



Chapter 5 Context-Free Languages

Context-Free Grammars (CFG)
Parsing and Ambiguity
CFG and Programming Languages



Needs for CFG

- RL is effective in describing certain simple pattern
- Regular expressions were inadequate to define all interesting languages, such as `(())`, nested structure in programming languages
- CFL has important applications in the design of programming languages as well as in the construction of efficient compilers



5.1 Context-Free Grammars

Definition 5.1

- A grammar $G = (V, T, S, P)$ is said to be **context-free** if all productions in P have the form
 - $A \rightarrow x$
 - where $A \in V$ and $x \in (V \cup T)^*$.
- A language L is said to be context-free if and only if there is a CFG G such that $L = L(G)$



Grammar Set Notations:

a concise way of design description

$$G = (V, T, S, P)$$

T – a set of terminal symbols

V – a set of nonterminal symbols

S – a element of V, starting symbol

P – a set of production rules with certain format restrictions

Write out the production rule format in set notation for each of the following grammars :

- RG
- CFG



Context Free Grammar (CFG): the representation defining CFL

- A grammar is said to be context-free if **all productions** in P have the form
 - $A \rightarrow x$
 - **English description:**
 - where A is a single nonterminal and x can be any string of terminals or/and nonterminals
 - Single nonterminal of the left handside of the production rule makes the grammar context free
 - **Set notation:**
 - where $A \in V$ and $x \in (V \cup T)^*$



Regular Grammars (RG):

one of 4 representations of RL

- A subset of CFG
- A CFG is said to be regular grammar if all productions in P have the form
 - $A \rightarrow xB, A \rightarrow x;$ or $A \rightarrow Bx, A \rightarrow x$
 - **English:**
 - Where A is a single nonterminal and x can be word (string of terminals) or semiword (a string of terminals ended or starting with a nonterminals)
 - **Set notation:**
 - Where $A, B \in V$, and $x \in T^*$



Example 5.1

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

- is context-free with productions
 - $S \rightarrow aSa$
 - $S \rightarrow bSb$
 - $S \rightarrow \lambda$
- A derivation in this grammar is
 - $S \Rightarrow aSa \Rightarrow aaSaa \Rightarrow aabSbaa \Rightarrow aabbbaa.$
 - We can derive that $L(G) = \{ww^R : w \in \{a, b\}^*\}.$
 - The language is context-free, but is not regular (as shown in Example 4.8)

Example 5.2

- The grammar G is context-free with productions
 - $S \rightarrow abB$
 - $A \rightarrow aaBb$
 - $B \rightarrow bbAa$
 - $A \rightarrow \lambda$
- Show that
 - $L(G) = \{ab(bbaa)^n bba(ba)^n : n \geq 0\}.$



Example 5.3 -- 1

To show $L = \{a^n b^m : n \neq m\}$ is context-free

- We need to produce a CFG for L
- We know how to produce a CFG with $n = m$
 - Write out
- Now we need to add 2 sets of productions
 - For the case $n > m$ we add the first set
 - $S \rightarrow AS_1$
 - $S_1 \rightarrow a S_1 b | \lambda$
 - $A \rightarrow aA | a$



Example 5.3 -- 2

To show $L = \{a^n b^m : n \neq m\}$ is context-free

- For the case $n < m$, we add another similar set, and we get the answer
 - $S \rightarrow AS_1 \mid S_1B$
 - $S_1 \rightarrow a S_1 b \mid \lambda$
 - $A \rightarrow aA \mid a$
 - $B \rightarrow bB \mid b$
- The resulting grammar is context-free, hence L is CFL.

Example 5.4

programming language L_2 include strings such as $(())$ and $()()()$

- CFL L can be generated by CFG with productions
 - $S \rightarrow aSb | SS | \lambda$
 - Substitute a with $($ and b with $)$ in the productions, we can generate L_2 above
 - $L = \{w \in \{a, b\}^* : n_a(w) = n_b(w) \text{ and } n_a(v) \geq n_b(v), \text{ where } v \text{ is any prefix of } w\}$



Leftmost and Rightmost Derivations

- Definition 5.2

- A derivation is said to be **leftmost** if in each step the leftmost variable in the sentential form is replaced. If in each step the right most variable is replaced, we call the derivation **rightmost**.

Example 5.5

Leftmost and Rightmost Derivations

- $S \rightarrow aAB$
- $A \rightarrow bBb$
- $B \rightarrow A|\lambda$
 - Leftmost derivation of abbbb
 - $S \Rightarrow aAB \Rightarrow abBbB \Rightarrow abAbB \Rightarrow abbBbbB \Rightarrow abbbbB \Rightarrow abbbb$
 - Rightmost derivation of abbbb
 - $S \Rightarrow aAB \Rightarrow aA \Rightarrow abBb \Rightarrow abAb \Rightarrow abbBbb \Rightarrow abbbb$

