

University of Sulaimani
College of Science
Department of Computer Science



Compiler

Second Phase of the Compiler

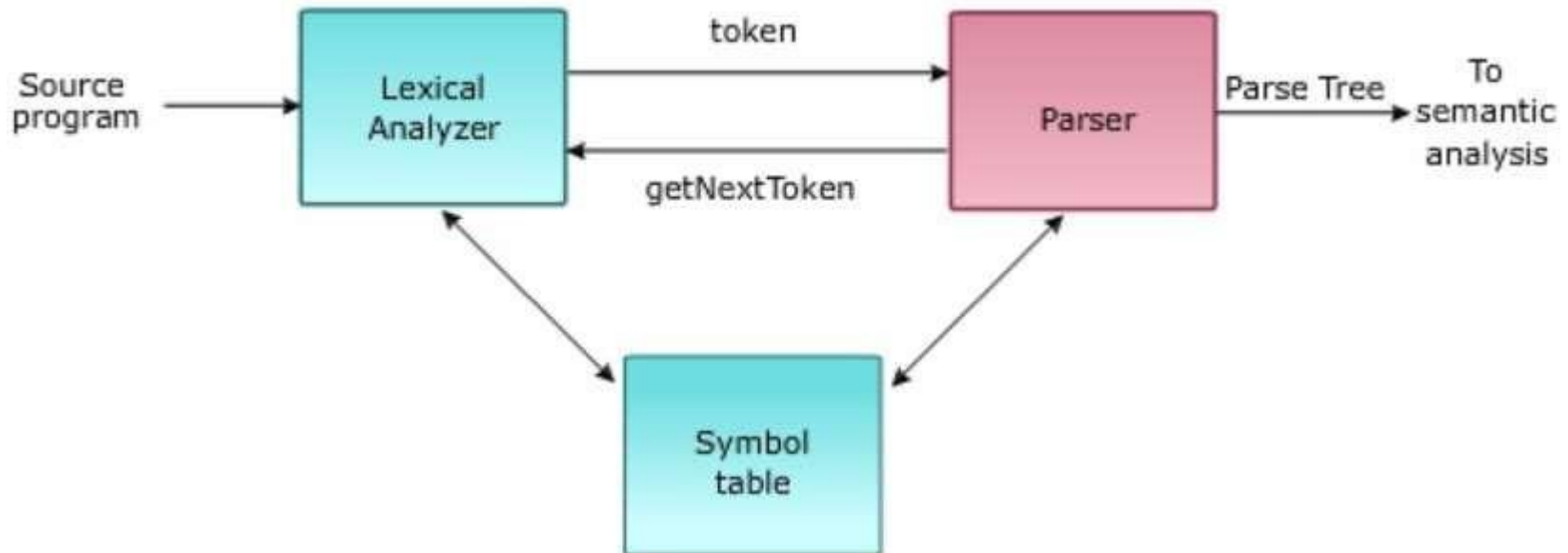
Syntax Analyzer

2024-2025

Mzhda Hiwa

Syntax Analyzer

The second phase of the compiler is syntax analysis or parsing. The parser uses the first components of the tokens produced by the lexical analyzer to create a tree-like intermediate representation that depicts the grammatical structure of the token stream.



Role of the Syntax Analyzer

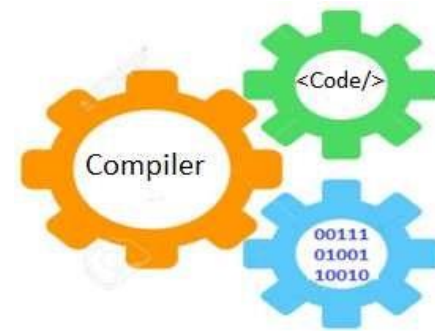
- Takes input from the lexical analyzer (a stream of tokens) and organizes them into a hierarchical structure.
- Builds a parse tree or abstract syntax tree (AST) that represents the syntactic structure of the code. This tree is used in later stages of the compilation process.
- Ensures the source code follows the correct grammatical structure.
- Detects syntax errors like missing semicolons, unmatched brackets, incorrect ordering of tokens, etc. Reports to user where any syntax error in the source code are.



Introduction to Parsing → Syntax Analysis

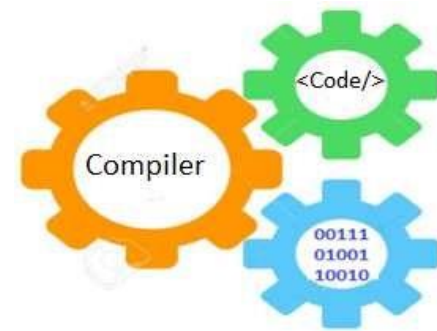
- **Parsing** is the process of determining how a string of terminals can be generated by a grammar. $F \rightarrow id \mid (E)$
- A parser must be capable of constructing the tree in principle, or else the translation cannot be guaranteed correct.
- A parser scans the input string from **left to right** and it makes use of production rules for choosing appropriate derivation.
- **Two parsing techniques:**
 - **Top-down parsing**
 - **Bottom-up parsing**

abbcde
→



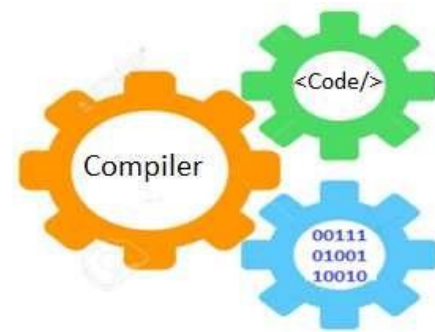
Top-down parsing

- A parser is top-down if it generates a parse tree starting from the root and precedes towards the leaves.
 - It is easier to understand and program manually
 - A leftmost derivation is applied at each derivation step
- Two kinds of top-down parsing techniques will be studied
1. Recursive-Decent parser
 2. Predictive parser



Bottom-Up Parsing

- A bottom-up parse corresponds to the construction of a parse tree for an input string beginning at the leaves (the bottom) and working up towards the root (the top).
- Bottom-up parsing is more general than top-down parsing.



