

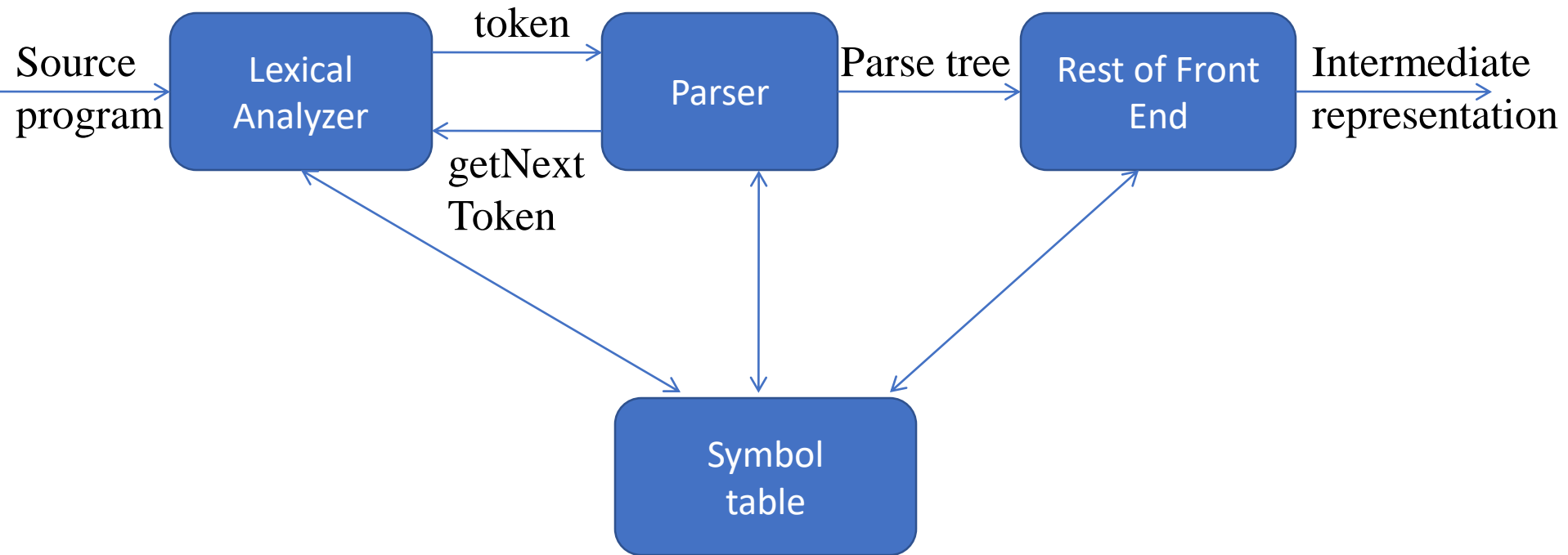
# Compiler course

Syntax Analysis

# Outline

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

# The role of parser



# SYNTAX ANALYSIS INTRODUCTION

- LEXICAL PHASE IS IMPLEMENTED ON FINITE AUTOMATA & FINITE AUTOMATA CAN REALLY ONLY EXPRESS THINGS WHERE YOU CAN COUNT MODULUS ON K.
- REGULAR LANGUAGES – THE WEAKEST FORMAL LANGUAGES WIDELY USED
  - – MANY APPLICATIONS
  - – CAN'T HANDLE ITERATION & NESTED LOOPS(NESTED IF ELSE ).
- TO SUMMARIZE, THE LEXER TAKES A STRING OF CHARACTER AS INPUT AND PRODUCES A STRING
  - OF TOKENS AS OUTPUT.
- THAT STRING OF TOKENS IS THE INPUT TO THE PARSER WHICH TAKES A STRING OF TOKENS AND PRODUCES A PARSE TREE OF THE PROGRAM.
- SOMETIMES THE PARSE TREE IS ONLY IMPLICIT. SO THE, A COMPILER MAY NEVER ACTUALLY BUILD THE FULL PARSE

# 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-recover 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**

$G=(\Sigma, T, P, S)$

$\Sigma$  – IS A FINITE SET OF TERMINALS

$T$ – IS A FINITE SET OF NON-

TERMINALS  $P$  – IS A FINITE

SUBSET OF PRODUCTION RULES

$S$  – START SYMBOL  
 IS THE START SYMBOL  
 OF THE GRAMMAR

A context-free grammar has four components:

- A set of **non-terminals** ( $V$ ). Non-terminals are syntactic variables that denote sets of strings. The non-terminals define sets of strings that help define the language generated by the grammar.
  - A set of tokens, known as **terminal symbols** ( $\Sigma$ ). Terminals are the basic symbols from which strings are formed.
  - A set of **productions** ( $P$ ). The productions of a grammar specify the manner in which the terminals and non-terminals can be combined to form strings. Each production consists of a **non-terminal** called the left side of the production, an arrow, and a sequence of tokens and/or **on- terminals**, called the right side of the production.
  - One of the non-terminals is designated as the start symbol ( $S$ ); from where the production begins.
- The strings are derived from the start symbol by repeatedly replacing a non-terminal (initially the start symbol) by the right side of a production, for that non-terminal.



# CONTEXT FREE GRAMMAR EXAMPLES

- ARITHMETIC EXPRESSIONS

$$\begin{aligned} E &::= T \mid E + T \mid E - T \\ T &::= F \mid T * F \mid T / F \mid F \\ &::= id \mid (E) \end{aligned}$$

Steps:

1. Begin with a string with only the start symbol S
2. Replace any non-terminal X in the string by the right-hand side of some production

- STATEMENTS

$$X \rightarrow Y_1 \dots Y_n$$

terminals

*If Statement* ::= if *E* then *Statement* else *Statement*

3. Repeat (2) until there are no non-

# Uses of grammars

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

# Derivations

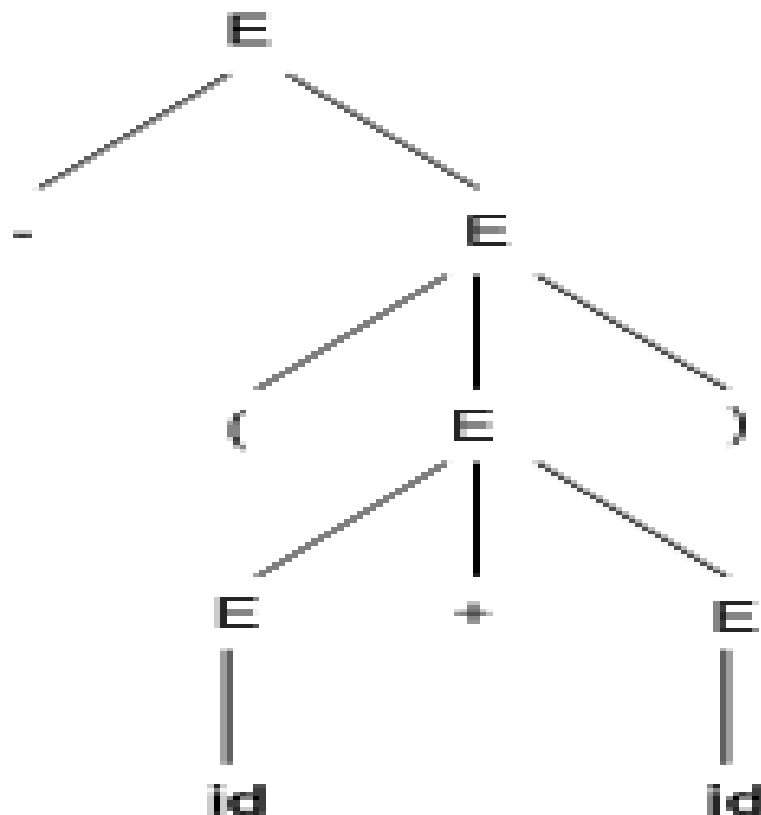
- Productions are treated as rewriting rules to generate a string
- Rightmost and leftmost derivations
  - $E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid \mathbf{id}$
  - Derivations for  $\mathbf{-(id+id)}$ 
    - $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow \mathbf{-(id+E)} \Rightarrow \mathbf{-(id+id)}$
- A derivation is basically a sequence of production rules, in order to get the input string. During parsing, we take two decisions for some sentential form of input:
- Deciding the non-terminal which is to be replaced.
- Deciding the production rule, by which, the non-terminal will be replaced.
- To decide which non-terminal to be replaced with production rule, we can have two options.

