Compiler course

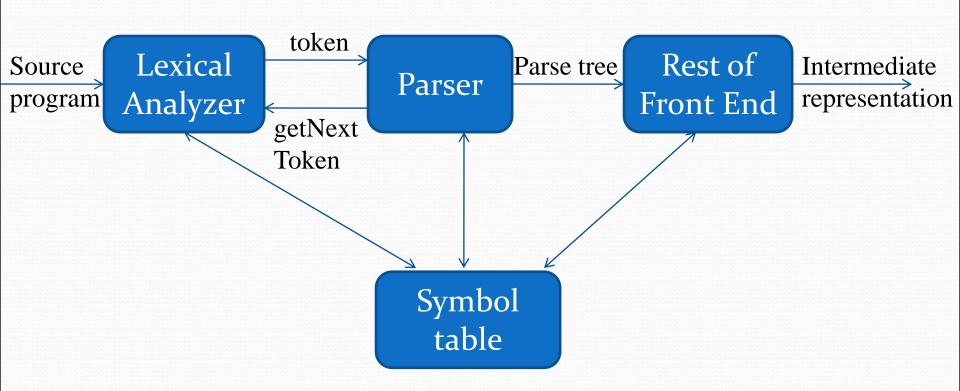
Chapter 4
Syntax Analysis

Outline

- Role of parser
- Context free grammars
- Top down parsing
- Bottom up parsing
- Parser generators

- Checking the structure of source program.
- Known as parsing
- 2nd phase of compiler
- Checks the syntax and produced IC code for the representation of source program.
- Why CFG used for syntax checking?

The role of parser



Uses of grammars

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

Error handling

- Common programming errors
 - Lexical errors
 - Syntactic errors
 - Semantic errors
- Error handler goals
 - Report the presence of errors clearly and accurately
 - Recover from each error quickly enough to detect subsequent errors
 - Add minimal overhead to the processing of correct programs

Error recovery:

Panic Mode Recovery

- In this method, successive characters from input are removed one at a time until a designated set of synchronizing tokens is found. Synchronizing tokens are deli-meters such as ; or }
- Advantage is that its easy to implement and guarantees not to go to infinite loop
- Disadvantage is that a considerable amount of input is skipped without checking it for additional errors

Statement Mode recovery

- In this method, when a parser encounters an error, it performs necessary correction on remaining input so that the rest of input statement allow the parser to parse ahead.
- The correction can be deletion of extra semicolons, replacing comma by semicolon or inserting missing semicolon.
- While performing correction, atmost care should be taken for not going in infinite loop.

Error production

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

Global correction

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

• For Given incorrect input X and grammar G, it will find parse tree for related string y such that no of insertions, deletions and changes of tokens required to transform x into y is as small as possible

Context free grammars

- Terminals
- Nonterminals
- Start symbol
- productions

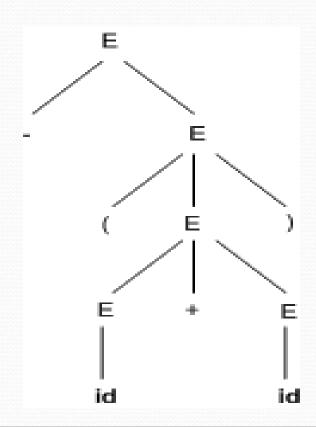
```
expression -> expression + term
expression -> expression - term
expression -> term
term -> term * factor
term -> term / factor
term -> factor
factor -> (expression)
factor -> id
```

Derivations

- Productions are treated as rewriting rules to generate a string
- Rightmost and leftmost derivations
 - $E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid id$
 - Derivations for –(id+id)
 - E => -E => -(E) => -(E+E) => -(id+E) => -(id+id)

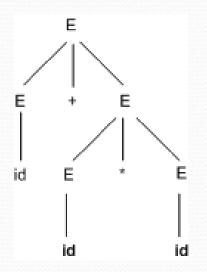
Parse trees

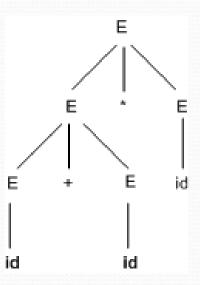
- -(id+id)
- E => -E => -(E) => -(E+E) => -(id+E) => -(id+id)



Ambiguity

- For some strings there exist more than one parse tree
- Or more than one leftmost derivation
- Or more than one rightmost derivation
- Example: id+id*id





