

BCSE307L – COMPILER DESIGN

Text Book(s)	
1.	A. V. Aho, Monica S. Lam, Ravi Sethi and Jeffrey D. Ullman, Compilers: Principles, techniques, & tools, 2007, Second Edition, Pearson Education, Boston.
Reference Books	
1.	Watson, Des. A Practical Approach to Compiler Construction. Germany, Springer International Publishing, 2017.
Mode of Evaluation: CAT, Quiz, Written assignment and FAT	

Module:2	SYNTAX ANALYSIS	8 hours
Role of Parser- Parse Tree - Elimination of Ambiguity – Top Down Parsing - Recursive Descent Parsing - LL (1) Grammars – Shift Reduce Parsers- Operator Precedence Parsing - LR Parsers, Construction of SLR Parser Tables and Parsing- CLR Parsing- LALR Parsing.		

Syntax analysis

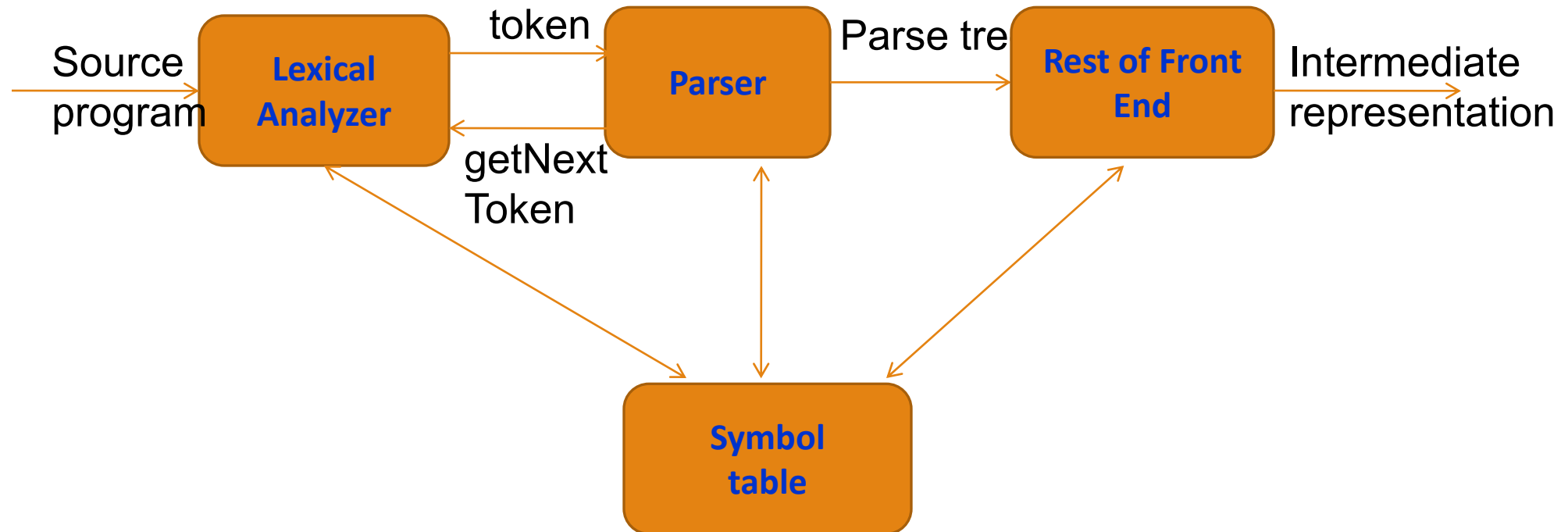
Introduction to Syntax Analysis

CFG

Types of Parser

The role of parser

- Parser works on a stream of tokens.
- The smallest item is a token.



Syntax Analyzer

Syntax Analyzer creates the syntactic structure of the given source program.

This syntactic structure is mostly a *parse tree*.

Syntax Analyzer is also known as *parser*.

The syntax of a programming is described by a *context-free grammar (CFG)*.

Syntax Analyzer(contd..)

The **syntax analyzer** (parser) checks whether a given source program satisfies the rules implied by a context-free grammar or not.

- If it satisfies, the parser creates the parse tree of that program.
- Otherwise the parser gives the error messages.

A **context-free grammar**

- gives a precise syntactic specification of a programming language.
- the design of the grammar is an initial phase of the design of a compiler.
- a grammar can be directly converted into a parser by some tools.

Expression Grammar

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

Parsers

We categorize the parsers into two groups:

1. Top-Down Parser

- the parse tree is created top to bottom, starting from the root.

2. Bottom-Up Parser

- the parse is created bottom to top; starting from the leaves

Both top-down and bottom-up parsers scan the input from left to right (one symbol at a time).

Efficient top-down and bottom-up parsers can be implemented only for sub-classes of context-free grammars.

- LL for top-down parsing
- LR for bottom-up parsing

Syntax Error Handling

- Lexical Errors
- Syntactic Errors
- Semantic Errors
- Logical Errors

Syntax Error Handling

Lexical Errors

- Misspellings of identifiers, Keywords, Operators....
- Ex: identifier `ellipseSize` instead of `ellipseSize`

Syntactic Errors

- Misplaced semicolons or missing braces

Semantic Errors

- Mismatches between operators and operands

Logical Errors

- It can be anything from incorrect reasoning on the part of the programmer

Error-recovery strategies

Panic mode recovery

- Discard input symbol one at a time until one of designated set of synchronization tokens is found

Phrase level recovery

- Replacing a prefix of remaining input by some string that allows the parser to continue

Error productions

- Augment the grammar with productions that generate the erroneous constructs

Global correction

- Choosing minimal sequence of changes to obtain a globally least-cost correction

Syntax Tree

A parse tree is a record of the rules (and tokens) used to match some input text

A syntax tree records the structure of the input and is insensitive to the grammar that produced it.

Syntax tree

Parse tree: interior nodes are non-terminals, leaves are terminals

Syntax tree: interior nodes are “operators”, leaves are operands

Parse tree: rarely constructed as a data structure

Syntax tree: when representing a program in a tree structure usually use a syntax tree

Parse tree: Represents the *concrete syntax* of a program

Syntax tree: Represents the *abstract syntax* of a program (the semantics)

A syntax tree is often called *abstract syntax tree* or **AST**

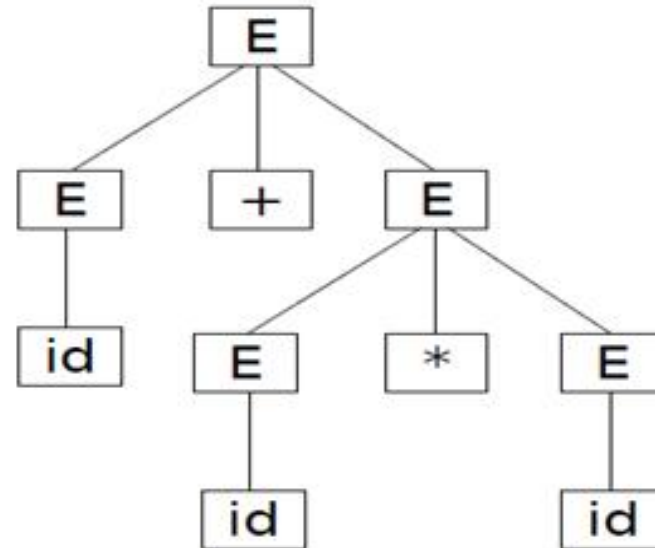
Syntax and parse tree examples

Grammar:

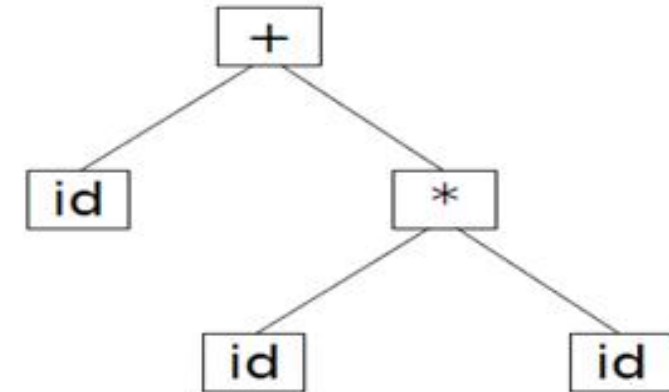
$E \rightarrow E * E$
 $\quad | E + E$
 $\quad | id$

Program: $a + b * c$

Parse tree



Syntax tree



CONTEXT FREE GRAMMAR (CFG)

- ❑ The formal definition of a context-free grammar
- ❑ Notational Conventions
- ❑ Derivations
- ❑ Parse Tree and Derivations
- ❑ Ambiguity
- ❑ Verifying the language generated by a Grammar
- ❑ Context-free grammar versus regular expression

Context-Free Grammars

Inherently recursive structures of a programming language are defined by a context-free grammar.

In a **context-free grammar** $G = (V, T, P, S)$, we have:

- **T** - A finite set of **terminals** (in our case, this will be the set of tokens)
- **V** - A finite set of **non-terminals** (syntactic-variables)
- **P** - A finite set of **productions** rules in the following form
 - $A \rightarrow \alpha$ where A is a non-terminal and α is a string of terminals and non-terminals (including the empty string)
- **S** - A **start symbol** (one of the non-terminal symbol)

Notational Conventions

1. These symbols are terminals:

- (a) Lowercase letters early in the alphabet, such as a , b , c .
- (b) Operator symbols such as $+$, $*$, and so on.
- (c) Punctuation symbols such as parentheses, comma, and so on.
- (d) The digits $0, 1, \dots, 9$.
- (e) Boldface strings such as **id** or **if**, each of which represents a single terminal symbol.

2. These symbols are nonterminals:

- (a) Uppercase letters early in the alphabet, such as A , B , C .
- (b) The letter S , which, when it appears, is usually the start symbol.
- (c) Lowercase, italic names such as *expr* or *stmt*.
- (d) When discussing programming constructs, uppercase letters may be used to represent nonterminals for the constructs. For example, nonterminals for expressions, terms, and factors are often represented by E , T , and F , respectively.

Notational Conventions

3. Uppercase letters late in the alphabet, such as X , Y , Z , represent *grammar symbols*; that is, either nonterminals or terminals.
4. Lowercase letters late in the alphabet, chiefly u, v, \dots, z , represent (possibly empty) strings of terminals.
5. Lowercase Greek letters, α, β, γ for example, represent (possibly empty) strings of grammar symbols. Thus, a generic production can be written as $A \rightarrow \alpha$, where A is the head and α the body.
6. A set of productions $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k$ with a common head A (call them *A-productions*), may be written $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$. Call $\alpha_1, \alpha_2, \dots, \alpha_k$ the *alternatives* for A .
7. Unless stated otherwise, the head of the first production is the start symbol.

Example

Using these conventions, the grammar of Example

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

The notational conventions tell us that E , T , and F are nonterminals, with E the start symbol. The remaining symbols are terminals.

Expression Grammar Example:

$E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid - E$

$E \rightarrow (E)$

$E \rightarrow \text{id}$

CFG - Terminology

$L(G)$ is *the language of G* (the language generated by G) which is a set of sentences.

A sentence of $L(G)$ is a string of terminal symbols of G.

If S is the start symbol of G then

ω is a sentence of $L(G)$ iff $S \xRightarrow{+} \omega$ where ω is a string of terminals of G.

If G is a context-free grammar, $L(G)$ is a *context-free language*.

Two grammars are *equivalent* if they produce the same language.

$S \xRightarrow{*} \alpha$

- If α contains **non-terminals**, it is called as a **sentential form** of G.
- If α **does not contain non-terminals**, it is called as a **sentence** of G.

Derivations

A sequence of replacements of non-terminal symbols is called a **derivation** of id+id from E.

$\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$ (α_n derives from α_1 or α_1 derives α_n)

\Rightarrow^* : derives in one step

\Rightarrow^+ : derives in zero or more steps

\Rightarrow : derives in one or more steps

Derivations

$$E \rightarrow E + E \mid E * E \mid - E \mid (E) \mid \mathbf{id} \quad \underline{\hspace{1cm}}$$

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(\mathbf{id})$$

$$\alpha A \beta \Rightarrow \alpha \gamma \beta$$

1. $\alpha \xRightarrow{*} \alpha$, for any string α , and
2. If $\alpha \xRightarrow{*} \beta$ and $\beta \Rightarrow \gamma$, then $\alpha \xRightarrow{*} \gamma$.

Derivations

$$E \Rightarrow E+E$$

$E+E$ derives from E

- we can replace E by $E+E$
- to be able to do this, we have to have a production rule $E \rightarrow E+E$ in our grammar.

$$E \Rightarrow E+E \Rightarrow id+E \Rightarrow id+id$$

Derivation Example

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(id+E) \Rightarrow -(id+id)$$

OR

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(E+id) \Rightarrow -(id+id)$$

At each derivation step, we can choose any of the non-terminal in the sentential form of G for the replacement.

If we always choose the left-most non-terminal in each derivation step, this derivation is called as **left-most derivation**.

If we always choose the right-most non-terminal in each derivation step, this derivation is called as **right-most derivation**.