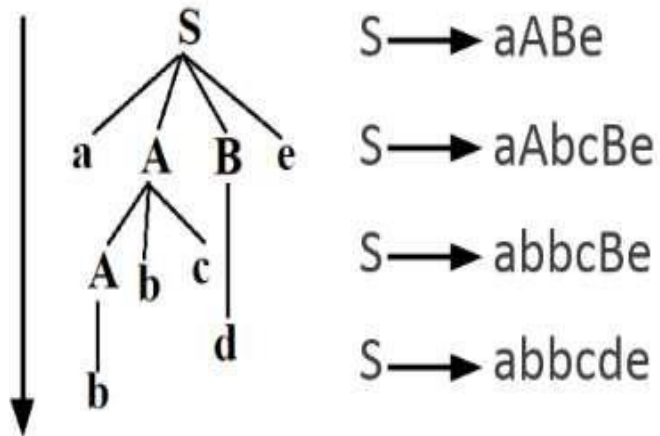
Example

$S \rightarrow aABe$ $w=abbcde$

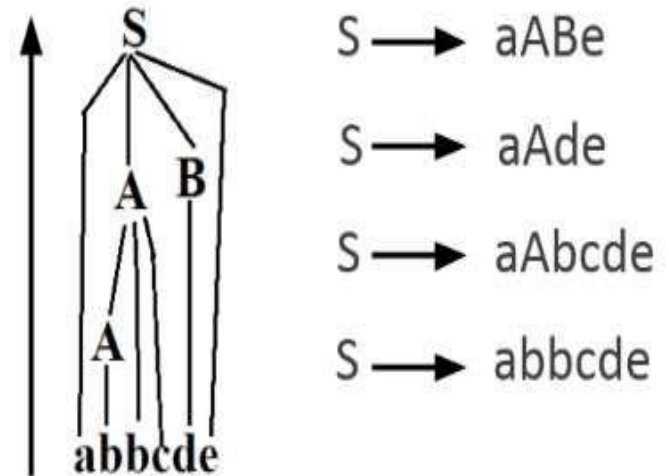
$A \rightarrow Abc \mid b$

$B \rightarrow d$

• Top-down LMD



• Bottom-up RMD

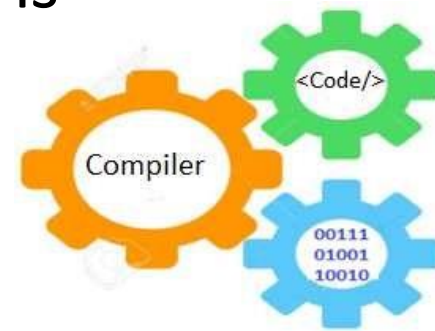


1- Top-Down Parsing by Recursive-Descent

- It reads characters from the input stream and matches them with terminals from the grammar.

The operation involved are : (↺)

- Start from the “Start non-terminal” and select a rule from the production rules (CFG)
- If it was not a correct rule, then backtrack and choose another rule.
- If every production is unsuitable for string match, then parse tree cannot be built, and syntax error is reported.



Example1

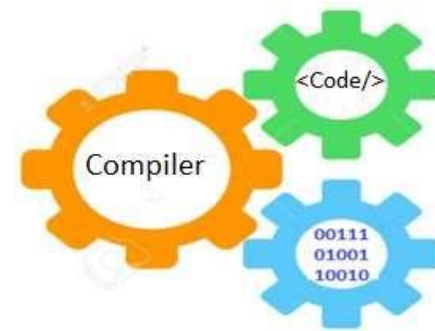
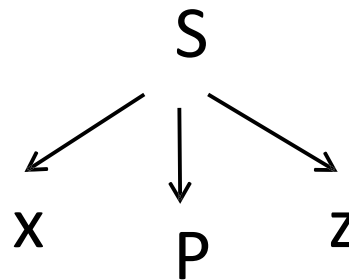
- Consider grammar