Elimination of left recursion

- A grammar is left recursive if it has a non-terminal A such that there is a derivation $A=>^+A \alpha$
- Top down parsing methods cant handle left-recursive grammars
- A simple rule for direct left recursion elimination:
 - For a rule like:
 - $A \rightarrow A \alpha \mid \beta$
 - We may replace it with
 - $A \rightarrow \beta A'$
 - A' -> α A' | ε

Eliminate left recursion

1.
$$A \rightarrow ABd / Aa / a$$

 $B \rightarrow Be / b$

2.
$$S \rightarrow (L) / a$$

 $L \rightarrow L, S / S$

$$3. A \rightarrow Ba/Aa/c$$

 $B \rightarrow Bb/Ab/d$

Left factoring

- Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive or top-down parsing.
- Consider following grammar:
 - Stmt -> if expr then stmt else stmt
 - | **if** expr **then** stmt
- On seeing input **if** it is not clear for the parser which production to use
- We can easily perform left factoring:
 - If we have A-> α β 1 | α β 2 then we replace it with
 - A $\rightarrow \alpha$ A'
 - A' -> $\beta 1 \mid \beta 2$

Left factoring (cont.)

• Example:

```
1.S \rightarrow IEtS | iEtSeS | a
E \rightarrow b
```

- 2. $A \rightarrow aAB / aBc / aAc$
- 3. $S \rightarrow bSSaaS / bSSaSb / bSb / a$

Top Down Parsing

Introduction

- A Top-down parser tries to create a parse tree from the root towards the leafs scanning input from left to right
- It can be also viewed as finding a leftmost derivation for an input string
- Example: id+id*id

$$E \xrightarrow{lm} E \xrightarrow{lm} E \xrightarrow{lm} E \xrightarrow{lm} E \xrightarrow{lm} E$$

$$T \quad E' \quad T \quad E' \quad T \quad E' \quad T \quad E'$$

$$F \quad T' \quad F \quad T' \quad F \quad T' \quad F \quad T' \quad + \quad T \quad E'$$

$$id \qquad id \qquad \epsilon \qquad id \qquad \epsilon$$

Recursive descent parsing

- Consists of a set of procedures, one for each nonterminal
- Execution begins with the procedure for start symbol
- A typical procedure for a non-terminal

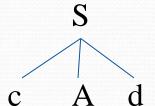
Recursive descent parsing (cont)

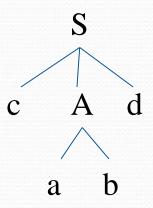
- General recursive descent may require backtracking
- The previous code needs to be modified to allow backtracking
- In general form it cant choose an A-production easily.
- So we need to try all alternatives
- If one failed the input pointer needs to be reset and another alternative should be tried
- Recursive descent parsers cant be used for left-recursive grammars

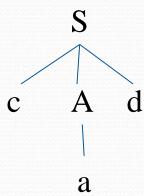
Example

S->cAd A->ab | a

Input: cad







$$2. S \rightarrow AB$$

$$A \rightarrow c|aB$$

$$B \rightarrow d$$
.

To understand first and follow

• Each time a predictive parser makes a decision, it needs to determine which production rule to apply to the leftmost non-terminal in an intermediate form, based on the next terminal

• FIRST?

- We saw the need of backtrack in the <u>previous article</u> of on Introduction to Syntax Analysis, which is really a complex process to implement. There can be easier way to sort out this problem:
- If the compiler would have come to know in advance, that what is the "first character of the string produced when a production rule is applied", and comparing it to the current character or token in the input string it sees, it can wisely take decision on which production rule to apply.

- $S \rightarrow aAb$
- A -> a | <epsilon>

• FOLLOW?

The parser faces one more problem. Let us consider below grammar to understand this problem.

