Chapter 4 Grammars and Parsing

- A context-free grammar G = (Vt, Vn, S, P)
 - A finite terminal vocabulary Vt
 - The token set produced by scanner
 - A finite set of nonterminal vacabulary Vn
 - Intermediate symbols
 - A start symbol S ∈ Vn that starts all derivations
 - Also called goal symbol
 - P, a finite set of productions (rewriting rules) of the form $A \rightarrow X_1 X_2 \dots X_m$
 - $A \in Vn$, $X_i \in Vn \cup Vt$, $1 \le i \le m$
 - A $\rightarrow \lambda$ is a valid production

Other notations

- Vacabulary V of G,
 - $V = V \cap \cup V t$
- L(G), the set of string s derivable from S
 - Context-free language of grammar G
- Notational conventions
 - a,b,c, ... denote symbols in Vt
 - A,B,C, ... denote symbols in Vn
 - U,V,W, ... denote symbols in V
 - $\alpha, \beta, \gamma, \dots$ denote strings in V^*
 - u,v,w, \dots denote strings in V_t^*

- Derivation
 - One step derivation
 - If $A \rightarrow \gamma$, then $\alpha A \beta \Rightarrow \alpha \gamma \beta$
 - One or more steps derivation ⇒⁺
 - Zero or more steps derivation ⇒*
- If S ⇒*β, then β is said to be sentential form of the CFG
 - SF(G) is the set of sentential forms of grammar G
- $L(G) = \{x \in V_t^* | S \Rightarrow^+ x\}$
 - $-L(G)=SF(G)\cap V_t^*$

Left-most derivation, a top-down parsers

$$\square \Rightarrow_{\operatorname{Im}} , \Rightarrow_{\operatorname{Im}} , \stackrel{+}{\Rightarrow}_{\operatorname{Im}} *$$

E.g. of leftmost derivation of F(V+V)

```
G_{0} \begin{cases} E \rightarrow Prefix(E) \\ E \rightarrow V \text{ Tail} \\ Prefix \rightarrow F \\ Prefix \rightarrow \lambda \\ Tail \rightarrow +E \\ Tail \rightarrow \lambda \end{cases} \qquad E \Rightarrow_{lm} Prefix(E) \\ \Rightarrow_{lm} F(E) \\ \Rightarrow_{lm} F(V \text{ Tail}) \\ \Rightarrow_{lm} F(V + V \text{ Tail}) \\ \Rightarrow_{lm} F(V + V \text{ Tail}) \\ \Rightarrow_{lm} F(V + V \text{ Tail}) \end{cases}
```

• Right-most derivation (canonical derivation)

```
\square \Rightarrow_{\mathsf{rm}} , \Rightarrow_{\mathsf{rm}} , \stackrel{\mathsf{+}}{\Rightarrow}_{\mathsf{rm}} \stackrel{\mathsf{*}}{\Rightarrow}
```

- Buttom-up parsers
- E.g. of leftmost derivation of F(V+V)

```
G_{0} \begin{cases} E \rightarrow Prefix(E) \\ E \rightarrow V \text{ Tail} \\ Prefix \rightarrow F \\ Prefix \rightarrow \lambda \\ Tail \rightarrow +E \\ Tail \rightarrow \lambda \end{cases} \qquad E \Rightarrow_{rm} Prefix(E) \\ \Rightarrow_{rm} Prefix(V \text{ Tail}) \\ \Rightarrow_{rm} Prefix(V+E) \\ \Rightarrow_{rm} Prefix(V+V) \\ \Rightarrow_{rm} Prefix(V+V) \\ \Rightarrow_{rm} F(V+V) \end{cases}
```

- A parse tree
 - rooted by the start symbol
 - Its leaves are grammar symbols or λ

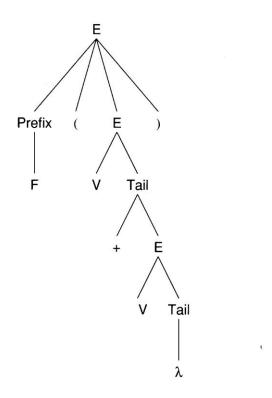


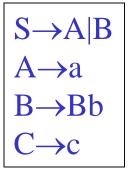
Figure 4.1 The Parse Tree Corresponding to F(V+V)