$S \rightarrow xPz$

$P \rightarrow yw \mid y$

- For Token stream is: xyz

1- Select rule $S \rightarrow xPz$



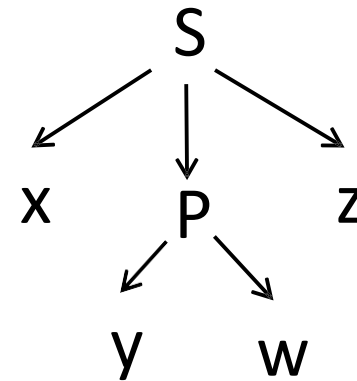
Example1...cont'd

$S \rightarrow xPz$

$P \rightarrow yw \mid y$

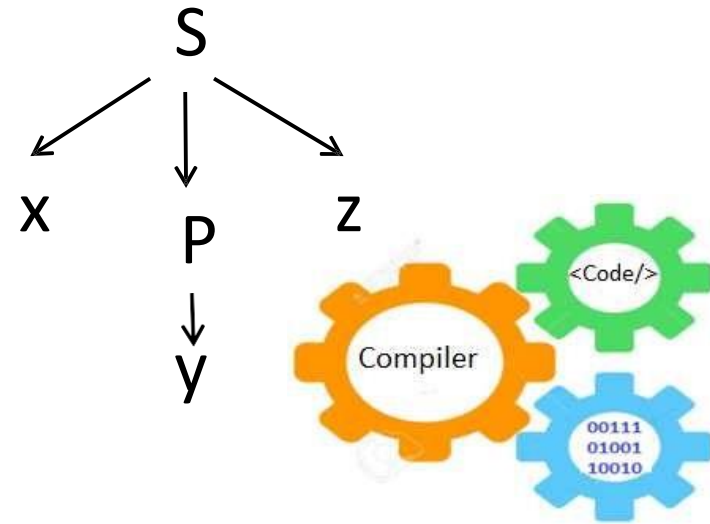
2- Select rule $P \rightarrow yw$

- Not correct



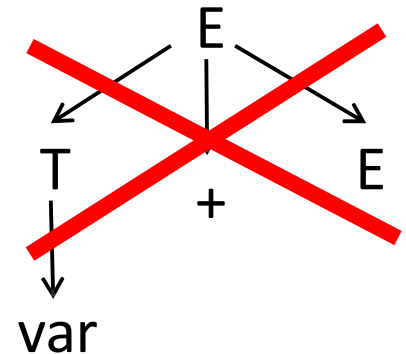
3- Select rule $P \rightarrow y$

- Correct



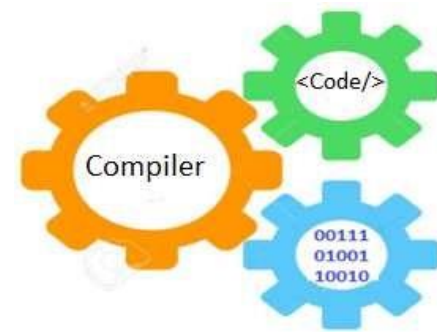
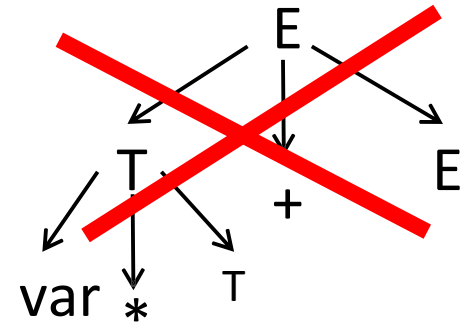
Example2

- Consider the grammar
 - $E \rightarrow T + E \mid T$
 - $T \rightarrow \text{var} \mid \text{var} * T$
- Token stream is: $\text{var} * \text{var}$
- Start with top-level non-terminal E
- Try the rules for E in order
 - Try $E \rightarrow T + E$
 - Then try a rule for $T \rightarrow \text{var}$
 - Token matches **var**
 - But **+** after **var** does not match input token *****



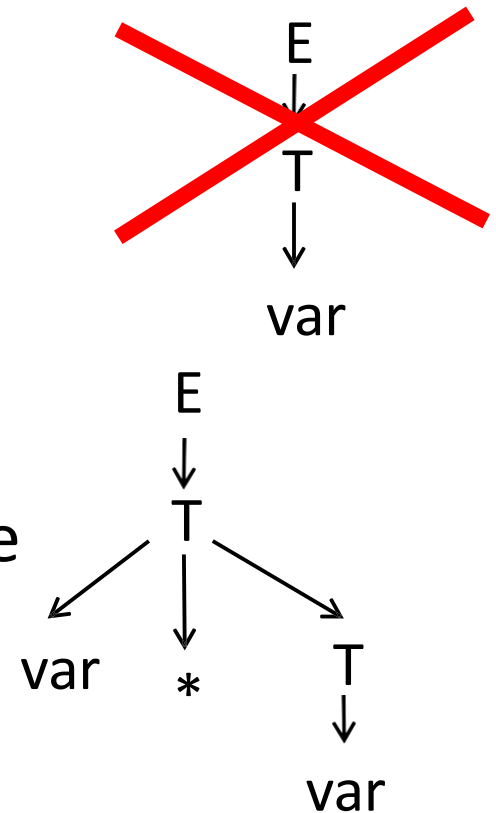
Example2...cont'd

- Try $T \rightarrow \text{var} * T$
 - Token matches.
 - This will match but $+$ after T will be unmatched
- Has exhausted the choices for $E \rightarrow T + E$
- Backtrack to choice for E



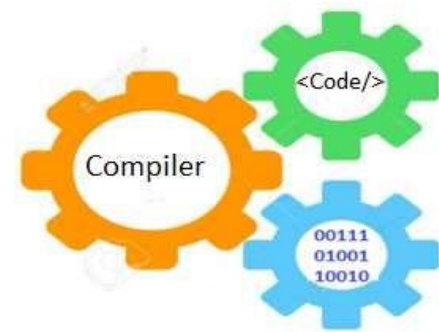
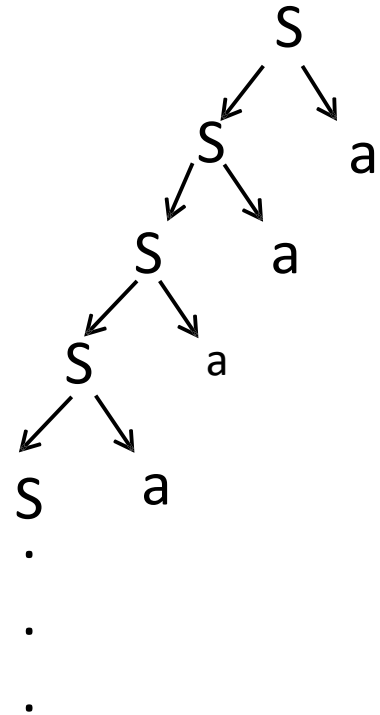
Example2...cont'd

