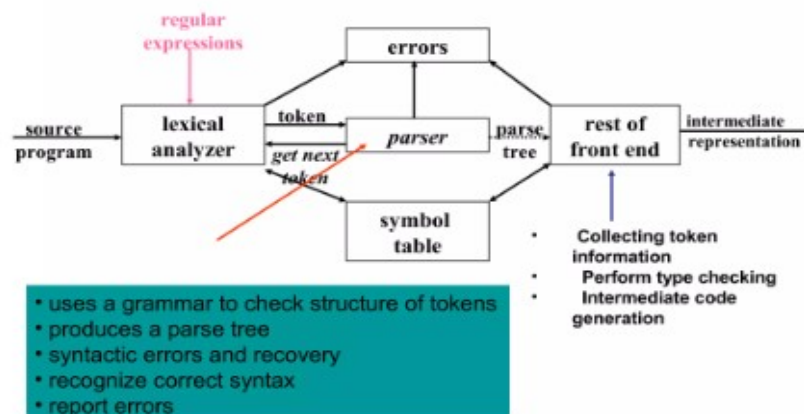


UNIT-2

- **Syntax Analysis:**
- **Introduction:** The Role of the parser – Representative Grammars – Syntax Error Handling – Error Recovery Strategies.
- **Context Free Grammars:** The formal definition of a CFG – Notational Conventions – Derivations – Parse trees and derivations – Ambiguity
- **Writing Grammar:** Lexical Versus Syntax Analysis – Eliminating Ambiguity – Elimination of Left Recursion – Left Factoring
- **Top Down Parsing:** Recursive Descent Parsing-FIRST and FOLLOW - LL(1) Grammars – Non recursive Predictive Parsing- Error Recovery in Predictive Parsing.

The Role of the parser



The Role of the parser

- The **parser** obtains a string of tokens **from lexical analyzer**.
- It then verifies whether the input string can be generated from the grammar of the source language.
- It has to report any syntactical errors and recover from commonly occurring error to continue the processing the remainder of the program
- The **parser constructs a parse tree** and passes it to the rest of the compiler for further processing.
- It should also **report syntactical errors** in an easily understandable way to the user.

Parser:

A parser is a program that generates a parse tree for the given string, if the string is generated from the underlying grammar.

```
id = id + id * id ;
```



Parser

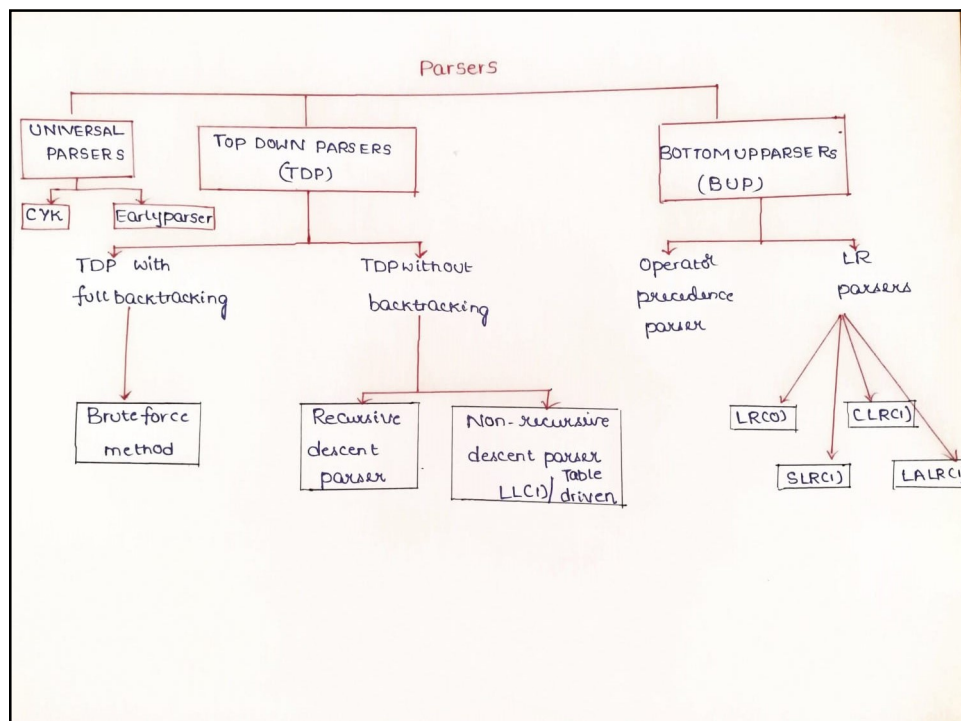


```
S → id = E ;
E → E + T | T
T → T * F | F
F → id
```



Categorization of parsers

- There are 03 general types of parsers for grammars: universal, top-down, and bottom-up.
 - Universal parsing methods such as the Cocke-Younger-Kasami algorithm and Earley's algorithm can parse any grammar. These general methods are, however, too inefficient to use in production compilers.
- The methods commonly used in compilers is either top-down or bottom-up.
- Top-down methods build parse trees from the top (root) to the bottom (leaves), while bottom-up methods start from the leaves and work their way up to the root.



Representative Grammars

$E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid \text{id}$

Useful for Bottom up parsing. Belongs to LR grammar Class.

$E \rightarrow T E'$
 $E' \rightarrow + T E' \mid \epsilon$
 $T \rightarrow F T'$
 $T' \rightarrow * F T' \mid \epsilon$
 $F \rightarrow (E) \mid \text{id}$

After removing left recursion used for top down parsing

$E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$

useful for illustrating handling ambiguities during parsing

Syntax Error Handling

Common programming errors can occur at many different levels

Types of Errors

- 1) Lexical
- 2) Syntactic
- 3) Semantic
- 4) Logical

- **Lexical errors** include misspellings of identifiers, keywords, or operators.
 - e.g., the use of an identifier *ellipseSize* instead of *ellipseSize*
 - missing quotes around text intended as a string
- **Syntactic errors** include misplaced semicolons or extra or missing braces; that is, “{” or “}”
 - Example: In C or Java, the appearance of a *case* statement without an enclosing *switch* is a syntactic error

Syntax Error Handling

Types of Errors

1) Lexical

2) Syntactic

3) Semantic

4) Logical

- **Semantic errors** include type mismatches between operators and operands,
 - e.g., the return of a value in a Java method with result type void.
- **Logical errors** occur when executed code does not produce the expected result.
 - incorrect reasoning on the part of the programmer
 - The use in a C program of the assignment operator = instead of the comparison operator ==

Syntax Error Handling

- The error handler in a parser has goals that are simple to state but challenging to realize:
 - Report the presence of errors clearly and accurately.
 - Recover from each error quickly enough to detect subsequent errors.
 - Add minimal overhead to the processing of correct programs.

Error-Recovery Strategies

1. Panic-Mode Recovery
2. Phrase-Level Recovery
3. Error Productions
4. Global Correction

Panic-Mode Recovery

- Once an error is found, the parser intends to find designated set of synchronizing tokens by discarding input symbols one at a time.
- Synchronizing tokens are *delimiters, semicolon or }* whose role in source program is clear.
- When parser finds an error in the statement, it ignores the rest of the statement by not processing the input.
- Advantage:
 - Simplicity
 - Never get into infinite loop
- Disadvantage:
 - Additional errors cannot be checked as some of the input symbols will be skipped.

Phrase-Level Recovery

- When a parser finds an error, it tries to take corrective measures so that the rest of inputs of statement allow the parser to parse ahead. The corrections may be
 - Replacing a prefix by some string.
 - Replacing comma by semicolon.
 - Deleting extraneous semicolon.
 - Inserting missing semicolon.

we must be careful to choose replacements that do not lead to infinite loops, as would be the case, for example, if we always inserted something on the input ahead of the current input symbol.