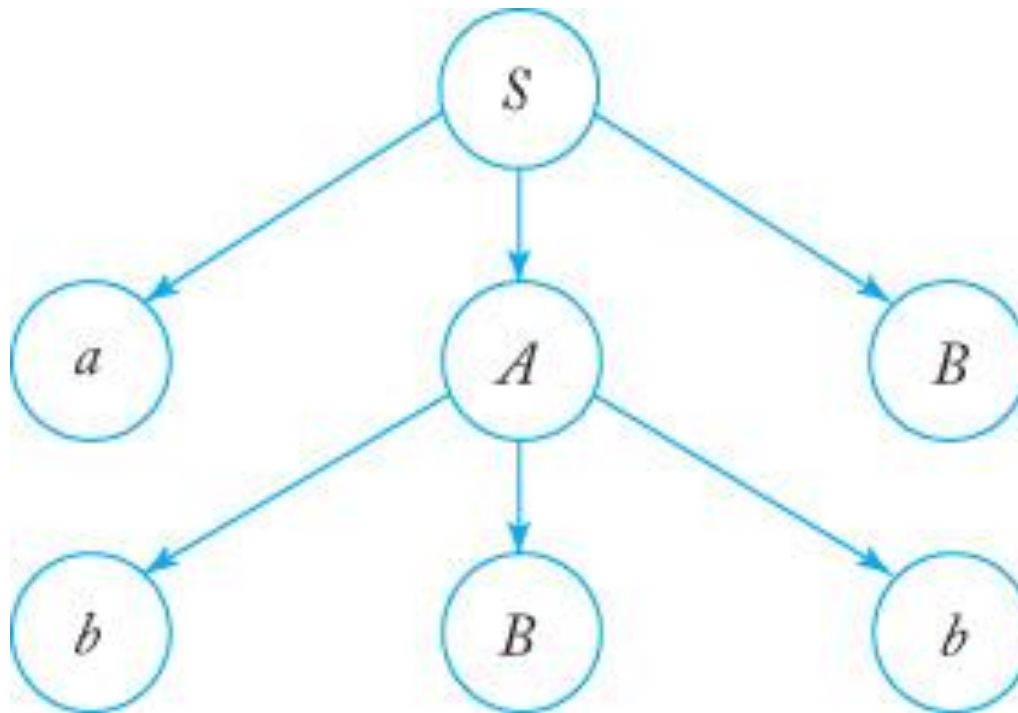


5.2 Derivation Trees/Parse Trees

- A second way of showing derivations, independent of the order in which productions are used
- Definition 5.3 (page 130, a formal definition of **derivation tree** and **partial derivation tree**)

