

- Syntax Analysis: Introduction,
- Role Of Parsers, Context Free Grammars,
- Writing a grammar,
- Top Down Parsers,
- Bottom-Up Parsers,
- Operator-Precedence Parsing

The role of parser

Uses of grammars

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

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

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

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

Error handling

- Common programming errors
 - Lexical errors
 - Syntactic errors
 - Semantic errors
 - Logical 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

Context free grammars

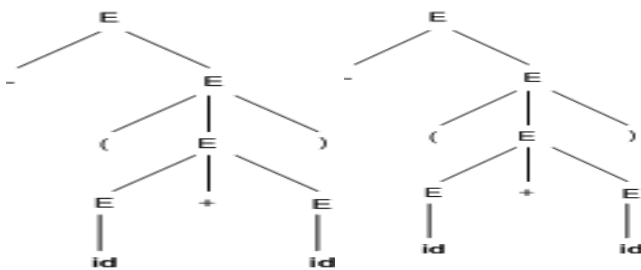
- Terminals
- Nonterminals
- Start symbol
- Productions

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 $-(id+id)$
 - $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(id+E) \Rightarrow -(id+id)$

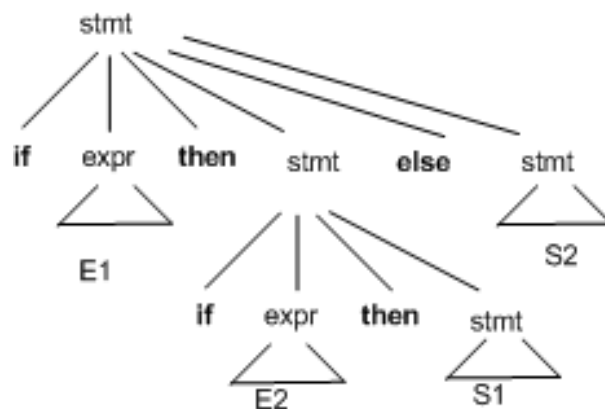
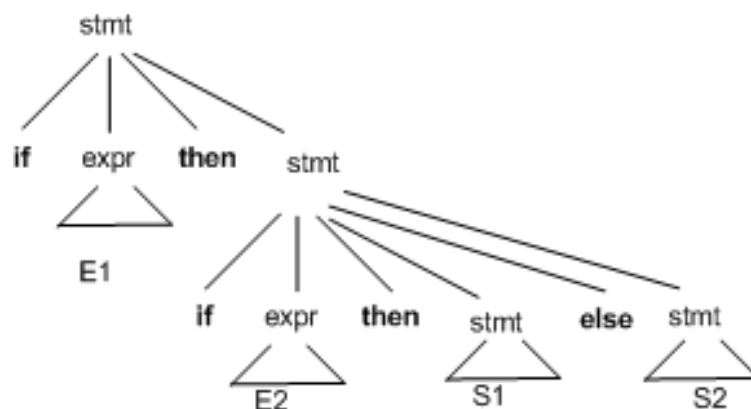
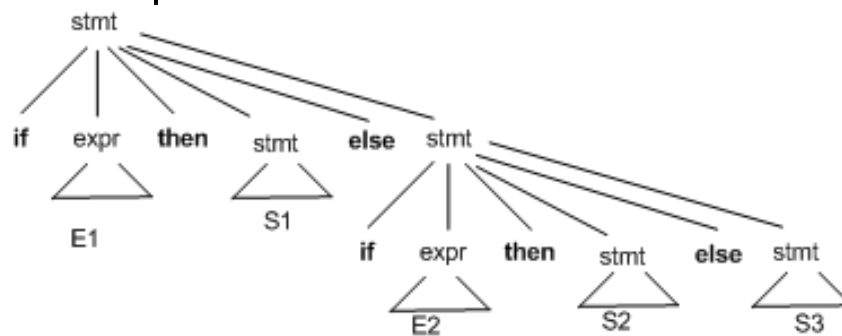
Parse trees

- $-(id+id)$
 - $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(id+E) \Rightarrow -(id+id)$



Elimination of ambiguity

$\text{stmt} \longrightarrow$ **If** expr **then** stmt
 | **If** expr **then** stmt **else** stmt
 | **other**



