

Grammar

- A 4-tuple $G = \langle V_N, V_T, P, S \rangle$ of a language $L(G)$
 - V_N is a set of nonterminal symbols used to write the grammar
 - V_T is the set of terminals (set of words in the language $L(G)$)
 - P is a set of production rules
 - S is a special symbol in V_N , called the start symbol of the grammar
- Strings in language $L(G)$ are those derived from S by applying the production rules from P
- Examples:

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

$E \rightarrow TE'$
 $E' \rightarrow +TE' \mid \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' \mid \epsilon$
 $F \rightarrow (E) \mid \text{id}$



Error handling

- Common programming errors
 - Lexical errors
 - Syntactic errors
 - Semantic errors
 - Lexical 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 strategies

- Panic mode recovery
 - Discard input symbol one at a time until one of designated set of synchronization tokens is found
- Phrase level recovery
 - Replacing a prefix of remaining input by some string that allows the parser to continue
- Error productions
 - Augment the grammar with productions that generate the erroneous constructs
- Global correction
 - Choosing minimal sequence of changes to obtain a globally least-cost correction



Context free grammars

- Terminals
- Nonterminals
- Start symbol
- Productions

expression \rightarrow expression + term

expression \rightarrow expression – term

expression \rightarrow term

term \rightarrow term * factor

term \rightarrow term / factor

term \rightarrow factor

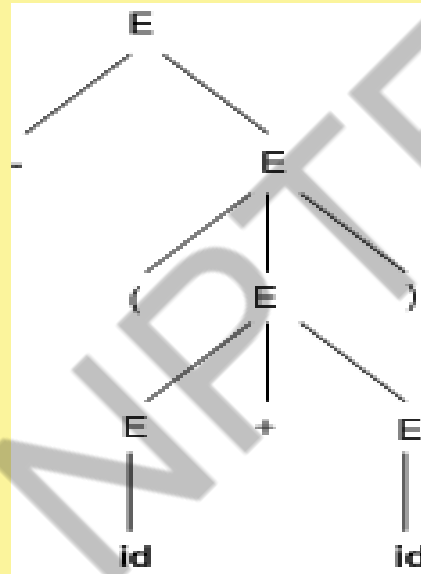
factor \rightarrow (expression)

factor \rightarrow **id**

Derivations

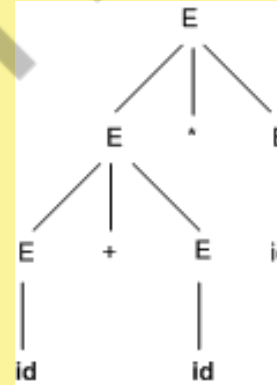
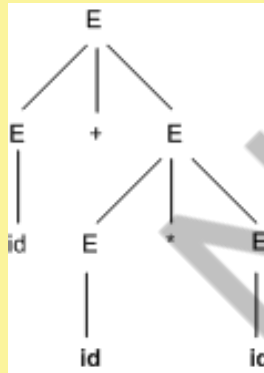
- Productions are treated as rewriting rules to generate a string
- Rightmost and leftmost derivations
 - $E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid \text{id}$
 - Derivations for $-(\text{id}+\text{id})$
 - $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(\text{id}+E) \Rightarrow -(\text{id}+\text{id})$
 - $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(E+\text{id}) \Rightarrow -(\text{id}+\text{id})$

Parse tree



Ambiguity

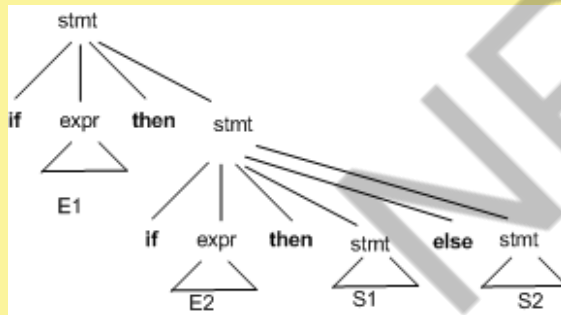
- For some strings there exist more than one parse tree
- Or more than one leftmost derivation
- Or more than one rightmost derivation
- Example: $\text{id}+\text{id}*\text{id}$