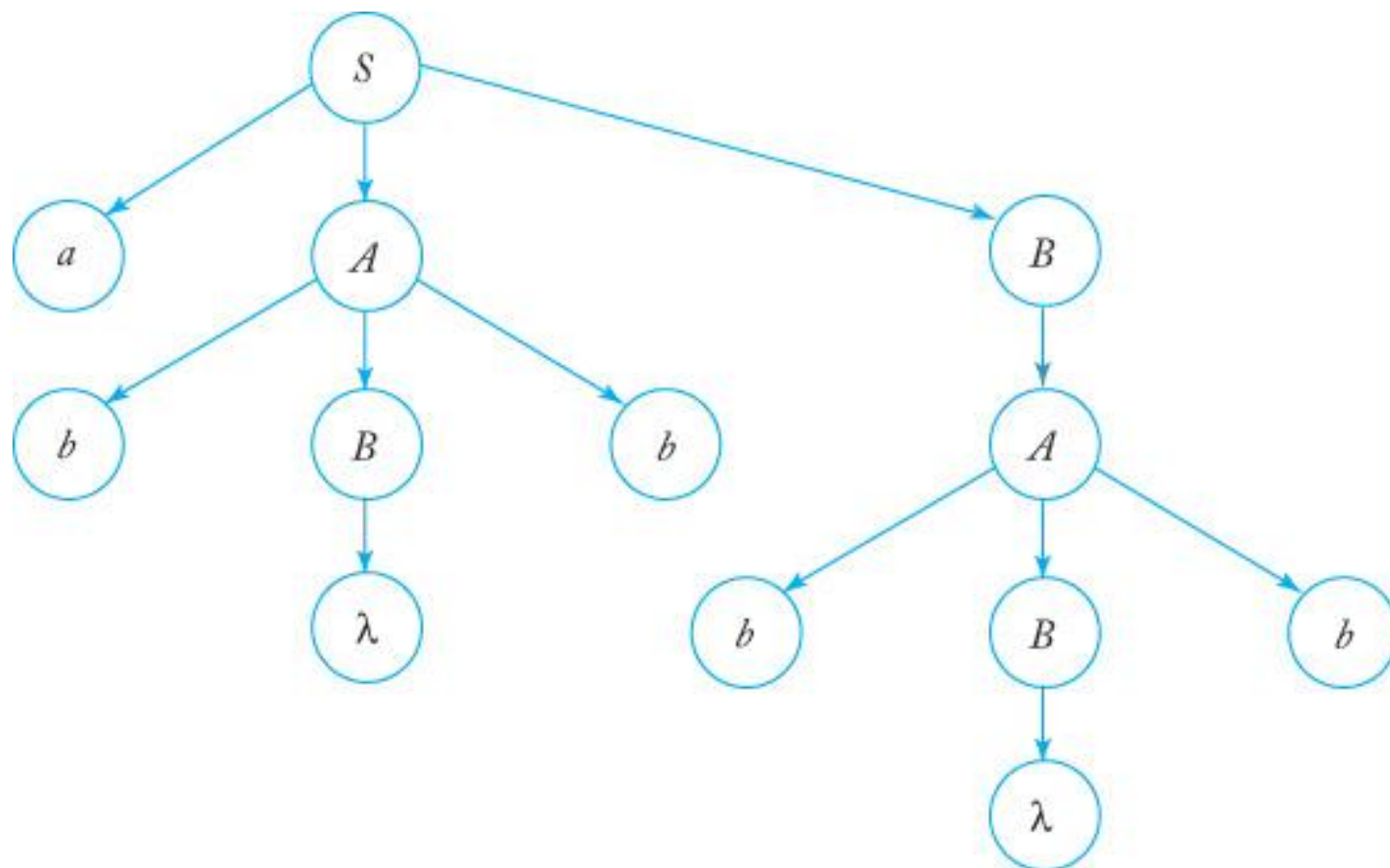
Example 5.6: Partial derivation tree

$S \rightarrow aAB, A \rightarrow bBb, B \rightarrow A|\lambda$



Example 5.6: a complete Derivation tree

$S \rightarrow aAB, A \rightarrow bBb, B \rightarrow A|\lambda$



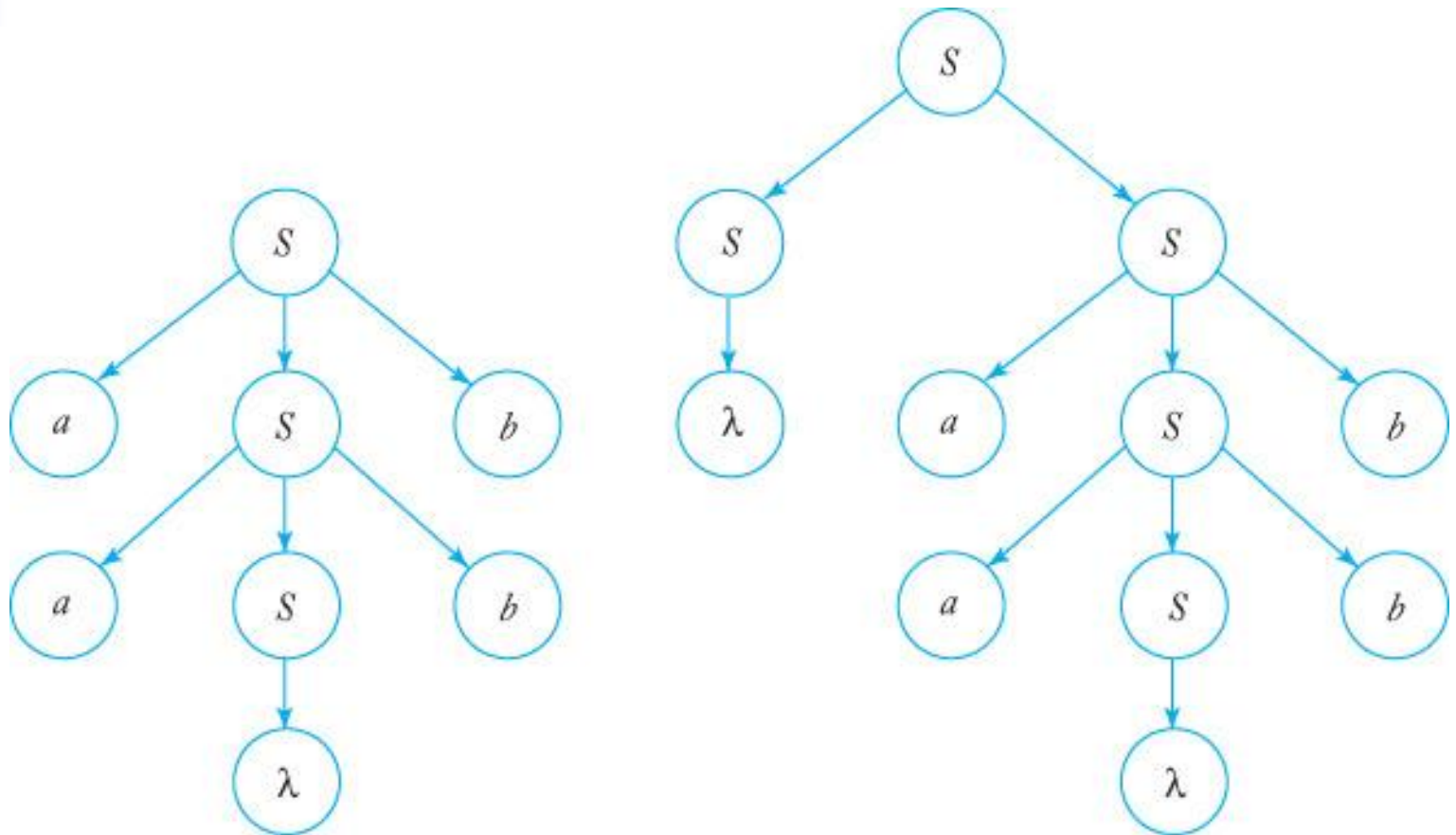


Parsing and Ambiguity

- **Parsing** – finding a sequence of productions by which a $w \in L(G)$ is derived
- Definition 5.5