- Try $E \rightarrow T$
- Follow same steps as before for T
 - Try a rule for $T \rightarrow \text{var}$
Token matches **var**
 - But there is no other token after **var**
- Try $T \rightarrow \text{var} * T$
 - Then try $T \rightarrow \text{var}$
 - Succeed with the following parse tree



Notes

- Easy to implement by hand.
- But does not always work ...
- Consider a production $S \rightarrow S a$
- S will get into an infinite loop.
- This case is called **left-recursion**.
- Recursive descent does not work in such cases.



Left Recursion

- A production of grammar is said to have **left recursion** if the leftmost variable of its RHS is same as variable of its LHS.
- A grammar containing a production having left recursion is called as **Left Recursive Grammar**.

Elimination of Left Recursion

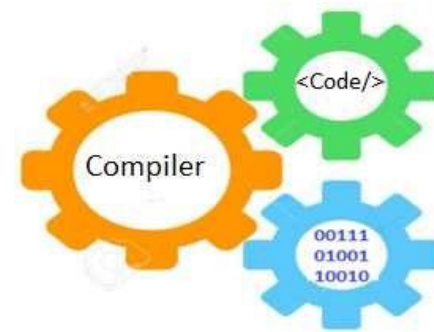
If we have the left-recursive pair of productions-

$$A \rightarrow A\alpha / \beta$$

Then, we can eliminate left recursion by replacing the pair of productions with-

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' / \epsilon$$



Example of Elimination of Left Recursion

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

Eliminate immediate left recursion from the Grammar

Solution

Comparing $E \rightarrow E + T \mid T$ with $A \rightarrow A \alpha \mid \beta$

$$A = E, \alpha = +T, \beta = T$$

$A \rightarrow A \alpha \mid \beta$ is changed to $A \rightarrow \beta A'$ and $A' \rightarrow \alpha A' \mid \epsilon$

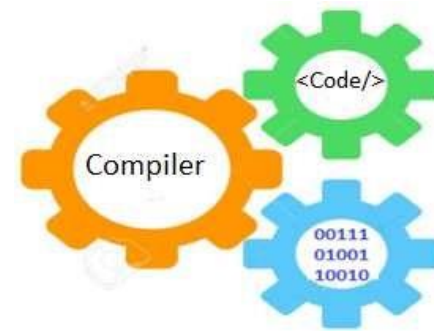
$A \rightarrow \beta A'$ means $E \rightarrow TE'$

$A' \rightarrow \alpha A' \mid \epsilon$ means $E' \rightarrow +TE' \mid \epsilon$

Result

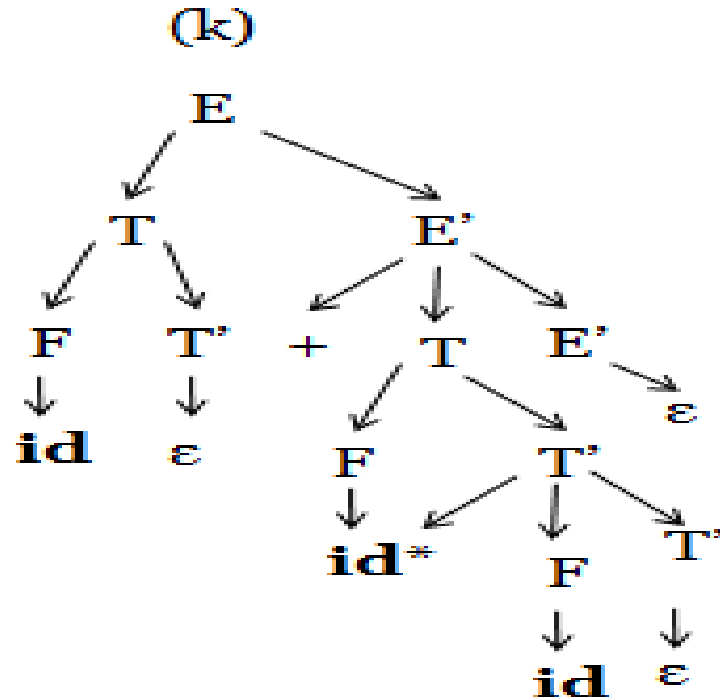
$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

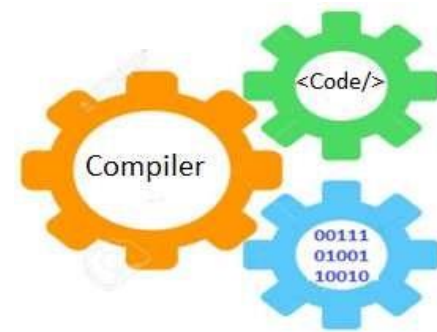


2- Top-Down Parsing by Predictive parser

- Consider Grammar

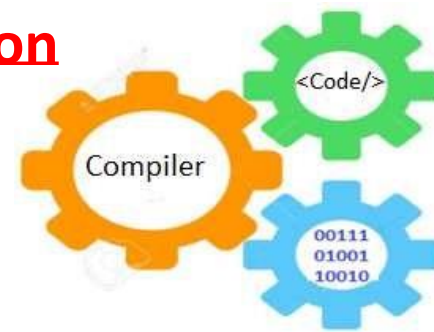
$$E \rightarrow T E'$$
$$E' \rightarrow + T E' \mid \varepsilon$$
$$T \rightarrow F T'$$
$$T' \rightarrow * F T' \mid \varepsilon$$
$$F \rightarrow (E) \mid \text{id}$$


