

PARSING

It is the process of analyzing a continuous stream of input in order to determine its grammatical structure with respect to a given formal grammar.

Parse tree:

Graphical representation of a derivation or deduction is called a parse tree. Each interior node of the parse tree is a non-terminal; the children of the node can be terminals or non-terminals.

Types of parsing:

1. Top down parsing
2. Bottom up parsing

Ø Top-down parsing : A parser can start with the start symbol and try to transform it to the input string. Example : LL Parsers.

Ø Bottom-up parsing : A parser can start with input and attempt to rewrite it into the start symbol. Example : LR Parsers.

TOP-DOWN PARSING

It can be viewed as an attempt to find a left-most derivation for an input string or an attempt to construct a parse tree for the input starting from the root to the leaves.

Types -downoftopparsing :

1. Recursive descent parsing
2. Predictive parsing

RECURSIVE DESCENT 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 lookahead.

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.

This parsing method may involve backtracking.

Example for :backtracking

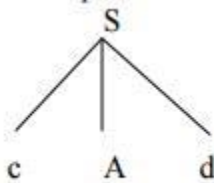
Consider the grammar $G : S \rightarrow cAd$
 $A \rightarrow ab|a$

and the input string $w=cad$.

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

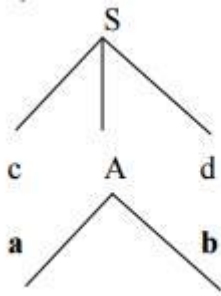
Step1:

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.

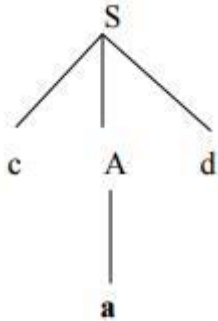


Step3:

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 **backtracking**.

Step4:

Now try the second alternative for A.



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

Predictive parsing

It is possible to build a nonrecursive predictive parser by maintaining a stack explicitly, rather than implicitly via recursive calls. The key problem during predictive parsing is that of determining the production to be applied for a nonterminal. The nonrecursive parser in figure looks up the production to be applied in parsing table. In what follows, we shall see how the table can be constructed directly from certain grammars.

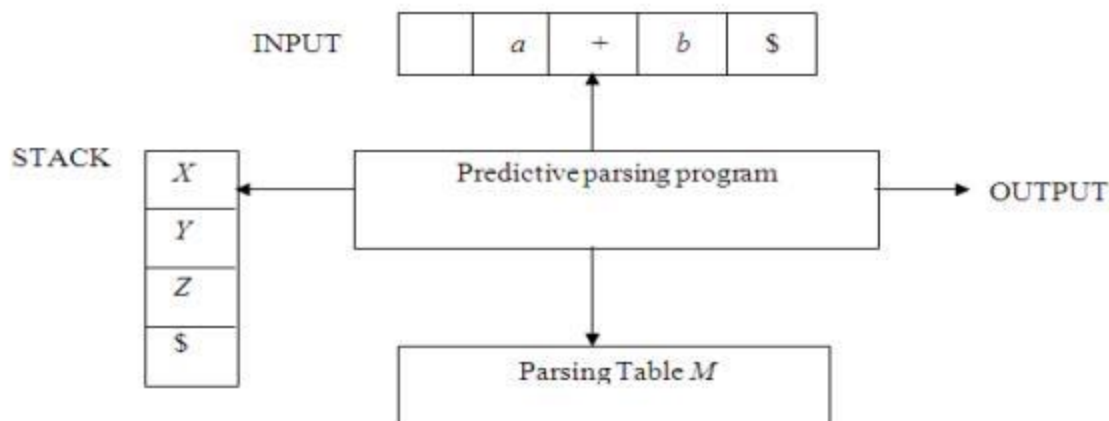


Fig. 2.4 Model of a nonrecursive predictive parser

Fig. 2.4 Model of a nonrecursive predictive parser

A table-driven predictive parser has an input buffer, a stack, a parsing table, and an output stream. The input buffer contains the string to be parsed, followed by \$, a symbol used as a right endmarker to indicate the end of the input string. The stack contains a sequence of grammar symbols with \$ on the bottom, indicating the bottom of the stack. Initially, the stack contains the start symbol of the grammar on top of \$. The parsing table is a two dimensional array $M[A,a]$ where A is a nonterminal, and a is a terminal or the symbol \$. The parser is controlled by a program that behaves as follows. The program considers X , the symbol on the top of the stack, and a , the current input symbol. These two symbols determine the action of the parser. There are three possibilities.