 - A CFG G is said to be **ambiguous** if there exists some $w \in L(G)$ that has at least two distinct derivation trees. Alternatively ambiguity implies the existence of two or more leftmost or rightmost derivations

Example 5.10

Grammar $S \rightarrow aSb \mid SS \mid \lambda$, is ambiguous since $aabb$ has the two derivation trees shown



Example 5.11 -- 1

$G = (V, T, E, P)$ with

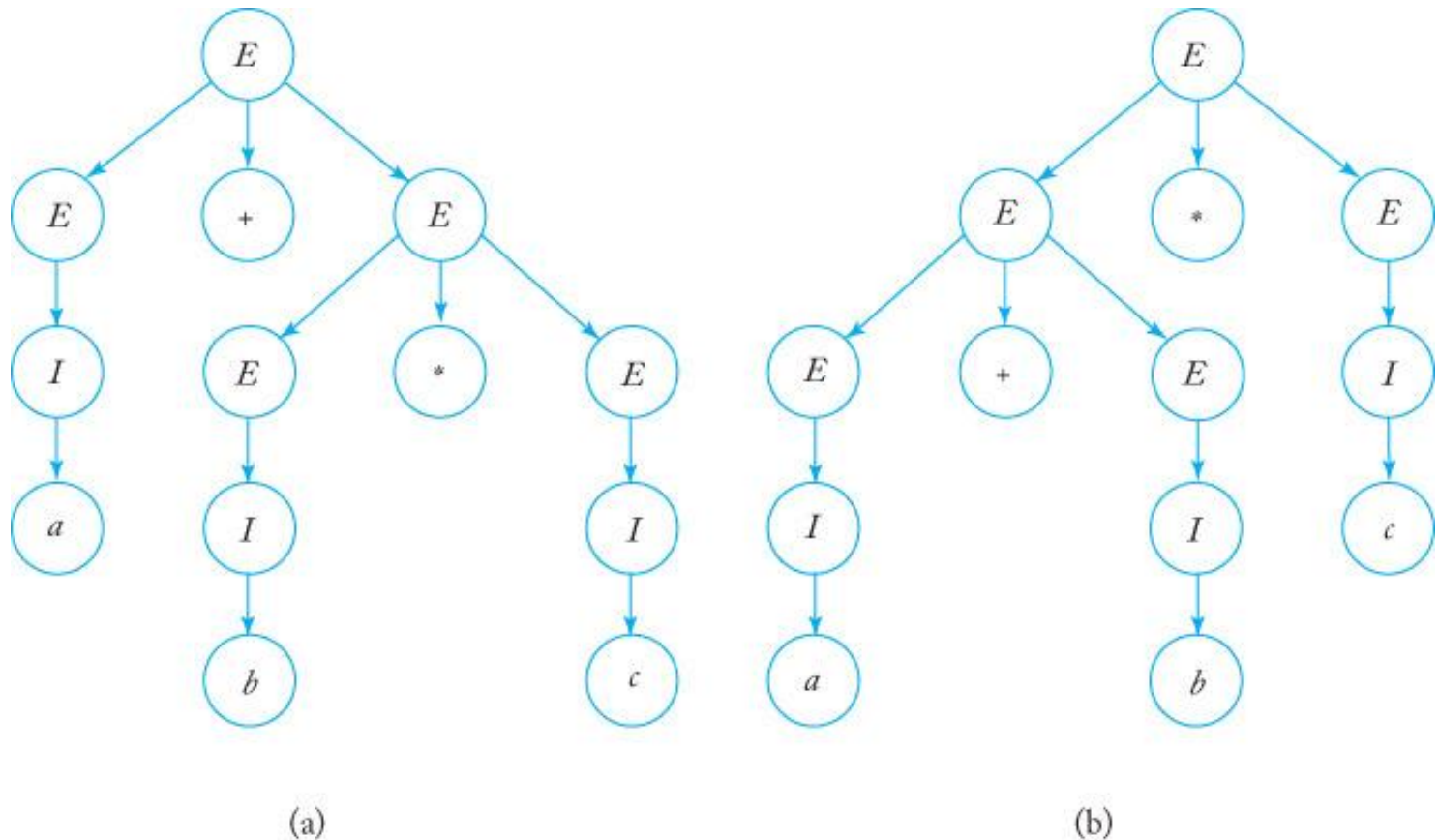
$V = \{E, I\}$, $T = \{a, b, c, +, *, (,)\}$

