

COP 3402 Systems Software

Predictive Parsing (First and Follow Sets)

Outline

1. First set
2. Nullable symbols
3. Follow set
4. Predictive parsing table
5. LL(1) parsing

First set

A recursive descent (or predictive) parser chooses the correct production looking ahead at the input string a fix number of symbols (typically one symbol or token).

First set:

Let **X** be a string of terminals and non-terminals.

First(**X**) is the set of terminals that can begin strings or sequence derivable from **X**.

which means that we are interested in knowing if a particular nonterminal **X** derives a string with a terminal **t**.

First set

Definition: $\text{FIRST}(X) = \{ t \mid X \Rightarrow^* t \omega \text{ for some } \omega \} \cup \{ \epsilon \mid X \Rightarrow^* \epsilon \}$

If $X \rightarrow A B C$ then $\text{FIRST}(X) = \text{FIRST}(A B C)$ and is computed as follows:

If A is a terminal,

$\text{FIRST}(X) = \text{FIRST}(A B C) = \{ A \}$ (for example, if $X \rightarrow t B C$, $\text{FIRST}(X) = \text{FIRST}(t B C) = \{ t \}$)

Otherwise, if X does not derive to an empty string,

$\text{FIRST}(X) = \text{FIRST}(A B C) = \text{FIRST}(A)$.

If $\text{FIRST}(A)$ contains the empty string then,

$\text{FIRST}(X) = \text{FIRST}(A B C) = \text{FIRST}(A) - \{ \epsilon \} \cup \text{FIRST}(BC)$

Similarly, for $\text{FIRST}(BC)$ we have:

$\text{FIRST}(BC) = \{ B \}$ if B is a terminal,

Otherwise, if B does not derive to an empty string,

$\text{FIRST}(BC) = \text{FIRST}(B)$

If $\text{FIRST}(B)$ contains the empty string then,

$\text{FIRST}(BC) = \text{FIRST}(B) - \{ \epsilon \} \cup \text{FIRST}(C)$

And so on...

First set

Example:

$S \rightarrow A B C \mid C b B \mid B a$

$A \rightarrow d a \mid B C$

$B \rightarrow g \mid \epsilon$

$C \rightarrow h \mid \epsilon$

$\text{FIRST}(S) = \text{FIRST}(A B C) \cup \text{FIRST}(C b B) \cup \text{FIRST}(B a)$

$\text{FIRST}(A) = \text{FIRST}(d a) \cup \text{First}(B C) = \{ d \} \cup \text{FIRST}(B C)$

$\text{FIRST}(B) = \text{FIRST}(g) \cup \text{First} \{ \epsilon \} = \{ g, \epsilon \}$

$\text{FIRST}(C) = \text{FIRST}(h) \cup \text{First} \{ \epsilon \} = \{ h, \epsilon \}$

Now we can compute:

$\text{FIRST}(BC) = \text{FIRST}(B) - \{ \epsilon \} \cup \{ h, \epsilon \} = \{ g, \epsilon \} - \{ \epsilon \} \cup \{ h, \epsilon \} = \{ g, h, \epsilon \}$

and

$\text{FIRST}(A) = \{ d \} \cup \{ g, h, \epsilon \} = \{ d, g, h, \epsilon \}$

Exercise: Compute $\text{FIRST}(C b B)$ and $\text{FIRST}(B a)$ in order to compute $\text{FIRST}(S)$

First set

Example: Given the following expression grammar:

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

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

$$\begin{aligned} \text{Because: } E + T &\rightarrow T + T \rightarrow F + T \rightarrow \text{id} + T \\ E + T &\rightarrow T + T \rightarrow F + T \rightarrow (E) + T \end{aligned}$$

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

$$\begin{aligned} \text{Because: } E &\rightarrow T \rightarrow F \rightarrow \text{id} \\ E &\rightarrow T \rightarrow F \rightarrow (E) \end{aligned}$$

Nullable Symbols

Nullable symbols are the ones that produce the empty (ϵ) string

Example: Given the following grammar, find the nullable symbols and the First set:

$Z \rightarrow d$ $Y \rightarrow \epsilon$ $X \rightarrow Y$

$Z \rightarrow X Y Z$ $Y \rightarrow c$ $X \rightarrow a$

Note that if X can derive the empty string, nullable(X) is true.

$X \rightarrow Y \rightarrow \epsilon$

$Y \rightarrow \epsilon$

$Z \rightarrow d$

$Z \rightarrow X Y Z$

	Nullable	First
X	Yes	{ a, c, ϵ }
Y	Yes	{ c, ϵ }
Z	No	{ a, c, d }

Follow set

$$\text{FOLLOW}(\mathbf{A}) = \{ \mathbf{t} \mid \mathbf{S} \Rightarrow^* \alpha \mathbf{A} \mathbf{t} \omega \text{ for some } \alpha, \omega \}$$

Given a production \mathbf{A} , $\text{Follow}(\mathbf{A})$ is the set of terminals symbols that can immediately follow \mathbf{A} .

Example 1: If there is a derivation containing $\mathbf{A} \mathbf{t}$, then $\text{Follow}(\mathbf{A}) = \mathbf{t}$.