Left-Most and Right-Most Derivations

Left-Most Derivation

$$\begin{array}{ccccccc} E & \Rightarrow & -E & \Rightarrow & -(E) & \Rightarrow & -(E+E) \Rightarrow -(id+E) \Rightarrow -(id+id) \\ \text{lm} & & \text{lm} & & \text{lm} & & \text{lm} \end{array}$$

Right-Most Derivation

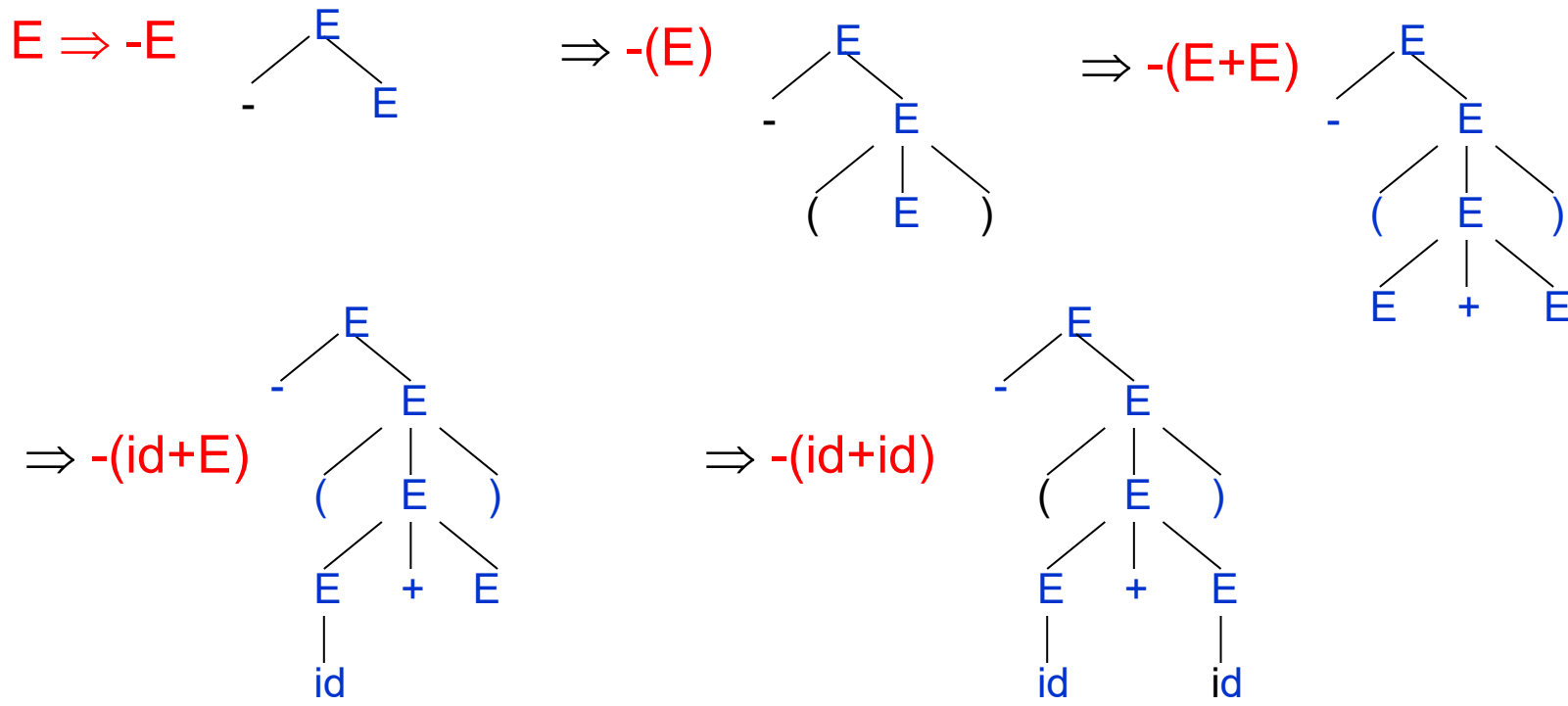
$$\begin{array}{ccccccc} E & \Rightarrow & -E & \Rightarrow & -(E) & \Rightarrow & -(E+E) \Rightarrow -(E+id) \Rightarrow -(id+id) \\ \text{rm} & & \text{rm} & & \text{rm} & & \text{rm} \end{array}$$

We will see that the **top-down parsers** try to find the **left-most derivation** of the given source program.

We will see that the **bottom-up parsers** try to find the **right-most derivation** of the given source program in the **reverse order**.

Parse Tree

- Inner nodes of a parse tree are non-terminal symbols.
 - The leaves of a parse tree are terminal symbols.
-
- A parse tree can be seen as a graphical representation of a derivation.



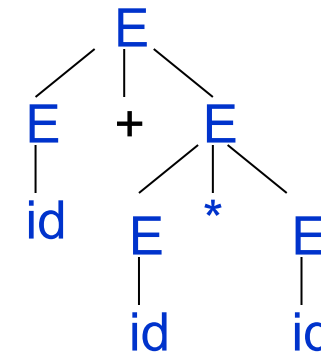
Ambiguity

- A grammar produces more than one parse tree for a sentence is called as an **ambiguous** grammar.

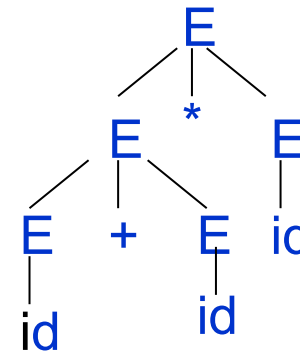
Ambiguous:

Two Different Parse Tree:

$E \Rightarrow E + E$
 $\Rightarrow id + E$
 $\Rightarrow id + E * E$
 $\Rightarrow id + id * E$
 $\Rightarrow id + id * id$



$E \Rightarrow E * E$
 $\Rightarrow E + E * E$
 $\Rightarrow id + E * E$
 $\Rightarrow id + id * E$
 $\Rightarrow id + id * id$

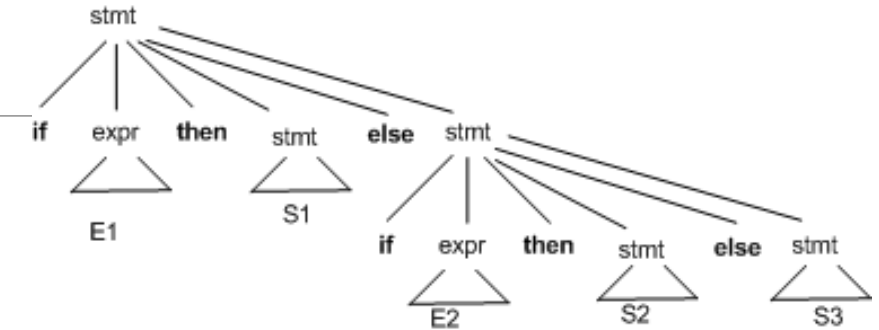


Ambiguity

stmt \longrightarrow If expr then stmt
| If expr then stmt else stmt
| other

One Parse Tree:

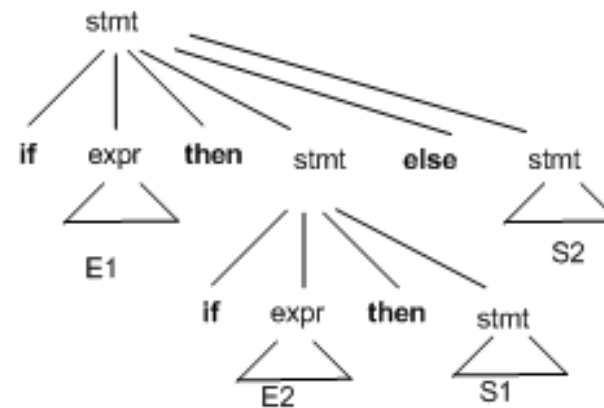
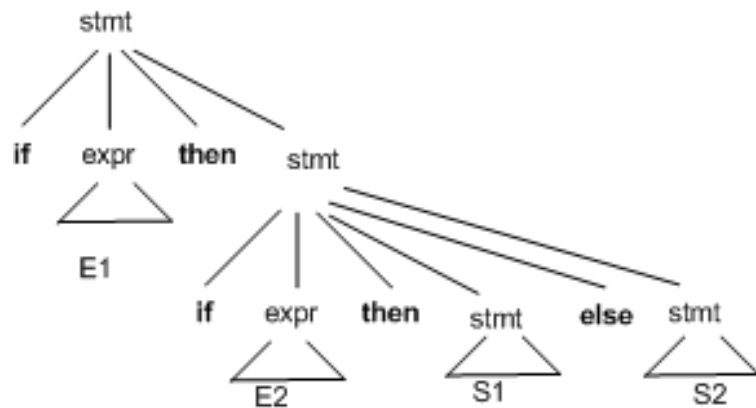
if E_1 then $S1$ else if E_2 then $S2$ else $S3$



Ambiguous:

Two Different Parse Tree:

if E_1 then if E_2 then S_1 else S_2



Ambiguity

For the most parsers, the grammar must be unambiguous.

unambiguous grammar

→ unique selection of the parse tree for a sentence

We should eliminate the ambiguity in the grammar during the design phase of the compiler.

An unambiguous grammar should be written to eliminate the ambiguity.

We have to prefer one of the parse trees of a sentence (generated by an ambiguous grammar) to disambiguate that grammar to restrict to this choice.

Verifying the language generated by a Grammar

Expression $S \rightarrow (S) S \mid \varepsilon$

$$S \underset{\text{lm}}{\Rightarrow} (S) \underset{\text{lm}}{S \Rightarrow} (x)S \underset{\text{lm}}{\Rightarrow} (x)y$$

Checking the string $(x)y$ is balanced with equal number of right and left parentheses and every prefix has at least as many left parentheses as right

$w \Rightarrow (x)y$ is also derivable from S .

Context-free grammar versus regular expression

R.E = $(a|b)^*abb$

CFG:

$S \rightarrow aS \mid bS \mid aA$

$A \rightarrow bB$

$B \rightarrow bC$

$C \rightarrow \epsilon$

Writing a Grammar

- ❑ Lexical versus Syntactic Analysis
- ❑ Eliminating Ambiguity
- ❑ Elimination of Left Recursion
- ❑ Elimination of left Factoring

Elimination of ambiguity

Idea:

- A statement appearing between a **then** and an **else** must be matched (Match each **else** with the closest previous unmatched **then**)

stmt \longrightarrow matched_stmt
|
open_stmt

matched_stmt \longrightarrow If expr **then** matched_stmt **else** matched_stmt
|
other

open_stmt \longrightarrow If expr **then** stmt
|
If expr **then** matched_stmt **else** open_stmt

Ambiguity

- We prefer the first parse tree (else matches with closest if).
 - So, we have to disambiguate our grammar to reflect this choice.
 - The unambiguous grammar will be:
-

$\text{stmt} \rightarrow \text{matched_stmt} \mid \text{unmatched_stmt}$

$\text{matchedstmt} \rightarrow \text{if expr then matched_stmt else matchedstmt} \mid$
 otherstmts

$\text{unmatchedstmt} \rightarrow \text{if expr then stmt} \mid$
 $\text{if expr then matched_stmt else unmatched_stmt}$

Expression Grammar Example:

$E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid - E$

$E \rightarrow (E)$

$E \rightarrow \text{id}$

Simple Arithmetic Expression (Unambiguous Expression) Grammar

$\text{expr} \rightarrow \text{expr} + \text{term}$

$\text{expr} \rightarrow \text{expr} - \text{term}$

$\text{expr} \rightarrow \text{term}$

$\text{term} \rightarrow \text{term} * \text{factor}$

$\text{term} \rightarrow \text{term} / \text{factor}$

$\text{term} \rightarrow \text{factor}$

$\text{factor} \rightarrow (\text{expr})$

$\text{factor} \rightarrow \text{id}$

Unambiguous Expression Grammar

$E \rightarrow E + T \mid E - T \mid T$

$T \rightarrow T * F \mid T / F \mid F$

$F \rightarrow (E)$

$F \rightarrow \text{id}$

Left Recursion

A grammar is **left recursive** if it has a non-terminal A such that there is a derivation.

$$A \xRightarrow{+} A\alpha \quad \text{for some string } \alpha$$

Top-down parsing techniques **cannot** handle left-recursive grammars.

So, we have to convert our left-recursive grammar into an equivalent grammar which is not left-recursive.

The left-recursion may appear in a single step of the derivation (*immediate left-recursion*), or may appear in more than one step of the derivation.

Immediate Left-Recursion

$A \rightarrow A \alpha \mid \beta$ where β does not start with A



eliminate immediate left recursion

$A \rightarrow \beta A'$

$A' \rightarrow \alpha A' \mid \varepsilon$ an equivalent grammar

In general,

$A \rightarrow A \alpha_1 \mid \dots \mid A \alpha_m \mid \beta_1 \mid \dots \mid \beta_n$ where $\beta_1 \dots \beta_n$ do not start with A



eliminate immediate left recursion

$A \rightarrow \beta_1 A' \mid \dots \mid \beta_n A'$

$A' \rightarrow \alpha_1 A' \mid \dots \mid \alpha_m A' \mid \varepsilon$ an equivalent grammar

Immediate Left-Recursion -- Example

$$E \rightarrow E+T \mid T$$

$$T \rightarrow T*F \mid F$$

$$F \rightarrow (E)$$

$$F \rightarrow \text{id}$$

eliminate immediate left recursion \Downarrow

$$E \rightarrow T E'$$

$$E' \rightarrow +T E' \mid \varepsilon$$

$$T \rightarrow F T'$$

$$T' \rightarrow *F T' \mid \varepsilon$$

$$F \rightarrow (E)$$

$$F \rightarrow \text{id}$$

Algorithm 4.19: Eliminating left recursion.

INPUT: Grammar G with no cycles or ϵ -productions.

OUTPUT: An equivalent grammar with no left recursion.