- Advantage:

- It can correct any input string.

- Disadvantage:

- It is difficult to cope up with actual error if it has occurred before the point of detection.

Error Productions

- The use of the error production method can be incorporated if the user is aware of common mistakes that are encountered in grammar in conjunction with errors that produce erroneous constructs.
 - Example: write $5x$ instead of $5*x$
- Advantage:
 - If this is used then, during parsing appropriate error messages can be generated and parsing can be continued.
- Disadvantage:
 - The disadvantage is that it's difficult to maintain.

Global Correction

- The parser considers the program in hand as a whole and tries to figure out what the program is intended to do and tries to find out a closest match for it, which is error-free.
- When an erroneous input statement X is fed, it creates a parse tree for some closest error-free statement Y.
- Advantage:
 - This may allow the parser to make minimal changes in the source code.
- Disadvantage:
 - Due to the complexity (time and space) of this strategy, it has not been implemented in practice yet.

Lexical analysis Vs syntax analysis

1. Separating the syntactic structure of a language into lexical and non-lexical parts provides a convenient way of modularizing the front end of a compiler into two manageable-sized components.
2. The lexical rules of a language are frequently quite simple, and to describe them we do not need a notation as powerful as grammars.
3. Regular expressions generally provide a more concise and easier-to-understand notation for tokens than grammars.
4. More efficient lexical analyzers can be constructed automatically from regular expressions than from arbitrary grammars.

Lexical Vs Syntactic analysis

Lexical analysis	Syntax analysis
It is responsible for converting a sequence of characters into a pattern of tokens.	Process of analyzing a string of symbols either in natural language or computer languages that satisfies the rules of a formal grammar.
Reads the program one character at a time, the output is meaningful lexemes.	Tokens are taken as input and a parse tree is generated as output.
It is the first phase of the compilation process.	It is the second phase of the compilation process.

Lexical analysis	Syntax analysis
It can also be referred to as lexing and tokenization.	It can also be referred to as syntactic analysis and parsing.
A lexical analyser is a pattern matcher.	A syntax analysis involves forming a tree to identify deformities in the syntax of the program.
Less complex approaches are often used for lexical analysis.	Syntax analysis requires a much more complex approach.

Context Free Grammars

- The formal definition of a CFG
- Notational Conventions
- Derivations
- Parse trees and derivations
- Ambiguity

The formal definition of a CFG

- Context-free grammars are named as such because **any of the production rules in the grammar can be applied regardless of context**—it does not depend on any other symbols that may or may not be around a given symbol that is having a rule applied to it.
- A context free grammar G is defined by four tuple format as

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

where,

G – Grammar

V – Set of variables

T – Set of terminals

P – Set of productions

S – Start symbol

A Context Free Grammar

- ▶ A context-free grammar has four components:
 - ▶ A set of terminal symbols, sometimes referred to as "tokens."
 - ▶ A set of nonterminals, sometimes called "syntactic variables."
 - ▶ A set of productions, where each production consists of a nonterminal, called the head or left side of the production, an arrow, and a sequence of terminals and/or nonterminals, called the body or right side of the production
 - ▶ A designation of one of the nonterminals as the start symbol.

Example 1:

<i>expression</i>	→	<i>expression</i> + <i>term</i>
<i>expression</i>	→	<i>expression</i> - <i>term</i>
<i>expression</i>	→	<i>term</i>
<i>term</i>	→	<i>term</i> * <i>factor</i>
<i>term</i>	→	<i>term</i> / <i>factor</i>
<i>term</i>	→	<i>factor</i>
<i>factor</i>	→	(<i>expression</i>)
<i>factor</i>	→	id

the terminal symbols are **id + - * / ()**

The nonterminal symbols are *expression*, *term* and *factor*,

expression is the start symbol

Example 2:

stmt → **if** (*expr*) *stmt* **else** *stmt*

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, ..., 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.

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

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

$$\mathbf{V} = \{E, T, F\}$$

$$\mathbf{T} = \{+, -, *, /, (,), id\}$$

$$\mathbf{S} = \{E\}$$

$$\mathbf{P} :$$

$$E \rightarrow E + T$$

$$E \rightarrow E - T$$

$$E \rightarrow T$$

$$T \rightarrow T * F$$

$$T \rightarrow T / F$$

$$T \rightarrow F$$

$$F \rightarrow (E)$$

$$F \rightarrow id$$

Derivations

- ▶ $E \Rightarrow E+E$: $E+E$ derives from E
 - ▶ $E \Rightarrow E+E \Rightarrow id+E \Rightarrow id+id$
 - ▶ A sequence of replacements of non-terminal symbols is called a **derivation** of $id+id$ from E .
 - ▶ $\alpha A \beta \Rightarrow \alpha \gamma \beta$ if there is a production rule $A \rightarrow \gamma$ in our grammar and α and β are arbitrary strings of terminal and non-terminal symbols
 - ▶ $\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$ (α_n derives from α_1 or α_1 derives α_n)
- | | |
|-------------------|---------------------------------|
| \Rightarrow | : derives in one step |
| $\xRightarrow{*}$ | : derives in zero or more steps |
| $\xRightarrow{+}$ | : derives in one or more steps |

Act

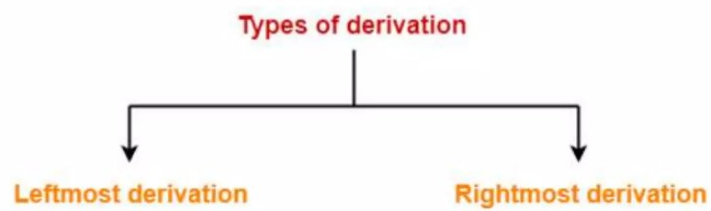
CFG - Terminology

- ▶ $L(G)$ 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

Parse Tree-

- The process of deriving a string is called as **derivation**.
- The geometrical representation of a derivation is called as a **parse tree** or **derivation tree**.



LEFT MOST DERIVATION

At each and every step the leftmost non-terminal is expanded by substituting its corresponding production to derive a string.

$$E \rightarrow E + E \mid E * E \mid \text{id}$$

$$w = \text{id} + \text{id} * \text{id}$$

$$E \xrightarrow{\text{lm}} E + E$$

$$E \xrightarrow{\text{lm}} \text{id} + E \quad [E \rightarrow \text{id}]$$

$$E \xrightarrow{\text{lm}} \text{id} + E * E \quad [E \rightarrow E * E]$$

$$E \xrightarrow{\text{lm}} \text{id} + \text{id} * E \quad [E \rightarrow \text{id}]$$

$$E \xrightarrow{\text{lm}} \text{id} + \text{id} * \text{id} \quad [E \rightarrow \text{id}]$$

RIGHT MOST DERIVATION

At each and every step the rightmost non-terminal is expanded by substituting its corresponding production to derive a string.

$$E \rightarrow E + E \mid E * E \mid \text{id}$$

$$w = \text{id} + \text{id} * \text{id}$$

$$E \xrightarrow{\text{rm}} E + E$$

$$E \xrightarrow{\text{rm}} E + E * E \quad [E \rightarrow E * E]$$

$$E \xrightarrow{\text{rm}} E + E * \text{id} \quad [E \rightarrow \text{id}]$$

$$E \xrightarrow{\text{rm}} E + \text{id} * \text{id} \quad [E \rightarrow \text{id}]$$

$$E \xrightarrow{\text{rm}} \text{id} + \text{id} * \text{id} \quad [E \rightarrow \text{id}]$$