## Parse trees

A parse tree is a graphical depiction of a derivation. It is convenient to see how strings are derived from the start symbol. The start symbol of the derivation becomes the root of the parse tree.

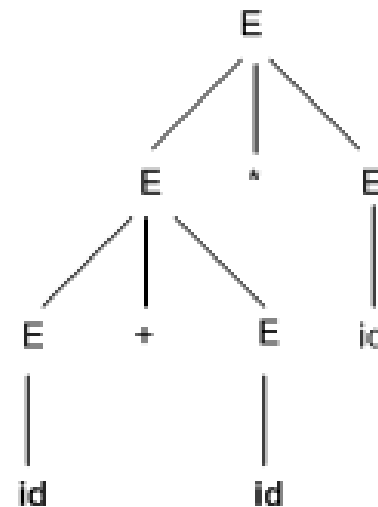
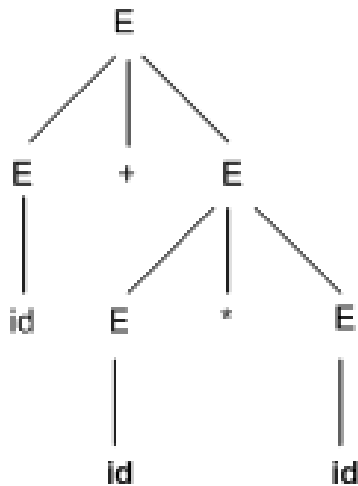
- **-(id+id)**

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



# Ambiguity

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



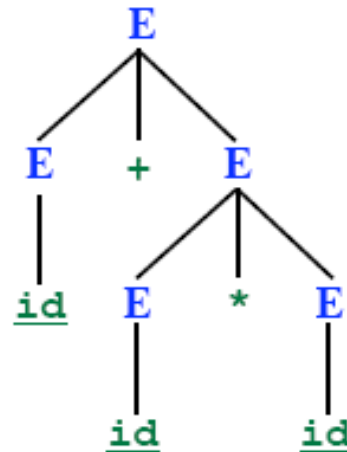
# AMBIGUOUS GRAMMAR

1.  $E \rightarrow E + E$
2.  $\rightarrow E * E$
3.  $\rightarrow ( E )$
4.  $\rightarrow - E$
5.  $\rightarrow ID$

Input:  $id+id*id$

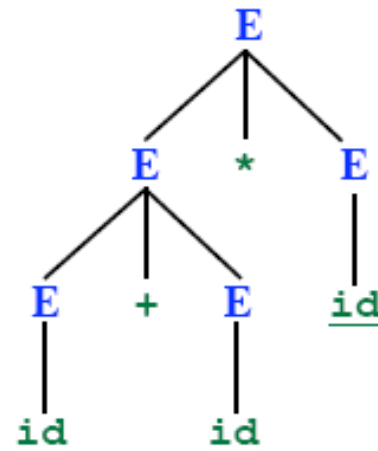
## Leftmost Derivation #1

$E$   
 $\Rightarrow E + E$   
 $\Rightarrow \underline{id} + E$   
 $\Rightarrow \underline{id} + E * E$   
 $\Rightarrow \underline{id} + \underline{id} * E$   
 $\Rightarrow \underline{id} + \underline{id} * \underline{id}$



## Leftmost Derivation #2

$E$   
 $\Rightarrow E * E$   
 $\Rightarrow E + E * E$   
 $\Rightarrow \underline{id} + E * E$   
 $\Rightarrow \underline{id} + \underline{id} * E$   
 $\Rightarrow \underline{id} + \underline{id} * \underline{id}$



# AMBIGUOUS GRAMMAR

- More than one Parse Tree for some sentence.
  - The grammar for a programming language may be ambiguous
  - Need to modify it for parsing.
- 
- Also: Grammar may be left recursive.
  - Need to modify it for parsing.

