



Program : **B.Tech**

Subject Name: **Compiler Design**

Subject Code: **CS-603**

Semester: **6th**



LIKE & FOLLOW US ON FACEBOOK

facebook.com/rgpvnotes.in

Department of Computer Science and Engineering
Subject Notes
CS 603(C) Compiler Design

UNIT- II:

Syntax analysis: CFGs, Top down parsing, Brute force approach, recursive descent parsing, transformation on the grammars, predictive parsing, bottom up parsing, operator precedence parsing, LR parsers (SLR, LALR, LR), Parser generation. Syntax directed definitions: Construction of Syntax trees, Bottom up evaluation of S-attributed definition, L-attribute definition, Top down translation, Bottom Up evaluation of inherited attributes Recursive Evaluation, Analysis of Syntax directed definition.

1 Introduction of Syntax Analysis

The second stage of translation is called Syntax analysis or parsing. In this phase expressions, statements, declarations etc... are identified by using the results of lexical analysis. Syntax analysis is aided by using techniques based on formal grammar of the programming language.

Where lexical analysis splits the input into tokens, the purpose of syntax analysis (also known as parsing) is to recombine these tokens. Not back into a list of characters, but into something that reflects the structure of the text. This “something” is typically a data structure called the syntax tree of the text. As the name indicates, this is a tree structure. The leaves of this tree are the tokens found by the lexical analysis, and if the leaves are read from left to right, the sequence is the same as in the input text. Hence, what is important in the syntax tree is how these leaves are combined to form the structure of the tree and how the interior nodes of the tree are labeled. In addition to finding the structure of the input text, the syntax analysis must also reject invalid texts by reporting syntax errors.

As syntax analysis is less local in nature than lexical analysis, more advanced methods are required. We, however, use the same basic strategy: A notation suitable for human understanding is transformed into a machine-like low-level notation suitable for efficient execution. This process is called parser generation.

A syntax analyzer or parser takes the input from a lexical analyzer in the form of token streams. The parser analyzes the source code (token stream) against the production rules to detect any errors in the code. The output of this phase is a **parse tree**.

This way, the parser accomplishes two tasks, i.e., parsing the code, looking for errors and generating a parse tree as the output of the phase.

1.1 Context Free Grammar

Definition – A context-free grammar (CFG) consisting of a finite set of grammar rules is a quadruple (N, T, P, S) , where

N is a set of non-terminal symbols.

T is a set of terminals where $N \cap T = \text{NULL}$.

P is a set of rules, $P: N \rightarrow (N \cup T)^*$, i.e., the left-hand side of the production rule **P** does not have any right context or left context.

S is the start symbol.

Example

The grammar $(\{A\}, \{a, b, c\}, P, A)$, $P: A \rightarrow aA, A \rightarrow abc$.

The grammar $(\{S, a, b\}, \{a, b\}, P, S)$, $P: S \rightarrow aSa, S \rightarrow bSb, S \rightarrow \epsilon$

Generation of Derivation Tree

A derivation tree or parse tree is an ordered rooted tree that graphically represents the semantic information a string derived from a context-free grammar.

Representation Technique

Root vertex – Must be labeled by the start symbol.

Vertex – Labeled by a non-terminal symbol.

Leaves – Labeled by a terminal symbol or ϵ .

If $S \rightarrow x_1x_2 \dots x_n$ is a production rule in a CFG, then the parse tree / derivation tree will be as follows –

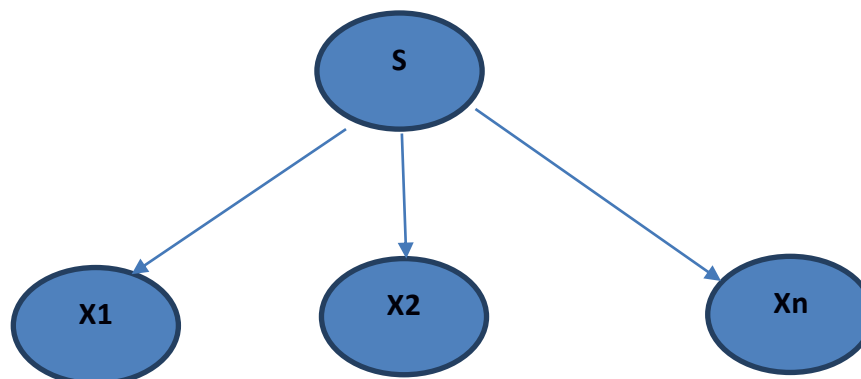


Figure 2.1: Derivation tree

There are two different approaches to draw a derivation tree –

Top-down Approach –

Starts with the starting symbol S

Goes down to tree leaves using productions

Bottom-up Approach –

Starts from tree leaves

Proceeds upward to the root which is the starting symbol S

Leftmost and Rightmost Derivation of a String

Leftmost derivation – A leftmost derivation is obtained by applying production to the leftmost variable in each step.

Rightmost derivation – A rightmost derivation is obtained by applying production to the rightmost variable in each step.

Example

Let any set of production rules in a CFG be

$X \rightarrow X+X \mid X*X \mid X \mid a$

over an alphabet $\{a\}$.

The leftmost derivation for the string " $a+a*a$ " may be –

$X \rightarrow X+X \rightarrow a+X \rightarrow a+X*X \rightarrow a+a*X \rightarrow a+a*a$

The stepwise derivation of the above string is shown as below –

The rightmost derivation for the above string " $a+a*a$ " may be –

$X \rightarrow X*X \rightarrow X*a \rightarrow X+X*a \rightarrow X+a*a \rightarrow a+a*a$

The stepwise derivation of the above string is shown as below –

2. Parsing

Syntax analyzers follow production rules defined by means of context-free grammar. The way the production rules are implemented (derivation) divides parsing into two types: top-down parsing and bottom-up parsing.

2.1 Top-Down Parsing

A program that performs syntax analysis is called a parser. A syntax analyzer takes tokens as input and output error message if the program syntax is wrong. The parser uses symbol-look-ahead and an approach called top-down parsing without backtracking. Top-down parsers check to see if a string can be generated by a grammar by creating a parse tree starting from the initial symbol and working down. Bottom-up parsers, however, check to see a string can be generated from a grammar by creating a parse tree from the

leaves, and working up. Early parser generators such as YACC creates bottom-up parsers whereas many of Java parser generators such as JavaCC create top-down parsers.

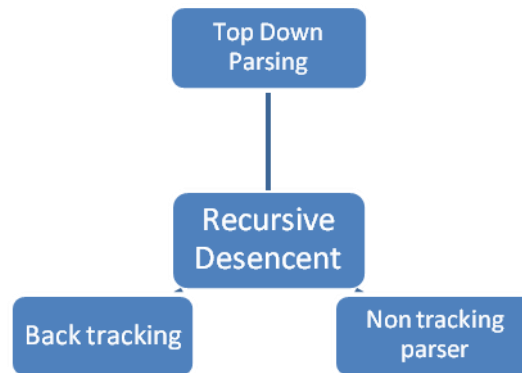


Figure 2.2: Top down parsing

2.2 Brute-Force Approach

A **top-down parse** moves from the goal symbol to a string of terminal symbols. In the terminology of trees, this is moving from the root of the tree to a set of the leaves in the syntax tree for a program. In using full backup we are willing to attempt to create a syntax tree by following branches until the correct set of terminals is reached. In the worst possible case, that of trying to parse a string which is not in the language, all possible combinations are attempted before the failure to parse is recognized.