- A -> aBb
- B -> c | ε

First and Follow

- First() is set of terminals that begins strings derived from
- If $\alpha \stackrel{*}{=} > \varepsilon$ then is also in First(ε)
- In predictive parsing when we have A-> $\alpha \mid \beta$, if First(α) and First(β) are disjoint sets then we can select appropriate A-production by looking at the next input
- Follow(A), for any nonterminal A, is set of terminals a that can appear immediately after A in some sentential form
 - If we have $S \stackrel{*}{=} > \alpha$ Aa β for some α and β then a is in Follow(A)
- If A can be the rightmost symbol in some sentential form, then \$ is in Follow(A)

Computing First

- To compute First(X) for all grammar symbols X, apply following rules until no more terminals or ε can be added to any First set:
 - 1. If X is a terminal then $First(X) = \{X\}$.
 - 2. If X is a nonterminal and X->Y1Y2...Yk is a production for some k>=1, then place a in First(X) if for some i, a is in First(Yi) and ε is in all of First(Y1),...,First(Yi-1) that is Y1...Yi- $\mathring{1}$ => ε . if ε is in First(Yj) for j=1,...,k then add ε to First(X).
 - 3. If $X \rightarrow \epsilon$ is a production then add ϵ to First(X)
- Example! *

Rules for First Sets

- 1. If X is a terminal then First(X) is just X!
- 2. If there is a Production $X \to \varepsilon$ then add ε to first(X)
- 3.If there is a Production $X \rightarrow Y1Y2...Yk$ then add first(Y1Y2...Yk) to first(X)
- 4.First(Y1Y2..Yk) is either
 - 1. First(Y1) (if First(Y1) doesn't contain ε)
 - 2.OR (if First(Y1) does contain ε) then First (Y1Y2..Yk) is
 - Everything in First(Y1) < except for ε > as well as everything in
 - First(Y2..Yk)
 - 3.If First(Y1) First(Y2)..First(Yk) all contain ε then add ε to
 - First(Y1Y2..Yk) as well.

Computing follow

- To compute First(A) for all nonterminals A, apply following rules until nothing can be added to any follow set:
 - 1. Place \$ in Follow(S) where S is the start symbol
 - 2. If there is a production $A -> \alpha B \beta$ then everything in First(β) except ϵ is in Follow(B).
 - 3. If there is a production A->B or a production A-> α B β where First(β) contains ϵ , then everything in Follow(A) is in Follow(B)
- Example!

Rules for Follow Sets

- 1.First put \$ (the end of input marker) in Follow(S) (S is the start symbol)
- 2. If there is a production $A \rightarrow aBb$, (where a can be a whole string) then everything in FIRST(b) except for ε is placed in FOLLOW(B).
- 3. If there is a production $A \rightarrow aB$, then everything in FOLLOW(A) is in FOLLOW(B)
- 4.If there is a production $A \rightarrow aBb$, where FIRST(b) contains ε , then everything in FOLLOW(A) is in FOLLOW(B)

Examples

```
1)
S \rightarrow aAB|bA|NULL
A-> aAb|NULL
B->bB|c
S->iCtSS'|a
S'-> eS|NULL
C->b
```

3.

 $S \rightarrow aBDh$

 $B \rightarrow cC$

 $C \rightarrow bC \mid NULL$

D -> EF

 $E \rightarrow g \mid NULL$

 $F \rightarrow f \mid NULL$

LL(1) Grammars

- Predictive parsers are those recursive descent parsers needing no backtracking
- Grammars for which we can create predictive parsers are called LL(1)
 - The first L means scanning input from left to right
 - The second L means leftmost derivation
 - And 1 stands for using one input symbol for lookahead
- A grammar G is LL(1) if and only if whenever A-> $\alpha \mid \beta$ are two distinct productions of G, the following conditions hold:
 - For no terminal a do α and β both derive strings beginning with a
 - At most one of α or β can derive empty string
 - If $\alpha => \varepsilon$ then β does not derive any string beginning with a terminal in Follow(A).

Construction of predictive parsing table

- For each production $A -> \alpha$ in grammar do the following:
 - 1. For each terminal a in First(α) add A-> in M[A,a]
 - 2. If ε is in First(α), then for each terminal b in Follow(A) add A-> ε to M[A,b]. If ε is in First(α) and \$ is in Follow(A), add A-> ε to M[A,\$] as well
- If after performing the above, there is no production in M[A,a] then set M[A,a] to error

Example

E -> TE' E' -> +TE' | ε T -> FT' T' -> *FT' | ε F -> (E) | id

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

Input Symbol

Non -	input Symbol						
terminal	id	+	*	()	\$	
Е	E -> TE'			E -> TE'			
E'		E'->+TE'			E'-> E	E'-> E	
T	T -> FT'			T -> FT'			
Т'		T'-> E	T'->*FT'		Τ'-> ε	Τ'-> ε	
F	F -> id			F -> (E)			

Another example

Non -	Input Symbol						
terminal	a	b	e	i	t	\$	
S	S -> a			S -> iEtSS'			
S'			$S' \rightarrow E$ $S' \rightarrow eS$			S'->ε	
E		E -> b					

Bottom-up Parsing

OPERATOR PRECEDENCE PARSING

- One of the bottom up parser
- It is mainly use to define the mathematical operators.
- What is operator grammar?
- A grammar which is use to define generally mathematical operations it is called as operator grammar.
- With some restrictions.?