Elimination of left recursion

- A grammar is left recursive if it has a non-terminal A such that there is a derivation $A \Rightarrow A\alpha$
- Top down parsing methods cant handle left-recursive grammars
- A simple rule for direct left recursion elimination:
 - For a rule like:

- $A \rightarrow A\alpha | \beta$
- We may replace it with
 - $A \rightarrow \beta A'$
 - $A' \rightarrow \alpha A' | \epsilon$

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:
 - $\text{Stmt} \rightarrow \text{if expr then stmt else stmt}$
 - $\quad \quad \quad | \text{if expr then stmt}$
- On seeing input if it is not clear for the parser which production to use
- We can easily perform left factoring:
 - If we have $A \rightarrow \alpha\beta_1 | \alpha\beta_2$ then we replace it with
 - $A \rightarrow \alpha A'$
 - $A' \rightarrow \beta_1 | \beta_2$

➤ TOP DOWN PARSING

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`

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

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

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

```
void A() {  
    choose an A-production, A->X1X2..Xk  
    for (i=1 to k) {  
        if (Xi is a nonterminal  
            call procedure Xi();  
        else if (Xi equals the current input symbol a)  
            advance the input to the next symbol;  
        else /* an error has occurred */  
    }  
}
```

Example

$S \rightarrow cAd$

$A \rightarrow ab \mid a$

Input: cad

First and Follow

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

If A can be the rightmost symbol in some sentential form, then $\$$ is in $\text{Follow}(A)$

Computing First

- To compute $\text{First}(X)$ 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 $\text{First}(X) = \{X\}$.
 2. If X is a nonterminal and $X \rightarrow Y_1 Y_2 \dots Y_k$ is a production for some $k \geq 1$, then place a in $\text{First}(X)$ if for some i a is in $\text{First}(Y_i)$ and ϵ is in all of $\text{First}(Y_1), \dots, \text{First}(Y_{i-1})$ that is $Y_1 \dots Y_{i-1} \Rightarrow \epsilon$. if ϵ is in $\text{First}(Y_j)$ for $j=1, \dots, k$ then add ϵ to $\text{First}(X)$.
 3. If $X \rightarrow \epsilon$ is a production then add ϵ to $\text{First}(X)$
- Example!

Computing follow

- To compute $\text{Follow}(A)$ for all nonterminals A, apply following rules until nothing can be added to any follow set:
 1. Place $\$$ in $\text{Follow}(S)$ where S is the start symbol

2. If there is a production $A \rightarrow \alpha B \beta$ then everything in $\text{First}(\beta)$ except ϵ is in $\text{Follow}(B)$.
3. If there is a production $A \rightarrow B$ or a production $A \rightarrow \alpha B \beta$ where $\text{First}(\beta)$ contains ϵ , then everything in $\text{Follow}(A)$ is in $\text{Follow}(B)$

- Example!

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 $A \rightarrow \alpha | \beta$ 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 $\alpha \Rightarrow \epsilon$ then β does not derive any string beginning with a terminal in $\text{Follow}(A)$.

Construction of predictive parsing table

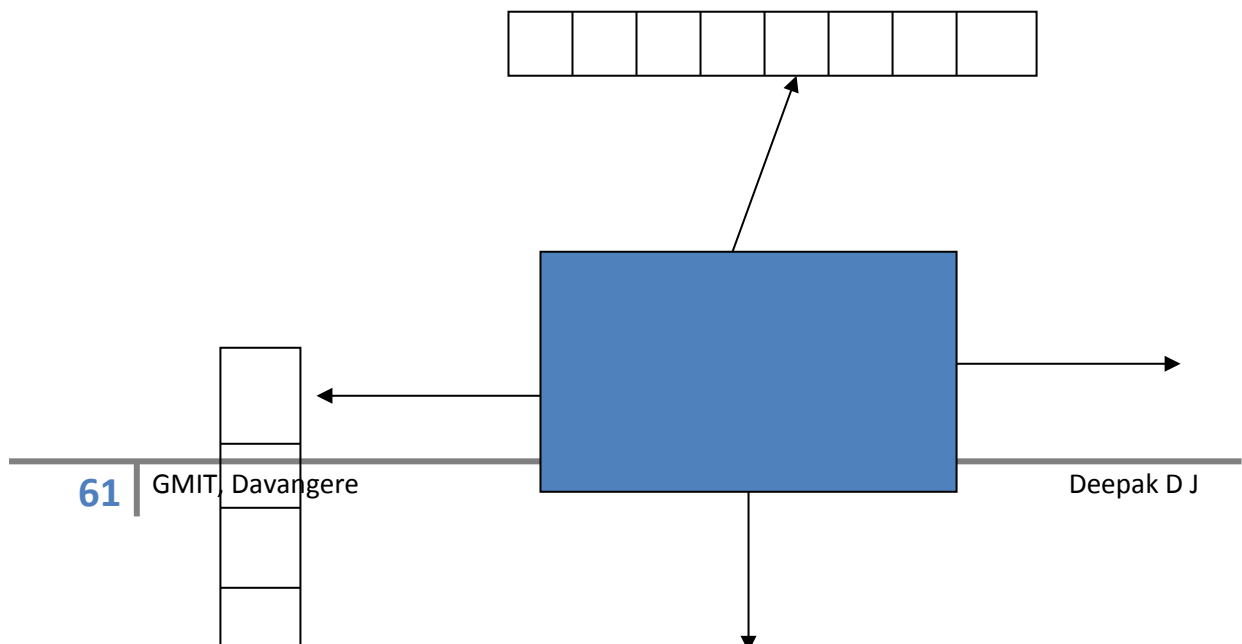
For each production $A \rightarrow \alpha$ in grammar do the following:

For each terminal a in $\text{First}(\alpha)$ add $A \rightarrow$ in $M[A, a]$

If ϵ is in $\text{First}(\alpha)$, then for each terminal b in $\text{Follow}(A)$ add $A \rightarrow \epsilon$ to $M[A, b]$. If ϵ is in $\text{First}(\alpha)$ and $\$$ is in $\text{Follow}(A)$, add $A \rightarrow \epsilon$ 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

Non-recursive predicting parsing

Predictive parsing algorithm

Set ip point to the first symbol of w;

Set X to the top stack symbol;

While (X<>\$) { /* stack is not empty */

 if (X is a) pop the stack and advance ip;

 else if (X is a terminal) error();

 else if (M[X,a] is an error entry) error();

 else if (M[X,a] = X->Y₁Y₂..Y_k) {

 output the production X->Y₁Y₂..Y_k;

 pop the stack;

 push Y_k,...,Y₂,Y₁ on to the stack with Y₁ on top;

 }

 set X to the top stack symbol;

}

Shift-reduce parser

The general idea is to shift some symbols of input to the stack until a reduction can be applied

At each reduction step, a specific substring matching the body of a production is replaced by the nonterminal at the head of the production

The key decisions during bottom-up parsing are about when to reduce and about what production to apply A reduction is a reverse of a step in a derivation

The goal of a bottom-up parser is to construct a derivation in reverse:
 $E \Rightarrow T \Rightarrow T * F \Rightarrow T * id \Rightarrow F * id \Rightarrow id * id$

Handle pruning

- A Handle is a substring that matches the body of a production and whose reduction represents one step along the reverse of a rightmost derivation

Shift reduce parsing (cont.)

Basic operations:

Shift, Reduce, Accept, Error Example: $id * id$

LR Parsing

The most prevalent type of bottom-up parsers

LR(k), mostly interested on parsers with $k \leq 1$

Why LR parsers?

Table driven

Can be constructed to recognize all programming language constructs

Most general non-backtracking shift-reduce parsing method

Can detect a syntactic error as soon as it is possible to do so

Class of grammars for which we can construct LR parsers are superset of those which we can construct LL parsers

States of an LR parser

States represent set of items

An LR(0) item of G is a production of G with the dot at some position of the body:

For $A \rightarrow XYZ$ we have following items

$A \rightarrow \cdot XYZ$

$A \rightarrow X \cdot YZ$

$A \rightarrow XY \cdot Z$

$A \rightarrow XYZ \cdot$

In a state having $A \rightarrow \cdot XYZ$ we hope to see a string derivable from XYZ next on the input.

What about $A \rightarrow X \cdot YZ$?

Constructing canonical LR(0) item sets

Augmented grammar:

G with addition of a production: $S' \rightarrow S$

Closure of item sets:

If I is a set of items, $\text{closure}(I)$ is a set of items constructed from I by the following rules:

Add every item in I to $\text{closure}(I)$

If $A \rightarrow \alpha \cdot B\beta$ is in $\text{closure}(I)$ and $B \rightarrow \gamma$ is a production then add the item $B \rightarrow \cdot \gamma$ to $\text{closure}(I)$.