- String `id + id * id`



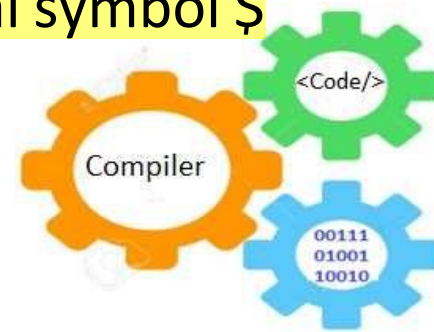
2- Top-Down Parsing by Predictive parser

- Predicts which production to use
 - By looking at the next few tokens, using “lookahead” variable
 - No backtracking
- Predictive parsers accept LL(k) grammars
 - L means “Left-to-Right” scan of input
 - L means “Leftmost derivation”
 - k means “predict based on k tokens of lookahead”
 - In practice, LL(1) is used
 - LL(k) grammar must be unambiguous
 - LL(k) grammar must not include any **left-recursion**



LL(1) Parser

- **input buffer** : The string to be parsed
- **Output**: A production rule representing a step of the derivation sequence (left-most derivation) of the string in the input buffer.
- **Stack**: keeps the grammar symbols, initially contains \$.
- The symbols in RHS of rule are pushed into the stack in reverse order i.e. from right to left
- **parsing table**:
 - a two-dimensional array
 - each row is a non-terminal symbol
 - each column is a terminal symbol or the special symbol \$
 - each entry holds a production rule.



Input token

A diagram of a stack structure. It consists of a vertical rectangle divided into three horizontal sections. The top section is shaded gray. The middle section is white and contains the text '\$'. The bottom section is shaded gray. To the left of the rectangle, the word 'top' is written, followed by an arrow pointing to the boundary between the top and middle sections.

LL(1) Parser

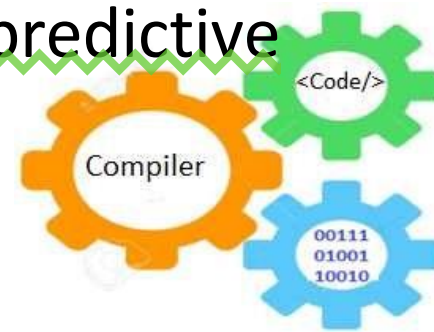
	a	+	b	\$
A				
B				
C				

Building Predictive Parser

Three steps

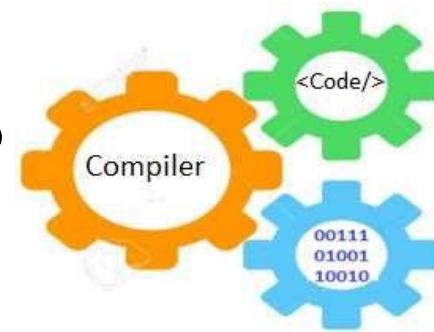
1. Compute FIRST and FOLLOW
2. Construct the predictive parsing table
3. Parse the input string

- **Note:** the grammar must be unambiguous and Left Recursion must be eliminated
- First and follow are used to construct the predictive parsing table



Computing First and Follow

- **FIRST(α)** is a set of the terminal symbols which occur as first symbols in strings derived from α . Where α is any string of grammar (terminals and non-terminals).
- if α derives to ϵ , then ϵ is also in FIRST(α).
- **FOLLOW(A)** is the set of the terminals which occur immediately after (follow) the *non-terminal* A. If the strings derived from the starting symbol.
_ \$ is in FOLLOW(A) if $S \Rightarrow \alpha A$
- a terminal **a** is in FOLLOW(A) if $S \Rightarrow \alpha A a \beta$

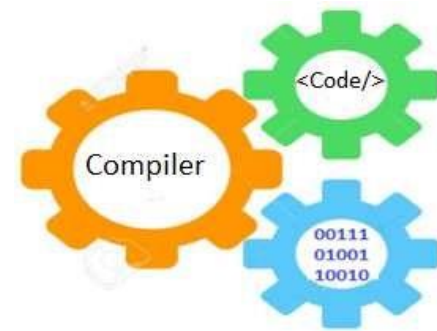


Computing FIRST

- $\text{FIRST}(a) = \{a\}$ if $a \in T$
- $\text{FIRST}(\epsilon) = \{\epsilon\}$
- $\text{FIRST}(X)$ for a non-terminal X

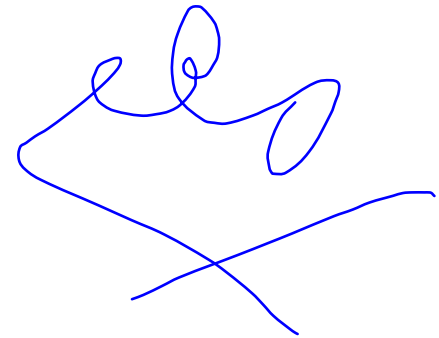
If there is production $X \rightarrow Y_1 Y_2 \dots Y_k$ then

