

Chapter 3

Syntactic Analysis I

- The Syntactic Analyzer, or Parser, is the heart of the front end of the compiler.
- The parser's main task is to analyze the structure of the program and its component statements.
- Our principle resource in Parser Design is the theory of Formal Languages.
- We will use and study context free grammars (They cannot handle definition before use, but we can get around this other ways)

1. Grammars

- **Informal Definition** -- a finite set of rules for generating an infinite set of sentences.
- **Def: Generative Grammar:** this type of grammar builds a sentence in a series of steps, refining each step, to go from an abstract to a concrete sentence.

- **Def: Parse Tree:** a tree that represents the analysis/structure of a sentence (following the refinement steps used by a generative grammar to build it).

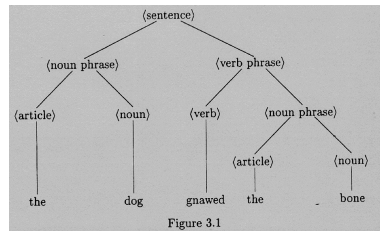


Figure 3.1

- **Def: Productions/Re-Write Rules:** rules that explain how to refine the steps of a generative grammar.

- **Def: Terminals:** the actual words in the language.

- **Def: Non-Terminals:** Symbols not in the language, but part of the refinement process.

1.1 Syntax and Semantics

- Syntax deals with the way a sentence is put together.
- Semantics deals with what the sentence means.
- There are sentences that are grammatically correct that do not make any sense.

- There are things that make sense that are not grammatically correct.
- The compiler will check for syntactical correctness, yet it is the programmers responsibility (usually during debugging) to make sure it makes sense.

1.2 Grammars: Formal Definition

- $G = (T, N, S, R)$
 - ◆ T = set of Terminals
 - ◆ N = set of Non-Terminals
 - ◆ S = Start Symbol (element of N)
 - ◆ R = Set of Rewrite Rules ($\alpha \rightarrow \beta$)

- In your rewrite rules, if α is a single non-terminal the language is Context-Free.
- BNF stands for Backus-Naur Form
 - ◆ $::=$ is used in place of \rightarrow
 - ◆ in extended BNF $\{ \}$ is equivalent to $()^*$

1.3 Parse Trees and Derivations

- $a1 \Rightarrow a2$ -- string $a1$ is changed to string $a2$ via 1 rewrite rule.
- $\alpha \Rightarrow^* \beta$ -- 0 or more re-write rules
- sentential forms -- the strings appearing in various derivation steps
- $L(G) = \{ w \mid S \Rightarrow_G^* w \}$

1.4 Rightmost and Leftmost Derivations

- Which non-terminal do you rewrite-expand when there is more than one to choose from.
 - ◆ If you always select the rightmost NonTerminal to expand, it is a Rightmost Derivation.
- Leftmost and Rightmost derivations are unique.

- **Def:** any sentential form occurring in a leftmost {rightmost} derivation is termed a left {right} sentential form.
- Some parsers construct leftmost derivations and others rightmost, so it is important to understand the difference.

■ Given (pg 72) $G_E = (T, N, S, R)$

◆ $T = \{ i, +, -, *, /, (,) \}$,

◆ $N = \{ E \}$

◆ $S = E$

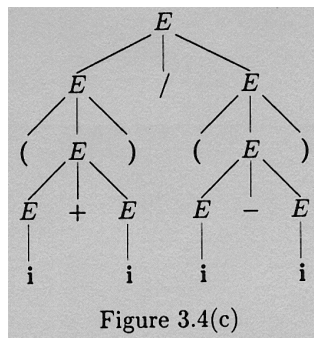
◆ $R = \{$

◆ $E \rightarrow E + E$ $E \rightarrow E - E$

◆ $E \rightarrow E * E$ $E \rightarrow E / E$

◆ $E \rightarrow (E)$ $E \rightarrow i \}$

■ consider: $(i+i)/(i-i)$



1.5 Ambiguous Grammars

■ Given (pg 72) $G_E = (T, N, S, R)$

◆ $T = \{ i, +, -, *, /, (,) \}$,

◆ $N = \{ E \}$

◆ $S = E$

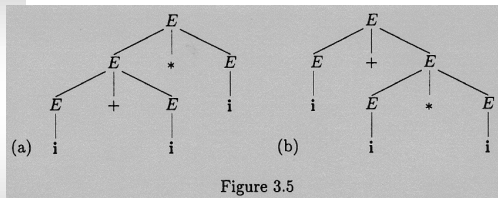
◆ $R = \{$

◆ $E \rightarrow E + E$ $E \rightarrow E - E$

◆ $E \rightarrow E * E$ $E \rightarrow E / E$

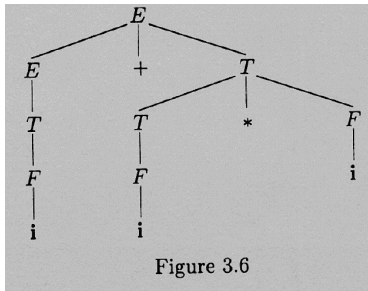
◆ $E \rightarrow (E)$ $E \rightarrow i \}$