# ELIMINATION OF AMBIGUITY

- Ambiguous
- A Grammar is ambiguous if there are multiple parse trees for the same sentence.
- Disambiguation
- Express Preference for one parse tree over others
  - Add disambiguating rule into the grammar



# RESOLVING PROBLEMS: AMBIGUOUS GRAMMARS

Consider the following grammar segment:

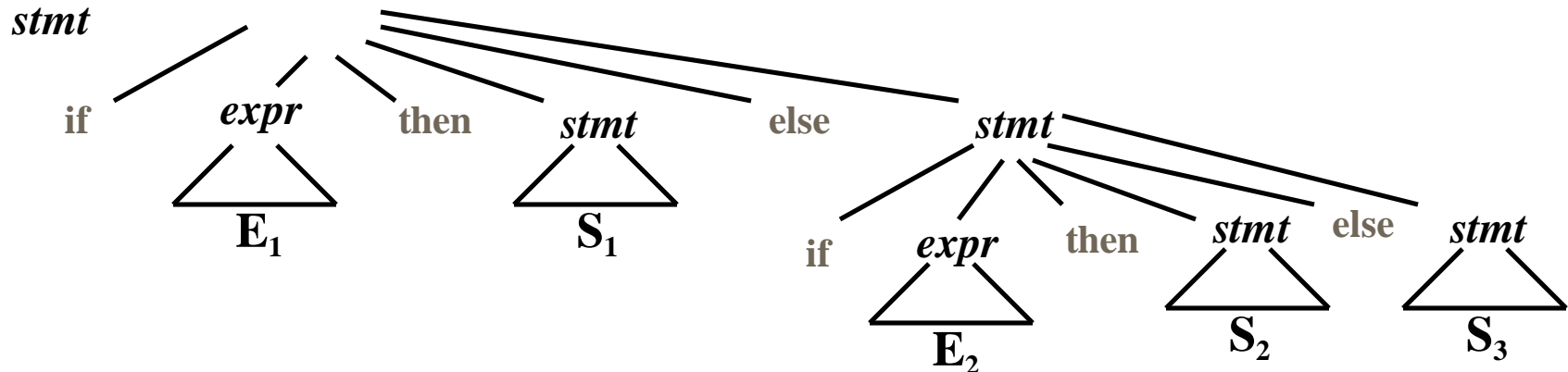
*stmt*  $\rightarrow$  **if** *expr* **then** *stmt*

| **if** *expr* **then** *stmt* **else** *stmt*

| **other** (any other statement)

If **E1** then **S1** else if **E2** then **S2** else **S3**

simple parse tree:



# EXAMPLE : WHAT HAPPENS WITH THIS STRING?

If  $E_1$  then if  $E_2$  then  $S_1$  else  $S_2$

How is this parsed ?

if  $E_1$  then  
  if  $E_2$  then  
     $S_1$   
  else  
     $S_2$

vs.