- $\text{FIRST}(X) = \text{FIRST}(Y_1) - \{\epsilon\}$.
- But, if $\epsilon \in \text{FIRST}(Y_1)$, then add $\text{FIRST}(Y_2) - \{\epsilon\}$
- And, if $\epsilon \in \text{FIRST}(Y_2), \dots$

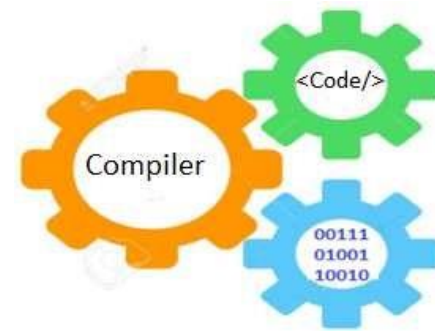


Example 1

- $E \rightarrow TE'$
 $E' \rightarrow +TE' \mid \varepsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' \mid \varepsilon$
 $F \rightarrow \text{id} \mid (E)$



- $\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{ \text{id}, (\}$
- $\text{FIRST}(E') = \{ +, \varepsilon \}$
- $\text{FIRST}(T) = \text{FIRST}(F) = \{ \text{id}, (\}$
- $\text{FIRST}(T') = \{ *, \varepsilon \}$
- $\text{FIRST}(F) = \{ \text{id}, (\}$



Example 2

- $type \rightarrow simple$

 - | $\wedge id$

 - | **array** [*simple*] **of type**

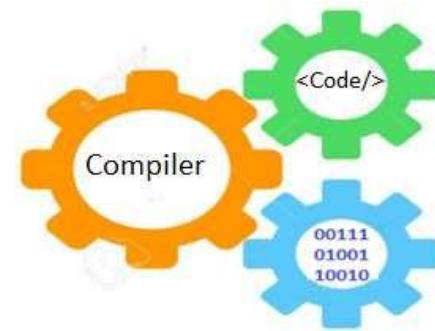
- $simple \rightarrow integer$

 - | **char**

 - | **num dot num**

- $FIRST(simple) = \{ integer, char, num \}$

- $FIRST(type) = \{ integer, char, num, \wedge, array \}$



Exercise

Find FIRST for the following grammar

• $S \rightarrow ACB \mid CbB \mid Ba$

$A \rightarrow da \mid BC$

$B \rightarrow g \mid \epsilon$

$C \rightarrow h \mid \epsilon$

• $S \rightarrow Aa$

$A \rightarrow bdZ \mid eZ$

$Z \rightarrow cZ \mid adZ \mid \epsilon$

$f(Z) = \{c, a, \epsilon\}$

$f(A) = \{b, e\}$

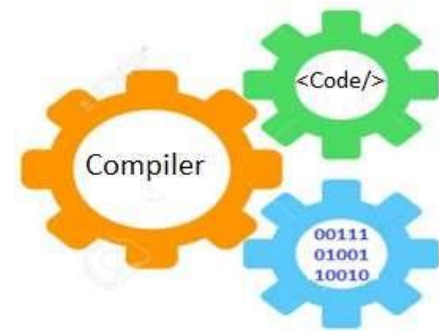
$f(S) = f(A) = \{b, e\}$

Handwritten solutions for the first grammar:

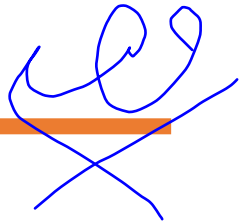
$$F(S) = \{d, g, h, \epsilon, b, a, A\}$$

$$F(A) = \{d, g, h, \epsilon\}$$

$$F(B) = \{g, \epsilon\}$$

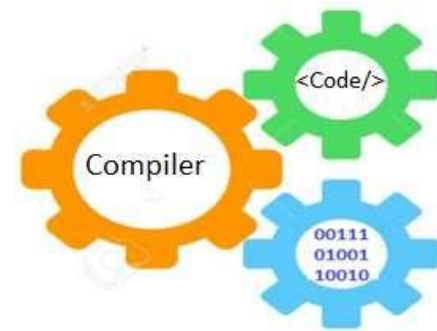
$$F(C) = \{h, \epsilon\}$$


Computing FOLLOW



- * If S is the start symbol $\rightarrow \$$ is in $\text{FOLLOW}(S)$
- * if $A \rightarrow \alpha B \beta$ is a production rule
 \rightarrow everything in $\text{FIRST}(\beta)$ is $\text{FOLLOW}(B)$ except ϵ
- * If ($A \rightarrow \alpha B$ is a production rule) or
($A \rightarrow \alpha B \beta$ is a production rule and ϵ is in $\text{FIRST}(\beta)$)
 \rightarrow everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(B)$.

We apply these rules until nothing more can be added to any follow set



Example 1

$E \rightarrow T E'$

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

$T \rightarrow F T'$

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

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

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

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

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

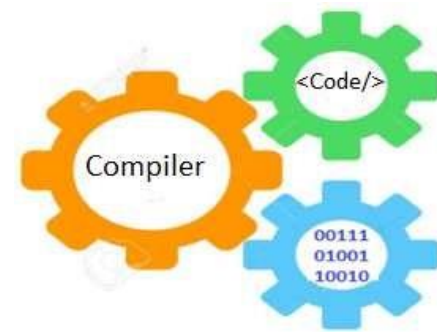
- $\text{FOLLOW}(E) = \{), \$ \}$

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

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

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

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



Example 2

$$E \rightarrow T E'$$

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

$$T \rightarrow F T'$$

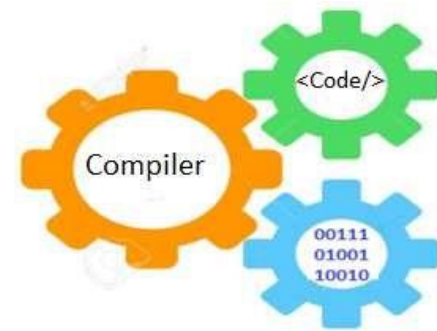
$$T' \rightarrow * F T' \mid / F T' \mid \varepsilon$$