Operator Grammar

- No E-transition.
- No two adjacent non-terminals.

Eg.

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

The above grammar is not an operator grammar but:

$$E \rightarrow E + E \mid E^* E \mid id$$

 $A \rightarrow ABA/a$

B->bAb/b

Operator Precedence

- If a has higher precedence over b; a .> b
- If a has lower precedence over b; a <. b
- If a and b have equal precedence; a =. b

Note:

- id has higher precedence than any other symbol
- \$ has lowest precedence.
- if two operators have equal precedence, then we check the **Associativity** of that particular operator.

Precedence Table

	id	+	*	\$
id		.>	.>	.>
+	v ·	.>	<.	.>
*	٧.	.>	.>	.>
\$	<.	<.	<.	

Example: w= \$id + id * id\$ \$<.id.>+<.id.>*<

BASIC PRINCIPLE

- Scan input string left to right, try to detect .> and put a pointer on its location.
- Now scan backwards till reaching <.
- String between <. And .> is our handle.
- Replace handle by the head of the respective production.
- REPEAT until reaching start symbol.

Operator-Precedence Parsing Algorithm

```
set p to point to the first symbol of w$;
repeat forever
if ($ is on top of the stack and p points to $) then return
else {
let a be the topmost terminal symbol on the stack and let b be the symbol
  pointed to by p;
if ( a <. b or a = \cdot b ) then { /* SHIFT */
push b onto the stack;
advance p to the next input symbol;
else if (a.>b) then /* REDUCE */
repeat pop stack
until (the top of stack terminal is related by <.
to the terminal most recently popped);
else error();
```

EXAMPLE

STACK	INPUT	ACTION/REMARK
\$	id + id * id\$	\$ <. Id
\$ id	+ id * id\$	id >. +
\$	+ id * id\$	\$ <. +
\$ +	id * id\$	+ <. Id
\$ + id	* id\$	id .> *
\$ +	* id\$	+ <. *
\$ + *	id\$	* <. Id
\$ + * id	\$	id .> \$
\$ + *	\$	* .> \$
\$ +	\$	+ .> \$
\$	\$	accept

PRECEDENCE FUNCTIONS

• Operator precedence parsers use **precedence functions** that map terminal symbols to integers.

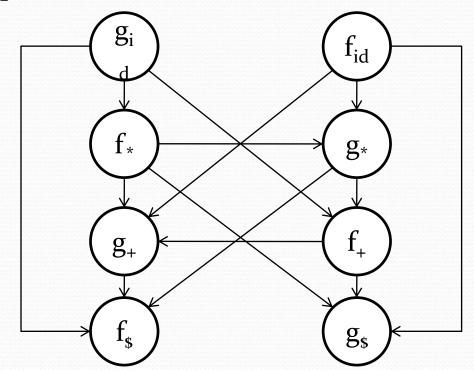
Algorithm for Constructing Precedence Functions

- 1. Create functions f_a for each grammar terminal a and for the end of string symbol.
- 2. Partition the symbols in groups so that f_a and g_b are in the same group if a = b (there can be symbols in the same group even if they are not connected by this relation).
- Create a directed graph whose nodes are in the groups, next for each symbols a and b do: place an edge from the group of g_b to the group of f_a if a < b, otherwise if a > b place an edge from the group of f_a to that of g_b .
- 4. If the constructed graph has a cycle then no precedence functions exist. When there are no cycles collect the length of the longest paths from the groups of f_a and g_b respectively.

Consider the following table:

	id	+	*	\$
id		.>	.>	.>
+	٧.	.>	<.	.>
*	۷.	.>	.>	.>
\$	<,	<.	<,	

• Resulting graph:



• From the previous graph we extract the following precedence functions:

	id	+	*	\$
f	4	2	4	О
id	5	1	3	О