METHOD: Apply the algorithm in Fig. 4.11 to G . Note that the resulting non-left-recursive grammar may have ϵ -productions. \square

- 1) arrange the nonterminals in some order A_1, A_2, \dots, A_n .
- 2) **for** (each i from 1 to n) {
- 3) **for** (each j from 1 to $i - 1$) {
- 4) replace each production of the form $A_i \rightarrow A_j \gamma$ by the
 productions $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$, where
 $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ are all current A_j -productions
- 5) }
- 6) eliminate the immediate left recursion among the A_i -productions
- 7) }

Apply Left recursion for this example:

$E \rightarrow E + T \mid E - T \mid T$

$T \rightarrow T * F \mid T / F \mid F$

$F \rightarrow (E)$

$F \rightarrow \text{id}$

Elimination of left recursion

LHS \rightarrow RHS

$S \rightarrow P1 \mid P2 \mid P3 \mid \dots \mid Pn$

Rule:

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

\Downarrow

$A \rightarrow \beta A'$

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

Example 2

$S \rightarrow Aa \mid b$

$A \rightarrow Ac \mid Sd \mid \varepsilon$

Example 2

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Sd \mid \varepsilon$$

$$A \rightarrow Ac \mid \underline{Aad} \mid \underline{bd} \mid \varepsilon$$

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \beta_1 \mid \beta_2$$

$$S \rightarrow Aa \mid b$$

$$A \rightarrow bdA' \mid A'$$

$$A' \rightarrow cA' \mid adA' \mid \varepsilon$$

Elimination of Left Factoring

A **predictive parser** (a top-down parser without backtracking) insists that the grammar must be *left-factored*.

grammar \rightarrow a new equivalent grammar suitable for predictive parsing

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

Elimination of Left-Factoring

In general,

$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$ where α is non-empty and the first symbols of β_1 and β_2 (if they have one) are different.

Elimination of Left-Factoring

But, if we re-write the grammar as follows

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2 \quad \text{so, we can immediately expand } A \text{ to } \alpha A'$$

Algorithm 4.21: Left factoring a grammar.

INPUT: Grammar G .

OUTPUT: An equivalent left-factored grammar.

METHOD: For each nonterminal A , find the longest prefix α common to two or more of its alternatives. If $\alpha \neq \epsilon$ — i.e., there is a nontrivial common prefix — replace all of the A -productions $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \cdots \mid \alpha\beta_n \mid \gamma$, where γ represents all alternatives that do not begin with α , by

$$\begin{array}{l} A \rightarrow \alpha A' \mid \gamma \\ A' \rightarrow \beta_1 \mid \beta_2 \mid \cdots \mid \beta_n \end{array}$$

Here A' is a new nonterminal. Repeatedly apply this transformation until no two alternatives for a nonterminal have a common prefix. \square

Left-Factoring

Algorithm

For each non-terminal A with two or more alternatives (production rules) with a common non-empty prefix, let say

$$A \rightarrow \alpha\beta_1 \mid \dots \mid \alpha\beta_n \mid \gamma_1 \mid \dots \mid \gamma_m$$

convert it into

$$A \rightarrow \alpha A' \mid \gamma_1 \mid \dots \mid \gamma_m$$

$$A' \rightarrow \beta_1 \mid \dots \mid \beta_n$$

Example 1

$$S \rightarrow iEtS \mid iEtSeS \mid a$$
$$E \rightarrow b$$

Removal of Left Factoring

$$S \rightarrow iEtSS' \mid a$$
$$S' \rightarrow eS \mid \varepsilon$$
$$E \rightarrow b$$

Example 2 – Removal of Left Factoring

$A \rightarrow ad \mid a \mid ab \mid abc \mid b$



$A \rightarrow aA' \mid b$

$A' \rightarrow d \mid \varepsilon \mid b \mid bc$



$A \rightarrow aA' \mid b$

$A' \rightarrow d \mid \varepsilon \mid bA''$

$A'' \rightarrow \varepsilon \mid c$

Example 3 - Removal of Left Factoring

$A \rightarrow \underline{a}bB \mid \underline{a}B \mid cdg \mid cdeB \mid cdfB$



$A \rightarrow aA' \mid \underline{cd}g \mid \underline{cd}eB \mid \underline{cd}fB$

$A' \rightarrow bB \mid B$



$A \rightarrow aA' \mid cdA''$

$A' \rightarrow bB \mid B$

$A'' \rightarrow g \mid eB \mid fB$

Parsing Technique

We categorize the parsers into two groups:

1. **Top-Down Parsing**

- the parse tree is created top to bottom, starting from the root.

2. **Bottom-Up Parsing**

- the parse is created bottom to top; starting from the leaves

Both top-down and bottom-up parsers scan the input from left to right (one symbol at a time).

Efficient top-down and bottom-up parsers can be implemented only for subclasses of context-free grammars.

- **LL** for top-down parsing
- **LR** for bottom-up parsing

Top Down Parsing

- ☐ Recursive-Descent Parsing
- ☐ Non-Recursive Predictive Parsing

Top Down Parser

A **Top-down parser** tries to create a parse tree from the root towards the leafs scanning input from left to right