- A *phase* of a sentential form is a sequence of symbols descended from a single nonterminal in the parse tree
 - Simple or prime phrase
- The *handle* of a sentential form is the leftmost simple phrase

- Regular grammars
 - is of CFGs
 - Limited to productions of the form
 - A→aB
 - $C \rightarrow \lambda$
 - See exercise 6
- The *handle* of a sentential form is the leftmost simple phrase

Errors in Context-Free Grammars

- CFGs are a definitional mechanism. They may have errors, just as programs may.
- Flawed CFG
 - 1. Useless nonterminals
 - Unreachable
 - Derive no terminal string



Nonterminal C cannot be reached form S Nonterminal B derives no terminal string

Errors in Context-Free Grammars

- Ambiguous:
 - Grammars that allow different parse trees for the same terminal string
- It is impossible to decide whether a given CFG is ambiguous

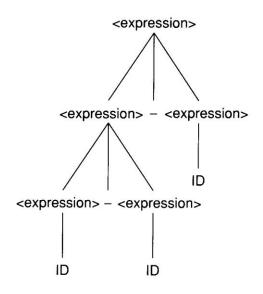


Figure 4.2 A Parse Tree for ID–ID

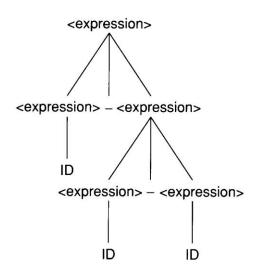


Figure 4.3 An Alternate Parse Tree for ID-ID-ID

Errors in Context-Free Grammars

- It is impossible to decide whether a given CFG is ambiguous
 - For certain grammar classes, we can prove that constituent grammars are unambiguous
- Wrong language
- A general comparison algorithm applicable to all CFGs is known to be impossible

Transforming Extened BNF Grammars

- - Extended BNF allows
 - Square bracket []
 - Optional list {}

```
for (each production P = A \rightarrow \alpha \ [X_1 \ \dots \ X_n] \ \beta) { Create a new nonterminal, N. Replace production P with P' = A \rightarrow \alpha \ N \ \beta Add the productions: N \rightarrow X_1 \ \dots \ X_n and N \rightarrow \lambda } for (each production Q = B \rightarrow \gamma \ \{Y_1 \ \dots \ Y_m\} \ \delta) { Create a new nonterminal, M. Replace production Q with Q' = B \rightarrow \gamma \ M \ \delta Add the productions: M \rightarrow Y_1 \ \dots \ Y_m \ M and M \rightarrow \lambda }
```

Figure 4.4 Algorithm to Transform Extended BNF Grammars into Standard Form

Parsers and Recognizers

- Recognizer
 - An algorithm that does boolean-valued test
 - "Is this input syntactically valid?
- Parser
 - Answers more general questions
 - Is this input valid?
 - And, if it is, what is its structure (parse tree)?

- Two general approaches to parsing
 - Top-down parser
 - Expanding the parse tree (via predictions) in a depth-first manner
 - Preorder traversal of the parse tree
 - *Predictive* in nature
 - lm
 - LL

- Buttom-down parser
 - Beginning at its bottom (the leaves of the tree, which are terminal symbols) and determining the productions used to generate the leaves
 - Postorder traversal of the parse tree
 - rm
 - LR

```
 \begin{array}{ll} <\text{Program}> & \rightarrow \text{begin} <\text{Stmts}> \text{end} \ \$ \\ <\text{Stmts}> & \rightarrow <\text{Stmt}> \ ; <\text{Stmts}> \\ <\text{Stmt}> & \rightarrow \lambda \\ <\text{Stmt}> & \rightarrow \text{SimpleStmt} \\ <\text{Stmt}> & \rightarrow \text{begin} <\text{Stmts}> \text{end} \\ \end{array}
```