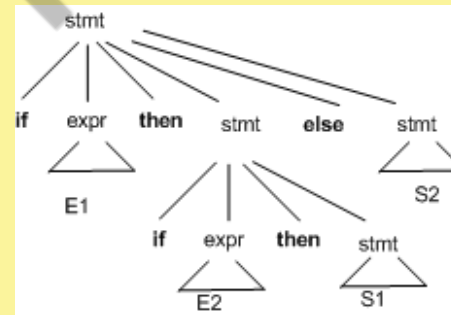
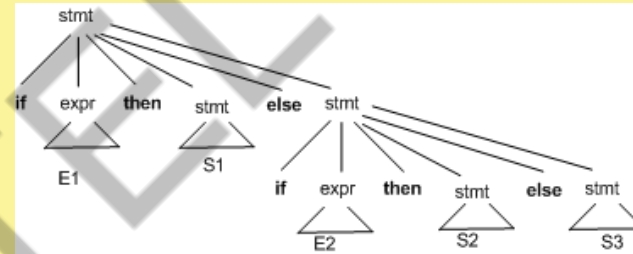
Elimination of ambiguity

if E1 then if E2 then S1 else S2

stmt \rightarrow If expr then stmt
| If expr then stmt else stmt
| other

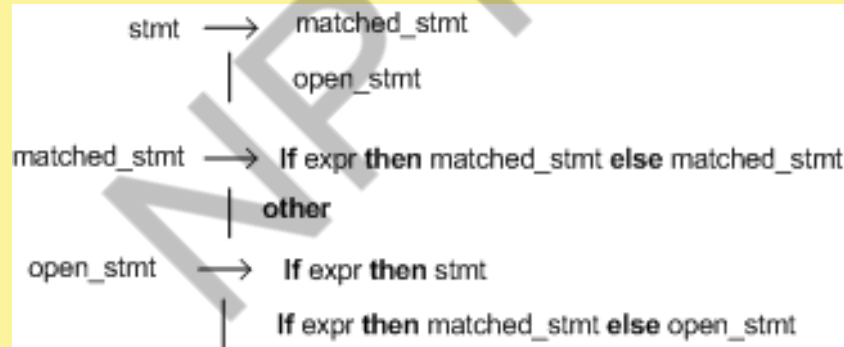


if E1 then S1 else if E2 then S2 else S3



Elimination of ambiguity (cont.)

- Idea:
 - A statement appearing between a **then** and an **else** must be matched



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



Left recursion elimination (cont.)

- There are cases like following
 - $S \rightarrow Aa \mid b$
 - $A \rightarrow Ac \mid Sd \mid \varepsilon$
- Left recursion elimination algorithm:
 - Arrange the nonterminals in some order A_1, A_2, \dots, A_n .
 - For (each i from 1 to n) {
 For (each j from 1 to $i-1$) {
 Replace each production of the form $A_i \rightarrow A_j \gamma$ by the production $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$ where $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ are all current A_j productions
 }
 Eliminate left recursion among the A_i -productions
 }



Left Recursion Elimination Example

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



$E \rightarrow TE'$
 $E' \rightarrow +TE' \mid \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' \mid \epsilon$
 $F \rightarrow (E) \mid \text{id}$

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 \rightarrow **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 $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$ then we replace it with
 - $A \rightarrow \alpha A'$
 - $A' \rightarrow \beta_1 \mid \beta_2$



Left factoring (cont.)

- Algorithm
 - For each non-terminal A , find the longest prefix α common to two or more of its alternatives. If $\alpha \neq \epsilon$, then replace all of A -productions $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$ by
 - $A \rightarrow \alpha A' \mid \gamma$
 - $A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$
- Example:
 - $S \rightarrow iEtS \mid iEtSeS \mid a$
 - $E \rightarrow b$
- Modifies to
 - $S \rightarrow iEtSS' \mid a$
 - $S' \rightarrow eS \mid \epsilon$
 - $E \rightarrow b$



TOP DOWN PARSING

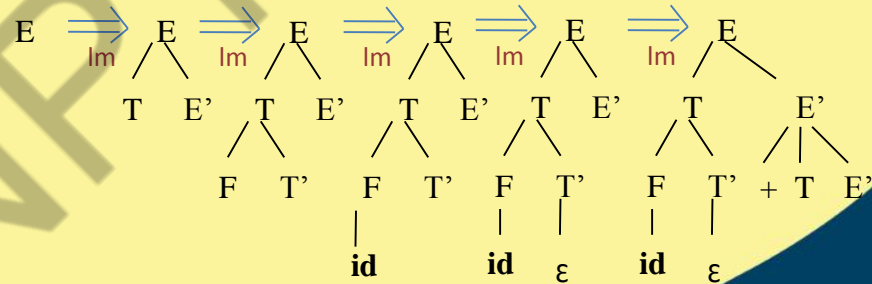


Introduction

- 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`

$$E \rightarrow TE'$$
$$E' \rightarrow +TE' \mid \varepsilon$$
$$T \rightarrow FT'$$
$$T' \rightarrow *FT' \mid \varepsilon$$