Example 2: If the derivation contains $\mathbf{A} \mathbf{B} \mathbf{C} \mathbf{t}$ and \mathbf{B} and \mathbf{C} are nullable, \mathbf{t} is in $\text{Follow}(\mathbf{A})$.

Example 3: Given the following grammar:

$\mathbf{Z} \rightarrow \mathbf{d} \mathbf{\$}$	$\mathbf{Y} \rightarrow \mathbf{\epsilon}$	$\mathbf{X} \rightarrow \mathbf{Y}$
$\mathbf{Z} \rightarrow \mathbf{X} \mathbf{Y} \mathbf{Z}$	$\mathbf{Y} \rightarrow \mathbf{c}$	$\mathbf{X} \rightarrow \mathbf{a}$

Compute First, Follow, and nullable.

	Nullable	First	Follow
X	Yes	{ a, c, $\mathbf{\epsilon}$ }	{ a, c, d }
Y	Yes	{ c, $\mathbf{\epsilon}$ }	{ a, c, d }
Z	No	{ a, c, d }	{ $\mathbf{\$}$ \leftarrow EOF }

Predictive parsing table

Method to construct the predictive parsing table

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

- 1.- For each terminal t in $\text{First}(A)$, add $A \rightarrow \alpha$ to $m[A, t]$, where m is the table.
- 2.- If $\text{nullable}(\alpha)$ is true, add the production $A \rightarrow \alpha$ in row A , column t , for each t in $\text{Follow}(A)$.

Example: Given the grammar:

$Z \rightarrow d$ $Y \rightarrow \epsilon$ $X \rightarrow Y$
 $Z \rightarrow XYZ$ $Y \rightarrow c$ $X \rightarrow a$

	a	c	d
X	$X \rightarrow a$ $X \rightarrow Y$	$X \rightarrow Y$	$X \rightarrow Y$
Y	$Y \rightarrow \epsilon$	$Y \rightarrow c$ $Y \rightarrow \epsilon$	$Y \rightarrow \epsilon$
Z	$Z \rightarrow XYZ$	$Z \rightarrow XYZ$	$Z \rightarrow d$ $Z \rightarrow XYZ$

$m[Y, d]$
 ← Table m

Predictive parsing table

Example: Given the grammar:

$S \rightarrow E\$$

$E \rightarrow E + T$

$E \rightarrow T$

$T \rightarrow T * F$

$T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow (E)$

We can rewrite the grammar to avoid left recursion obtaining thus:

$S \rightarrow E\$$

$E \rightarrow T E'$

$E' \rightarrow + T E'$

$E' \rightarrow \epsilon$

$T \rightarrow F T'$

$T' \rightarrow * F T'$

$T' \rightarrow \epsilon$

$F \rightarrow \text{id}$

$F \rightarrow (E)$

Compute First, Follow, and nullable.

	Nullable	First	Follow
E	No	{ id , (}	{) , \$ }
E'	Yes	{ + , ϵ }	{) , \$ }
T	No	{ id , (}	{) , + , \$ }
T'	Yes	{ * , ϵ }	{) , + , \$ }
F	No	{ id , (}	{) , * , + , \$ }

Predictive parsing table

Parsing table for the expression grammar:

	+	*	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			$F \rightarrow id$	$F \rightarrow (E)$		

Predictive parsing table

Using the predictive parsing table, it is easy to write a recursive-descent parser:

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

```
Void Tprime (void) { swith (token)
    { case PLUS:      break ;
      case TIMES:    accept (TIMES) ; F ( ) ; Tprime ( ) ; break ;
      case RPAREN : break ;
      default:      error ( ) ;
    }
}
```

Left factoring

Another problem that we must avoid in predictive parsers is when two productions for the same non-terminal start with the same symbol.

Example: $S \rightarrow \text{if } E \text{ then } S$
 $S \rightarrow \text{If } E \text{ then } S \text{ else } S$

Solution: Left-factor the grammar. Take allowable ending “**else S**” and **e**, and make a new production (new non-terminal) for them:

$$\begin{aligned} S &\rightarrow \text{if } E \text{ then } S X \\ X &\rightarrow \text{else } S \\ X &\rightarrow \varepsilon \end{aligned}$$

Grammars whose predictive parsing tables contain no multiples entries are called LL(1).

The first L stands for left-to-right parse of input string. (input string scanned from left to right)

The second L stands for leftmost derivation of the grammar

The “1” stands for one symbol lookahead

Nonrecursive predictive parsing

Example: Given the grammar:

$$\begin{array}{lll} S \rightarrow E\$ & & \\ E \rightarrow TE' & T \rightarrow FT' & F \rightarrow \text{id} \\ E' \rightarrow +TE' & T' \rightarrow *FT' & F \rightarrow (E) \\ E' \rightarrow \epsilon & T' \rightarrow \epsilon & \end{array}$$

With the following First, Follow, and nullable.