To parse

begin SimpleStmt; SimpleStmt; end \$

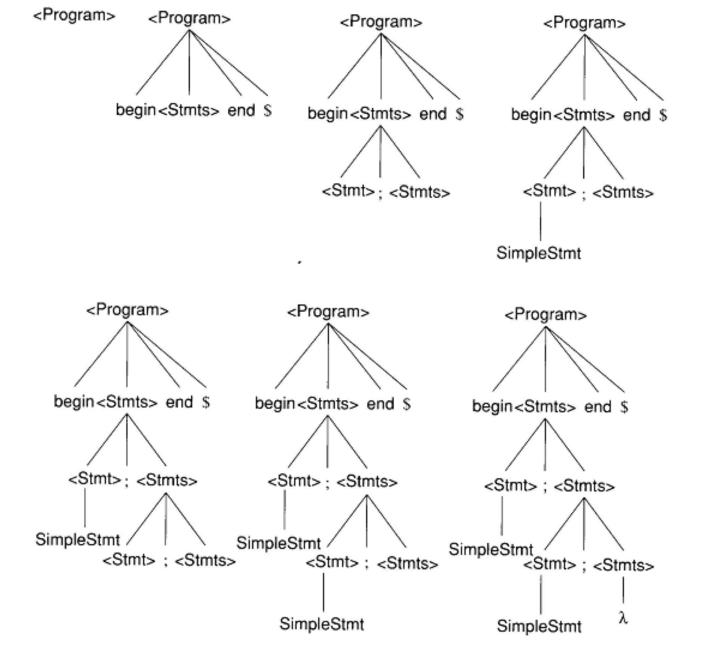


Figure 4.5 A Top-Down Parse

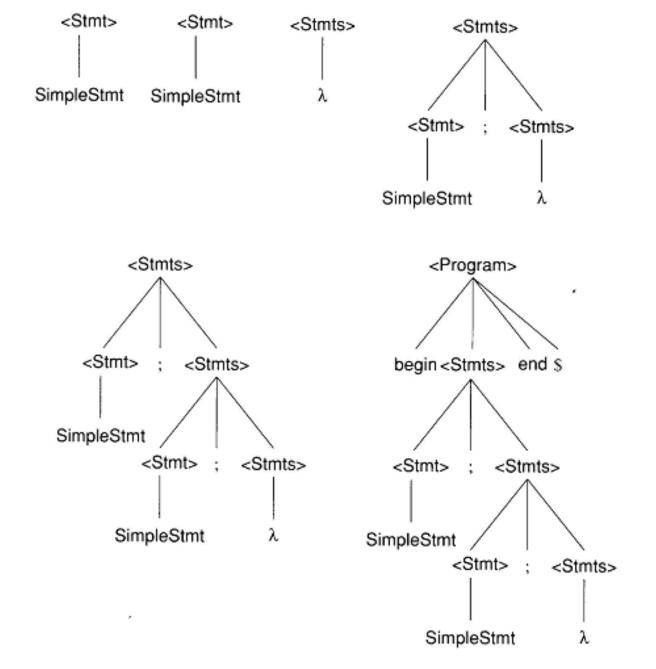
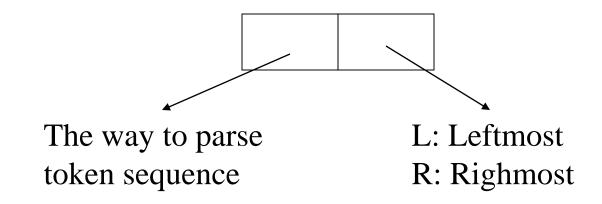


Figure 4.6 A Bottom-Up Parse

Naming of parsing techniques



- Top-down
 - > LL
- Bottom-up
 - > LR

Grammar Analysis Algorithms

- Goal of this section:
 - Discuss a number of important analysis algorithms for Grammars

• The data structure of a grammar G

```
typedef int symbol; /* a symbol in the grammar */
* The symbolic constants used below, NUM TERMINALS,
* NUM NONTERMINALS, and NUM PRODUCTIONS are
* determined by the grammar. MAX RHS LENGTH should
 * simply be "big enough."
 */
#define VOCABULARY (NUM_NONTERMINALS + NUM_TERMINALS)
typedef struct gram {
   symbol terminals[NUM_TERMINALS];
   symbol nonterminals[NUM_NONTERMINALS];
   symbol start symbol;
   int num productions;
   struct prod {
     symbol lhs;
     int rhs length;
     symbol rhs[MAX_RHS_LENGTH];
   } productions[NUM_PRODUCTIONS];
   symbol vocabulary[VOCABULARY];
 } grammar;
typedef struct prod production;
typedef symbol terminal;
typedef symbol nonterminal;
```

• What nonterminals can derive λ ?