■ consider: $i + i * i$



- a grammar in which it is possible to parse even one sentence in two or more different ways is ambiguous
- A language for which no unambiguous grammar exists is said to be inherently ambiguous

- The previous example is "fixed" by operator-precedence rules,
- or re-write the grammar
 - ◆ $E \rightarrow E + T \mid E - T \mid T$
 - ◆ $T \rightarrow T * F \mid T / F \mid F$
 - ◆ $F \rightarrow (E) \mid i$
- Try: $i+i*i$



1.6 The Chomsky Hierarchy (from the outside in)

■ Type 0 grammars

- ◆ $\gamma A \delta \rightarrow \gamma \beta \delta$
- ◆ these are called phrase structured, or unrestricted grammars.
- ◆ It takes a Turing Machine to recognize these types of languages.

■ Type 1 grammars

- ◆ $\gamma A \delta \rightarrow \gamma \beta \delta$
 $\beta \neq \epsilon$
- ◆ therefore the sentential form never gets shorter.
- ◆ Context Sensitive Grammars.
- ◆ Recognized by a simpler Turing machine
 [linear bounded automata (lba)]

■ Type 2 grammars:

- ◆ $A \rightarrow \beta$
- ◆ Context Free Grammars
- ◆ it takes a stack automaton to recognize CFG's (FSA with temporary storage)
- ◆ Nondeterministic Stack Automaton cannot be mapped to a DSA, but all the languages we will look at will be DSA's

■ Type 3 grammars

- ◆ The Right Hand Side may be
 - ◆ a single terminal
 - ◆ a single non-terminal followed by a single terminal.
- ◆ Regular Grammars
- ◆ Recognized by FSA's

1.7 Some Context-Free and Non-Context-Free Languages

■ Example 1:

- ◆ $S \rightarrow S S$
- ◆ $| (S)$
- ◆ $| ()$
- ◆ This is Context Free.

■ Example 2:

◆ $a^n b^n c^n$

■ this is NOT Context Free.

Chapter 3 -- Syntactic Analysis I25

■ Example 3:

◆ $S \rightarrow aSBC$

◆ $S \rightarrow abC$

◆ $CB \rightarrow BC$

◆ $bB \rightarrow bb$

◆ $bC \rightarrow bc$

◆ $cC \rightarrow cc$

■ This is a Context Sensitive Grammar

Chapter 3 -- Syntactic Analysis I26

■ $L_2 = \{wcw \mid w \in (T-c)^*\}$ is NOT a Context Free Grammar.

Chapter 3 -- Syntactic Analysis I27

1.8 More about the Chomsky Hierarchy

- There is a close relationship between the productions in a CFG and the corresponding computations to be carried out by the program being parsed.
- This is the basis of Syntax-directed translation which we use to generate intermediate code.

2. Top-Down parsers

- The top-down parser must start at the root of the tree and determine, from the token stream, how to grow the parse tree that results in the observed tokens.
- This approach runs into several problems, which we will now deal with.

2.1 Left Recursion

- Productions of the form $A \rightarrow A\alpha$ are left recursive.
- No Top-Down parser can handle left recursive grammars.
- Therefore we must re-write the grammar and eliminate both direct and indirect left recursion.

■ How to eliminate Left Recursion (direct)

◆ Given:

- ◆ $A \rightarrow A\alpha_1 \mid A\alpha_2 \mid A\alpha_3 \mid \dots$
- ◆ $A \rightarrow \delta_1 \mid \delta_2 \mid \delta_3 \mid \dots$

◆ Introduce A'

- ◆ $A \rightarrow \delta_1 A' \mid \delta_2 A' \mid \delta_3 A' \mid \dots$
- ◆ $A' \rightarrow \epsilon \mid \alpha_1 A' \mid \alpha_2 A' \mid \alpha_3 A' \mid \dots$

◆ Example:

- ◆ $S \rightarrow Sa \mid b$
- ◆ **Becomes**
- ◆ $S \rightarrow bS'$
- ◆ $S' \rightarrow \epsilon \mid aS'$

■ How to remove ALL Left Recursion.