EXAMPLE

□ Leftmost

$$S \rightarrow SS + \mid SS * \mid a$$

$$w = aa + a*$$

$$S \xrightarrow{\text{lm}} SS*$$

$$S \xrightarrow{\text{lm}} SS + S* \quad [S \rightarrow SS +]$$

$$S \xrightarrow{\text{lm}} aS + S* \quad [S \rightarrow a]$$

$$S \xrightarrow{\text{lm}} aa + S* \quad [S \rightarrow a]$$

$$S \xrightarrow{\text{lm}} aa + a* \quad [S \rightarrow a]$$

□ Rightmost

$$S \rightarrow SS + \mid SS * \mid a$$

$$w = aa + a*$$

$$S \xrightarrow{\text{rm}} SS*$$

$$S \xrightarrow{\text{rm}} Sa* \quad [S \rightarrow a]$$

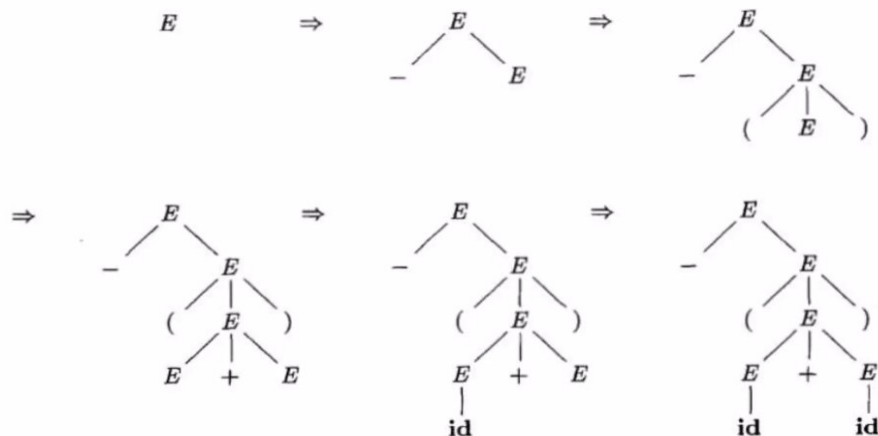
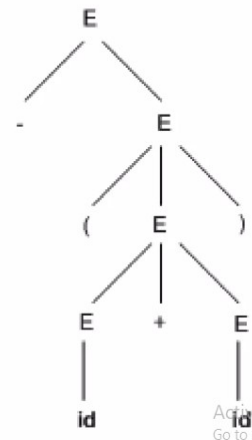
$$S \xrightarrow{\text{rm}} SS + a* \quad [S \rightarrow SS +]$$

$$S \xrightarrow{\text{rm}} Sa + a* \quad [S \rightarrow a]$$

$$S \xrightarrow{\text{rm}} aa + a* \quad [S \rightarrow a]$$

Parse Tree

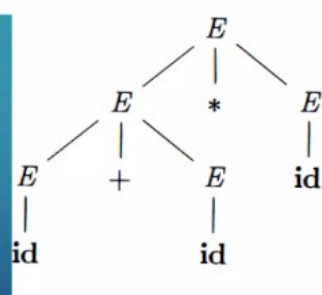
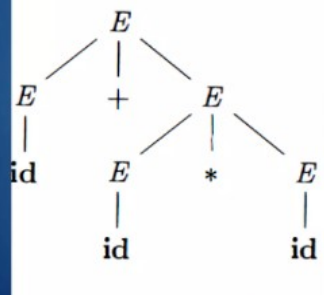
- Parse tree is a hierarchical structure which represents the derivation of the grammar to yield input strings.
- Root node of parse tree has the start symbol of the given grammar from where the derivation proceeds.
- Leaves of parse tree represent terminals.
- Each interior node represents productions of grammar.
- If $G : E \rightarrow E+E \mid E * E \mid (E) \mid - E \mid id$ is the grammar, then Parse tree for the input string $-(id + id)$ is shown.



Ambiguity

- a grammar that produces more than one parse tree for some sentence is said to be ambiguous

$E \Rightarrow E + E$	$E \Rightarrow E * E$
$\Rightarrow \text{id} + E$	$\Rightarrow E + E * E$
$\Rightarrow \text{id} + E * E$	$\Rightarrow \text{id} + E * E$
$\Rightarrow \text{id} + \text{id} * E$	$\Rightarrow \text{id} + \text{id} * E$
$\Rightarrow \text{id} + \text{id} * \text{id}$	$\Rightarrow \text{id} + \text{id} * \text{id}$



Assignment 1

Exercise 4.2.1: Consider the context-free grammar:

$$S \rightarrow S S + \mid S S * \mid a$$

and the string $aa + a*$.

- Give a leftmost derivation for the string.
- Give a rightmost derivation for the string.
- Give a parse tree for the string.
- ! d) Is the grammar ambiguous or unambiguous? Justify your answer.
- ! e) Describe the language generated by this grammar.

- 1 Give a leftmost derivation for the string.
- 2 Give a rightmost derivation for the string.
- 3 Give a parse tree for the string.
- 4 Is the grammar ambiguous or unambiguous?

a $S \rightarrow 0 S 1 \mid 0 1$ with string 000111.

b $S \rightarrow + S S \mid * S S \mid a$ with string $+ * a a a$.

c $S \rightarrow S (S) S \mid \epsilon$ with string $((())())$.

d $S \rightarrow S + S \mid S S \mid (S) \mid S * \mid a$ with string $(a + a) * a$.

e $S \rightarrow (L) \mid a$ and $L \rightarrow L , S \mid S$ with string $((a, a), a, (a))$.

f $S \rightarrow a S b S \mid b S a S \mid \epsilon$ with string $a a b b a b$.

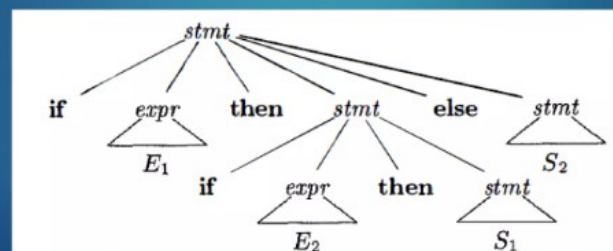
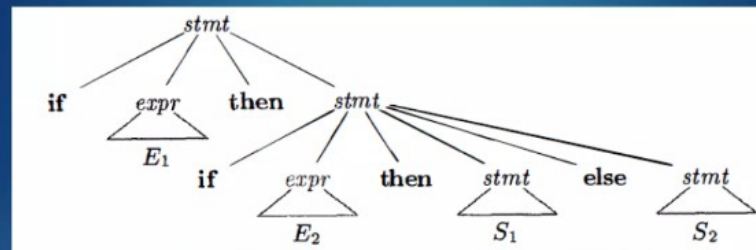
Writing a Grammar

- ▶ Grammars are capable of describing most, of the syntax of programming languages .
- ▶ Grammar should be unambiguous.
- ▶ Left-recursion elimination and left factoring - are useful for rewriting grammars .
- ▶ From the resulting grammar we can create top down parsers without backtracking.
- ▶ Such parsers are called predictive parsers or recursive-descent parser

Eliminating Ambiguity

- ▶ ambiguous grammar can be rewritten to eliminate the ambiguity.
- ▶ $stmt \rightarrow if\ expr\ then\ stmt$
 $\quad | if\ expr\ then\ stmt\ else\ stmt$
 $\quad | other$
- ▶ is ambiguous since the string
 $\quad \triangleright if\ E_1\ then\ if\ E_2\ then\ S_1\ else\ S_2$ has the two parse trees

Two parse trees for an ambiguous sentence



Eliminating Ambiguity

- ▶ The general rule is, "Match each else with the closest unmatched then."