1 If $X = a = \$$, the parser halts and announces successful completion of parsing.

2 If $X \neq a \neq \$$, the parser pops X off the stack and advances the input pointer to the next input symbol.

3 If X is a nonterminal, the program consults entry $M[X,a]$ of the parsing table M . This entry will be either an X -production of the grammar or an error entry. If, for example, $M[X,a] = \{X \rightarrow UVW\}$, the parser replaces X on top of the stack by WVU (with U on top). As output, we shall assume that the parser just prints the production used; any other code could be executed here. If $M[X,a] = \text{error}$, the parser calls an error recovery routine

Algorithm for Nonrecursive predictive parsing.

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

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

Method. Initially, the parser is in a configuration in which it has SS on the stack with S , the start symbol of G on top, and $w\$$ in the input buffer. The program that utilizes the predictive parsing table M to produce a parse for the input is shown in Fig.

set ip to point to the first symbol of $w\$$. repeat

let X be the top stack symbol and a the symbol pointed to by ip . if X is a terminal of $\$$ then

if $X = a$ then

pop X from the stack and advance ip else error()

else

if $M[X,a]=X \rightarrow Y_1Y_2...Y_k$ then begin pop X from the stack;

push $Y_k, Y_{k-1}...Y_1$ onto the stack, with Y_1 on top; output the production $X \rightarrow Y_1Y_2...Y_k$

end

else error()

until $X=\$$ /* stack is empty */

Predictive parsing table construction:

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

3. FIRST

4. 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 $X \rightarrow Y_1 Y_2...Y_k$ 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 $\text{FIRST}(Y_1), ..., \text{FIRST}(Y_{i-1})$; that is, $Y_1, ..., Y_{i-1} \Rightarrow \epsilon$. If ϵ is in $\text{FIRST}(Y_j)$ for all $j=1,2,...,k$, then add ϵ to $\text{FIRST}(X)$.

Rules for follow():

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

Algorithm for construction of predictive parsing table:

Input : Grammar G

Output : Parsing table M

Method :

1. For each production $A \rightarrow \alpha$ of the grammar, do steps 2 and 3.
2. For each terminal a in $\text{FIRST}(\alpha)$, add $A \rightarrow \alpha$ to $M[A, a]$.
3. If ε is in $\text{FIRST}(\alpha)$, add $A \rightarrow \alpha$ to $M[A, b]$ for each terminal b in $\text{FOLLOW}(A)$. If ε is in $\text{FIRST}(\alpha)$ and $\$$ is in $\text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, \$]$.
4. Make each undefined entry of M be error.

Example:

Consider the following grammar :

$E \rightarrow E+T \mid T$

$T \rightarrow T * F \mid F$

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

After eliminating left-recursion the grammar is

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

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

First() :

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

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

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

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

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

Follow():

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

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

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

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

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

Predictive parsing Table

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

Stack Implementation

<u>stack</u>	Input	Output
SE	<u>id+id*id</u> S	
SE'T	<u>id+id*id</u> S	$E \rightarrow TE'$
SE'T'F	<u>id+id*id</u> S	$T \rightarrow FT'$
<u>SE'T'id</u>	<u>id+id*id</u> S	<u>$F \rightarrow id$</u>
SE'T'	+id*id S	
SE'	+id*id S	$T' \rightarrow \epsilon$
SE'T+	+id*id S	$E' \rightarrow +TE'$
SE'T	id*id S	
SE'T'F	id*id S	$T \rightarrow FT'$
<u>SE'T'id</u>	id*id S	<u>$F \rightarrow id$</u>
SE'T'	*id S	
SE'T'F*	*id S	$T' \rightarrow *FT'$
SE'T'F	id S	
<u>SE'T'id</u>	id S	<u>$F \rightarrow id$</u>
SE'T'	S	
SE'	S	$T' \rightarrow \epsilon$
S	S	$E' \rightarrow \epsilon$

LL(1) grammar:

The parsing table entries are single entries. So each location has not more than one entry. This type of grammar is called LL(1) grammar.