◆ 1.Sort the nonterminals

◆ 2.for each nonterminal

- ◆ if $B \rightarrow A\beta$
- ◆ and $A \rightarrow \gamma_1 \mid \gamma_2 \mid \gamma_3 \mid \dots$
- ◆ then $B \rightarrow \gamma_1\beta \mid \gamma_2\beta \mid \gamma_3\beta \mid \dots$

◆ 3.After all done, remove immediate left recursion.

■ Example:

- ◆ $S \rightarrow aA \mid b \mid cS$
- ◆ $A \rightarrow Sd \mid e$

becomes

- ◆ $S \rightarrow aA \mid b \mid cS$
- ◆ $A \rightarrow aAd \mid bd \mid cSd \mid e$

■ note: the S in A(3) \rightarrow but it is NOT left recursion

2.2 Backtracking

- One way to carry out a top-down parse is simply to have the parser try all applicable productions exhaustively until it finds a tree.
- This is sometimes called the brute force method.
- It is similar to depth-first search of a graph
- Tokens may have to be put back on the input stream

- Given a grammar:

- ◆ $S \rightarrow ee \mid bAc \mid bAe$
- ◆ $A \rightarrow d \mid cA$

- A Backtracking algorithm will not work properly with this grammar.

- Example: input string is bcde

- ◆ When you see a b you select $S \rightarrow bAc$
- ◆ This is wrong since the last letter is e not c

- The solution is Left Factorization.

- ◆ **Def: Left Factorization:** -- create a new non-terminal for a unique right part of a left factorable production.

- Left Factor the grammar given previously.

- ◆ $S \rightarrow ee \mid bAQ$
- ◆ $Q \rightarrow c \mid e$
- ◆ $A \rightarrow d \mid cA$

3. Recursive-Descent Parsing

- There is one function for each non-terminal these functions try each production and call other functions for non-terminals.
- The stack is invisible for CFG's
- The problem is -- a new grammar requires new code.

■ Example:

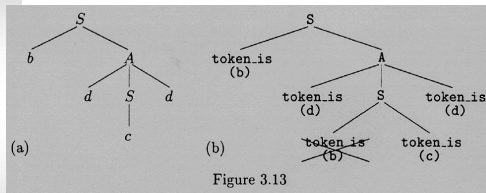
- ◆ $S \rightarrow bA \mid c$
- ◆ $A \rightarrow dSd \mid e$

■ Code:

- | | |
|-------------------------|-------------------------|
| ■ function S: boolean; | ■ else |
| ■ begin | ■ if token_is('c') then |
| ■ S := true; | ■ writeln('S --> c') |
| ■ if token_is('b') then | ■ else |
| ■ if A then | ■ begin |
| ■ writeln('S --> bA') | ■ error('S'); |
| ■ else | ■ S := false |
| ■ S := false; | ■ end |
| | ■ end; { S } |

- | | |
|-------------------------|-------------------------|
| ■ function A: boolean | ■ else |
| ■ begin | ■ A := false |
| ■ A := true; | ■ end |
| ■ if token_is('d') then | ■ else |
| ■ begin | ■ if token_is('e') then |
| ■ if S then | ■ writeln('A --> e') |
| ■ if token_is('d') then | ■ else |
| ■ writeln('A --> dSd'); | ■ begin |
| ■ else | ■ error('A'); |
| ■ begin | ■ A := false |
| ■ error('A'); | ■ end |
| ■ A := false | ■ end; { A } |
| ■ end | |

■ Input String: bcdcd



4. Predictive Parsers

- The goal of a predictive parser is to know which characters on the input string trigger which productions in building the parse tree.
- Backtracking can be avoided if the parser had the ability to look ahead in the grammar so as to anticipate what terminals are derivable (by leftmost derivations) from each of the various nonterminals on the RHS.

■ First (α)

(you construct first() of RHS's)

- ◆ 1. if α begins with a terminal x ,
 - ◆ then $\text{first}(\alpha) = x$.
- ◆ 2. if $\alpha \Rightarrow^* \epsilon$,
 - ◆ then $\text{first}(\alpha)$ includes ϵ .
- ◆ 3. $\text{First}(\epsilon) = \{\epsilon\}$.
- ◆ 4. if α begins with a nonterminal A ,
 - ◆ then $\text{first}(\alpha)$ includes $\text{first}(A) - \{\epsilon\}$

■ Follow(α)

- ◆ 1. if A is the start symbol,
 - ◆ then put the end marker $\$$ into $\text{follow}(A)$.
- ◆ 2. for each production with A on the right hand side

$Q \rightarrow xAy$

- ◆ 1. if y begins with a terminal q ,
 - q is in $\text{follow}(A)$.
- ◆ 2. else $\text{follow}(A)$ includes $\text{first}(y) - \epsilon$.
- ◆ 3. if $y = \epsilon$, or y is nullable ($y \Rightarrow^* \epsilon$)
 - then add $\text{follow}(Q)$ to $\text{follow}(A)$.