<i>stmt</i>	→	<i>matched_stmt</i>
		<i>open_stmt</i>
<i>matched_stmt</i>	→	if <i>expr</i> then <i>matched_stmt</i> else <i>matched_stmt</i>
		other
<i>open_stmt</i>	→	if <i>expr</i> then <i>stmt</i>
		if <i>expr</i> then <i>matched_stmt</i> else <i>open_stmt</i>

Left Recursion

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

$A \Rightarrow A\alpha_*$ for some string α

- ▶ Top-down parsing techniques cannot handle left-recursive grammars.
- ▶ 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
 \Downarrow eliminate immediate left recursion

$A \rightarrow \beta A'$
 $A' \rightarrow \alpha A' \mid \varepsilon$

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
 \Downarrow 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

- Here is an example of a (directly) left-recursive grammar:

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow (E) \mid \text{id}$

- This grammar can be re-written as the following non left-recursive grammar:

$E \rightarrow T E' \quad E' \rightarrow + T E' \mid \epsilon$

$T \rightarrow F T' \quad T' \rightarrow * F T' \mid \epsilon$

$F \rightarrow (E) \mid \text{id}$

Left Factoring

- ▶ Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive, or top-down, parsing.
- ▶ Stmt \rightarrow if expr then stmt else stmt
| if expr then stmt
- ▶ $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$
- ▶ So it should be left factored as

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

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \cdots \mid \alpha \beta_n \mid \gamma$$

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

$$\begin{array}{l} S \rightarrow i E t S \mid i E t S e S \mid a \\ E \rightarrow b \end{array}$$

$$\begin{array}{l} S \rightarrow i E t S S' \mid a \\ S' \rightarrow e S \mid \epsilon \\ E \rightarrow b \end{array}$$

Left-Factoring -- Algorithm

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

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

convert it into

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

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

Left-Factoring – Example1

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



$$A \rightarrow aA' \mid \underline{cdg} \mid \underline{cdeB} \mid \underline{cdfB}$$

$$A' \rightarrow bB \mid B$$



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

$$A' \rightarrow bB \mid B$$

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

Left-Factoring – Example2

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



$$A \rightarrow aA' \mid b$$

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



$$A \rightarrow aA' \mid b$$

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

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

Left-Recursion -- Problem

- A grammar cannot be immediately left-recursive, but it still can be left-recursive.

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Sc \mid d$$

$$\underline{S} \Rightarrow Aa \Rightarrow \underline{S}ca$$

$$\underline{A} \Rightarrow Sc \Rightarrow \underline{A}ac \quad \text{causes to a left-recursion}$$

Removing left recursion: Example 1

- $S \rightarrow Sca/b/da$
- $A \rightarrow Aac/d/bc$

$S \rightarrow bS'/daS'$

$S' \rightarrow caS'/\epsilon$

$A \rightarrow dA'/bcA'$

$A' \rightarrow acA'/\epsilon$

Example 2

- $E \rightarrow EZM/M/2M$
- $M \rightarrow MUF/F$

• $E \rightarrow ME'/2ME'$

• $E' \rightarrow ZME'/\epsilon$

• $M \rightarrow FM'$

• $M' \rightarrow UFM'/\epsilon$

Example 3

- $S \rightarrow S; S / id := E / print(L)$
 - $S \rightarrow id := ES' / print(L) S'$
 - $S' \rightarrow ; SS' / \epsilon$
- $E \rightarrow E + E / id / num / (S, E)$
 - $E \rightarrow id E' / num E' / (S, E) E'$
 - $E' \rightarrow + EE' / \epsilon$
- $L \rightarrow E / L, E$
 - $L \rightarrow EL'$
 - $L' \rightarrow , EL' / \epsilon$

Remove left factoring/left recursion

- $D \rightarrow X / X, D$
- $X \rightarrow id / id[c]$
- $S \rightarrow T / T; S$
- $A \rightarrow id = E / id(E) = E$
- $F \rightarrow C / id / id[E]$
- $C \rightarrow G / GC$
- $S \rightarrow iEtS / iEtSeS$

Types of parsing techniques

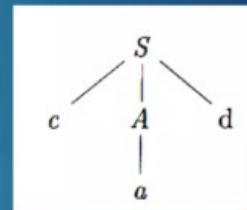
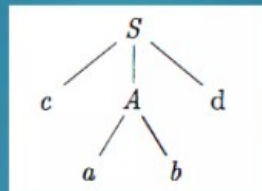
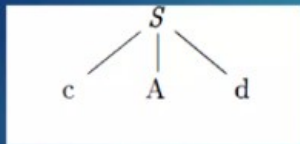
- Top down parsing and bottom up parsing
- Top down parsing:
 - Recursive Descent parsing(or) brute force method
 - Predictive parsing(or) non recursive descent parsing

Top-Down Parsing

- The parse tree is created top to bottom.
- Top-down parser
 - ▶ Recursive-descent parsing
 - ▶ Backtracking is needed
 - ▶ It is a general parsing technique, but not widely used.
 - ▶ Not efficient
 - ▶ Predictive parsing
 - ▶ No backtracking
 - ▶ Efficient
 - ▶ Needs a special form of grammars - (LL(1) grammars).
 - ▶ Recursive predictive parsing is a special form of recursive descent parsing without backtracking.
 - ▶ Non-recursive (table driven) predictive parser is also known as LL(1) parser.

Recursive Descent parsing : string “cad”

$$\begin{array}{lcl} S & \rightarrow & c A d \\ A & \rightarrow & a b \mid a \end{array}$$



Examples for recursive descent parsing

- $S \rightarrow abA$ derive string “ab”
- $A \rightarrow cd/c/\epsilon$

- $S \rightarrow abA$ derive string “abb”
- $A \rightarrow bc/b$

Predictive parsing

- This is also called as **recursive predictive parsing**
- No need of backtracking
- To construct predictive parsing, we need to compute the **First and Follow** for all Non terminals

- ▶ Predictive parsing
 - ▶ No backtracking
 - ▶ Efficient
 - ▶ Needs a special form of grammars - (LL(1) grammars).
 - ▶ Recursive predictive parsing is a special form of recursive descent parsing without backtracking.
 - ▶ Non-recursive (table driven) predictive parser is also known as LL(1) parser.

RULES FOR FIRST SET

1. If X is terminal, then $\text{FIRST}(X)$ is $\{X\}$.
2. If $X \Rightarrow \epsilon$ is a production, then add ϵ to $\text{FIRST}(X)$.
3. If X is nonterminal and $X \Rightarrow Y_1 Y_2 \dots Y_k$ is a production, then place a in $\text{FIRST}(X)$ if for some i , a is in $\text{FIRST}(Y_i)$, and ϵ is in all of $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$; that is, $Y_1, \dots, Y_{i-1} \Rightarrow \epsilon$. If ϵ is in $\text{FIRST}(Y_j)$ for all $j = 1, 2, \dots, k$, then add ϵ to $\text{FIRST}(X)$.

For example, everything in $\text{FIRST}(Y_1)$ is surely in $\text{FIRST}(X)$. If Y_1 does not derive ϵ , then we add nothing more to $\text{FIRST}(X)$, but if $Y_1 \Rightarrow \epsilon$, then we add $\text{FIRST}(Y_2)$ and so on.

RULES FOR FOLLOW SET

1. Place \$ in FOLLOW(S), where S is the start symbol and \$ is the input right end marker.
2. If there is a production $A \Rightarrow a B\beta$, then everything in $\text{FIRST}(\beta)$, except for ϵ , is placed in FOLLOW(B).
3. If there is a production $A \Rightarrow a B$, or a production $A \Rightarrow a B\beta$ where $\text{FIRST}(\beta)$ contains ϵ (i.e., $\beta \Rightarrow \epsilon$), then everything in FOLLOW(A) is in FOLLOW(B).