It can be also viewed as finding a **leftmost derivation** for an input string

Example: **id+id*id**

Top-Down Parsing

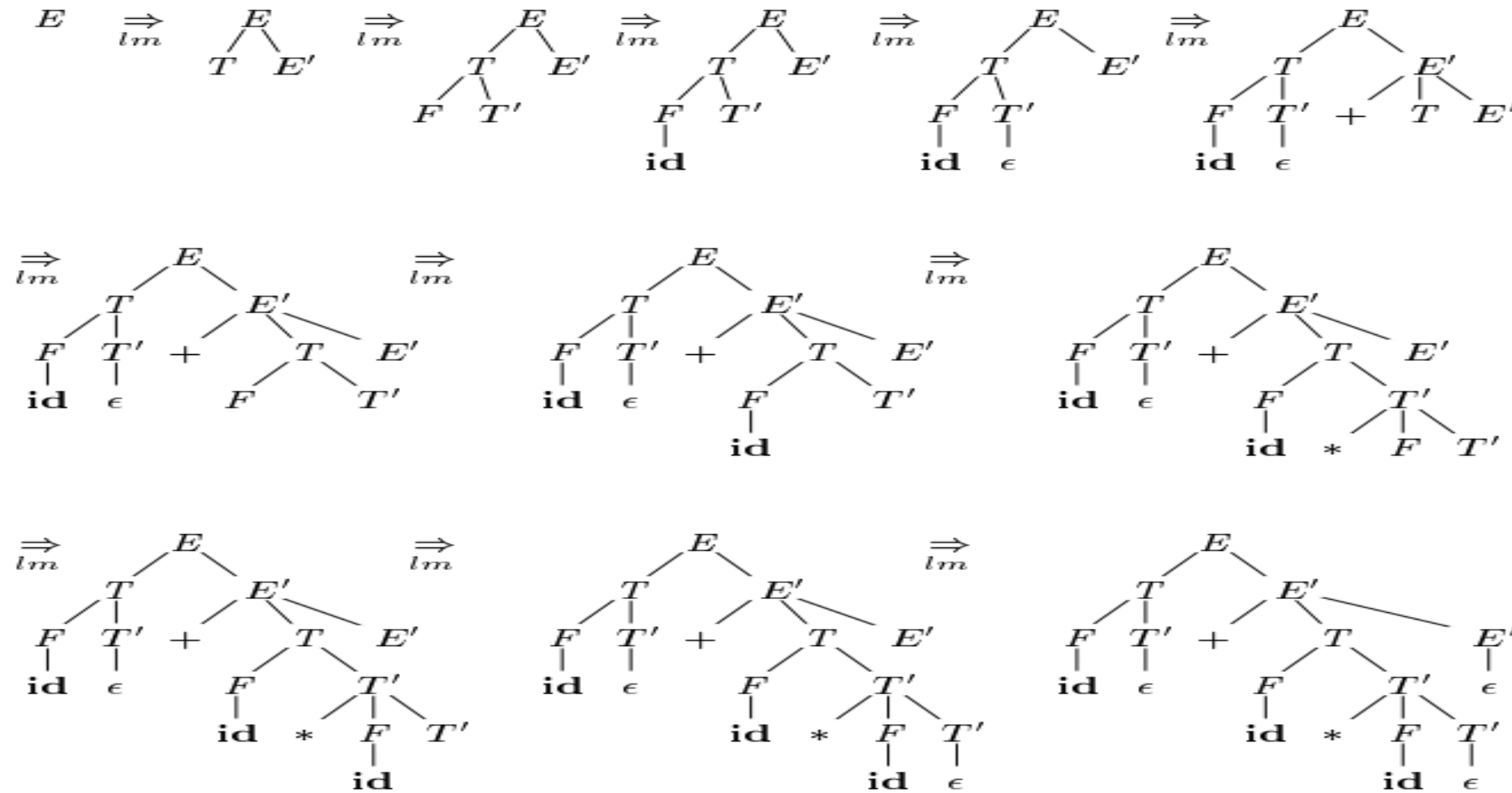


Figure 4.12: Top-down parse for **id + id * id**

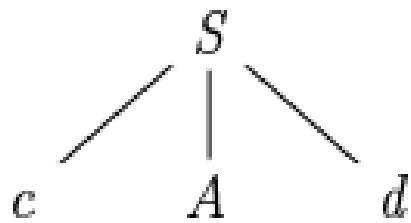
Recursive-Descent Parsing

```
void A() {  
1)    Choose an  $A$ -production,  $A \rightarrow X_1 X_2 \cdots X_k$ ;  
2)    for (  $i = 1$  to  $k$  ) {  
3)        if (  $X_i$  is a nonterminal )  
4)            call procedure  $X_i()$ ;  
5)        else if (  $X_i$  equals the current input symbol  $a$  )  
6)            advance the input to the next symbol;  
7)        else /* an error has occurred */;  
    }  
}
```

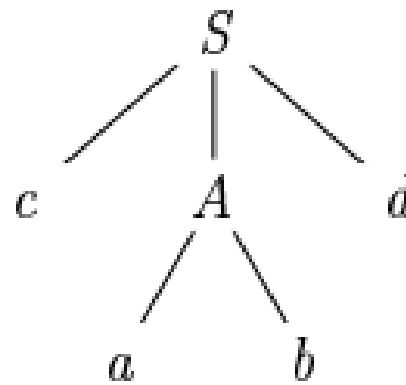
Recursive-Descent Parsing

$$S \rightarrow c A d$$
$$A \rightarrow a b \mid a$$

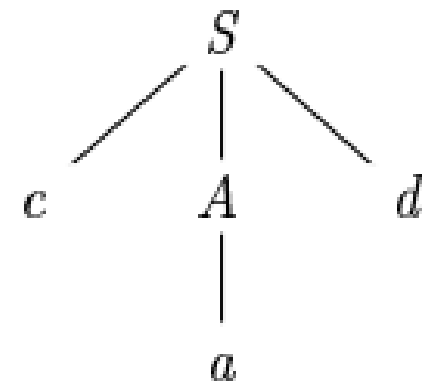
Recursive-Descent Parsing



(a)



(b)



(c)

Figure 4.14: Steps in a top-down parse

Non-Recursive Predictive Parsing

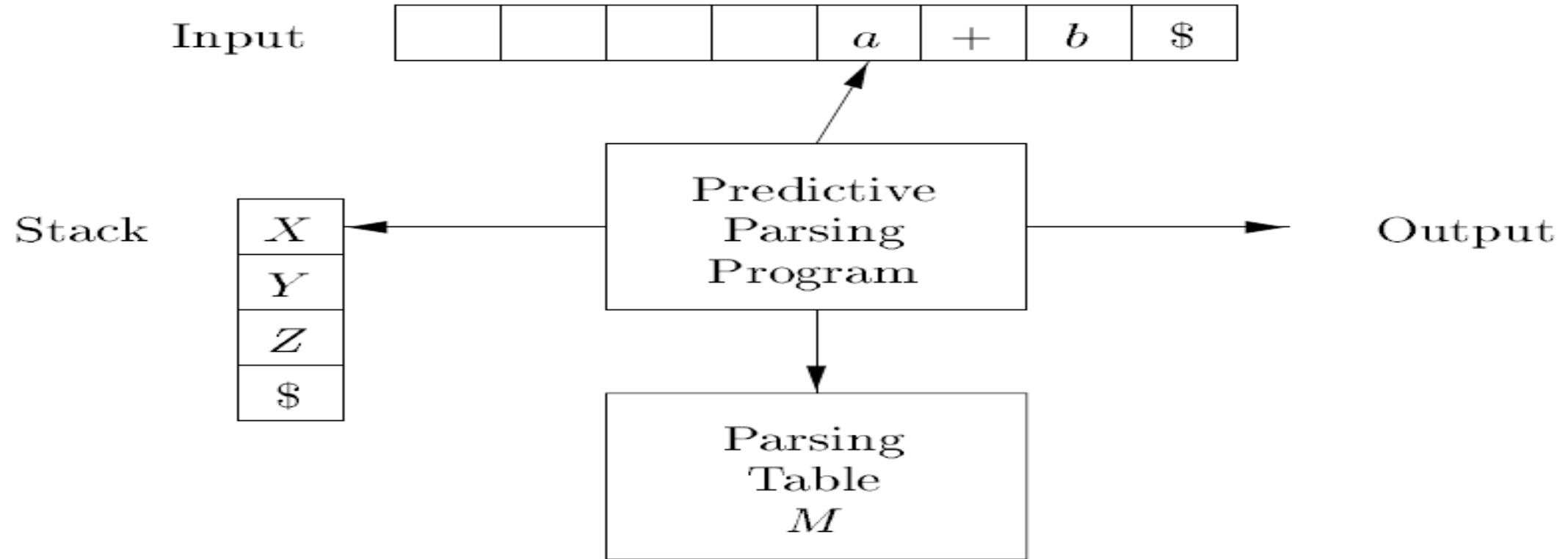


Figure 4.19: Model of a table-driven predictive parser

Non-Recursive Predictive Parsing / Top Down Parser / LL(1)