■ Grammar:

- ◆ $E \rightarrow T Q$
- ◆ $Q \rightarrow + T Q \mid - T Q \mid \epsilon$
- ◆ $T \rightarrow F R$
- ◆ $R \rightarrow * F R \mid / F R \mid \epsilon$
- ◆ $F \rightarrow (E) \mid I$

■ Construction of First and Follow Sets:

- $\text{First}(E) = \text{First}(T) = \text{First}(F) = \{i, (\}$
- $\text{First}(Q) = \{+, -, \epsilon\}$
- $\text{First}(R) = \{*, /, \epsilon\}$
- $\text{Follow}(E) = \{\$,)\}$
- $\text{Follow}(Q) = \text{Follow}(E) = \{\$,)\}$
- $\text{Follow}(T) = \text{First}(Q) - \epsilon + \text{Follow}(E)$
 $\{+, -\} + \{), \$\}$
- $\text{Follow}(R) = \text{Follow}(T)$
- $\text{Follow}(F) = \text{First}(R) - \epsilon + \text{Follow}(T)$
 $\{*, /\} + \{+, -,), \$\}$

■ LL(1) Grammars

- ◆ In a predictive parser, Follow tells us when to use the epsilon productions.
- ◆ **Def:** LL(1) -- Left to Right Scan of the tokens, Leftmost derivation, 1 token lookahead.

■ For a grammar to be LL(1), we require that for every pair of productions $A \rightarrow \alpha | \beta$

- ◆ 1. $\text{First}(\alpha) - \epsilon$ and $\text{First}(\beta) - \epsilon$ must be disjoint.
- ◆ 2. if α is nullable, then $\text{First}(\beta)$ and $\text{Follow}(A)$ must be disjoint.
- ◆ if rule 1 is violated, we may not know which right hand side to choose
- ◆ if rule 2 is violated, we may not know when to choose β or ϵ .

4.1 A Predictive Recursive-Descent Parser

- The book builds a predictive recursive-descent parser for

- ◆ $E \rightarrow E + T \mid T$
- ◆ $T \rightarrow T * F \mid F$
- ◆ $F \rightarrow (E) \mid I$

- First step is -- Remove Left Recursion

4.2 Table-Driven Predictive Parsers

- Grammar

- ◆ $E \rightarrow E + T \mid E - T \mid T$
- ◆ $T \rightarrow T * F \mid T / F \mid F$
- ◆ $F \rightarrow (E) \mid I$

- Step 1: Eliminate Left Recursion.

- Grammar without left recursion

- ◆ $E \rightarrow T Q$
- ◆ $Q \rightarrow + T Q \mid - T Q \mid \text{epsilon}$
- ◆ $T \rightarrow F R$
- ◆ $R \rightarrow * F R \mid / F R \mid \text{epsilon}$
- ◆ $F \rightarrow (E) \mid I$

- It is easier to show you the table, and how it is used first, and to show how the table is constructed afterward.

■ Table

	i	+	-	*	/	()	\$
E	TQ					TQ		
Q		+TQ	-TQ				ε	ε
T	FR					FR		
R		ε	ε	*FR	/FR		ε	ε
F	i					(E)		

■ Driver Algorithm:

- ◆ Push \$ onto the stack
- ◆ Put a similar end marker on the end of the string.
- ◆ Push the start symbol onto the stack.

◆ While (stack not empty do)

- ◆ Let x = top of stack and a = incoming token.
- ◆ If x is in T (a terminal)
 - ◆ if x == a then pop x and goto next input token
 - ◆ else error
- ◆ else (nonterminal)
 - ◆ if Table[x,a]
 - ◆ pop x
 - ◆ push Table[x,a] onto stack in reverse order
 - ◆ else error

- It is a successful parse if the stack is empty and the input is used up.