F -> (E) | 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,  $A \rightarrow X_1X_2..X_k$   
    for (i = 1 to k) {  
        if ( $X_i$  is a nonterminal  
            call procedure  $X_i()$ ;  
        else if ( $X_i$  equals the current input symbol a)  
            advance the input to the next symbol;  
        else /* an error has occurred */  
    }  
}
```



Recursive descent parsing (cont)

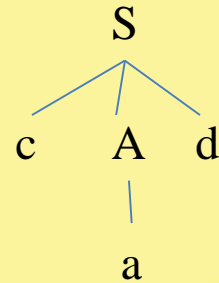
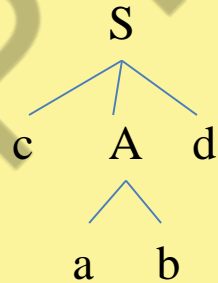
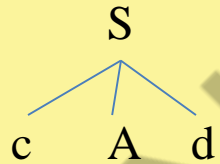
- General recursive descent may require backtracking
- The previous code needs to be modified to allow backtracking
- In general form it cannot choose an appropriate production easily.
- So we need to try all alternatives
- If one fails, the input pointer needs to be reset and another alternative has to be tried
- Recursive descent parsers cannot be used for left-recursive grammars



Example

$S \rightarrow cAd$
 $A \rightarrow ab \mid a$

Input: cad



Predictive parser

- It is a recursive-descent parser that needs no backtracking
- Suppose $A \rightarrow A1 \mid A2 \mid \dots \mid An$
- If the non-terminal to be expanded next is 'A', then the choice of rule is made on the basis of the current input symbol 'a'.



Procedure

- Make a **transition diagram** (like dfa/nfa) for every rule of the grammar.
- **Optimize** the dfa by reducing the number of states, yielding the final transition diagram
- To parse a string, **simulate** the string on the transition diagram
- If after consuming the input the transition diagram reaches an **accept state**, it is parsed.

Example

Consider the grammar:

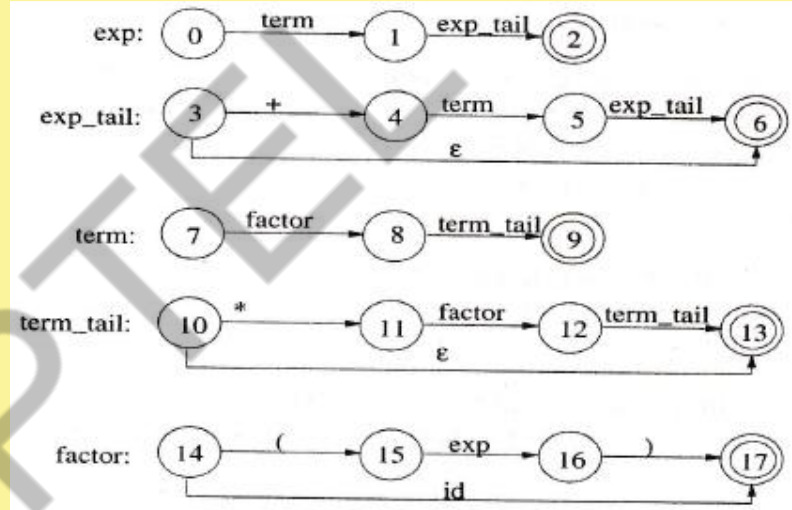
$\text{exp} \rightarrow \text{term exp_tail}$

$\text{exp_tail} \rightarrow + \text{term exp_tail} \mid \epsilon$

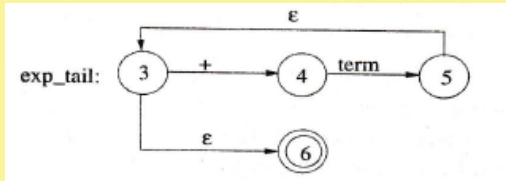
$\text{term} \rightarrow \text{factor term_tail}$

$\text{term_tail} \rightarrow * \text{factor term_tail} \mid \epsilon$

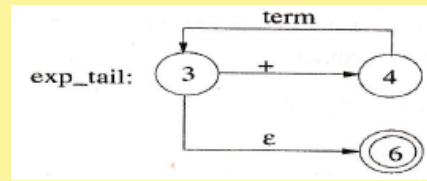
$\text{factor} \rightarrow (\text{exp}) \mid \text{id}$



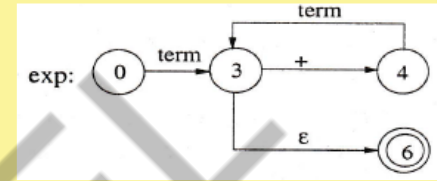
Example – Simplification



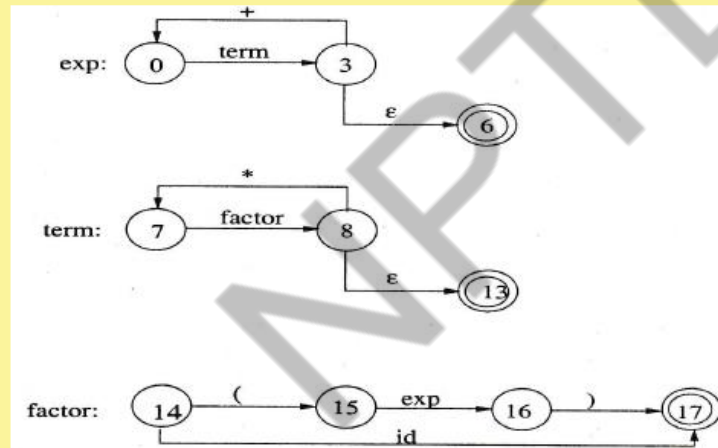
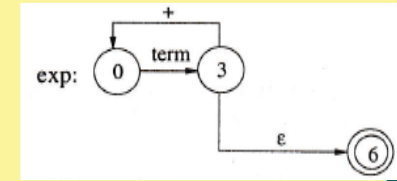
Eliminate self-recursion



Remove redundant ϵ edge



Substituting exp_tail into exp



Final set of diagrams

SIMULATION METHOD

- Start from the start state
- If a **terminal** comes **consume** it, move to next state
- If a **non – terminal** comes go to the state of the “dfa” of the non-term and return on reaching the final state