Top-down parsing with full backup is a "brute-force" method of parsing. In general terms, this method operates as follows:

1. Given a particular non-terminal that is to be expanded, the first production for this non-terminal is applied.
2. Then, within this newly expanded string, the next (leftmost) non-terminal is selected for expansion and its first production is applied.
3. This process (step 2) of applying productions is repeated for all subsequent non-terminals that are selected until such time as the process cannot or should not be continued. This termination (if it ever occurs) may be due to two causes. First, no more non-terminals may be present, in which case the string has been successfully parsed. Second, it may result from an incorrect expansion which would be indicated by the production of a substring of terminals which does not match the appropriate segment of the source string. In the case of such an incorrect expansion, the process is "backed up" by undoing the most recently applied production. Instead of using the particular expansion that caused the error, the next production of this non-terminal is used as the next expansion, and then the process of production application continues as before.

If, on the other hand, no further productions are available to replace the production that caused the error, this error-causing expansion is replaced by the non-terminal itself, and the process is backed up again to undo the next most recently applied production. This backing up continues either until we are able to resume normal application of productions to selected non-terminals or until we have backed up to the goal symbol and there are no further productions to be tried. In the latter case, the given string must be unappeasable because it is not part of the language determined by this particular grammar.

As an example of this brute-force parsing technique, let us consider the simple grammar

S->aAd/aB

A->b/c

B->ccd/ddc

Where S is the goal or start symbol. Figure 6-1 illustrates the working of this brute-force parsing technique by showing the sequence of syntax trees generated during the parse of the string 'accd'.

2.3 Recursive Decent Parsing

Typically, top-down parsers are implemented as a set of recursive functions that descent through a parse tree for a string. This approach is known as recursive descent parsing, also known as LL(k) parsing where the first L stands for left-to-right, the second L stands for leftmost-derivation, and k indicates k-symbol look ahead. Therefore, a parser using the single symbol look-ahead method and top-down parsing without backtracking is called LL(1) parser. In the following sections, we will also use an extended BNF notation in which some regulation expression operators are to be incorporated.

A syntax expression defines sentences of the form, or. A syntax of the form defines sentences that consist of a sentence of the form followed by a sentence of the form followed by a sentence of the form. A syntax of the form defines zero or one occurrence of the form. A syntax of the form defines zero or more occurrences of the form.

A usual implementation of an LL(1) parser is.

- initialize its data structures,
- get the lookahead token by calling scanner routines, and
- call the routine that implements the start symbol.

Here is an example.

```
proc syntax Analysis()
begin
initialize(); // initialize global data and structures
nextToken(); // get the lookahead token
program(); // parser routine that implements the start symbol
end;
```

Example for backtracking:

Consider the grammar G :

$S \rightarrow cAd$

$A \rightarrow ab \mid a$

And the input string $w=cad$.

The parse tree can be constructed using the following top-down approach:

Step1:

Initially create a tree with single node labeled S. An input pointer points to 'c', the first symbol of w. Expand the tree with the production of S.

Step2:

The leftmost leaf 'c' matches the first symbol of w, so advance the input pointer to the second symbol of w 'a' and consider the next leaf 'A'. Expand A using the first alternative.

Step3:

The second symbol 'a' of w also matches with second leaf of tree. So advance the input pointer to third symbol of w 'd'. But the third leaf of tree is b which does not match with the input symbol d.

Hence discard the chosen production and reset the pointer to second position. This is called backtracking.

Step4:

Now try the second alternative for A.

Now we can halt and announce the successful completion of parsing.

Example for recursive decent parsing:

A left-recursive grammar can cause a recursive-descent parser to go into an infinite loop. Hence, elimination of left-recursion must be done before parsing.

2.4 Predictive Parsing

Predictive parsing is a special case of recursive descent parsing where no backtracking is required.

The key problem of predictive parsing is to determine the production to be applied for a non-terminal in case of alternatives.

Non-recursive predictive parser

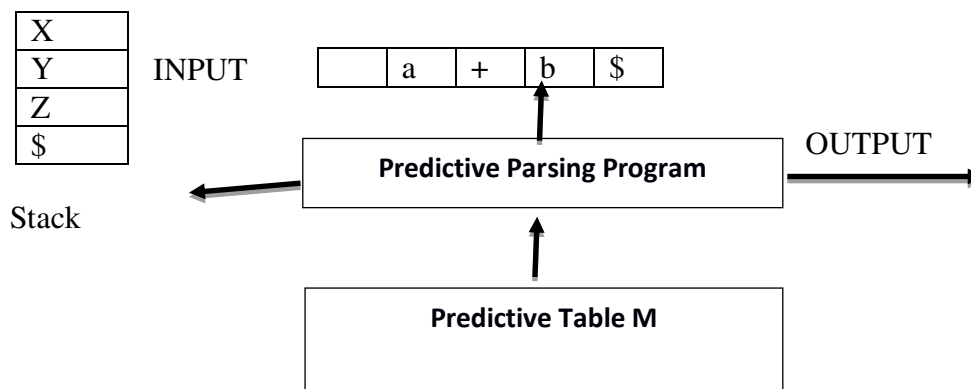


Figure 2.3: Predictive parsing

The table-driven predictive parser has an input buffer, stack, a parsing table and an outputstream.

Input buffer:

It consists of strings to be parsed, followed by \$ to indicate the end of the input string.

Stack:

It contains a sequence of grammar symbols preceded by \$ to indicate the bottom of the stack. Initially, the stack contains the start symbol on top of \$.

Parsing table:

It is a two-dimensional array $M[A, a]$, where 'A' is a non-terminal and 'a' is a terminal. Predictive parsing program. The parser is controlled by a program that considers X, the symbol on top of stack, and a, the current input symbol. These two symbols determine the parser action. There are three possibilities:

- If $X = a = \$$, the parser halts and announces successful completion of parsing.
- If $X = a \neq \$$, the parser pops X off the stack and advances the input pointer to the next input symbol.
- If X is a non-terminal, the program consults entry $M[X, a]$ of the parsing table M. This entry will either be an X-production of the grammar or an error entry.
- If $M[X, a] = \{X \rightarrow UVW\}$, the parser replaces X on top of the stack by UVW
- If $M[X, a] = \text{error}$, the parser calls an error recovery routine.

Predictive parsing table construction:

The construction of a predictive parser is aided by two functions associated with a grammar G :

1. FIRST
2. FOLLOW

Rules for first ():

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 non-terminal and $X \rightarrow a\alpha$ is a production then add a to $\text{FIRST}(X)$.

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

Rules for follow ():

1. If S is a start symbol, then $FOLLOW(S)$ contains $\$$.
2. If there is a production $A \rightarrow \alpha B \beta$, then everything in $FIRST(\beta)$ except ϵ is placed in $FOLLOW(B)$.
3. If there is a production $A \rightarrow \alpha B$, or a production $A \rightarrow \alpha B \beta$ where $FIRST(\beta)$ contains ϵ , then everything in $FOLLOW(A)$ is in $FOLLOW(B)$.