LL(1) Grammar

- ❑ Elimination of Left Recursion / Left Factoring
- ❑ FIRST and FOLLOW
- ❑ Predictive parsing table
- ❑ Stack Implementation

LL(1) Grammar

Predictive Parser, i.e recursive-descent parsers without backtracking is constructed for a class of grammar called LL(1)

- L – Leftmost derivation
- L – scans the input from left to right
- 1 – one input symbol of lookahead at each step to make parsing action

Elimination of Left Recursion

$E \rightarrow E+T \mid T$

$T \rightarrow T*F \mid F$

$F \rightarrow (E)$

$F \rightarrow \text{id}$



$E \rightarrow T E'$

$E' \rightarrow +T E' \mid \varepsilon$

$T \rightarrow F T'$

$T' \rightarrow *F T' \mid \varepsilon$

$F \rightarrow (E)$

$F \rightarrow \text{id}$

FIRST and FOLLOW

FIRST(X):

Rules:

1. X is a Terminal, $\text{FIRST}(X) = \{ X \}$
2. X is a Non-terminal, $X \rightarrow Y_1, Y_2, Y_3 \dots Y_k, K \geq 1$,
 $\text{FIRST}(X) = \text{FIRST}(Y_1)$
3. $X \rightarrow \epsilon$ is a production, $\text{FIRST}(X) = \{ \epsilon \}$

FIRST(x):

$$1) E \rightarrow T E'$$

$$\text{FIRST}(E) = \text{FIRST}(T)$$

$$2) E' \rightarrow +T E' \mid \varepsilon$$

$$E' \rightarrow +T E'$$

$$E' \rightarrow \varepsilon$$

$$\text{FIRST}(E') = +$$

$$\text{FIRST}(E') = \varepsilon$$

FIRST(x):

3) $T \rightarrow F T'$

$$\text{FIRST}(T) = \text{FIRST}(F)$$

4) $T' \rightarrow *F T' \mid \varepsilon$

$T' \rightarrow *F T'$

$$\text{FIRST}(T') = *$$

$T' \rightarrow \varepsilon$

$$\text{FIRST}(T') = \varepsilon$$

FIRST(x):

5) $F \rightarrow (E)$

$\text{FIRST}(E) = ($

6) $F \rightarrow \text{id}$

$\text{FIRST}(E) = \text{id}$

FIRST (X)

Non-terminal	FIRST
E	(, id
E'	+ , ϵ
T	(, id
T'	* , ϵ
F	(, id

FIRST and FOLLOW

FOLLOW(X):

Rules:

1. S is a Start Symbol, $\text{FOLLOW}(S) = \$$
2. If a production $A \rightarrow \alpha B \beta$, $\text{FOLLOW}(B) = \text{FIRST}(\beta) - \epsilon$
3. If a production $A \rightarrow \alpha B$ or $A \rightarrow \alpha B \beta$ i.e, $\text{FIRST}(\beta) = \epsilon$
 $\text{FOLLOW}(B) = \text{FOLLOW}(A)$

FOLLOW(X)

1) $E \rightarrow T E'$

$\text{FOLLOW}(T) = \text{FIRST}(E') - \epsilon$

$= +, \epsilon - \epsilon$

$= +$

$\text{FOLLOW}(T) = \text{FOLLOW}(E)$

$\text{FOLLOW}(T) = +, \text{FOLLOW}(E)$

FOLLOW(X)

1) $E \rightarrow T E'$

$\text{FOLLOW}(E') = \text{FOLLOW}(E)$

$\text{FOLLOW}(E') = \text{FOLLOW}(E)$

FOLLOW(X)

2) $E' \rightarrow +T E' \mid \varepsilon$

$$\begin{aligned}\text{FOLLOW}(T) &= \text{FIRST}(E') - \varepsilon \\ &= +, \varepsilon - \varepsilon \\ &= +\end{aligned}$$

$$\text{FOLLOW}(T) = \text{FOLLOW}(E')$$

$$\text{FOLLOW}(T) = +, \text{FOLLOW}(E')$$

FOLLOW(X)

$$2) E' \rightarrow +T E' \mid \varepsilon$$

$$\text{FOLLOW}(E') = \text{FOLLOW}(E')$$

$$\text{FOLLOW}(E') = \text{FOLLOW}(E')$$

FOLLOW(X)

3) $T \rightarrow F T'$

$$\begin{aligned}\text{FOLLOW}(F) &= \text{FIRST}(T') - \varepsilon \\ &= *, \varepsilon - \varepsilon \\ &= *\end{aligned}$$

$$\text{FOLLOW}(F) = \text{FOLLOW}(T)$$

$$\text{FOLLOW}(F) = *, \text{FOLLOW}(T)$$

FOLLOW(X)

3) $T \rightarrow F T'$

$\text{FOLLOW}(T') = \text{FOLLOW}(T)$

$\text{FOLLOW}(T') = \text{FOLLOW}(T)$

FOLLOW(X)

$$4) T' \rightarrow * F T' \mid \varepsilon$$

$$\begin{aligned} \text{FOLLOW}(F) &= \text{FIRST}(T') - \varepsilon \\ &= *, \varepsilon - \varepsilon \\ &= * \end{aligned}$$

$$\text{FOLLOW}(F) = \text{FOLLOW}(T')$$

$$\text{FOLLOW}(F) = *, \text{FOLLOW}(T')$$

FOLLOW(X)

4) $T' \rightarrow *F T' \mid \varepsilon$

$\text{FOLLOW}(T') = \text{FOLLOW}(T')$

$\text{FOLLOW}(T') = \text{FOLLOW}(T')$

FOLLOW(X)

5) $F \rightarrow (E)$

$\text{FOLLOW}(E) = \text{FIRST}(') ') =)$

$\text{FOLLOW}(E) = \$$

$\text{FOLLOW}(E) = \{) , \$ \}$

FOLLOW(X)

6) $F \rightarrow id$

FOLLOW

Non-terminal	FOLLOW
E) , \$
E'	$\text{FOLL}(E') = \text{FOLL}(E)$ $\text{FOLL}(E') = \text{FOLL}(E')$
T	$\text{FOLL}(T) = + , \text{FOLL}(E)$ $\text{FOLL}(T) = + , \text{FOLL}(E')$
T'	$\text{FOLL}(T') = \text{FOLL}(T)$ $\text{FOLL}(T') = \text{FOLL}(T')$
F	$\text{FOLL}(F) = * , \text{FOLL}(T)$ $\text{FOLL}(F) = * , \text{FOLL}(T')$

FOLLOW

Non-terminal	FOLLOW	FOLLOW
E) , \$) , \$
E'	FOLLOW(E) , FOLLOW(E')) , \$
T	+ , FOLLOW(E) , FOLLOW(E')	+,) , \$
T'	FOLLOW(T), FOLLOW(T')	+,) , \$
F	* , FOLLOW(T), FOLLOW(T')	*, +,) , \$

FIRST and FOLLOW

Non-terminal	FIRST	FOLLOW
E	(, id) , \$
E'	+ , ϵ) , \$
T	(, id	+,) , \$
T'	* , ϵ	+,) , \$
F	(, id	* , +,) , \$

Non-Recursive Predictive Parsing

Algorithm 4.31: Construction of a predictive parsing table.