- and productions
 - $E \rightarrow I \mid E + E \mid E * E \mid (E)$
 - $I \rightarrow a \mid b \mid c$
- The grammar is ambiguous since $a + b * c$ has two different derivation trees
 - Try to draw the two different derivation trees

Example 5.11 – 2

$E \rightarrow I \mid E+E \mid E^*E \mid (E), I \rightarrow a \mid b \mid c$

$a+b*c$ has two different derivation trees



Example 5.12 -- 1

Rewrite the grammar

- One way to resolve the ambiguity is to associate precedence rules with the operators + and *
- * has higher precedence than +
 - + is one level closer to root of the derivation tree, E
- We introduce new variables
 - $V = \{E, T, F, I\}$
 - Replace productions in such a way that + is closer to E than * in the derivation tree

Example 5.12 – 2

Rewrite the productions

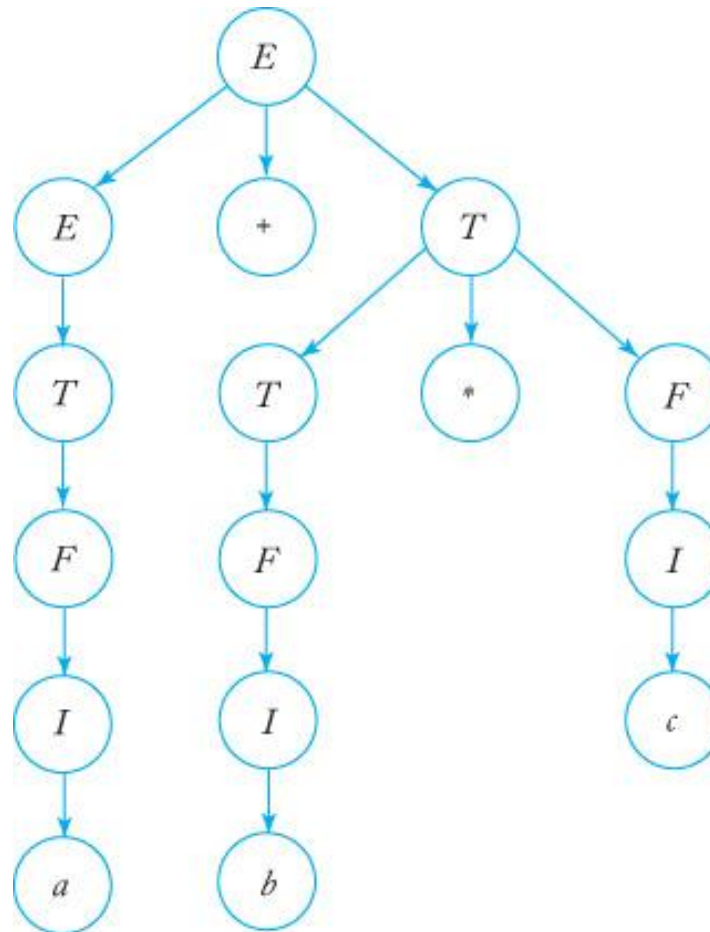
- ambiguous
- $E \rightarrow I \mid E + E \mid E * E \mid (E)$
- $I \rightarrow a \mid b \mid c$

Rewrite to remove ambiguity

- $E \rightarrow E + T \mid T$
- $T \rightarrow T * F \mid F$
- $F \rightarrow (E) \mid I$
- $I \rightarrow a \mid b \mid c$

Example 5.12 – 3:

a unique derivation tree of $a+b*c$ (using new grammar)



A Practical Problem

Assume the following rules of associativity and precedence for expressions

- Precedence:

- Highest *****, **/**, **not**
- **+**, **-**, **&**, **mod**
- **-** (unary)
- **=**, **≠**, **<**, **≤**, **≥**, **>**
- **and**
- Lowest **or**, **xor**

- Associativity: left to right

- Write a CFG for the expression. Assume the only operands are the names a, b, c, d, and e.

Recursion to specify associativity

Precedence cascade to specify precedence

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle \text{ or } \langle e1 \rangle \mid \langle \text{expr} \rangle \text{ xor } \langle e1 \rangle \mid \langle e1 \rangle$
 $\langle e1 \rangle ::= \langle e1 \rangle \text{ and } \langle e2 \rangle \mid \langle e2 \rangle$
 $\langle e2 \rangle ::= \langle e2 \rangle = \langle e3 \rangle \mid \langle e2 \rangle \neq \langle e3 \rangle \mid \langle e2 \rangle < \langle e3 \rangle$
 $\mid \langle e2 \rangle \leq \langle e3 \rangle \mid \langle e2 \rangle > \langle e3 \rangle \mid \langle e2 \rangle \geq \langle e3 \rangle \mid \langle e3 \rangle$
 $\langle e3 \rangle ::= \langle e4 \rangle \mid -\langle e4 \rangle$
 $\langle e4 \rangle ::= \langle e4 \rangle + \langle e5 \rangle \mid \langle e4 \rangle - \langle e5 \rangle \mid \langle e4 \rangle \& \langle e5 \rangle \mid \langle e4 \rangle$
 $\text{mod } \langle e5 \rangle \mid \langle e5 \rangle$
 $\langle e5 \rangle ::= \langle e5 \rangle * \langle e6 \rangle \mid \langle e5 \rangle / \langle e6 \rangle \mid \text{not } \langle e5 \rangle \mid \langle e6 \rangle$
 $\langle e6 \rangle ::= a \mid b \mid c \mid d \mid e \mid \text{const} \mid (\langle \text{expr} \rangle)$