EXAMPLES

$E \Rightarrow T E'$	$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{ (, \text{id} \}$
$E' \Rightarrow + T E' \mid \epsilon$	$\text{FIRST}(E') = \{ +, \epsilon \}$
$T \Rightarrow F T'$	$\text{FIRST}(T') = \{ *, \epsilon \}$
$T' \Rightarrow * F T' \mid \epsilon$	
$F \Rightarrow (E) \mid \text{id}$	$\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{ \}, \$\}$ $\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{ +,), \$\}$ $\text{FOLLOW}(F) = \{ +, *,), \$\}$

EXAMPLE

$S \Rightarrow A a$	$\text{First}(S) = \{b, d, a\}$
$A \Rightarrow B D$	$\text{First}(A) = \{b, d, \epsilon\}$
$B \Rightarrow b \mid \epsilon$	$\text{First}(B) = \{b, \epsilon\}$
$D \Rightarrow d \mid \epsilon$	$\text{First}(D) = \{d, \epsilon\}$
	$\text{Follow}(S) = \{\$ \}$
	$\text{Follow}(A) = \{a\}$
	$\text{Follow}(B) = \{d, a\}$
	$\text{Follow}(D) = \{a\}$

1. $C \Rightarrow P F \text{ class id } X Y$	$\text{First}(C) = \{\text{public, final, class}\}$	$\text{Follow}(C) = \{\$ \}$
2. $P \Rightarrow \text{public}$		$\text{Follow}(P) = \{\text{final, class}\}$
3. $P \Rightarrow \epsilon$	$\text{First}(P) = \{\text{public, } \epsilon \}$	$\text{Follow}(F) = \{\text{class}\}$
4. $F \Rightarrow \text{final}$	$\text{First}(F) = \{\text{final, } \epsilon \}$	$\text{Follow}(X) = \{\text{implements, } \$ \}$
5. $F \Rightarrow \epsilon$	$\text{First}(X) = \{\text{extends, } \epsilon \}$	$\text{Follow}(Y) = \{\$ \}$
6. $X \Rightarrow \text{extends id}$	$\text{First}(Y) = \{\text{implements, } \epsilon \}$	$\text{Follow}(I) = \{\$ \}$
7. $X \Rightarrow \epsilon$	$\text{First}(I) = \{\text{id}\}$	$\text{Follow}(J) = \{\$ \}$
8. $Y \Rightarrow \text{implements I}$	$\text{First}(J) = \{', \epsilon \}$	
9. $Y \Rightarrow \epsilon$		
10. $I \Rightarrow \text{id J}$		
11. $J \Rightarrow , I$		
12. $J \Rightarrow \epsilon$		

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.

If, after performing the above, there is no production at all in $M[A, a]$, then set $M[A, a]$ to **error** (which we normally represent by an empty entry in the table). \square

Production	Symbol	FOLLOW
$E \rightarrow TE'$	{ (, id }	{ \$,) }
$E' \rightarrow +TE' \epsilon$	{ +, ϵ }	{ \$,) }
$T \rightarrow FT'$	{ (, id }	{ +, \$,) }
$T' \rightarrow *FT' \epsilon$	{ *, ϵ }	{ +, \$,) }
$F \rightarrow (E) id$	{ (, id }	{ *, +, \$,) }

TABLE:
FIRST & FOLLOW

TABLE:
PARSING
TABLE

Non Terminal	INPUT SYMBOLS					
	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	$F \rightarrow id$			$F \rightarrow (E)$		

SYMBOL	FIRST	FOLLOW
$S \rightarrow iEtSS' \mid a$	a, i	\$, e
$S' \rightarrow eS \mid \epsilon$	e, ϵ	\$, e
$E \rightarrow b$	b	t

TABLE:
FAST & FOLLOW

TABLE:
PARSING TABLE

Non Terminal	INPUT SYMBOLS					
	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow eS$ $S' \rightarrow \epsilon$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

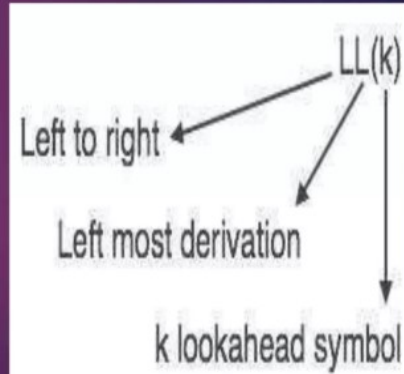
AMBIGUITY

- ❑ LL(1) grammars have distinct properties.
 - No ambiguous grammar or left recursive grammar can be LL(1).
- ❑ Thus , the given grammar is not LL(1).

□ An LL parser is called an $LL(k)$ parser if it uses k tokens of look ahead when parsing a sentence.

□ LL grammars, particularly $LL(1)$ grammars, as parsers are easy to construct, and many computer languages are designed to be $LL(1)$ for this reason.

□ The 1 stands for using **one** input symbol of look ahead at each step to make parsing action decision.



Properties of $LL(1)$

- No ambiguous or left recursive grammar can be $LL(1)$.
- Grammar G is $LL(1)$ **iff** whenever $A \rightarrow \alpha / \beta$ are two distinct productions of G and:
 - For no terminal a do both α and β derive strings beginning with a .
 - At most one of α and β can derive the empty string.
 - If $\beta \Rightarrow \epsilon$, the α does not derive any string beginning with a terminal in $FOLLOW(A)$.

$$FIRST(\alpha) \cap FIRST(\beta) = \emptyset$$

$$FIRST(\alpha FOLLOW(A)) \cap FIRST(\beta FOLLOW(A)) = \emptyset$$

Nonrecursive Predictive Parsing

- It is possible to build a **nonrecursive** predictive parser
- This is done by maintaining an **explicit stack**

2

Table-driven Parsers

- The nonrecursive LL(1) parser looks up the production to apply by looking up a **parsing table**

3

Table-driven Parsers

LL(1) table:

- One dimension for current **non-terminal** to expand
- One dimension for **next token**
- Table entry contains one **production**

Consider the expression grammar

1	E	\rightarrow	$T E'$
2	E'	\rightarrow	$+ T E'$
3		$ $	ε
4	T	\rightarrow	$F T'$
5	T'	\rightarrow	$* F T'$
6		$ $	ε
7	F	\rightarrow	(E)
8		$ $	<u>id</u>

Predictive Parsing Table

◆◆	id	+	*	()	\$
<i>E</i>	$E \rightarrow TE'$			$E \rightarrow TE'$		
<i>E'</i>		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
<i>T</i>	$T \rightarrow FT'$			$T \rightarrow FT'$		
<i>T'</i>		$T' \rightarrow \epsilon$	$T \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
<i>F</i>	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

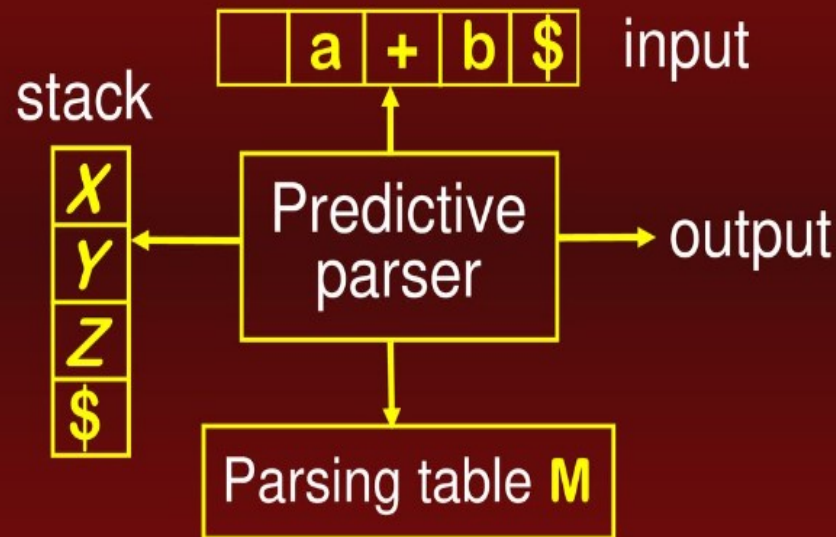
Columns for next token

Blank entries are errors

Predictive Parsers

- The predictive parser uses an **explicit stack** to keep track of pending non-terminals
- It can thus be implemented **without recursion**.

