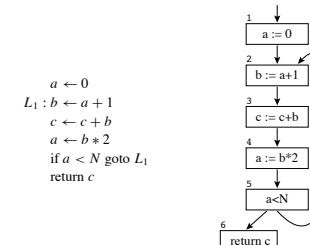
EQUATIONS 10.3. Dataflow equations for liveness analysis.

From *Modern Compiler Implementation in Java*,
 Cambridge University Press, ©1998 Andrew W. Appel

(in[n] = live-in at n , out[n] = live-out at n)



Computation of Liveness Example



GRAPH 10.1. Control-flow graph of a program.
 From *Modern Compiler Implementation in Java*,
 Cambridge University Press, ©1998 Andrew W. Appel

```

for each  $n$ 
   $in[n] \leftarrow \{ \}; out[n] \leftarrow \{ \}$ 
repeat
  for each  $n$ 
     $in'[n] \leftarrow in[n]; out'[n] \leftarrow out[n]$ 
     $in[n] \leftarrow use[n] \cup (out[n] - def[n])$ 
     $out[n] \leftarrow \bigcup_{s \in succ[n]} in[s]$ 
  until  $in'[n] = in[n]$  and  $out'[n] = out[n]$  for all  $n$ 

```

ALGORITHM 10.4. Computation of liveness by iteration.
 From *Modern Compiler Implementation in Java*,
 Cambridge University Press, ©1998 Andrew W. Appel



Computation of Liveness Example continued

Calculations using example graph:

	use	def	1st in	out	2nd in	out	3rd in	out	4th in	out	5th in	out	6th in	out	7th in	out
1		a			a		a		ac	c	ac	c	ac	c	ac	
2	a	b	a		bc	ac	bc	ac	bc	ac	bc	bc	ac	bc	ac	
3	bc	c	bc		b	bc	b	bc	c	bc	b	bc	bc	bc	bc	
4	b	a	b		a	b	a	b	ac	bc	ac	bc	ac	bc	ac	
5	a		a	a	ac	ac	ac									
6	c		c	c	c	c	c	c	c	c	c	c	c	c	c	

TABLE 10.5. Liveness calculation following forward control-flow edges.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel

	use	def	1st out	in	2nd out	in	3rd out	in
6	c		c		c		c	
5	a		c	ac	ac	ac	ac	ac
4	b	a	ac	bc	ac	bc	ac	bc
3	bc	c	bc	bc	bc	bc	bc	bc
2	a	b	bc	ac	bc	ac	bc	ac
1		a	ac	c	ac	c	ac	c

TABLE 10.6. Liveness calculation following reverse control-flow edges.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel



Conservative Approximation

- If a value a is needed when execution reaches node n , then we can be assured that a is live-out in any solution of the equations
- However, the converse is *not* true: any variable d that is live-out at node n does not mean that it will really be used
- A conservative approximation of liveness is one that may erroneously believe a variable is alive, but it will never erroneously believe it is dead (why is this important?)
- A dataflow equation used for compiler optimization should be set up so that any solution to it provides conservative information to the optimizer; imprecise information may lead to sub-optimal but never incorrect programs

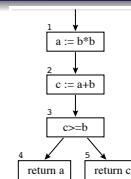


Time Complexity of Liveness Analysis

- Worst case: $O(N^4)$!
- However, ordering the nodes using depth-first search usually brings the number of repeat-loop iterations to two or three, plus the live-sets are usually sparse, so the algorithm runs between $O(N)$ and $O(N^2)$
- Other, better methods exist



Static vs. Dynamic Liveness



GRAPH 10.8. Standard static dataflow analysis will not take advantage of the fact that node 4 can never be reached.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel

Dynamic liveness A variable a is dynamically live at node n if some execution of the program goes from n to a use of a without going through any definition of a

Static liveness A variable a is statically live at node n if there is some path of control-flow edges from n to some use of a that does not go through a definition of a

- Clearly if a is dynamically live it is also statically live
- An optimizing compiler must perform optimizations on the basis of static liveness