■ Example 1: $(i+i)*i$ (pg 108)

	i	+	-	*	/	()	\$
E	TQ					TQ		
Q		+TQ	-TQ				ε	ε
T	FR					FR		
R		ε	ε	*FR	/FR		ε	ε
F	i					(E)		

\$E	(i + i)*i\$	$E \rightarrow TQ$	$E \Rightarrow TQ$
\$QT	(i + i)*i\$		

\$QT	(i + i)*i\$	$T \rightarrow FR$	$\Rightarrow FRQ$
\$QRF	(i + i)*i\$	$F \rightarrow (E)$	$\Rightarrow (E)RQ$
\$QR)E((i + i)*i\$		

\$QR)E	i + i) * i\$	$E \rightarrow TQ$	$\Rightarrow (TQ)RQ$
\$QR)QT	i + i) * i\$	$T \rightarrow FR$	$\Rightarrow (FRQ)RQ$
\$QR)QRF	i + i) * i\$	$F \rightarrow ($	$\Rightarrow ((RQ)RQ$
\$QR)QRi	i + i) * i\$	[pop and go to next token]	
\$QR)QR	i + i) * i\$	$R \rightarrow \epsilon$	$\Rightarrow (iQ)RQ$

\$QR)Q	+ i) * i\$	$Q \rightarrow +TQ$	$\Rightarrow (i+TQ)RQ$
\$QR)QT+	+ i) * i\$	[pop and go to next token]	
\$QR)QT	i) * i\$	$T \rightarrow FR$	$\Rightarrow (i+FRQ)RQ$
\$QR)QRF	i) * i\$	$F \rightarrow ($	$\Rightarrow ((i+RQ)RQ$
\$QR)QRi	i) * i\$	[pop and go to next token]	
\$QR)QR) * i\$	$R \rightarrow \epsilon$	$\Rightarrow ((i+iQ)RQ$
\$QR)Q) * i\$	$Q \rightarrow \epsilon$	$\Rightarrow ((i+i)RQ$
\$QR)) * i\$	[pop and go to next token]	
\$QR	* i\$	$R \rightarrow +FR$	$\Rightarrow ((i+i)*FRQ$
\$QRF*	* i\$	[pop and go to next token]	
\$QRF	i\$	$F \rightarrow ($	$\Rightarrow ((i+i)*iRQ$
\$QRi	\$	[pop and go to next token]	
\$QR	\$	$R \rightarrow \epsilon$	$\Rightarrow ((i+i)*iQ$
\$Q	\$	$Q \rightarrow \epsilon$	$\Rightarrow ((i+i)*i$
\$	\$	[pop and go to next token]	

■ Example 2: (i*) (pg 109)

	i	+	-	*	/	()	\$
E	TQ					TQ		
Q		+TQ	-TQ				ε	ε
T	FR					FR		
R		ε	ε	*FR	/FR		ε	ε
F	i					(E)		

Stack	Input	Production	Derivation
\$E	(i*\$	$E \rightarrow TQ$	$E \Rightarrow TQ$
\$QT	(i*\$	$T \rightarrow FR$	$\Rightarrow FRQ$
\$QRF	(i*\$	$F \rightarrow (E)$	$\Rightarrow (E)RQ$
\$QR)E((i*\$	[pop and go to next token]	
\$QR)E	i*\$	$E \rightarrow TQ$	$\Rightarrow (TQ)RQ$
\$QR)QT	i*\$	$T \rightarrow FR$	$\Rightarrow (FRQ)RQ$
\$QR)QRF	i*\$	$F \rightarrow i$	$\Rightarrow (iRQ)RQ$
\$QR)QRi	i*\$	[pop and go to next token]	
\$QR)QR	*)\$	$R \rightarrow *FR$	$\Rightarrow (i*FRQ)RQ$
\$QR)QRF*	*)\$	[pop and go to next token]	
\$QR)QRF)\$	**Error: no table entry for [F,)].	

4.3 Constructing the Predictive Parser Table

- Go through all the productions.
 $X \rightarrow \beta$ is your typical production.
 - ◆ 1. For all terminals a in $\text{First}(\beta)$, except ϵ ,
 $\text{Table}[X, a] = \beta$.
 - ◆ 2. If $\beta = \epsilon$, or if ϵ is in $\text{first}(\beta)$ then For ALL a in $\text{Follow}(X)$, $\text{Table}[X, a] = \epsilon$.
- So, Construct First and Follow for all Left and right hand sides.

4.4 Conflicts

- A conflict occurs if there is more than 1 entry in a table slot. This can sometimes be fixed by Left Factoring, ...
- If a grammar is LL(1) there will not be multiple entries.

5. Summary

- Left Recursion
- Left Factorization
- First (A)
- Follow (A)
- Predictive Parsers (table driven)