Predictive Parsers



LL(1) Parsing Algorithm

- The *input buffer* contains the string to be parsed; **\$** is the end-of-input marker
- The *stack* contains a sequence of grammar symbols
- **Initially**, the *stack* contains the **start symbol** of the grammar on the top of **\$**.

The parser is controlled by a program that behaves as follows:

- The program considers X , the symbol on **top** of the **stack**, and **a**, the current **input** symbol.
- These two symbols, X and **a** **determine** the **action** of the parser.
- There are **three** possibilities.

1. $X = a = \$$,
the parser **halts** and announces successful completion.
2. $X = a \neq \$$
the parser pops X off the stack and advances input pointer to next input symbol
3. If X is a nonterminal, the program consults entry **$M[X, a]$** of parsing table **M** .

If the entry is a production
 $M[X,a] = \{X \rightarrow UVW\}$
 the parser replaces X on
 top of the stack by WVU
 (with U on top).

If $M[X,a] = \text{error}$, the parser
 calls an error recovery
 routine

Example:

Let's parse the input string

id+id*id

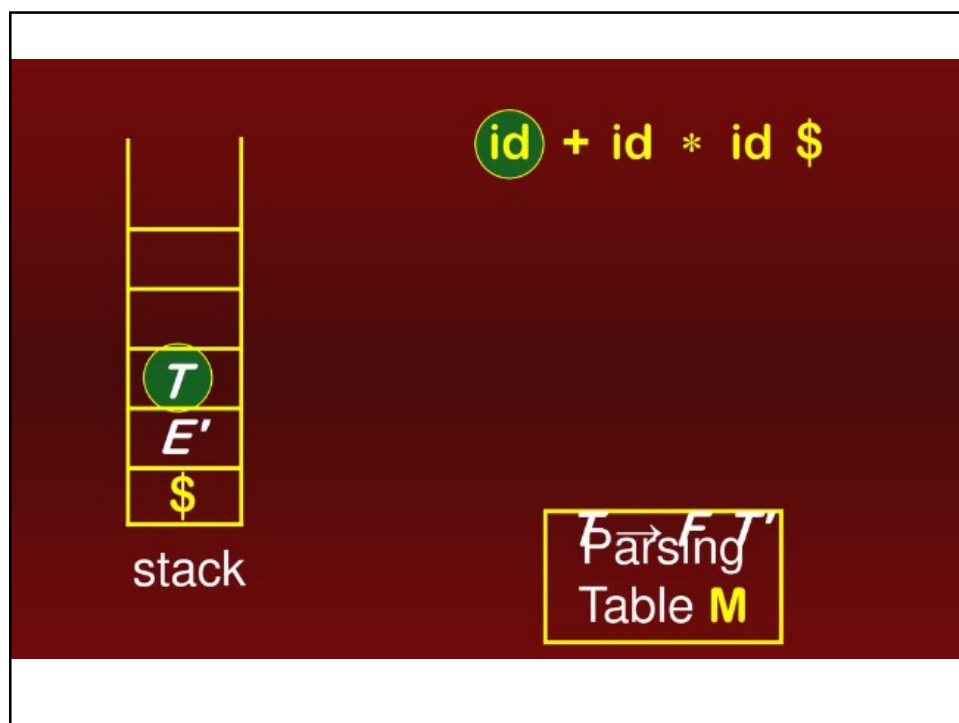
using the nonrecursive
 LL(1) parser

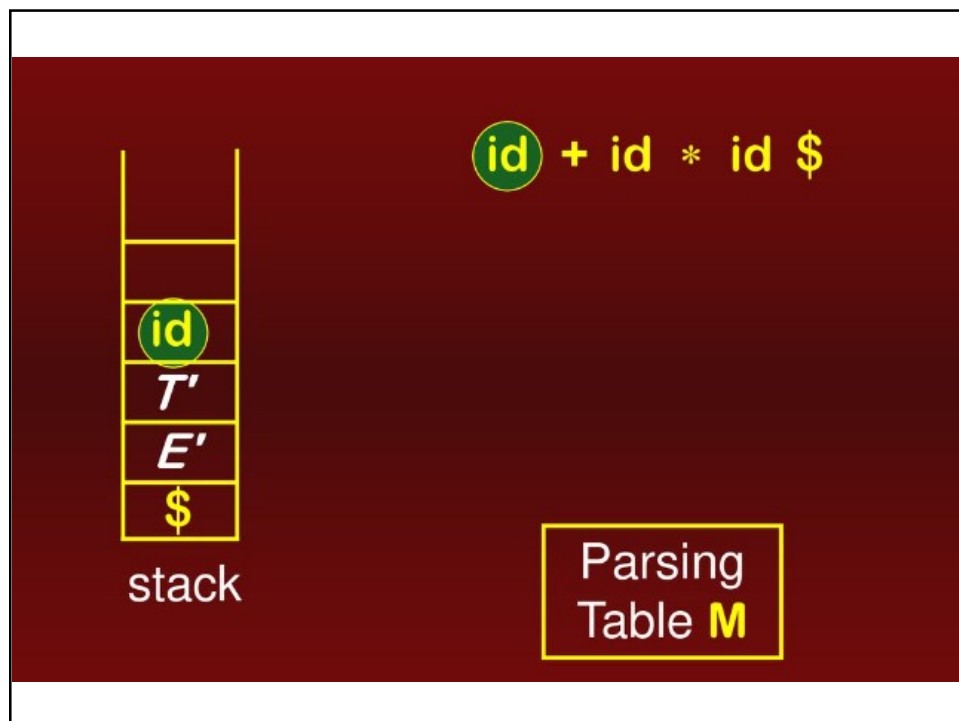
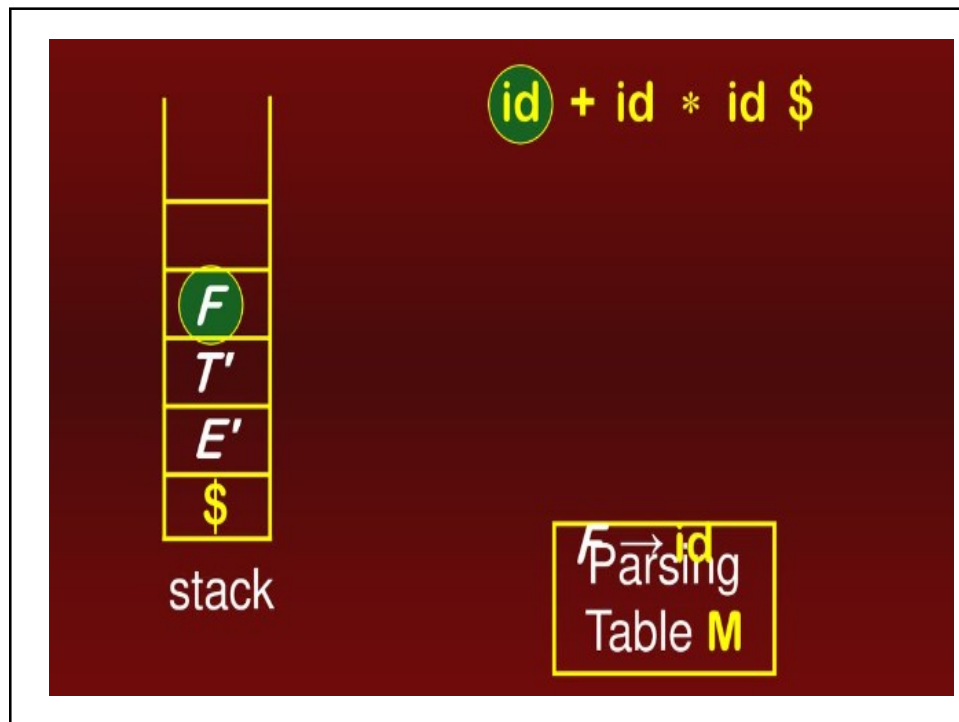