if  $E_1$  then  
  if  $E_2$  then  
     $S_1$   
  else  
     $S_2$



MD

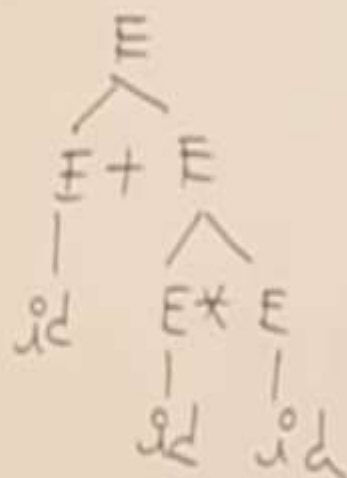
$$E \Rightarrow E * E$$

$$\Rightarrow E + E * E$$

$$\Rightarrow id + E * E$$

$$\Rightarrow id + id * E$$

$$\Rightarrow id + id * id$$



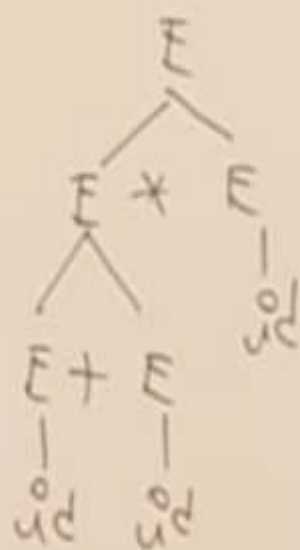
RMD:  $E \Rightarrow E * E$

$$\Rightarrow E * id$$

$$\Rightarrow E + E * id$$

$$\Rightarrow E + id * id$$

$$\Rightarrow id + id * id$$



MD

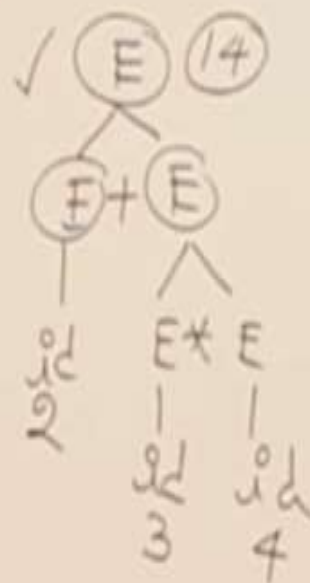
$$E \Rightarrow E * E$$

$$\Rightarrow E + E * E$$

$$\Rightarrow id + E * E$$

$$\Rightarrow id + id * E$$

$$\Rightarrow id + id * id$$



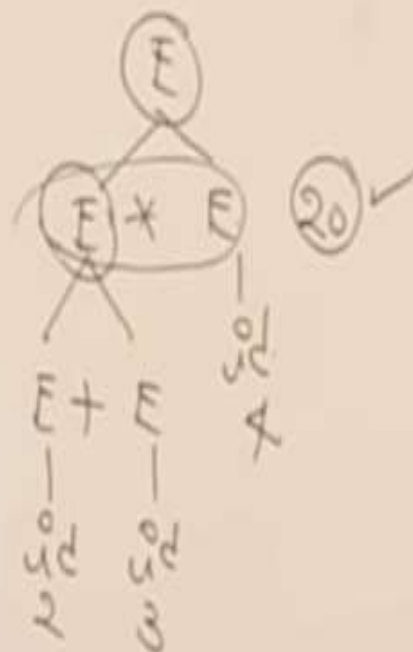