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

	First	Follow
E	{num , id}	{ \$ }
E'	{ + , - , ε }	{ \$ }
T	{num , id}	{ + , - , \$ }
T'	{ * , / , ε }	{ + , - , \$ }
F	{num , id}	{ * , / , + , - , \$ }

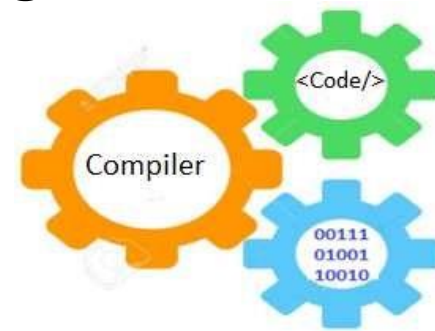
Notes

- To compute $\text{FIRST}(A)$ you must look for A on a production's **left-hand** side.
- To compute $\text{FOLLOW}(A)$ you must look for A on a production's **right-hand** side.
- FIRST sets are always sets of terminals (plus, perhaps epsilon).
- FOLLOW sets are always sets of terminals (plus, perhaps \$).
- Nonterminals are *never* in a FIRST or a FOLLOW set.
- epsilon is *never* in a FOLLOW set.



Constructing the Parse Table

- Parse table summarizes the applicable RHS for each terminal/non-terminal combination.
- Construct a parsing table T for CFG :
- For each production $X \rightarrow \alpha$
 - Add $\rightarrow \alpha$ to the X row for each symbol in $\text{FIRST}(\alpha)$
 - If α is nullable, add $\rightarrow \alpha$ for each symbol in $\text{FOLLOW}(X)$
 - Entry for $[S, \$]$ is ACCEPT
 - All other undefined entries of the parsing table are error entries.



Parse table Example (1)

$E \rightarrow T E'$

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

$T \rightarrow F T'$

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

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

	First	Follow
E	{(, id}	{), \$}
E'	{+, ε }	{), \$}
T	{(, id}	{+,), \$}
T'	{*, ε }	{+,), \$}
F	{(, id}	{*, +,), \$}

Parse table Example 1

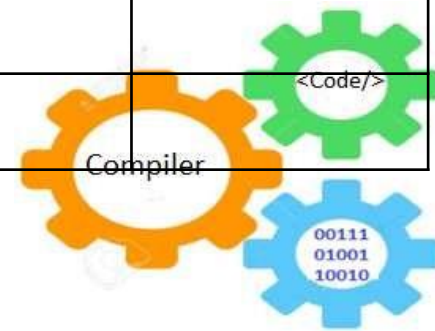
- Create table with:
 - Put each terminals in the columns
 - Put each non-terminals to rows

	+	*	()	id	\$
E						
E'						
T						
T'						
F						

Parse table Example 1

- For production $E \rightarrow T E'$
 - Add $E \rightarrow T E'$ to the E row for each symbol in $\text{FIRST}(E)$

	+	*	()	id	\$
E			$E \rightarrow T E'$		$E \rightarrow T E'$	
E'						
T						
T'						
F						



Parse table Example 1

- For production $E' \rightarrow + T E' \mid \varepsilon$
 - $E' \rightarrow + T E'$, Add $\rightarrow + T E'$ to the E' row for each symbol in $\text{FIRST}(E')$
 - $E' \rightarrow \varepsilon$, add $\rightarrow \varepsilon$ for each symbol in $\text{FOLLOW}(E')$

	+	*	()	id	\$
E			$E \rightarrow T E'$		$E \rightarrow T E'$	
E'	$E' \rightarrow + T E'$			$E' \rightarrow \varepsilon$		$E' \rightarrow \varepsilon$
T						
T'						
F						

Parse table Example 1

- For production $T \rightarrow F T'$
 - Add $T \rightarrow F T'$ to the T row for each symbol in $\text{FIRST}(T)$

	+	*	()	id	\$
E			$E \rightarrow T E'$		$E \rightarrow T E'$	
E'	$E' \rightarrow + T E'$			$E' \rightarrow \epsilon$		$E' \rightarrow \epsilon$
T			$T \rightarrow F T'$		$T \rightarrow F T'$	
T'						
F						

Parse table Example 1

- For production $T' \rightarrow * F T' \mid \varepsilon$
 - Add $T' \rightarrow * F T'$ to the T' row for each symbol in $\text{FIRST}(T')$
 - Add $T' \rightarrow \varepsilon$ for each symbol in $\text{FOLLOW}(T')$

	+	*	()	id	\$
E			$E \rightarrow T E'$		$E \rightarrow T E'$	
E'	$E' \rightarrow + T E'$			$E' \rightarrow \varepsilon$		$E' \rightarrow \varepsilon$
T			$T \rightarrow F T'$		$T \rightarrow F T'$	
T'	$T' \rightarrow \varepsilon$	$T' \rightarrow * F T'$		$T' \rightarrow \varepsilon$		$T' \rightarrow \varepsilon$
F						