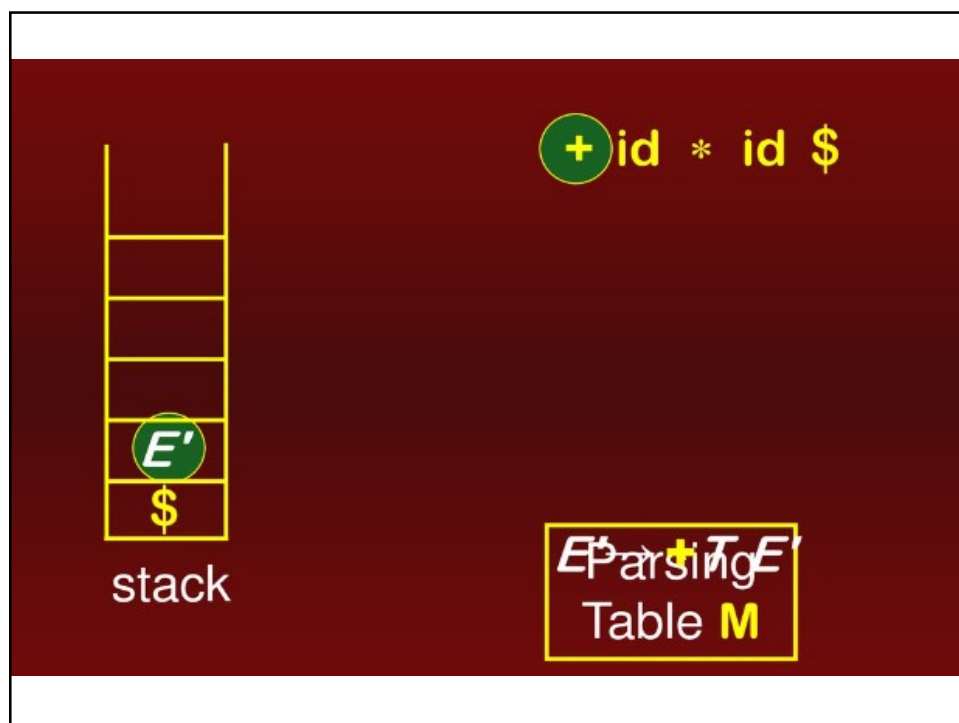
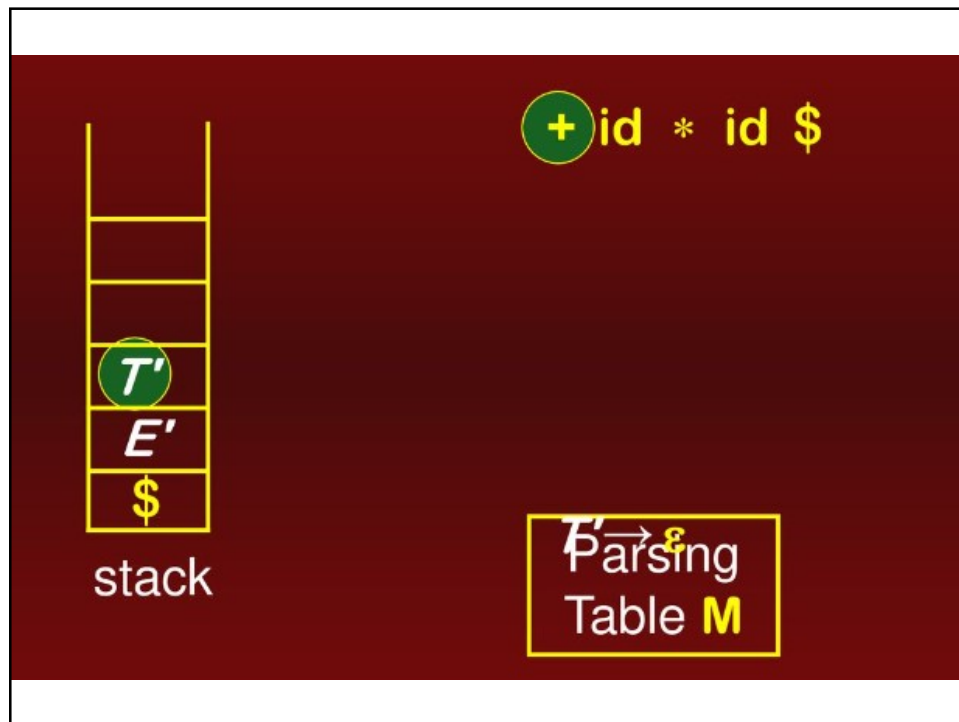
stack

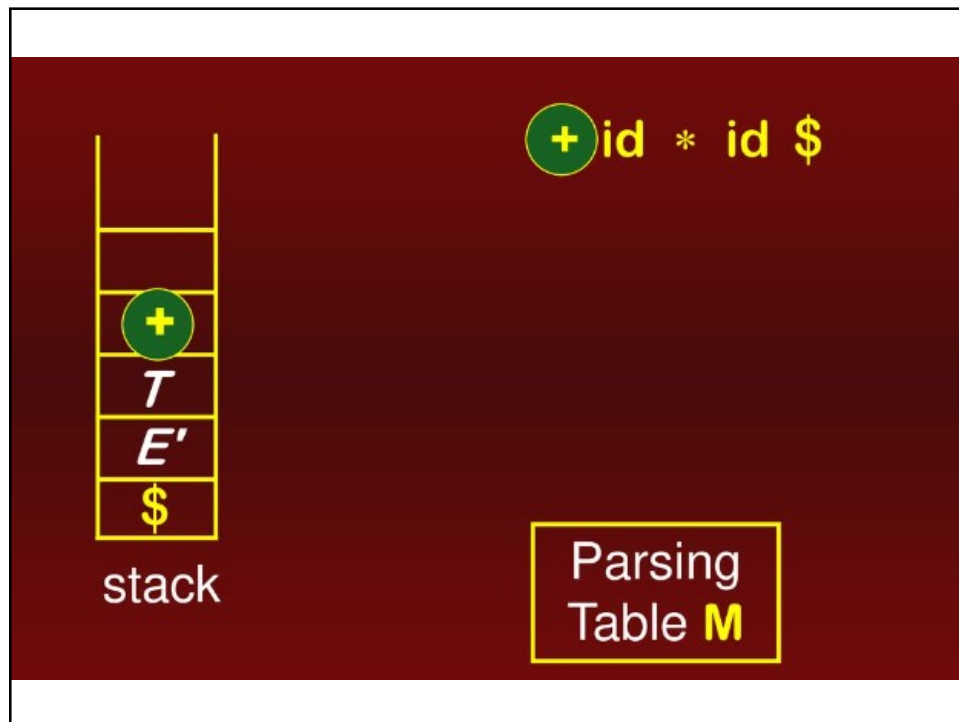
id + id * id \$

Parsing
 Table M









Stack	Input	Output
\$E	id+id*id\$	
\$E' T	id+id*id\$	$E \rightarrow TE'$
\$E' T' F	id+id*id\$	$T \rightarrow FT'$
\$E' T' id	id+id*id\$	$F \rightarrow id$
\$E' T'	+id*id\$	
\$E'	+id*id\$	$T' \rightarrow \epsilon$
\$E' T+	+id*id\$	$E' \rightarrow +TE'$