Example:

Consider the following grammar:

$E \rightarrow E+T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow (E) \mid id$

After eliminating left-recursion the grammar is

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

First () :

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

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

$FIRST(T) = \{ (, id \}$

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

$FIRST(F) = \{ (, id \}$

Follow ():

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

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

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

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

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

Non terminal	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'		$T \rightarrow \epsilon$			$T \rightarrow \epsilon$	$T \rightarrow \epsilon$
F	$F \rightarrow id$			$F \rightarrow (E)$		

Table 2.1: Predictive parsing table

3. Bottom-Up Parsing

Constructing a parse tree for an input string beginning at the leaves and going towards the root is called bottom-up parsing.

A general type of bottom-up parser is a shift-reduce parser.

3.1 Shift-Reduce Parsing

Shift-reduce parsing uses two unique steps for bottom-up parsing. These steps are known as shift-step and reduce-step.

- **Shift step:** The shift step refers to the advancement of the input pointer to the next input symbol, which is called the shifted symbol. This symbol is pushed onto the stack. The shifted symbol is treated as a single node of the parse tree.
- **Reduce step:** When the parser finds a complete grammar rule (RHS) and replaces it to (LHS), it is known as reduce-step. This occurs when the top of the stack contains a handle. To reduce, a POP function is performed on the stack which pops off the handle and replaces it with LHS non-terminal symbol.

Example:

Consider the grammar:

$S \rightarrow aABe$

$A \rightarrow Abc \mid b$

$B \rightarrow d$

The sentence to be recognized is abbcd

REDUCTION (LEFT MOST)	RIGHTMOST DERIVATION
abcde ($A \rightarrow b$)	$S \rightarrow aABe$
aAcde ($A \rightarrow Avc$)	$\rightarrow aAde$
aAde ($B \rightarrow d$)	$\rightarrow aAbcde$
aABe ($S \rightarrow aABe$)	$\rightarrow abbcd$

Table 2.2: Shift Reduce Parser

3.2 Operator Precedence Parsing

Bottom-up parsers for a large class of context-free grammars can be easily developed using operator grammars. Operator grammars have the property that no production right side is empty or has two adjacent non-terminals. This property enables the implementation of efficient operator-precedence parsers. These parser rely on the following three precedence relations:

Relation Meaning

$a < \cdot b$ a yields precedence to b

$a = \cdot b$ a has the same precedence as b

$a \cdot > b$ a takes precedence over b

These operator precedence relations allow to delimit the handles in the right sentential forms: $< \cdot$ marks the left end, $= \cdot$ appears in the interior of the handle, and $\cdot >$ marks the right end.

Example: The input string:

id1 + id2 * id3

After inserting precedence relations becomes

$\$ < \cdot id1 \cdot > + < \cdot id2 \cdot > * < \cdot id3 \cdot > \$$

Having precedence relations allows to identify handles as follows:

Scan the string from left until seeing $\cdot >$

Scan backwards the string from right to left until seeing $< \cdot$

Everything between the two relations $< \cdot$ and $\cdot >$ forms the handle

	Id	+	*	\$
Id		$\cdot >$	$\cdot >$	$\cdot >$
+	$< \cdot$	$\cdot >$	$< \cdot$	$\cdot >$
*	$< \cdot$	$\cdot >$	$\cdot >$	$\cdot >$
\$	$< \cdot$	$< \cdot$	$< \cdot$	$\cdot >$

Table 2.3: Operator precedence parsing

3.3 LR Parsing Introduction

The "L" is for left-to-right scanning of the input and the "R" is for constructing a rightmost derivation in reverse.

Why LR Parsing:

LR parsers can be constructed to recognize virtually all programming-language constructs for which context-free grammars can be written.

The LR parsing method is the most general non-backtracking shift-reduce parsing method known, yet it can be implemented as efficiently as other shift-reduce methods.

The class of grammars that can be parsed using LR methods is a proper subset of the class of grammars that can be parsed with predictive parsers.

An LR parser can detect a syntactic error as soon as it is possible to do so on a left-to-right scan of the input.

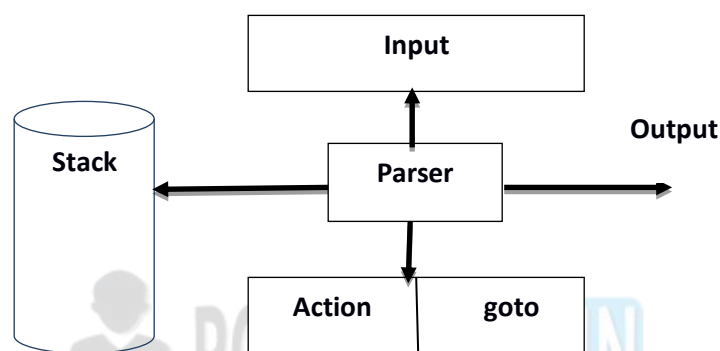


Figure 2.4: LR Parser

The disadvantage is that it takes too much work to construct an LR parser by hand for a typical programming-language grammar. But there are lots of LR parser generators available to make this task easy.

There are three widely used algorithms available for constructing an LR parser:

SLR (1) – Simple LR Parser:

- Works on smallest class of grammar
- Few number of states, hence very small table
- Simple and fast construction

LR (1) – LR Parser:

- Works on complete set of LR(1) Grammar
- Generates large table and large number of states
- Slow construction

LALR (1) – Look-Ahead LR Parser:

- Works on intermediate size of grammar
- Number of states are same as in SLR(1)

3.4 SLR (1) – Simple LR Parser:

Shift-reduce parsing attempts to construct a parse tree for an input string beginning at the leaves and working up towards the root. In other words, it is a process of “reducing” (opposite of deriving a symbol using a production rule) a string w to the start symbol of a grammar. At every (reduction) step, a particular substring matching the RHS of a production rule is replaced by the symbol on the LHS of the production.

A general form of shift-reduce parsing is LR (scanning from Left to right and using Right-most derivation in reverse) parsing, which is used in a number of automatic parser generators like Yacc, Bison, etc.

A convenient way to implement a shift-reduce parser is to use a stack to hold grammar symbols and an input buffer to hold the string w to be parsed. The symbol $\$$ is used to mark the bottom of the stack and also the right-end of the input.

Notation ally, the top of the stack is identified through a separator symbol $|$, and the input string to be parsed appears on the right side of $|$. The stack content appears on the left of $|$.

For example, an intermediate stage of parsing can be shown as follows:

$\$id1 \mid + id2 * id3\$$ (1)

Here “ $\$id1$ ” is in the stack, while the input yet to be seen is “ $+ id2 * id3\$$ ”

In shift-reduce parser, there are two fundamental operations: shift and reduce.

Shift operation: The next input symbol is shifted onto the top of the stack.

After shifting $+$ into the stack, the above state captured in (1) would change into:

$\$id1 + \mid id2 * id3\$$

Reduce operation: Replaces a set of grammar symbols on the top of the stack with the LHS of a production rule.