RMD:  $E \Rightarrow E * E$

$$\Rightarrow E * id$$

$$\Rightarrow E + E * id$$

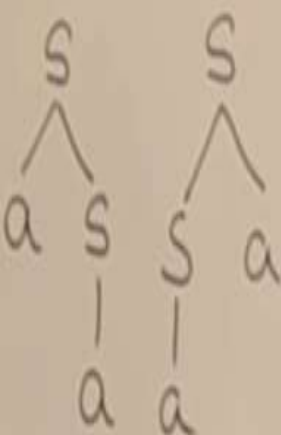
$$\Rightarrow E + id * id$$

$$\Rightarrow id + id * id$$



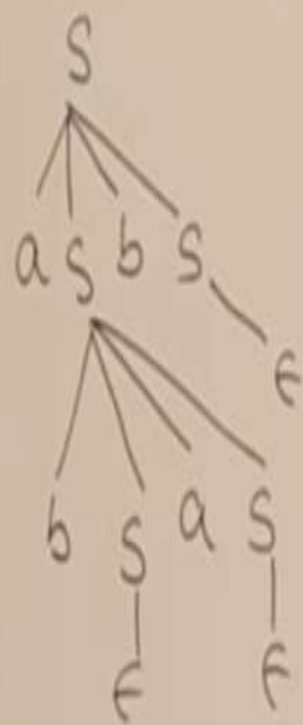
$$S \rightarrow aS / Sa / a$$

$$w = aa$$



$$S \rightarrow aSbS / bSaS / \epsilon$$

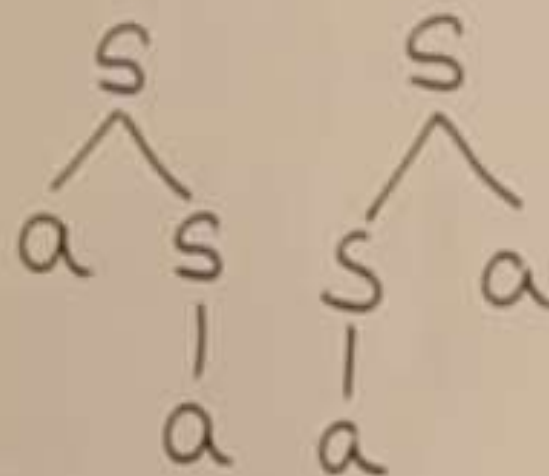
$$w = abab$$



$$R \rightarrow R+R / RR / R^* / a / b / c$$

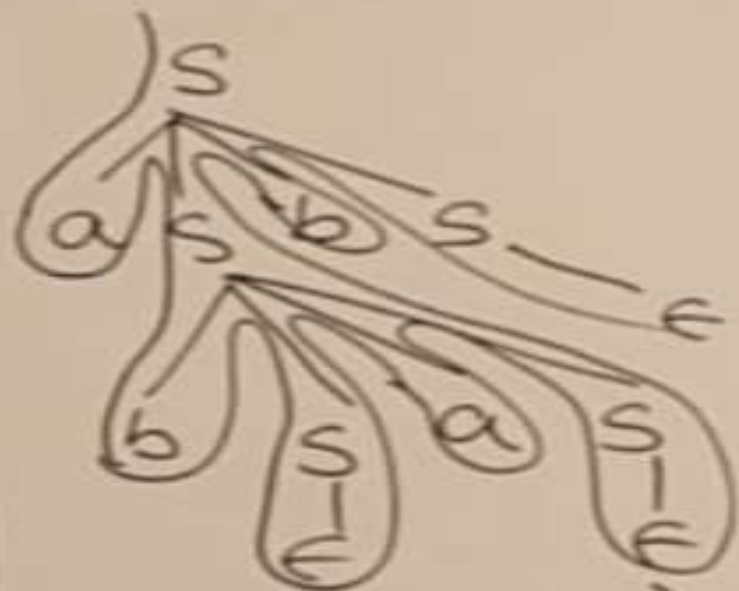
$$\checkmark S \rightarrow aS / Sa / a$$

$$w = aa$$



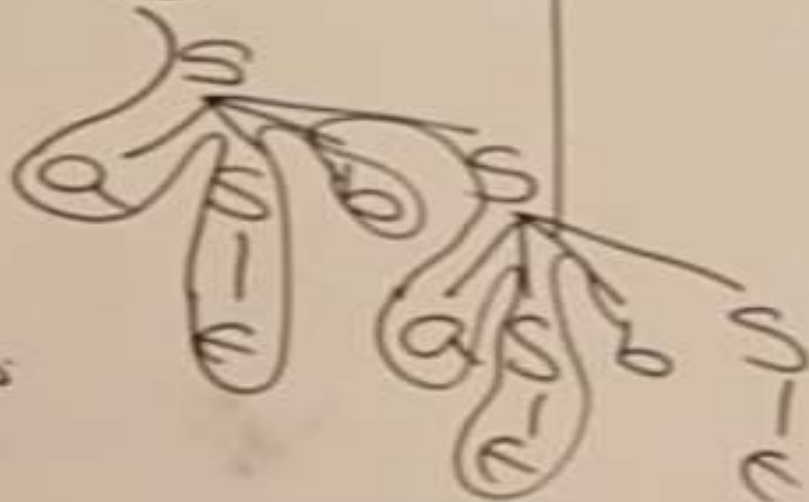
$S \rightarrow aSbS / bSaS / \epsilon$  |  $R \rightarrow R+R / RR / R^* / \epsilon$