<i>Stack</i>	<i>Input</i>	<i>Ouput</i>
$\$E'T$	$\underline{id}*\underline{id}\$$	
$\$E'T'F$	$\underline{id}*\underline{id}\$$	$T \rightarrow FT'$
$\$E'T'\underline{id}$	$\underline{id}*\underline{id}\$$	$F \rightarrow \underline{id}$
$\$E'T'$	$*\underline{id}\$$	
$\$E'T'F*$	$*\underline{id}\$$	$T \rightarrow *FT'$
$\$E'T'F$	$\underline{id}\$$	
$\$E'T'\underline{id}$	$\underline{id}\$$	$F \rightarrow \underline{id}$

<i>Stack</i>	<i>Input</i>	<i>Ouput</i>
$\$E'T'$	$\$$	
$\$E'$	$\$$	$T' \rightarrow \epsilon$
$\$$	$\$$	$E' \rightarrow \epsilon$

Error Recovery in Predictive Parsing

- An error is detected during the predictive parsing when the terminal on top of the stack does not match the next input symbol, or when nonterminal A on top of the stack, a is the next input symbol, and parsing table entry $M[A,a]$ is empty.
- The parser design should be able to provide an **error message** (an error message which depicts as much possible information as it can).
- It should be **recovering from that error case**, and it should be able to **continue parsing** with the rest of the input.

4 Types of error recovery strategies

1. Panic-Mode Recovery
2. Phrase-Level Recovery
3. Error Productions
4. Global Correction

Panic mode error recovery

- Panic-mode error recovery says that all the input symbols are skipped until a synchronizing token is found from the string.
- **What is the synchronizing token?**
- The terminal symbols which lie in follow set of non-terminal they can be used as a synchronizing token set for that non-terminal

How to select synchronizing set?

- Place all symbols in **FOLLOW(A)** into the **synchronizing set for non terminal A**. If we skip tokens until an element of FOLLOW(A) is seen and **pop A from the stack**, it likely that parsing can continue.
- We might add keywords that begins statements to the synchronizing sets for the non terminals generating expressions.

- If a **non terminal can generate the empty string**, then the **production deriving ϵ can be used as a default**. This may postpone some error detection, but cannot cause an error to be missed. This approach reduces the number of non terminals that have to be considered during error recovery.
- If **a terminal on top of stack cannot be matched**, a simple idea is to **pop the terminal**, issue a message saying that the terminal was inserted.

Example: error recovery

“synch” indicating synchronizing tokens obtained from FOLLOW set of the nonterminal in question.

If the parser looks up entry $M[A,a]$ and finds that it is blank, the input symbol a is skipped.

If the entry is synch, the nonterminal on top of the stack is popped.

If a token on top of the stack does not match the input symbol, then we pop the token from the stack.

$$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{ (, \text{id} \}.$$

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

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

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

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

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

NONTERMINAL	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

Fig. 4.18. Synchronizing tokens added to parsing table of Fig. 4.15.

Example: error recovery (II)

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

Fig. 4.19. Parsing and error recovery moves made by predictive parser.

Example 2

$S \rightarrow AbS|E|\epsilon$
 $A \rightarrow a|cAd$
 $\text{Follow}(S) = \{\$ \}$
 $\text{Follow}(A) = \{b, d\}$

	a	b	c	d	e	\$
S	$S \rightarrow Abs$	sync	$S \rightarrow Abs$	sync	$S \rightarrow e$	$S \rightarrow \epsilon$
A	$A \rightarrow A$	sync	$A \rightarrow cAd$	sync	sync	sync

stack	input	output	stack	input	output
\$S	aab\$	$S \rightarrow AbS$	\$S	ceadb\$	$S \rightarrow AbS$
\$SbA	aab\$	$A \rightarrow a$	\$SbA	ceadb\$	$A \rightarrow cAd$
\$Sba	aab\$		\$SbdAc	ceadb\$	
\$Sb	ab\$	Error: missing b, inserted	\$SbdA	ceadb\$	Error: unexpected
\$S	ab\$	$S \rightarrow AbS$	e		(illegal A)
\$SbA	ab\$	$A \rightarrow a$	(Remove all input tokens until first b or d .pop A)		
\$Sba	ab\$		\$Sbd	db\$	
\$Sb	b\$		\$Sb	b\$	
\$S	S	$S \rightarrow \epsilon$	\$S	\$	$S \rightarrow \epsilon$
\$	\$	accept	\$	\$	accept

Phrase-level recovery

- Phrase level recovery is implemented by filling in the **blank entries** in the predictive parsing table with **pointers to error routines**.
- At each unfilled entry in the parsing table, it is filled by a pointer to a special error routine that will take care of that error case specifically.
- These error routines can be of different types like :
 - change, insert, or delete input symbols.
 - issue appropriate error messages.
 - pop items from the stack.
- We should be careful while we design these error routines because we may put the parser into an infinite loop

Error productions

- If a user has knowledge of common errors that can be encountered then, these errors can be incorporated by augmenting the grammar with error productions that generate erroneous constructs.
- If this is used then, during parsing appropriate error messages can be generated and parsing can be continued.

Example: Suppose the input string is abcd.

Grammar: $S \rightarrow A$

$A \rightarrow aA \mid bA \mid a \mid b$

$B \rightarrow cd$

The input string is not obtainable by the above grammar, so we need to add Augmented Grammar.

Grammar: $E \rightarrow SB$ // AUGMENT THE GRAMMAR

$S \rightarrow A$

$A \rightarrow aA \mid bA \mid a \mid b$

$B \rightarrow cd$

Now, string abcd is possible to obtain.

Global corrections

- The parser examines the whole program and tries to find out the closest match for it which is error-free.
- The closest match program has less number of insertions, deletions, and changes of tokens to recover from erroneous input.
- Due to high time and space complexity, this method is not implemented practically.

1. Universal parsers:-

- These parsers perform parsing. Cocke - Kalamy - Younger (CKY) algorithm. It uses the Chomsky Normal Form (CNF) of the CFG.
- It is efficient and not used in Commercial Compilers.

2. Top down parsers: -

- These parsers build parse trees starting from the root node and work up to the Leaves.
- These parsers Satisfy LL grammars.
- These Scanners scan from left to right and use left-most derivation.
- These parsers are Categorized into
 - back track parsing.
 - Predictive parsing

- In back track parsing, if a sequence of erroneous expansions are made and a mismatch is discovered, the input pointer rolled back to the initial position and all effects of parsing are undone.
- In predictive parsing, using current non terminal to be expanded and the next input symbol, the parser makes a decision of which production needs to be expanded.
- **Two ways of implementing predictive parsers:**
 - Recursive descent parsing have Procedure for Parsing each non-terminal using the recursive nature of The procedures.
 - Table driven parsing. determines the next production to be applied by using a parsing table. Indexed by the Current non-terminal and the next input symbol.

3. Bottom up Parsers:

- These Parsers build parse trees starting from the leaves and work up to the root node.
- They satisfy LR grammars.
- They scan from left to right and use right most derivation.
- Here the input is reduced to the start symbol.
- These parsers are implemented using stacks.
- Here the next input symbol is shifted onto the stack
- Then the top few elements are popped out of the stack if a match occurs to right side of the production. It is then reduced to the left side of the production.
- Bottom up parsing is also called as shift reduce parsing.
- **These parsers are categorized into**
 - Operator precedence parsers
 - LR parsers

- In operator precedence parsers, the choice of whether to shift the next symbol. Or reduce the processed input is decided using a precedence relation table.
- In LR Parsers, parsing table derived out of grammar is used for deciding whether to shift the next input Symbol or reduce the processed input.

These parsers are categorized as

- Simple LR (SLR)
- Canonical LR (CLR) and
- Look Ahead LR (LALR) based on the way the parser tables are derived.