After reducing $id1$ using $E \rightarrow id$, the state (1) would change into:

$\$E \mid + id2 * id3\$$

In every example, we introduce a new start symbol (S'), and define a new production from this new start symbol to the original start symbol of the grammar.

Consider the following grammar (putting an explicit end-marker $\$$ at the end of the first production:

(1) $S' \rightarrow S\$$

(2) $S \rightarrow Sa$

(3) $S \rightarrow b$

For this example, the NFA for the stack can be shown as follows:

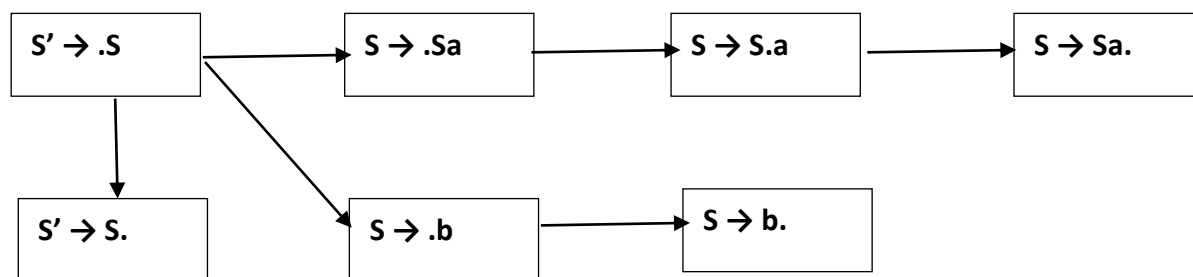


Figure 2.5: Shift Operation

After doing ϵ -closure, the resulting DFA is as follows:

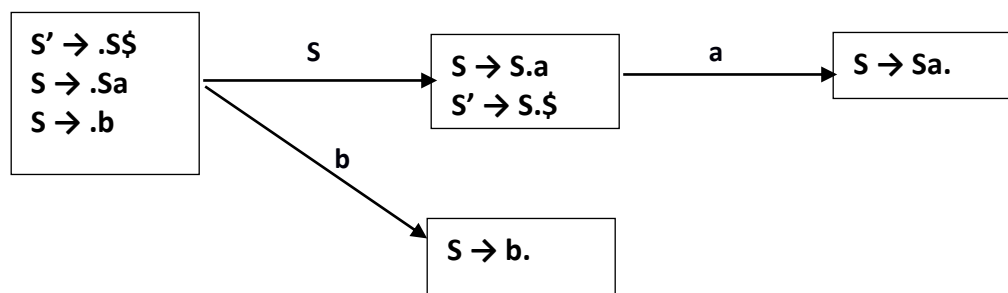


Figure 2.6: Reduce Operation

The states of DFA are also called “Canonical Collection of Items”. Using the above notation, the ACTION-GOTO table can be shown as follows:

State	A	B	\$	<u>s</u>
1		S3		
2	S4,r1	R1	R1	
3	R3	R3	R3	
4	R2	R2	R2	

Table 2.4: Simple LR Parser

3.5. Canonical LR Parsing

- Carry extra information in the state so that wrong reductions by $A \rightarrow \alpha$ will be ruled out.
- Redefine LR items to include a terminal symbol as a second component (look ahead symbol).
- The general form of the item becomes $[A \rightarrow \alpha . \square, a]$ which is called LR(1) item.
- Item $[A \rightarrow \alpha ., a]$ calls for reduction only if next input is a . The set of symbols

Canonical LR parsers solve this problem by storing extra information in the state itself. The problem we have with SLR parsers is because it does reduction even for those symbols of follow (A) for which it is invalid. So LR items are redefined to store 1 terminal (look ahead symbol) along with state and thus, the items now are LR(1) items.

An LR(1) item has the form : $[A \rightarrow a . \square, a]$ and reduction is done using this rule only if input is 'a'. Clearly the symbols a's form a subset of follow (A).

Closure (I)

repeat

for each item $[A \rightarrow \alpha . B \square, a]$ in I

for each production $B \rightarrow \gamma$ in G'

and for each terminal b in $\text{First}(\square a)$

add item $[B \rightarrow . \gamma, b]$ to I

until no more additions to I

To find closure for Canonical LR parsers:

Repeat

for each item $[A \rightarrow \alpha . B \square, a]$ in I

for each production $B \rightarrow \gamma$ in G'

and for each terminal b in $\text{First}(\square a)$

add item $[B \rightarrow . \gamma, b]$ to I

until no more items can be added to I

Example

Consider the following grammar

$S' \rightarrow S$

$S \rightarrow CC$

$C \rightarrow cC \mid d$

Compute closure (I) where $I = \{[S' \rightarrow .S, \$]\}$

$S' \rightarrow .S,$ \$

$S \rightarrow .CC,$ \$

$C \rightarrow .cC,$ c

$C \rightarrow .cC,$ d

$C \rightarrow .d,$ c

$C \rightarrow .d,$ d

For the given grammar:

$S' \rightarrow S$

$S \rightarrow CC$

$C \rightarrow cC \mid d$

$I : \text{closure}([S' \rightarrow S, \$])$

$S' \rightarrow .S$	\$	as first(e \$) = { \$ }
$S \rightarrow .CC$	\$	as first(C\$) = first(C) = { c, d }
$C \rightarrow .cC$	c	as first(Cc) = first(C) = { c, d }
$C \rightarrow .cC$	d	as first(Cd) = first(C) = { c, d }
$C \rightarrow .d$	c	as first(e c) = { c }
$C \rightarrow .d$	d	as first(e d) = { d }

Table 2.5: Canonical LR Parsing

Example

Construct sets of LR(1) items for the grammar on previous slide

$I_0 :$

$S' \rightarrow .S,$	\$
$S \rightarrow .CC,$	\$
$C \rightarrow .cC,$	c / d
$C \rightarrow .d,$	c / d

$I_1 :$

goto(I_0, S)	
$S' \rightarrow S.,$	\$

$I_2 :$	goto(I_0, C)	
	$S \rightarrow C.C,$	\$
	$C \rightarrow .cC,$	\$
	$C \rightarrow .d,$	\$
$I_3 :$	goto(I_0, c)	
	$C \rightarrow c.C,$	c/d
	$C \rightarrow .cC,$	c/d
	$C \rightarrow .d,$	c/d
$I_4 :$	goto(I_0, d)	
	$C \rightarrow d.,$	c/d
$I_5 :$	goto(I_2, C)	
	$S \rightarrow CC.,$	\$
$I_6 :$	goto(I_2, c)	
	$C \rightarrow c.C,$	\$
	$C \rightarrow .cC,$	\$
	$C \rightarrow .d,$	\$
$I_7 :$	goto(I_2, d)	
	$C \rightarrow d.,$	\$
$I_8 :$	goto(I_3, C)	

	$C \rightarrow cC.,$	c/d
$I_9 :$	goto(I 6 ,C)	
	$C \rightarrow cC.,$	\$

To construct sets of LR (1) items for the grammar given in previous slide we will begin by computing closure of $\{[S' \rightarrow .S, \$]\}$.

To compute closure we use the function given previously.

In this case $\alpha = \epsilon$, $B = S$, $\beta = \epsilon$ and $a = \$$. So add item $[S \rightarrow .CC, \$]$.