Interference Graphs

- Liveness analysis can be used for register allocation and code optimizations
- If we have a set of temporaries a, b, c, \dots that must be allocated to registers r_1, \dots, r_k , a condition that prevents a and b from being allocated to the same register is called **interference**
- Interference information can be represented as a matrix or graph:

	a	b	c
a		x	
b		x	
c	x	x	

(a) Matrix

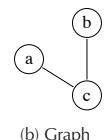


FIGURE 10.9. Representations of interference.
From *Modern Compiler Implementation in Java*, Cambridge University Press, ©1998 Andrew W. Appel



Spilling Example

- Given the following code for expression $(b * c) + a - (c * a)$:

```
t1 ← b
t2 ← c
t3 ← t1 * t2
t4 ← a
t5 ← t3 + t4
t6 ← t2 * t4
t7 ← t5 - t6
```

- If $K = 2$ (with registers $r1, r2$), where operands must be registers:

```
r2 = b
r1 = c
r2 = r2 * r1
spill1 = r1
r1 = a
r2 = r2 + r1
spill2 = r2
r2 = spill1
r1 = r2 * r1
r2 = spill2
r1 = r2 - r1
```



Register Allocation using Graph Coloring

- register allocation assigns intermediate-code temporary variables to actual machine registers
- Method:
 - ① construct an *interference graph*
 - ② *color* the graph *spilling* registers where needed
- The graph coloring problem derives from the old mapmakers' rule that no adjacent countries on a map should have the same color; in our case we "color" the interference graph with colors corresponding to registers
- The goal: if our target machine has K registers we wish to color the graph with K (or less) colors
- If a K -coloring is not possible, some variables will have to be spilled to memory



Coloring by Simplification

- General register allocation and graph coloring are *NP*-complete problems
- Linear-time approximation algorithm:
 - ① **Build** – construct the interference graph
 - ② **Simplify** – color the graph using a simple heuristic:
 - ★ suppose graph G contains a node m with fewer than K neighbors (where $K = \#$ of registers)
 - ★ let G' be the graph $G - \{m\}$ obtained by removing node m ; we therefore know that if G' can be colored, so then can G
 - ★ repeatedly remove (and push onto a stack) nodes of degree less than K until all nodes have been removed
 - ③ **Spill** – if at some point the heuristic fails (i.e. a node n is encountered with degree $\geq K$):
 - ★ mark node n for spilling, remove it from the graph and push it onto the stack
 - ★ continue with the simplification
 - ④ **Select** – assign colors (registers) to nodes of the graph:
 - ★ rebuild the graph from nodes from the stack, assigning a non-interfering register to each node when popped
 - ★ add code for spills when spill nodes are popped from stack



Graph Coloring Example

```
live-in: k j
g := mem[j+12]
h := k - 1
f := g * h
e := mem[j+8]
m := mem[j+16]
b := mem[f]
c := e + 8
d := c
k := m + 4
j := b
live-out: d k j
```



Coalescing

- Redundant move instructions can be removed from an interference graph; if there is no edge in the graph between the source and destination of a move, then the move can be eliminated
- The source and destination nodes are *coalesced* (combined) into a new node whose edges are the union of the edges being replaced
- Algorithm can be found in text

e	1
m	2
f	3
j&b	4
c&d	1
k	2
h	2
g	1
stack	
coloring	

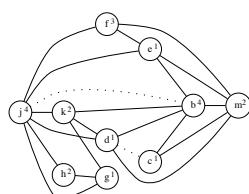
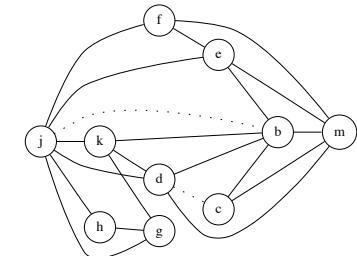


FIGURE 11.6. A coloring, with coalescing, for Graph 11.1.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel



Graph Coloring Example continued

```
live-in: k j
g := mem[j+12]
h := k - 1
f := g * h
e := mem[j+8]
m := mem[j+16]
b := mem[f]
c := e + 8
d := c
k := m + 4
j := b
live-out: d k j
```



GRAPH 11.1.

Interference graph for a program. Dotted lines are not interference edges but indicate move instructions.
From *Modern Compiler Implementation in Java*,
Cambridge University Press, ©1998 Andrew W. Appel