$$A \Rightarrow BCD \Rightarrow BC \Rightarrow B \Rightarrow \lambda$$

An iterative marking algorithm

```
typedef short boolean;
typedef boolean marked vocabulary[VOCABULARY];
/*
 * Mark those vocabulary symbols found to
 * derive \lambda (directly or indirectly).
marked vocabulary mark lambda(const grammar g)
   static marked vocabulary derives lambda;
  boolean changes;
                  /* any changes during last iteration? */
  boolean rhs derives lambda;
                  /* does the RHS derive λ? */
   symbol v;
                  /* a word in the vocabulary */
  production p; /* a production in the grammar */
                 /* loop variables */
   int i, j;
   for (v = 0; v < VOCABULARY; v++)
     derives lambda[v] = FALSE;
      /* initially, nothing is marked */
  do {
     changes = FALSE;
      for (i = 0; i < q.num productions; i++) {
         p = g.productions[i];
        if (! derives lambda[p.lhs]) {
            if (p.rhs length == 0) {
               /* derives λ directly */
               changes = derives lambda[p.lhs] = TRUE;
               continue;
            }
            /* does each part of RHS derive λ? */
            rhs derives lambda = derives lambda[p.rhs[0]];
            for (j = 1; j < p.rhs length; j++)
               rhs derives lambda = rhs derives lambda
                           && derives lambda[p.rhs[j]];
            if (rhs derives lambda)
               changes = TRUE;
               derives lambda[p.lhs] = TRUE;
         }
   } while (changes);
   return derives lambda;
}
```

Figure 4.7 Algorithm for Determining If a Nonterminal Can Derive λ

- Follow(A)
 - A is any nonterminal
 - Follow(A) is the set of terminals that my follow A in some sentential form

```
Follow(A)=\{a \in V_t | S \Rightarrow^* \dots Aa \dots \} \cup \{if S \Rightarrow^+ \alpha A \text{ then } \{\lambda\} \text{ else } \emptyset \}
```

- First(α)
 - The set of all the terminal symbols that can begin a sentential form derivable from α
 - If α is the right-hand side of a production, then First(α) contains terminal symbols that begin strings derivable from α

```
First(\alpha)={a \in V_t | \alpha \Rightarrow^* a\beta}\cup {if \alpha \Rightarrow^* \lambda then {\lambda} else \emptyset}
```

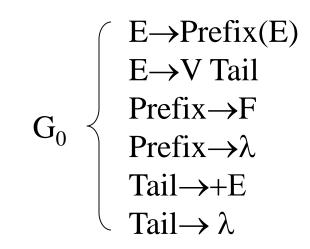
- Definition of C data structures and subroutines
 - first_set[X]
 - contains terminal symbols and λ
 - X is any single vocabulary symbol
 - follow_set[A]
 - contains terminal symbols and λ
 - A is a nonterminal symbol

```
typedef set of terminal or lambda termset;
termset follow set[NUM NONTERMINAL];
termset first set[SYMBOL];
marked vocabulary derives lambda = mark lambda (g);
/* mark lambda(g) as defined above */
termset compute first(string_of_symbols alpha)
 int i, k;
 termset result;
                                       It is a subroutine of
 k = length(alpha);
                                       fill first set()
 if (k == 0)
    result = SET OF(\lambda);
 else {
   result = first set[alpha[0]];
   for (i = 1; i < k && \lambda \in first set[alpha[i-1]]; i++)
    result = result () (first set[alpha[i]] - SET_OF(\lambda));
   if (i == k &   \lambda \in first set[alpha[k-1]])
     result = result () SET_OF(\lambda);
 return result;
```

Figure 4.8 Algorithm to Compute First(alpha)

```
extern grammar g;
void fill first set(void)
  nonterminal A;
  terminal
               a;
  production p;
  boolean
             changes;
  int
              i, j;
  for (i = 0; i < NUM NONTERMINAL; i++) {
     A = q.nonterminals[i];
     if (derives lambda[A])
        first set[A] = SET OF(\lambda);
     else
        first set[A] = \emptyset;
  for (i = 0; i < NUM TERMINAL; i++) {
     a = g.terminals[i];
     first set[a] = SET OF( a );
     for (j = 0; j < NUM NONTERMINAL; j++) {
        A = g.nonterminals[j];
      if (there exists a production A \rightarrow a\beta)
          first set[A] = first set[A] () SET OF( a );
     }
   }
  do {
     changes = FALSE;
     for (i = 0; i < g.num productions; i++) {</pre>
        p = g.productions[i];
        first set[p.lhs] = first set[p.lhs] ()
           compute first(p.rhs);
        if (first set changed)
           changes = TRUE;
  } while (changes);
```