Now $\text{first}(C\$)$ contains c and d so we add following items

We have $A = S$, $\alpha = \epsilon$, $B = C$, $\beta = C$ and $a = \$$

Now $\text{first}(C\$) = \text{first}(C)$ contains c and d

So we add the items $[C \rightarrow .cC, c]$, $[C \rightarrow .cC, d]$, $[C \rightarrow .dC, c]$, $[C \rightarrow .dC, d]$.

Similarly we use this function and construct all sets of LR (1) items.

Construction of Canonical LR parse table

Construct $C = \{I_0, \dots, I_n\}$ the sets of LR(1) items.

If $[A \rightarrow \alpha .a \square, b]$ is in I_i and $\text{goto}(I_i, a) = I_j$ then $\text{action}[i, a] = \text{shift } j$

If $[A \rightarrow \alpha ., a]$ is in I_i then $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha$

If $[S' \rightarrow S., \$]$ is in I_i then $\text{action}[i, \$] = \text{accept}$

If $\text{goto}(I_i, A) = I_j$ then $\text{goto}[i, A] = j$ for all non

We are representing shift j as sj and reduction by rule number j as rj. Note that entries corresponding to [state, terminal] are related to action table and [state, non-terminal] related to goto table. We have [1, \$] as accept because $[S' \rightarrow S., \$] \in I_1$.

Parse table

state	id	+	*	()	\$		E	T	F
0	s5			s4				1	2	3
1		s6				acc				
2		r2	s7		r2	r2				
3		r4	r4		r4	r4				
4	s5			s4				8	2	3
5		r6	r6		r6	r6				
6	s5			s4					9	3
7	s5			s4						10
8		s6			s11					
9		r1	s7		r1	r1				

10		r3	r3		r3	r3				
11		r5	r5		r5	r5				

Table 2.6: Parser Table

We are representing shift j as s_j and reduction by rule number j as r_j . Note that entries corresponding to [state, terminal] are related to action table and [state, non-terminal] related to goto table. We have [1,\$] as accept because $[S' \rightarrow S., \$] \in I_1$

LALR Parse table

Look Ahead LR parsers

Consider a pair of similar looking states (same kernel and different lookaheads) in the set of LR(1) items

$I_4 : C \rightarrow d., c/d$ $I_7 : C \rightarrow d., \$$

Replace I_4 and I_7 by a new state I_{47} consisting of $(C \rightarrow d., c/d/\$)$

Similarly I_3 & I_6 and I_8 & I_9 form pairs

Merge LR(1) items having the same core

We will combine I_i and I_j to construct new I_{ij} if I_i and I_j have the same core and the difference is only in look ahead symbols. After merging the sets of LR(1) items for previous example will be as follows:

$I_0 : S' \rightarrow S \$$

$S \rightarrow .CC \$$

$C \rightarrow .cC c/d$

$C \rightarrow .d c/d$

$I_1 : \text{goto}(I_0, S)$

$S' \rightarrow S. \$$

$I_2 : \text{goto}(I_0, C)$

$S \rightarrow C.C \$$

$C \rightarrow .cC \$$

$C \rightarrow .d \$$

$I_{36} : \text{goto}(I_2, c)$

$C \rightarrow c.C c/d/\$$

$C \rightarrow .cC c/d/\$$

$C \rightarrow .d c/d/\$$

$I_4 : \text{goto}(I_0, d)$

$C \rightarrow d. c/d$

$I_5 : \text{goto}(I_2, C)$

$S \rightarrow CC. \$$

$I_7 : \text{goto}(I_2, d)$

$C \rightarrow d. \$$

$I_{89} : \text{goto}(I_{36}, C)$

$C \rightarrow cC. c/d/\$$

Construct LALR parse table

Construct $C = \{I_0, \dots, I_n\}$ set of LR(1) items

For each core present in LR(1) items find all sets having the same core and replace these sets by their union

Let $C' = \{J_0, \dots, J_m\}$ be the resulting set of items

Construct action table as was done earlier

Let $J = I_1 \cup I_2 \dots \cup I_k$

since $I_1, I_2 \dots, I_k$ have same core, $\text{goto}(J, X)$ will have the same core

Let $K = \text{goto}(I_1, X) \cup \text{goto}(I_2, X) \dots \text{goto}(I_k, X)$ the $\text{goto}(J, X) = K$

The construction rules for LALR parse table are similar to construction of LR(1) parse table.

LALR parse table

state	id	+	*	()	\$		E	T	F
0	s5			s4				1	2	3
1		s6				acc				
2		r2	s7		r2	r2				
3,6		r4	r4		r4	r4			9	3
4,7	s5,7			s4,7				8	2	3,10
5,8	s5	r6	r6	s4,8	r6	r6				
9,10		r1	s7		r1	r1				
10		r3	r3		r3	r3				
11		r5	r5		r5	r5				

Table 2.7: LALR parse table.

The construction rules for LALR parse table are similar to construction of LR(1) parse table.

Notes on LALR parse table

Modified parser behaves as original except that it will reduce $C \rightarrow d$ on inputs like ccd. The error will eventually be caught before any more symbols are shifted.

In general core is a set of LR(0) items and LR(1) grammar may produce more than one set of items with the same core.

4 Parser Generators

Some common parser generators

YACC: Yet Another Compiler Compiler

Bison: GNU Software

ANTLR: AN other Tool for Language Recognition

Yacc/Bison source program specification (accept LALR grammars)

declaration

%%

translation rules

%%

supporting C routines

Yacc and Lex schema

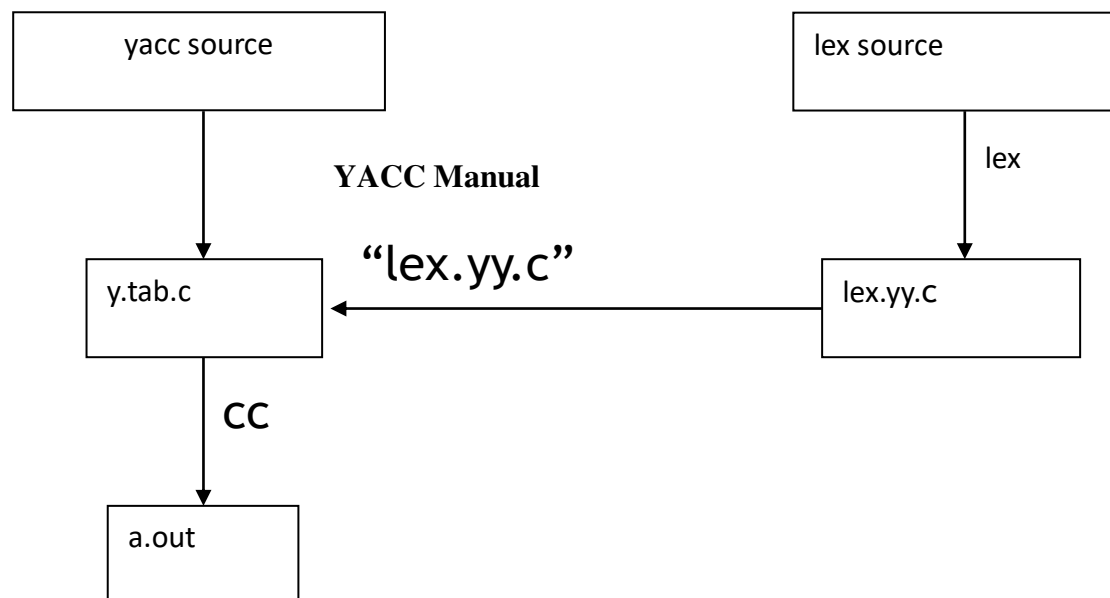


Figure 2.7: YACC

5. Syntax Directed Definition

Specifies the values of attributes by associating semantic rules with the grammar productions.