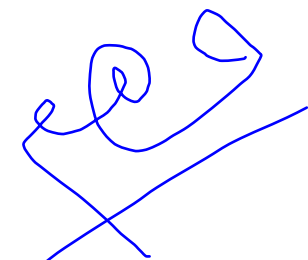
Parse table Example 1

- For production $F \rightarrow (E) \mid \text{id}$
 - Add $F \rightarrow (E)$ to the F row and symbol $($
 - Add $F \rightarrow \text{id}$ to the F row and symbol id

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

Parser Actions

- The symbol at the top of the stack (say X) and the current symbol in the input string (say a) determine the parser action.
- There are four possible parser actions.
 1. If X and a are $\$$ \rightarrow parser halts (successful completion)
 2. If X and a are the same terminal symbol (different from $\$$) \rightarrow parser pops X from the stack, and moves the next symbol in the input buffer.
 3. If X is a non-terminal \rightarrow parser looks at the parsing table entry $M[X,a]$. If $M[X,a]$ holds a production rule $X \rightarrow Y_1 Y_2 \dots Y_k$, it pops X from the stack and pushes Y_k, Y_{k-1}, \dots, Y_1 into the stack. The parser also outputs the production rule $X \rightarrow Y_1 Y_2 \dots Y_k$ to represent a step of the derivation.
 4. none of the above \rightarrow error
 - all empty entries in the parsing table are errors.
 - If X is a terminal symbol different from a , this is also an error case.



LL(1)Parser Example

$$S \rightarrow aBa$$
$$B \rightarrow bB \mid \varepsilon$$

	a	b	\$
S	$S \rightarrow aBa$		
B	$B \rightarrow \varepsilon$	$B \rightarrow bB$	

stack

input

output

\$S

abba\$

$$S \rightarrow aBa$$

\$aBa

abba\$

\$aB

bba\$
$$B \rightarrow bB$$

\$aB**b**

c bba\$ \rightarrow Syntactical error

\$aB

ba\$

$$B \rightarrow bB$$

\$aB**b**

ba\$

\$aB

a\$

$$B \rightarrow \varepsilon$$

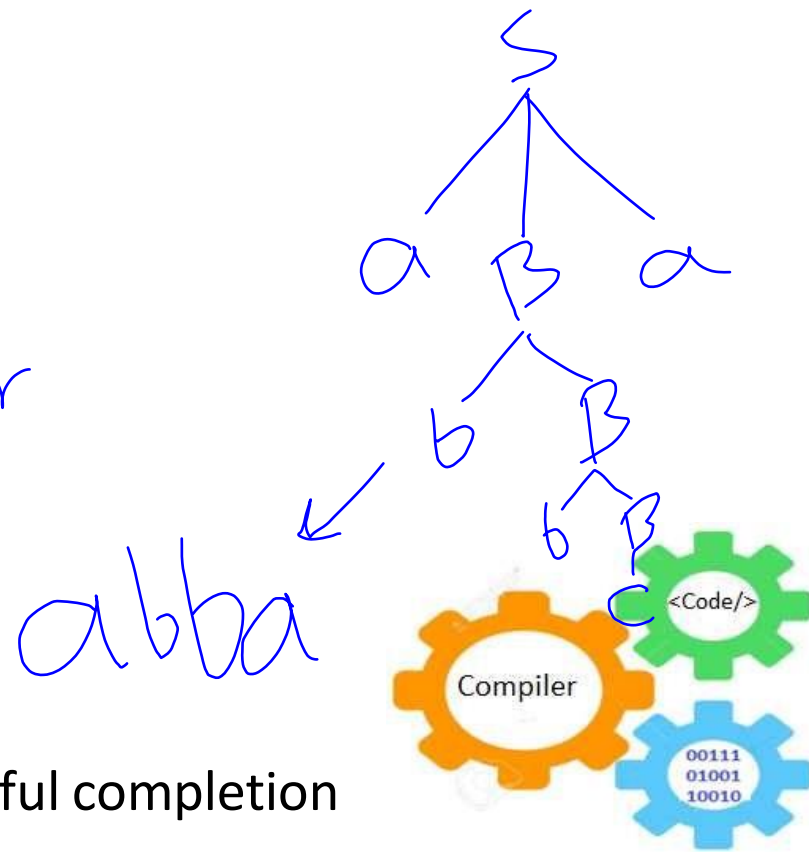
\$a

a\$

\$

\$

accept, successful completion



LL(1) Parser – Example2

$E \rightarrow TE'$

$E' \rightarrow +TE' \mid \varepsilon$

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \varepsilon$	$E' \rightarrow \varepsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \varepsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \varepsilon$	$T' \rightarrow \varepsilon$
F	$F \rightarrow id$			$F \rightarrow (E)$		

LL(1) Parser – Example2...cont'd

Input is id+id

<u>stack</u>	<u>input</u>	<u>output</u>
\$E	id+id\$	$E \rightarrow TE'$
\$E'T	id+id\$	$T \rightarrow FT'$
\$E'T'F	id+id\$	$F \rightarrow id$
\$E'T'id	id+id\$	
\$E'T'	+id\$	$T' \rightarrow \epsilon$
\$E'	+id\$	$E' \rightarrow +TE'$
\$E'T+	+id\$	
\$E'T	id\$	$T \rightarrow FT'$
\$E'T'F	id\$	$F \rightarrow id$
\$E'T'id	id\$	
\$E'T'	\$	$T' \rightarrow \epsilon$
\$E'	\$	$E' \rightarrow \epsilon$
\$	\$	accept