INPUT: Grammar G .

OUTPUT: Parsing table M .

METHOD: For each production $A \rightarrow \alpha$ of the grammar, do the following:

1. For each terminal a in $\text{FIRST}(\alpha)$, add $A \rightarrow \alpha$ to $M[A, a]$.
2. If ϵ is in $\text{FIRST}(\alpha)$, then for each terminal b in $\text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, b]$. If ϵ is in $\text{FIRST}(\alpha)$ and $\$$ is in $\text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, \$]$ as well.

Predicting parsing Table

NON - TERMINAL	INPUT SYMBOL					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

Figure 4.17: Parsing table M for Example 4.32

Stack Implementation

Algorithm 4.34: Table-driven predictive parsing.

INPUT: A string w and a parsing table M for grammar G .

OUTPUT: If w is in $L(G)$, a leftmost derivation of w ; otherwise, an error indication.

Stack Implementation

METHOD: Initially, the parser is in a configuration with $w\$$ in the input buffer and the start symbol S of G on top of the stack, above $\$$. The program in Fig. 4.20 uses the predictive parsing table M to produce a predictive parse for the input. \square

```
let  $a$  be the first symbol of  $w$ ;  
let  $X$  be the top stack symbol;  
while (  $X \neq \$$  ) { /* stack is not empty */  
    if (  $X = a$  ) pop the stack and let  $a$  be the next symbol of  $w$ ;  
    else if (  $X$  is a terminal ) error();  
    else if (  $M[X, a]$  is an error entry ) error();  
    else if (  $M[X, a] = X \rightarrow Y_1 Y_2 \cdots Y_k$  ) {  
        output the production  $X \rightarrow Y_1 Y_2 \cdots Y_k$ ;  
        pop the stack;  
        push  $Y_k, Y_{k-1}, \dots, Y_1$  onto the stack, with  $Y_1$  on top;  
    }  
    let  $X$  be the top stack symbol;  
}
```


Table-driven Predictive Parsing

Stack	Input	Action
\$E	id + id * id \$	push E \rightarrow TE'
\$E'T	id + id * id \$	push T \rightarrow FT'
\$E'T'F	id + id * id \$	push F \rightarrow id
\$E'T'id	id + id * id \$	pop (match id)
\$E'T'	+ id * id \$	push T' \rightarrow ϵ
\$E'	+ id * id \$	push E' \rightarrow +TE'
\$E'T+	+ id * id \$	pop (match +)
\$E'T	id * id \$	push T \rightarrow FT'

Table-driven Predictive Parsing

Stack	Input	Action
\$E'T'F	id * id \$	push F \rightarrow id
\$E'T'id	id * id \$	pop (match id)
\$E'T'	* id\$	push T' \rightarrow *FT'
\$E'T'F*	*id\$	pop (match *)
\$E'T'F	id\$	push F \rightarrow id
\$E'T'id	id\$	pop (match id)
\$E'T'	\$	push T' \rightarrow ϵ
\$E'	\$	push E' \rightarrow ϵ
\$	\$	Accept

Table-driven Predictive Parsing

MATCHED	STACK	INPUT	ACTION
	$E\$$	id + id * id \$	
	$TE' \$$	id + id * id \$	output $E \rightarrow TE'$
	$FT'E' \$$	id + id * id \$	output $T \rightarrow FT'$
	id $T'E' \$$	id + id * id \$	output $F \rightarrow \mathbf{id}$
id	$T'E' \$$	+ id * id \$	match id
id	$E' \$$	+ id * id \$	output $T' \rightarrow \epsilon$
id	+ $TE' \$$	+ id * id \$	output $E' \rightarrow + TE'$
id +	$TE' \$$	id * id \$	match +
id +	$FT'E' \$$	id * id \$	output $T \rightarrow FT'$
id +	id $T'E' \$$	id * id \$	output $F \rightarrow \mathbf{id}$
id + id	$T'E' \$$	* id \$	match id
id + id	* $FT'E' \$$	* id \$	output $T' \rightarrow * FT'$
id + id *	$FT'E' \$$	id \$	match *
id + id *	id $T'E' \$$	id \$	output $F \rightarrow \mathbf{id}$
id + id * id	$T'E' \$$	\$	match id
id + id * id	$E' \$$	\$	output $T' \rightarrow \epsilon$
id + id * id	\$	\$	output $E' \rightarrow \epsilon$

Error Recovery in Predictive Parsing

Error is detected in predictive parsing, $M[A,a]$ is error

- Panic mode
- Phrase-level recovery

Panic mode

- Skipping symbols on the input until a token in a selected set of synchronizing tokens appears.
- FOLLOW(A)
 - The synchronizing set for nonterminal A
 - Eg: if semicolons terminate statements as in C
- FIRST(A)
 - Add to the synchronizing set for nonterminal A
 - NT produces the empty string
 - Terminal cannot be matched, pop the terminal

Predictive Parsing Table

NON - TERMINAL	INPUT SYMBOL					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$	synch	synch
E'		$E \rightarrow +TE'$			$E \rightarrow \epsilon$	$E \rightarrow \epsilon$
T	$T \rightarrow FT'$	synch		$T \rightarrow FT'$	synch	synch
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow \text{id}$	synch	synch	$F \rightarrow (E)$	synch	synch

Figure 4.22: Synchronizing tokens added to the parsing table of Fig. 4.17

Parsing and Error Recovery

STACK	INPUT	REMARK
$E \$$) $\text{id} * + \text{id} \$$	error, skip)
$E \$$	$\text{id} * + \text{id} \$$	id is in $\text{FIRST}(E)$
$TE' \$$	$\text{id} * + \text{id} \$$	
$FT'E' \$$	$\text{id} * + \text{id} \$$	
$\text{id } T'E' \$$	$\text{id} * + \text{id} \$$	
$T'E' \$$	$* + \text{id} \$$	
$* FT'E' \$$	$* + \text{id} \$$	
$FT'E' \$$	$+ \text{id} \$$	error, $M[F, +] = \text{synch}$
$T'E' \$$	$+ \text{id} \$$	F has been popped
$E' \$$	$+ \text{id} \$$	
$+ TE' \$$	$+ \text{id} \$$	
$TE' \$$	$\text{id} \$$	
$FT'E' \$$	$\text{id} \$$	
$\text{id } T'E' \$$	$\text{id} \$$	
$T'E' \$$	$\$$	
$E' \$$	$\$$	
$\$$	$\$$	

Phrase-level recovery

- It is implemented by filling in the blank entries in the predictive parsing table with pointers to error routines.
- These routines may change, insert, or delete symbols on the input and issue appropriate error messages.

Thank You