It is a context free grammar with attributes and rules together which are associated with grammar symbols and productions respectively.

The process of syntax directed translation is two-fold:

- Construction of syntax tree and
- Computing values of attributes at each node by visiting the nodes of syntax tree.

Semantic actions

Semantic actions are fragments of code which are embedded within production bodies by syntax directed translation.

They are usually enclosed within curly braces ({ }).

It can occur anywhere in a production but usually at the end of production.

$E \rightarrow E_1 + T \{ \text{print '+'} \}$

Types of translation

• L-attributed translation

It performs translation during parsing itself.

No need of explicit tree construction.

L represents 'left to right'.

• S-attributed translation

It is performed in connection with bottom up parsing.

'S' represents synthesized.

Types of attributes

• Inherited attributes

It is defined by the semantic rule associated with the production at the parent of node.

Attributes values are confined to the parent of node, its siblings and by itself.

The non-terminal concerned must be in the body of the production.

• Synthesized attributes

It is defined by the semantic rule associated with the production at the node.

Attributes values are confined to the children of node and by itself.

The non-terminal concerned must be in the head of production.

Terminals have synthesized attributes which are the lexical values (denoted by *lex val*) generated by the lexical analyzer.

Syntax directed definition of simple desk calculator

Production	Semantic rules
$L \rightarrow E_n$	$L.val = E.val$
$E \rightarrow E_1 + T$	$E.val = E_1.val + T.val$
$E \rightarrow T$	$E.val = T.val$
$T \rightarrow T_i * F$	$T.val = T_i.val \times F.val$
$T \rightarrow F$	$T.val = F.val$
$F \rightarrow (E)$	$F.val = E.val$
$F \rightarrow \text{digit}$	$F.val = \text{digit.lexval}$

Table 2.8: Syntax directed

Syntax-directed definition-inherited attributes

Production	Semantic Rules
$D \rightarrow TL$	$L.inh = T.type$
$T \rightarrow \text{int}$	$T.type = \text{integer}$
$T \rightarrow \text{float}$	$T.type = \text{float}$
$L \rightarrow L_1, \text{id}$	$L_1.inh = L.inh$ $\text{addType}(\text{id.entry}, L.inh)$
$L \rightarrow \text{id}$	$\text{addType}(\text{id.entry}, L.inh)$

Table 2.9: Syntax directed

- Symbol T is associated with a synthesized attribute *type*.
- Symbol L is associated with an inherited attribute *inh*,

Types of Syntax Directed Definitions

5.1 S-Attributed Definitions

Syntax directed definition that involves only synthesized attributes is called S-attributed. Attribute values for the non-terminal at the head is computed from the attribute values of the symbols at the body of the production.

The attributes of a S-attributed SDD can be evaluated in bottom up order of nodes of the parse tree. i.e., by performing post order traversal of the parse tree and evaluating the attributes at a node when the traversal leaves that node for the last time.

Production	Semantic rules
$L \rightarrow E_n$	$L.val = E.val$
$E \rightarrow E_1 + T$	$E.val = E_1.val + T.val$
$E \rightarrow T$	$E.val = T.val$
$T \rightarrow T_i * F$	$T.val = T_i.val \times F.val$
$T \rightarrow F$	$T.val = F.val$
$F \rightarrow (E)$	$F.val = E.val$
$F \rightarrow \text{digit}$	$F.val = \text{digit.lexval}$

Table 2.10: S-attributed Definitions

L-Attributed Definitions

The syntax directed definition in which the edges of dependency graph for the attributes in production body, can go from left to right and not from right to left is called L-attributed definitions. Attributes of L-attributed definitions may either be synthesized or inherited.

If the attributes are inherited, it must be computed from:

- Inherited attribute associated with the production head.
- Either by inherited or synthesized attribute associated with the production located to the left of the attribute which is being computed.
- Either by inherited or synthesized attribute associated with the attribute under consideration in such a way that no cycles can be formed by it in dependency graph.

Production	Semantic Rules
$T \rightarrow FT'$	$T'.inh = F.val$
$T' \rightarrow *FT_1'$	$T'_1.inh = T'.inh \times F.val$

Table 2.11: L-attributed Definitions

In production 1, the inherited attribute T' is computed from the value of F which is to its left. In production 2, the inherited attributed T'_1 is computed from $T'.inh$ associated with its head and the value of F which appears to its left in the production. i.e., for computing inherited attribute it must either use *from the above* or *from the left* information of SDD

Construction of Syntax Trees

SDDs are useful for is construction of syntax trees. A syntax tree is a condensed form of parse tree.

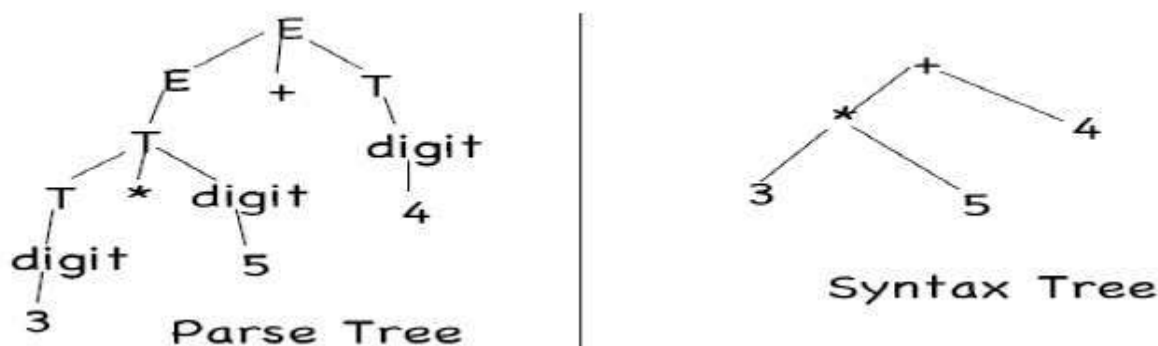


Figure 2.8 : Syntax Trees

- Syntax trees are useful for representing programming language constructs like expressions and statements.
- They help compiler design by decoupling parsing from translation.
- Each node of a syntax tree represents a construct; the children of the node represent the meaningful components of the construct.
- e.g. a syntax-tree node representing an expression $E_1 + E_2$ has label $+$ and two children representing the sub expressions E_1 and E_2
- Each node is implemented by objects with suitable number of fields; each object will have an op field that is the label of the node with additional fields as follows:

_____ If the node is a leaf, an additional field holds the lexical value for the leaf.

This is created by function Leaf (op, val)

_____ If the node is an interior node, there are as many fields as the node has children in the syntax tree. This is created by function Node (op, c1, c2,...,ck) .