Consider this following grammar:

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

$$E \rightarrow b$$

After eliminating left factoring, we have

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow eS \mid \varepsilon$$

$$E \rightarrow b$$

To construct a parsing table, we need FIRST() and FOLLOW() for all the non-terminals.

$$\text{FIRST}(S) = \{ i, a \}$$

$$\text{FIRST}(S') = \{ e, \varepsilon \}$$

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

$$\text{FOLLOW}(S) = \{ \$, e \}$$

$$\text{FOLLOW}(S') = \{ \$, e \}$$

$$\text{FOLLOW}(E) = \{ t \}$$
Parsing table:

NON- TERMINAL	a	b	e	i	t	\$
S	<u>$S \rightarrow a$</u>			<u>$S \rightarrow iEtSS'$</u>		
S'			<u>$S' \rightarrow eS$</u> <u>$S' \rightarrow \varepsilon$</u>			<u>$S' \rightarrow \varepsilon$</u>
E		<u>$E \rightarrow b$</u>				

Since there are more than one production, the grammar is not LL(1) grammar.

Actions performed in predictive parsing:

1. Shift
2. Reduce

3. Accept
4. Error

Implementation of predictive parser:

1. Elimination of left recursion, left factoring and ambiguous grammar.
2. Construct FIRST() and FOLLOW() for all non-terminals.
3. Construct predictive parsing table.
4. Parse the given input string using stack and parsing table

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.

SHIFT-REDUCE PARSING

Shift-reduce parsing is a type of bottom-up parsing that attempts to construct a parse tree for an input string beginning at the leaves (the bottom) and working up towards the root (the top).

Example:

Consider the grammar:

$$S \rightarrow aABe$$

$$A \rightarrow Abc \mid b$$

$$B \rightarrow d$$

The sentence to be recognized is abbcd.

REDUCTION (LEFTMOST) RIGHTMOST DERIVATION

abbcd (A \rightarrow b) S \rightarrow aABe

aAbcd (A \rightarrow Abc) \rightarrow aAde

aAde (B \rightarrow d) \rightarrow aAbcd

aABe (S \rightarrow aABe) \rightarrow abbcd

S

The reductions trace out the right-most derivation in reverse.

Handles:

A handle of a string is a substring that matches the right side of a production, and whose reduction to the non-terminal on the left side of the production represents one step along the reverse of a rightmost derivation.

Example:

Consider the grammar:

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

$E \rightarrow \text{id}$

And the input string $\text{id1} + \text{id2} * \text{id3}$

The rightmost derivation is :

$E \rightarrow E + E$

$\rightarrow E + E * E$

$\rightarrow E + E * \text{id3}$

$\rightarrow E + \text{id2} * \text{id3}$

$\rightarrow \text{id1} + \text{id2} *$

In the above derivation the underlined substrings are called handles.

Handle pruning:

A rightmost derivation in reverse can be obtained by “handle pruning”. (i.e.) if w is a sentence or string of the grammar at hand, then $w = \gamma_n$, where γ_n is the n th rightsentinel form of some rightmost derivation.

Actions in shift-reduce parser:

- shift - The next input symbol is shifted onto the top of the stack.
- reduce - The parser replaces the handle within a stack with a non-terminal.
- accept - The parser announces successful completion of parsing.
- error - The parser discovers that a syntax error has occurred and calls an error recovery routine.

Conflicts in shift-reduce parsing:

There are two conflicts that occur in shift-reduce parsing:

1. Shift-reduce conflict: The parser cannot decide whether to shift or to reduce.
2. Reduce-reduce conflict: The parser cannot decide which of several reductions to make.

Stack implementation of shift-reduce parsing :

Stack	Input	Action
\$	$id_1 + id_2 * id_3$ \$	shift
\$ id_1	$+ id_2 * id_3$ \$	reduce by $E \rightarrow id$
\$ E	$+ id_2 * id_3$ \$	shift
\$ E +	$id_2 * id_3$ \$	shift
\$ E + id_2	$* id_3$ \$	reduce by $E \rightarrow id$
\$ E + E	$* id_3$ \$	shift
\$ E + E *	id_3 \$	shift
\$ E + E * id_3	\$	reduce by $E \rightarrow id$
\$ E + E * E	\$	reduce by $E \rightarrow E * E$
\$ E + E	\$	reduce by $E \rightarrow E + E$
\$ E	\$	accept