CFG and Programming Languages

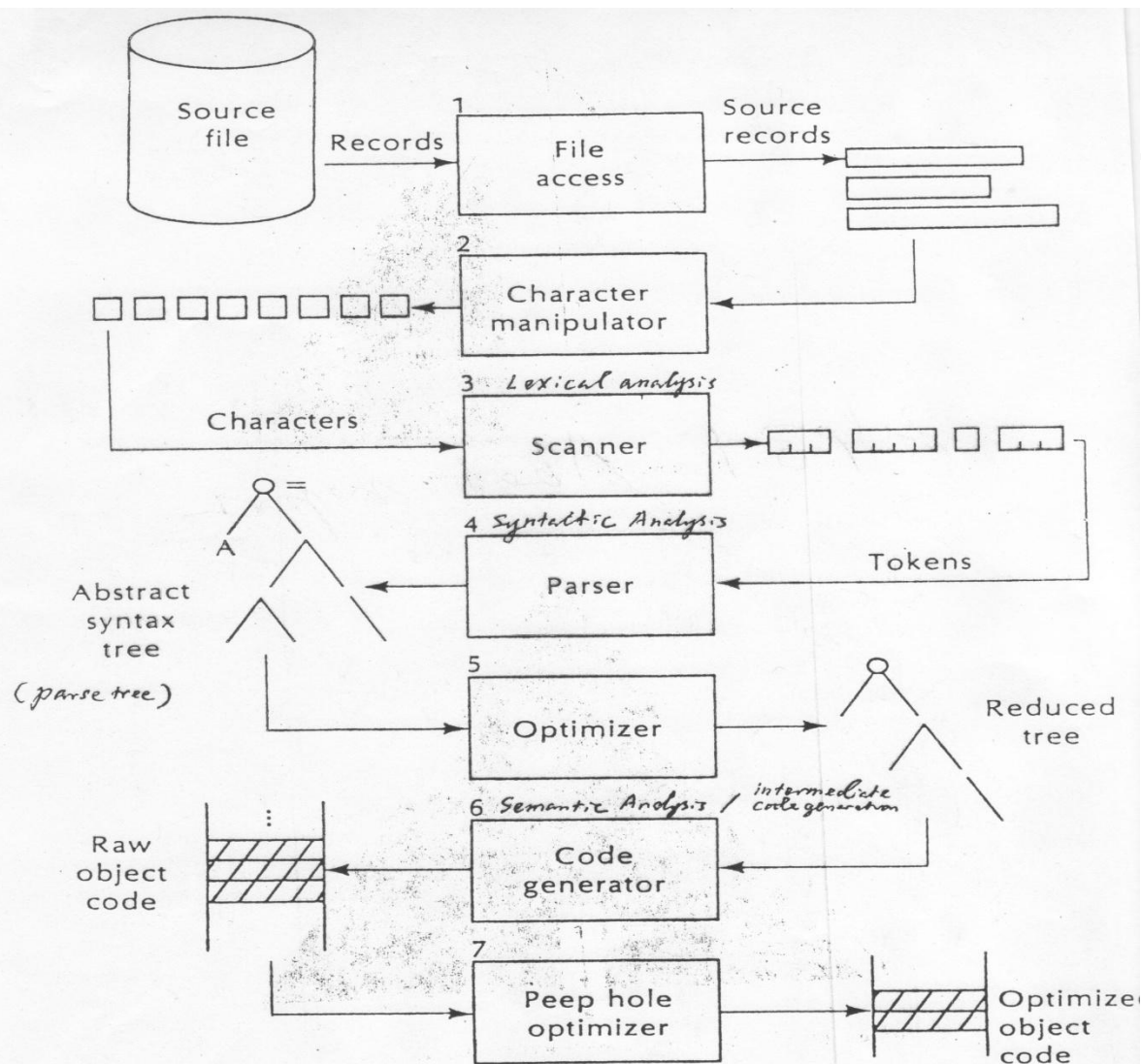
- One of the most important uses of the computing theory is in
 - the definition of programming language
 - the construction of interpreters/compiler
 - Both RL and CFL are used in model all aspects
 - We can define a programming language by a grammar