Example: The S-attributed definition in figure below constructs syntax trees for a simple expression grammar involving only the binary operators $+$ and $-$. As usual, these operators are at the same precedence level and are jointly left associative. All non-terminals have one synthesized attribute node, which

represents a node of the syntax tree.

S.no	Production	Semantic Rules
1	$E \rightarrow E_1 + E$	$E.\text{node} = \text{new Node}('+', E_1.\text{node}, T.\text{node})$
2	$E \rightarrow E_1 - T$	$E.\text{node} = \text{new Node}('-', E_1.\text{node}, T.\text{node})$
3	$E \rightarrow T$	$E.\text{node} = T.\text{node}$
4	$E \rightarrow (E)$	$E.\text{node} = T.\text{node}$
5	$T \rightarrow \text{id}$	$T.\text{node} = \text{new leaf}(\text{id}, \text{id}.\text{entry})$
6	$T \rightarrow \text{num}$	$T.\text{node} = \text{new Leaf}(\text{num}, \text{num}.\text{val})$

Syntax tree for $a-4+c$ using the above SDD is shown below.

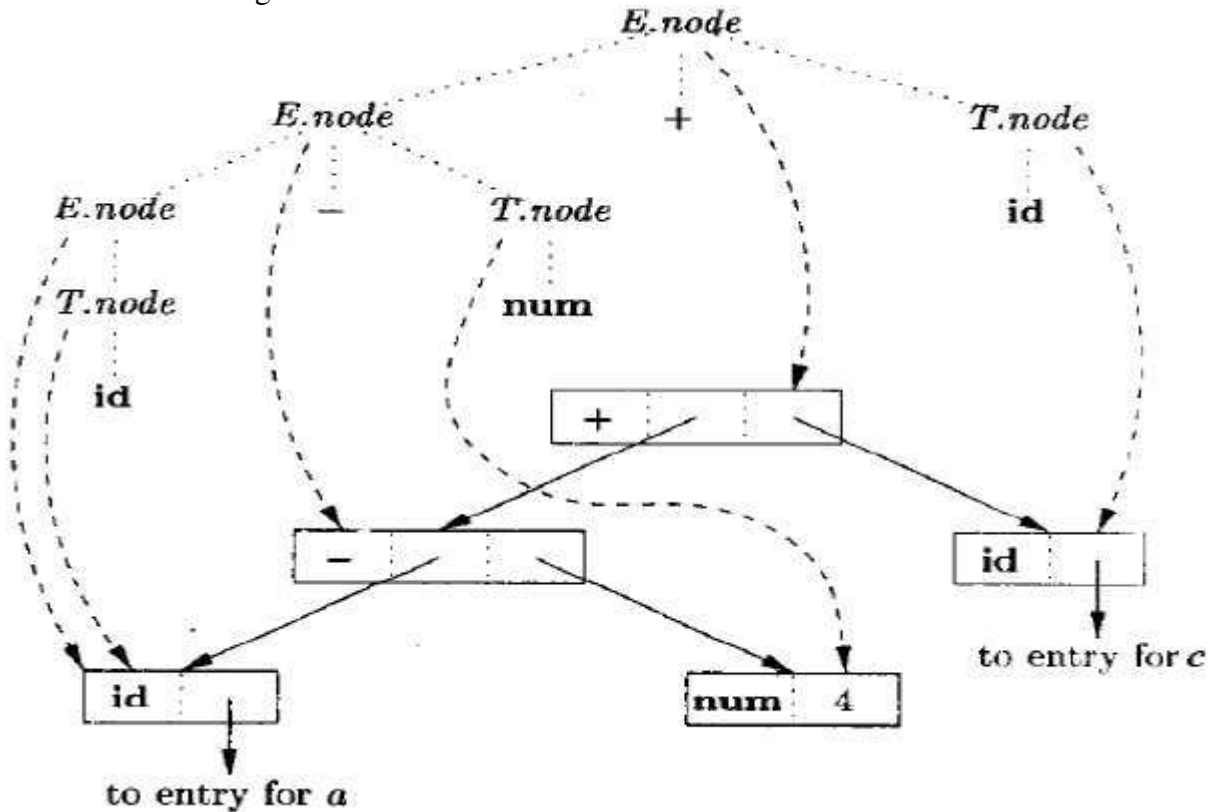


Figure 2.9:L-Attribute

5.2 Bottom-Up Evaluation Of SDT

Given an SDT, we can evaluate attributes even during bottom-up parsing. To carry out the semantic actions, parser stack is extended with semantic stack. The set of actions performed on semantic stack are mirror reflections of parser stack. Maintaining semantic stack is very easy.

During shift action, the parser pushes grammar symbols on the parser stack, whereas attributes are pushed on to semantic stack.

During reduce action, parser reduces handle, whereas in semantic stack, attributes are evaluated by the corresponding semantic action and are replaced by the result.

For example, consider the SDT

$A \rightarrow X Y Z$	$\{A \cdot a := f(X \cdot x, Y \cdot y, Z \cdot z);\}$
-----------------------	--

Strictly speaking, attributes are evaluated as follows

$A \rightarrow X Y Z$	$\{\text{val}[\text{ntop}] := f(\text{val}[\text{top} - 2], \text{val}[\text{top} - 1], \text{val}[\text{top}]);\}$
-----------------------	---

Evaluation of Synthesized Attributes

- Whenever a token is shifted onto the stack, then it is shifted along with its attribute value placed in val[top].
- Just before a reduction takes place the semantic rules are executed.
- If there is a synthesized attribute with the left-hand side non-terminal, then carrying out semantic rules will place the value of the synthesized attribute in val[ntop].

Let us understand this with an example:

$E \rightarrow E_1 \text{ “+” } T$	$\{ \text{val[ntop]} := \text{val[top-2]} + \text{val[top]}; \}$
$E \rightarrow T$	$\{ \text{val[top]} := \text{val[top]}; \}$
$T \rightarrow T_1 \text{ “*” } F$	$\{ \text{val[ntop]} := \text{val[top-2]} * \text{val[top]}; \}$
$T \rightarrow F$	$\{ \text{val[top]} := \text{val[top]}; \}$
$F \rightarrow \text{“(” } E \text{ “)”}$	$\{ \text{val[ntop]} := \text{val[top-1]}; \}$
$F \rightarrow \text{num}$	$\{ \text{val[top]} := \text{num.lvalue}; \}$

Table 2.12:

Figure shows the result of shift action. Now after performing reduce action by $E \rightarrow E * T$ resulting stack is as shown in figure.

Along with bottom-up parsing, this is how attributes can be evaluated using shift action/reduce action.

5.3 L-Attributed Definition

It allows both types, that is, synthesized as well as inherited. But if at all an inherited attribute is present, there is a restriction. The restriction is that each inherited attribute is restricted to inherit either from parent or from left sibling only.

For example, consider the rule

$A \rightarrow XYZPQ$ assume that there is an inherited attribute, “i” is present with each non-terminal. Then, $Z \cdot i = f(A \cdot i | X \cdot i | Y \cdot i)$ but $Z \cdot i = f(P \cdot i | Q \cdot i)$ is wrong as they are right siblings.