- Return to the original “dfa” and continue parsing
- If on completion(**reading input string completely**), you reach a final state, string is successfully parsed.



Disadvantage

- It is inherently a recursive parser, so it consumes a lot of memory as the stack grows.
- To remove this recursion, we use LL-parser, which uses a table for lookup.

First and Follow

- **First**(α) is set of terminals that begins strings derived from α
- If $\alpha \xRightarrow{*} \varepsilon$ then ε is also in First(α)
- In predictive parsing when we have $A \rightarrow \alpha \mid \beta$, 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 $S \xRightarrow{*} \alpha A a \beta$ 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 $\text{First}(X)$, apply following rules until no more terminals or ϵ can be added to any First set:

1. If X is a terminal then $\text{First}(X) = \{X\}$.

1. 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 $\text{First}(X)$. *

1. If $X \rightarrow \epsilon$ is a production then add ϵ to $\text{First}(X)$



Computing Follow

- To compute Follow(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 $A \rightarrow \alpha B \beta$ then everything in First(β) except ϵ is in Follow(B).
 3. If there is a production $A \rightarrow \alpha B$ or a production $A \rightarrow \alpha B \beta$ where First(β) contains ϵ , then everything in Follow(A) is in Follow(B)



Example of First and Follow Sets

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

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

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

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
 - More general one is LL(k), with k symbol lookahead
- A grammar G is LL(1) if and only if whenever $A \rightarrow \alpha \mid \beta$ 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 $\alpha \xRightarrow{*} \varepsilon$ then β does not derive any string beginning with a terminal in Follow(A)



Construction of predictive parsing table

- For each production $A \rightarrow \alpha$ in grammar do the following:
 1. For each terminal a in $\text{First}(\alpha)$ add $A \rightarrow \alpha$ in $M[A,a]$
 2. 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

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

Non - terminal	Input Symbol					
	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	$F \rightarrow id$			$F \rightarrow (E)$		

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