	Nullable	First	Follow
S	No	{ id }	
E	No	{ id , (}	{), \$ }
E'	Yes	{ + }	{), \$ }
T	No	{ id , (}	{), +, \$ }
T'	Yes	{ * }	{), +, \$ }
F	No	{ id , (}	{), *, +, \$ }

Nonrecursive predictive parsing

	+	*	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			$F \rightarrow id$	$F \rightarrow (E)$		

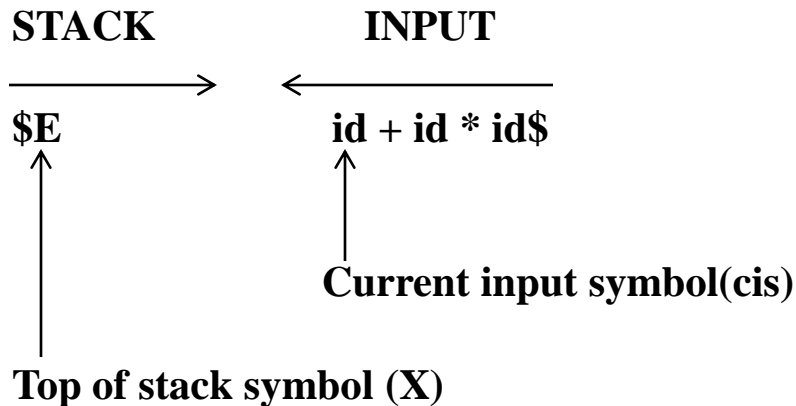
A nonrecursive predictive parser can be implemented using a stack instead of via recursive procedure calls. This approach is called table driven.

To implement it we need:

- 1) As input a string “w”.
- 2) A parsing table.
- 3) A stack.

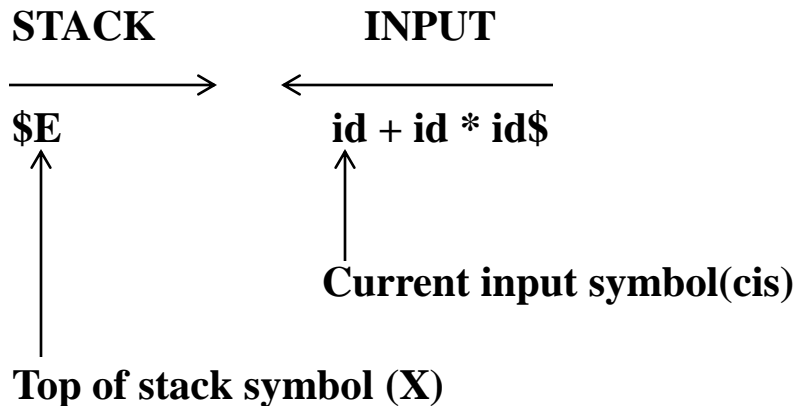
Initial configuration:

- 1) The string $w\$$ in the input buffer
- 2) The start symbol S on top of the stack, above the end of file symbol \$.



Nonrecursive predictive parsing

	+	*	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			$F \rightarrow id$	$F \rightarrow (E)$		



Algorithm:

Push \$ onto the stack

Push start symbol E onto the stack

Repeat { /*stack not empty */

 If ($X = cis$) { pop the stack;
 advance cis to next symbol;
 }

 elseif (X is a terminal) error();

 elseif ($M[X, cis]$ is an error entry) error();

 elseif ($M[X, cis] = \text{nonterminal}$) {
 pop the stack;
 push the right hand side of
 the production in reverse order ;

 }

 Let X point to the top of the stack.

}

until ($X == \$$) {accept}

Nonrecursive predictive parsing

Stack	Input	Production	Algorithm:
\$E	id + id * id\$		push \$ onto the stack
\$E'T	id + id * id\$	$E \rightarrow TE'$	push start symbol E onto the stack
\$E'T'F	id + id * id\$	$T \rightarrow FT'$	repeat ($X \neq \$$) { /*stack not empty */
\$E'T'id	id + id * id\$	$F \rightarrow id$	If ($X = cis$) { pop the stack;
\$E'T'	+ id * id\$	match id	advance cis to next symbol;
\$E'	+ id * id\$	$T' \rightarrow \epsilon$	}
\$E'T+	+ id * id\$	$E' \rightarrow +TE'$	elseif (X is a terminal) error();
\$E'T	id * id\$	match +	elseif ($M[X, cis]$ is an error entry) error();
\$E'T'F	id * id\$	$T \rightarrow FT'$	elseif ($M[X, cis] = \text{nonterminal}$) {
\$E'T'id	id * id\$	$F \rightarrow id$	pop the stack;
\$E'T'	* id\$	match id	push the right hand side of
\$E'T'F*	* id\$	$T' \rightarrow *FT'$	the production in reverse order ;
\$E'T'F	id\$	match *	}
\$E'T'id	id\$	$F \rightarrow id$	let X point to the top of the stack.
\$E'T'	\$	match id	}
\$E'	\$	$T' \rightarrow \epsilon$	until ($X == \$$) {accept}
\$	\$	$E' \rightarrow \epsilon$	else {error() }

COP 3402 Systems Software

Predictive Parsing (First and Follow Sets)

The End