Figure 4.9 Algorithm to Compute First Sets for V

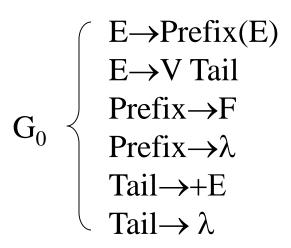


The execution of fill_first_set() using grammar G₀

Step	first_set							
	E	Prefix	Tail	()	V	F	+_
(1) First loop	Ø	{λ}	{λ}					
(2) Second (nested) loop	{V}	{F,λ}	$\{+,\lambda\}$	{(}	{)}	{V}	{ F }	{+}
(3) Third loop, production 1	{V,F,(}	{F,λ}	{+,λ}	{(}	{)}	{ V }	{ F }	{+}

```
void fill follow set (void)
  nonterminal A, B;
   int i;
  boolean
             changes;
   for (i = 0; i < NUM NONTERMINAL; i++) {
     A = q.nonterminals[i];
      follow set[A] = \emptyset;
   follow_set[g.start_symbol] = SET_OF(\lambda);
   do {
      changes = FALSE;
      for (each production A\to\alpha B\beta) {
         /*
          * I.e. for each production and each occurrence
          * of a nonterminal in its right-hand side.
          */
         follow set[B] = follow set[B] ( )
                 (compute first(\beta) - SET_OF(\lambda));
         if ( \lambda \in \mathtt{compute\_first}(\beta) )
            follow_set[B] = follow_set[B] \( \cup \) follow_set[A];
         if ( follow_set[B] changed )
            changes = TRUE;
    } while (changes);
```

Figure 4.10 Algorithm to Compute Follow Sets for All Nonterminals



The execution of fill_follow_set() using grammar G₀

Step	follow_se		et	
	E	Prefix	Tail	
(1) Initialization	{λ}	Ø	Ø	
(2) Process Prefix in production 1	{λ}	{(}	Ø	
(3) Process E in production 1	{λ,)}	{(}	Ø	
(4) Process Tail in production 2	{λ,)}	{(}	<i>{λ,)}</i>	

$$S \rightarrow aSe$$

$S \rightarrow B$

$$B \rightarrow bBe$$

$$B \rightarrow C$$

$$C \rightarrow cCe$$

$$C \rightarrow d$$

More examples

The execution of fill_first_set()

Step	first set							
	S	В	С	a	b	С	d	е
(1) First loop	Ø	Ø	Ø		3777	2		
(2) Second (nested) loop	{a}	{b}	{c,d}	{a}	{b}	{c}	{d}	{e}
(3) Third loop, production 2	{a,b}	{b}	{c,d}	{a}	{b}	{c}	{d}	{e}
(4) Third loop, production 4	{a,b}	{b,c,d}	{c,d}	{a}	{b}	{c}	{d}	{e}
(5) Third loop, production 2	{a,b,c,d}	{b,c,d}	{c,d}	{a}	{b}	{c}	{d}	{e}

The execution of fill_follow_set()

Step	follow_set				
	S	В	С		
(1) Initialization	{λ}	Ø	Ø		
(2) Process S in production 1	{e,λ}	Ø	Ø		
(3) Process B in production 2	{e,λ}	{e,λ}	Ø		
(4) Process B in production 3	no changes				
(5) Process C in production 4	{e,λ}	{e,λ}	{e,λ}		
(6) Process C in production 5	no changes				

 $S \rightarrow ABc$

 $A \rightarrow a$

 $A \rightarrow \lambda$

 $B \rightarrow b$

 $B \rightarrow \lambda$

More examples

The execution of fill_first_set()

Step	first_set						
	S	Α	В	а	b	С	
(1) First loop	Ø	{λ}	{λ}				
(2) Second (nested) loop	Ø	{a,λ}	{b,λ}	{a}	{b}	{c}	
(3) Third loop, production 1	{a,b,c}	{a,λ}	{b,λ}	{a}	{b}	{c}	

The execution of fill_follow_set()

Step	follow_set				
	S	Α	В		
(1) Initialization	{λ}	Ø	Ø		
(2) Process A in production 1	{λ}	{b,c}	Ø		
(3) Process B in production 1	{λ}	{b,c}	{c}		