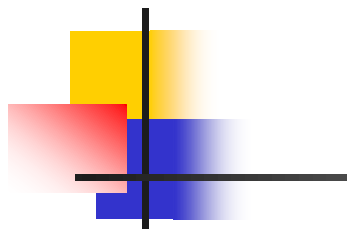
BNF and CFG

- A traditional notation in writing on programming languages is called the Backus-Naur form or BNF
- BNF is in essence the same as notation we use for CFG here
 - $\langle \text{expression} \rangle ::= \langle \text{term} \rangle \mid \langle \text{expression} \rangle + \langle \text{term} \rangle$
 - $\langle \text{term} \rangle ::= \langle \text{factor} \rangle \mid \langle \text{term} \rangle * \langle \text{factor} \rangle$
 - $\langle \text{while_statement} \rangle ::= \text{while} \langle \text{expression} \rangle \langle \text{statement} \rangle$

Major Operations in a Compiler



The Phases of a Compiler



SYMBOL TABLE

1	position	...
2	initial	...
3	rate	...
4		

position := initial + rate * 60

lexical analyzer

id₁ := id₂ + id₃ * 60

syntax analyzer

id₁ := id₂ + id₃ * 60

semantic analyzer

id₁ := id₂ + id₃ * inttoreal
60

intermediate code generator

temp1 := inttoreal(60)
temp2 := id3 * temp1
temp3 := id2 + temp2
id1 := temp3

code optimizer

temp1 := id3 * 60.0
id1 := id2 + temp1

code generator

MOVF id3, R2
MULF #60.0, R2
MOVF id2, R1

Parsing



In-class exercise:

for given grammar, show
BNF, EBNF, Syntax Diagram, first and follow sets

- Write G_3 in EBNF
- Draw syntax diagrams
- What are the 2 requirements for a predictive parser to be able to distinguish between choices in the grammar rule?
- Give the first sets/follow sets of G_3 .
- Do we need to compute follow set for G_3 ?



G_3 in BNF \rightarrow EBNF -- 1

- To remove recursion and ε
- $V \rightarrow SR\$$
- $S \rightarrow + \mid - \mid \varepsilon$
- $R \rightarrow .dN \mid dN.N$
- $N \rightarrow dN \mid \varepsilon$



BNF \rightarrow EBNF -- 2

- To remove recursion and ε
 - $V \rightarrow SR\$$ $V ::= SR\$$
 - $S \rightarrow + \mid - \mid \varepsilon$ $S ::= [+ \mid -]$
 - $R \rightarrow .dN \mid dN.N$ $R ::= .dN \mid dN.N$
 - $N \rightarrow dN \mid \varepsilon$ $N ::= \{d\}$
-
- Draw syntax diagrams for each variable
 - Can we simplify EBNF more?



What are the 2 requirements for a predictive parser to be able to distinguish between choices in the grammar rule? (Hint: use first and follow sets)

1. Given $A \rightarrow B \mid C$
2. Given $A \rightarrow D[E]F$



Two Requirements

(1) Given $A \rightarrow B \mid C$

$$\text{First}(B) \cap \text{first}(C) = \emptyset$$

(2) Given $A \rightarrow D[E]F$

$$\text{First}(E) \cap \text{follow}(E) = \emptyset$$



EBNF to Syntax Diagram

first sets and follow sets -1

- $V ::= SR\$$
- $S ::= [+ \mid -]$
- $R ::= .dN \mid dN.N$
- $N ::= \{d\}$



first sets and follow sets - 2

- $V ::= SR\$$
 $\text{first}(V) = \text{first}(S) \cup \text{first}(R)$
 $= \{+, -, ., d\}$
- $S ::= [+ \mid -]$
 $\text{first}(S) = \{+, -\}$
- $R ::= .dN \mid dN.N$
 $\text{first}(R) = \{., d\}$
- $N ::= \{d\}$
 $\text{first}(N) = \{d\}$

- $\text{Follow}(V) = \{\}, \text{Follow}(R) = \{\$\}$
- $\text{Follow}(S) = \text{first}(R) = \{., d\}$
- $\text{Follow}(N) = \{.\} \cup \text{follow}(R) = \{., \$\}$



EBNF

using substitution to simplify

- $V ::= SR\$$
 - $S ::= [+ \mid -]$
 - $R ::= .dN \mid dN.N$
 - $N ::= \{d\}$
-
- $V ::= [+|-] R \$$
 - $R ::= .d\{d\} \mid d\{d\}.\{d\}$

A Top-down Parser - 1

$G_3 = (\{V, S, R, N\}, \{+, -, ., d, \perp\}, P, V)$, where P is:

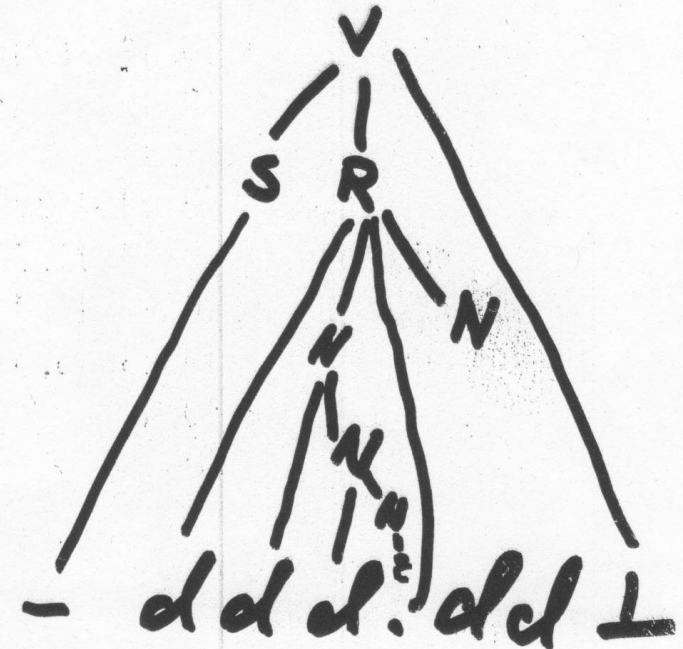
1. $V \rightarrow SR\perp$ $\{\perp \text{ is a stop symbol}\}$
2. $S \rightarrow +$
3. $S \rightarrow -$
4. $S \rightarrow \epsilon$ $\{\epsilon \text{ is the empty string}\}$
5. $R \rightarrow .dN$ $\{d \text{ is a decimal digit}\}$
6. $R \rightarrow dN.N$
7. $N \rightarrow dN$
8. $N \rightarrow \epsilon$

		Next token				
		+	-	.	d	\perp
Left-most Exposed Nonterminal	V	1	1	1	1	×
	S	2	3	4	4	×
	R	×	×	5	6	×
	N	×	×	8	7	8

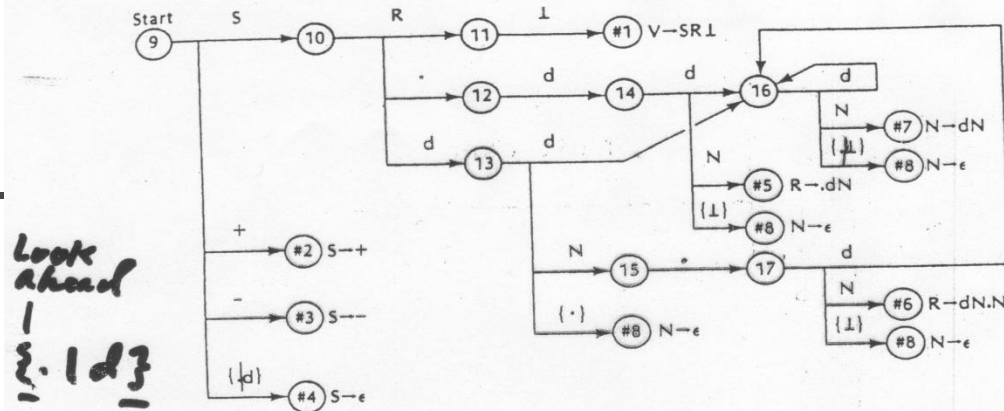
Figure 2.16. A top-down LL(1) parsing table for grammar G_3 .

A Top-down Parser - 2

Frontier	Remaining Input	Production
V	- ddd.dd ⊥	1
SR ⊥	- ddd.dd ⊥	3
- R ⊥	- ddd.dd ⊥	(match, drop -)
R ⊥	ddd.dd ⊥	6
dN.N ⊥	ddd.dd ⊥	(match)
N.N ⊥	dd.dd ⊥	7
dN.N ⊥	dd.dd ⊥	(match)
N.N ⊥	d.dd ⊥	7
dN.N ⊥	d.dd ⊥	(match)
N.N ⊥	.dd ⊥	8
.N ⊥	.dd ⊥	(match)
N ⊥	dd ⊥	7
dN ⊥	dd ⊥	(match)
N ⊥	d ⊥	7
dN ⊥	d ⊥	(match)
N ⊥	⊥	8
⊥	⊥	(match and halt)



A Bottom-up Parser



$\{.1L\}$

Look ahead
|
 $\{.1d\}$

Figure 2.17. A bottom-up parsing machine for grammar G_3 .

$LR(1)$

Rootline	State Path	Production
ddd.dd┐	9	4, $S \rightarrow \epsilon$
Sddd.dd┐	9, 10, 13, 16, 16	8, $N \rightarrow \epsilon$
SdddN.dd┐	9, 10, 13, 16, 16	7, $N \rightarrow dN$
SddN.dd┐	9, 10, 13, 16	7, $N \rightarrow dN$
SdN.dd┐	9, 10, 13, 15, 17, 16, 16	8, $N \rightarrow \epsilon$
SdN.ddN┐	9, 10, 13, 15, 17, 16, 16	7, $N \rightarrow dN$
SdN.dN┐	9, 10, 13, 15, 17, 16	7, $N \rightarrow dN$
SdN.N┐	9, 10, 13, 15, 17	6, $R \rightarrow dN.N$
SR┐	9, 10, 11	1, $V \rightarrow SR┐$
V	(halt)	

It should be clear that we have reproduced a right-most derivation of the sentence, in reverse order.

HW
- d. d┐ , + d. d┐