Semantic actions can be placed anywhere on the right hand side.

Attributes are generally evaluated by traversing the parse tree depth first and left to right. It is possible to rewrite any L-attributes to S-attributed definition.

L-attributed definition for converting infix to post fix form.

$$\begin{aligned}
 E &\rightarrow TE'' \\
 E'' &\rightarrow +T \#1 E'' \mid \epsilon \\
 T &\rightarrow FT'' \\
 T'' &\rightarrow *F \#2 T'' \mid \epsilon \\
 F &\rightarrow id \quad \#3
 \end{aligned}$$

Where #1 corresponds to printing “+” operator, #2 corresponds to printing “*,” and # 3 corresponds to printing id.val.

Look at the above SDT; there are no attributes, it is L-attributed definition as the semantic actions are in between grammar symbols. This is a simple example of L-attributed definition. Let us analyze this L-attributed definition and understand how to evaluate attributes with depth first left to right traversal. Take

the parse tree for the input string “a + b*c” and perform *Depth first left to right traversal*, i.e. at each node traverse the left sub tree depth wise completely then right sub tree completely.

Follow the traversal in. During the traversal whenever any dummy non-terminal is seen, carry out the translation.

Converting L-Attributed to S-Attributed Definition

Now that we understand that S-attributed is simple compared to L-attributed definition, let us see how to convert an L-attributed to an equivalent S-attributed.

Consider an L-attributed with semantic actions in between the grammar symbols. Suppose we have an L-attributed as follows:

$$S \rightarrow A \{\} B$$

How to convert it to an equivalent S-attributed definition? It is very simple!!

Replace actions by non-terminal as follows:

$$\begin{aligned} S &\rightarrow A M B \\ M &\rightarrow \varepsilon \quad \{\} \end{aligned}$$

Convert the following L-attributed definition to equivalent S-attributed definition.

	$E \rightarrow TE''$	
	$E'' \rightarrow +T$	#1 $E'' \mid \varepsilon$
	$T \rightarrow F T''$	
	$T'' \rightarrow *F$	#2 $T'' \mid \varepsilon$
	$F \rightarrow id$	#3

Table 2.13: Solution

Solution:

Replace dummy non-terminals that is, actions by non-terminals.

	$E \rightarrow TE''$	
	$E'' \rightarrow +T A E'' \mid \varepsilon$	
	$A \rightarrow$	{ print(“+”); }
	$T \rightarrow F T''$	
	$T'' \rightarrow *F B T'' \mid \varepsilon$	
	$B \rightarrow$	{ print(“*”); }
	$F \rightarrow id$	{ print(“id”); }

Table 2.14: Solution

6. YACC

YACC—Yet Another Compiler Compiler—is a tool for construction of automatic LALR parser generator.

Using Yacc

Yacc specifications are prepared in a file with extension “.y” For example, “test.y.” Then run this file with the Yacc command as “\$yacc test.y.” This translates yacc specifications into C-specifications under the default file name “y.tab.c,” where all the translations are under a function name called yyparse(); Now compile “y.tab.c” with C-compiler and test the program. The steps to be performed are given below:

Commands to execute

```
$yacc test.y
```

This gives an output “y.tab.c,” which is a parser in c under a function name yyparse().

With -v option (\$yacc -v test.y), produces file y.output, which gives complete information about the LALR parser like DFA states, conflicts, number of terminals used, etc.

```
$cc y.tab.c
```

```
$/a.out
```

Preparing the Yacc specification file

Every yacc specification file consists of three sections: the *declarations*, *grammar rules*, and *supporting subroutines*. The sections are separated by double percent “%%” marks.

```
declarations
```

```
% %
```

```
Translation rules
```

```
% %
```

```
supporting subroutines
```

The declaration section is optional. In case if there are no supporting subroutines, then the second %% can also be skipped; thus, the smallest legal Yacc specification is

```
% %
```

```
Translation rules
```

Declarations section

Declaration part contains two types of declarations—Yacc declarations or C-declarations. To distinguish between the two, C-declarations are enclosed within %{ and %}. Here we can have C-declarations like global variable declarations (int x=1;), header files (#include....), and macro definitions(#define...). This may be used for defining subroutines in the last section or action part in grammar rules.

Yacc declarations are nothing but tokens or terminals. We can define tokens by %token in the declaration part. For example, “num” is a terminal in grammar, then we define

% token num in the declaration part. In grammar rules, symbols within single quotes are also taken as terminals.

We can define the precedence and associativity of the tokens in the declarations section. This is done using %left, %right, followed by a list of tokens. The tokens defined on the same line will have the same precedence and associativity; the lines are listed in the order of increasing precedence. Thus,

```
%left '+' '-'
```

```
%left '*' '/'
```

are used to define the associativity and precedence of the four basic arithmetic operators ‘+’, ‘-’, ‘/’, ‘*’. Operators ‘*’ and ‘/’ have higher precedence than ‘+’ and both are left associative. The keyword %left is used to define left associativity and %right is used to define right associativity.

Translation rules section

This section is the heart of the yacc specification file. Here we write grammar. With each grammar rule, the user may associate actions to be performed each time the rule is recognized in the input process. These actions may return values, and may obtain the values returned by previous actions. Moreover, the lexical analyzer can return values when a token is matched. An action is defined with a set of C statements. Action can do input and output, call subprograms, and alter external vectors and variables. The action normally sets the pseudo variable “\$\$” to some value to return a value. For example, an action that does nothing but return the value 1 is { \$\$ = 1; }

To obtain the values returned by previous actions, we use the pseudo-variables \$1, \$2, . . . , which refer to the values returned by the grammar symbol on the right side of a rule, reading from left to right.

7. Analysis Syntax Directed Definition

Are a generalizations of context-free grammars in which:

1. Grammar symbols have an associated set of Attributes;
 2. Productions are associated with Semantic Rules for computing the values of attributes.
- Such formalism generates Annotated Parse-Trees where each node of the tree is a record with a field for each attribute (e.g., X.a indicates the attribute a of the grammar symbol X).
 - The value of an attribute of a grammar symbol at a given parse-tree node is defined by a semantic rule associated with the production used at that node.
- We distinguish between two kinds of attributes:
 1. **Synthesized Attributes:** They are computed from the values of the attributes of the children nodes.
 2. **Inherited Attributes:** They are computed from the values of the attributes of both the siblings and the parent nodes.

Form of Syntax Directed Definitions

Each production, $A \rightarrow \alpha$, is associated with a set of semantic rules: $b := f(c_1, c_2, \dots, c_k)$, where f is a function and either.

b is a synthesized attribute of A , and c_1, c_2, \dots, c_k are attributes of the grammar symbols of the production, or

b is an inherited attribute of a grammar symbol in α , and c_1, c_2, \dots, c_k are attributes of grammar symbols in α or attributes of A .



RGPVNOTES.IN

We hope you find these notes useful.

You can get previous year question papers at
<https://qp.rgpvnotes.in> .

If you have any queries or you want to submit your
study notes please write us at
rgpvnotes.in@gmail.com



LIKE & FOLLOW US ON FACEBOOK
facebook.com/rgpvnotes.in