Example: $E' \rightarrow E$

$E \rightarrow E + T \mid T$

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

Closure algorithm

SetOfItems CLOSURE(I) {

J=I;

repeat

 for (each item $A \rightarrow \alpha.B\beta$ in J) for (each production $B \rightarrow \gamma$ of G) if ($B \rightarrow \gamma$ is not in J) add $B \rightarrow \gamma$ to J;

until no more items are added to J on one round;

return J;

GOTO Algorithm

SetOfItems GOTO(I,X) {

J=empty;

 if ($A \rightarrow \alpha.X\beta$ is in I) add CLOSURE($A \rightarrow \alpha.X.\beta$) to J;

```
        return J;
    }
```

Canonical LR(0) items

```
Void items(G') {
    C= CLOSURE({[S'-.S]});
    repeat
        for (each set of items I in C)
            for (each grammar symbol X)
                if (GOTO(I,X) is not empty and not in C)
                    add GOTO(I,X) to C;
    until no new set of items are added to C on a round;
}
```

Line	Stack	Symbols	Input	Action
(1)	0	\$	id*id\$	Shift to 5
(2)	05	\$id	*id\$	Reduce by F- \rightarrow
(3)	03	\$F	*id\$	Reduce by T- \rightarrow
(4)	02	\$T	*id\$	Shift to 7
(5)	027	\$T*	id\$	Shift to 5
(6)	0275	\$T*id	\$	Reduce by F- \rightarrow
(7)	02710	\$T*F	\$	Reduce by T- \rightarrow
(8)	02	\$T	\$	Reduce by E- \rightarrow
(9)	01	\$E	\$	accept

a1	...	ai	...	an	\$
----	-----	----	-----	----	----

LR parsing algorithm

let a be the first symbol of w\$;

while(1) { /*repeat forever */

Sm

Sm-1

...

\$



```

let s be the state on top of the stack;
if (ACTION[s,a] = shift t) {
    push t onto the stack;
    let a be the next input symbol;
} else if (ACTION[s,a] = reduce A-> $\beta$ ) {
    pop  $|\beta|$  symbols of the stack;
    let state t now be on top of the stack;
    push GOTO[t,A] onto the stack;
    output the production A-> $\beta$ ;
} else if (ACTION[s,a]=accept) break; /* parsing is done */
else call error-recovery routine;
}

```

STATE	ACTON						GOTO		
	id	+	*	()	\$	E	T	F
0	S ₅			S ₄			1	2	3
1		S ₆				Acc			
2		R ₂	S ₇		R ₂	R ₂			
3		R ₄	R ₇		R ₄	R ₄			
4	S ₅			S ₄			8	2	3
5		R ₆	R ₆		R ₆	R ₆			
6	S ₅			S ₄				9	3
7	S ₅			S ₄					10
8		S ₆			S ₁₁				
9	Method	R ₁	S ₇		R ₁	R ₁			
10		R ₃	R ₃		R ₃	R ₃			
11		R ₅	R ₅		R ₅	R ₅			

Constructing SLR parsing table

Line	Stack	Symbols
(1)	0	
(2)	05	id
(3)	03	F
(4)	02	T
(5)	027	T*
(6)	0275	T*id
(7)	02710	T*F
(8)	02	T
(9)	01	E
(10)	016	E+
(11)	0165	E+id
(12)	0163	E+F

Construct $C=\{I_0, I_1, \dots, I_n\}$, the collection of LR(0) items for G'

State i is constructed from state li :

If $[A \rightarrow \alpha.a\beta]$ is in li and $Goto(li, a) = lj$, then set $ACTION[i, a]$ to "shift j "

If $[A \rightarrow \alpha.]$ is in li , then set $ACTION[i, a]$ to "reduce $A \rightarrow \alpha$ " for all a in $follow(A)$

If $\{S' \rightarrow .S\}$ is in li , then set $ACTION[i, \$]$ to "Accept"

If any conflicts appears then we say that the grammar is not SLR(1).

If $GOTO(li, A) = lj$ then $GOTO[i, A] = j$

All entries not defined by above rules are made "error"

The initial state of the parser is the one constructed from the set of items containing $[S' \rightarrow .S]$