$w = abab$



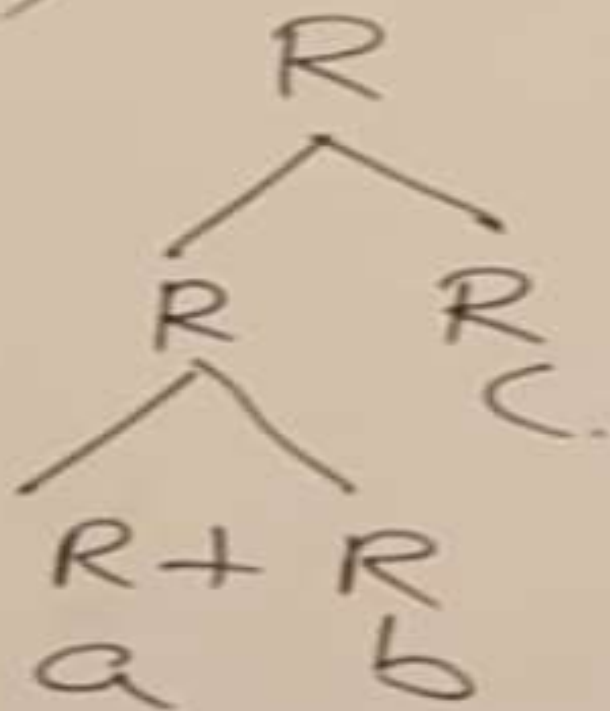
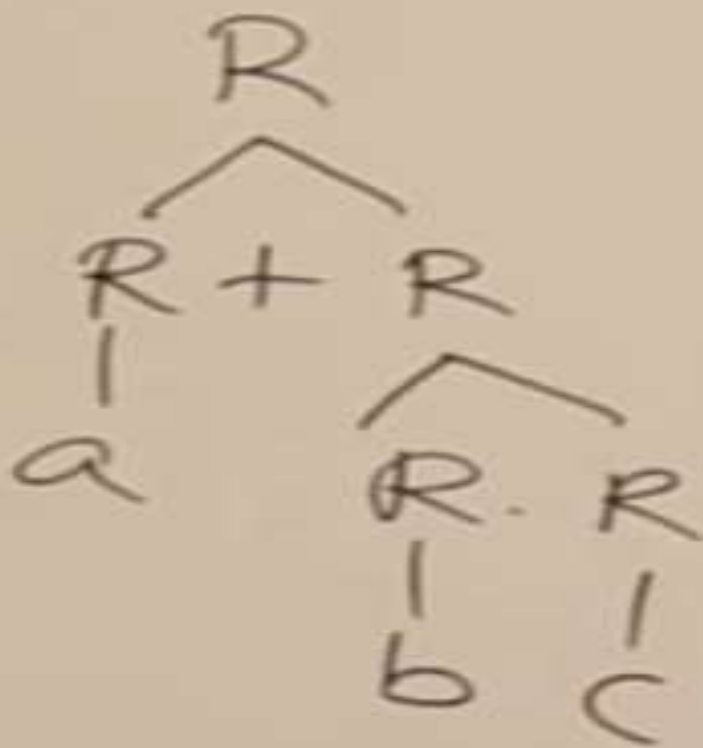
abab

abab



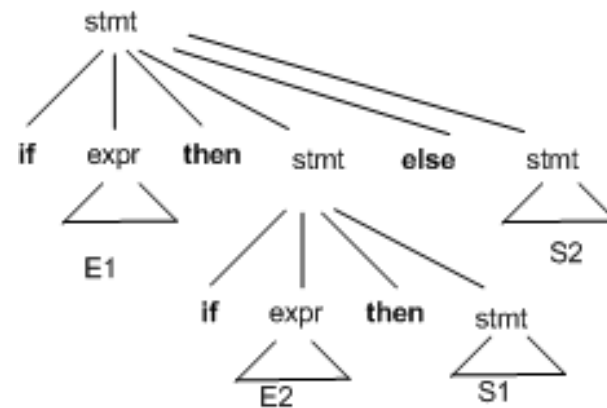
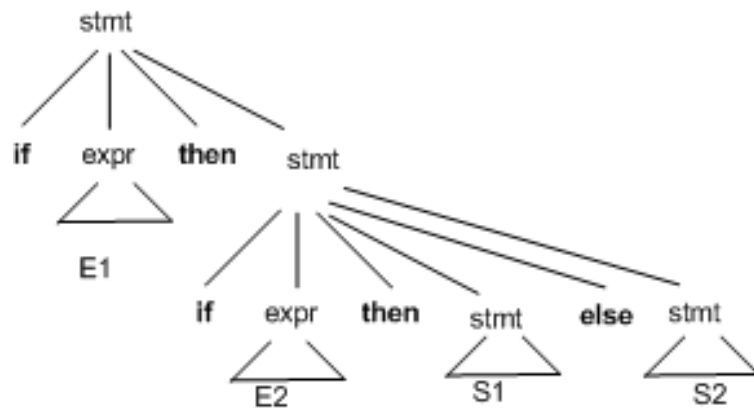
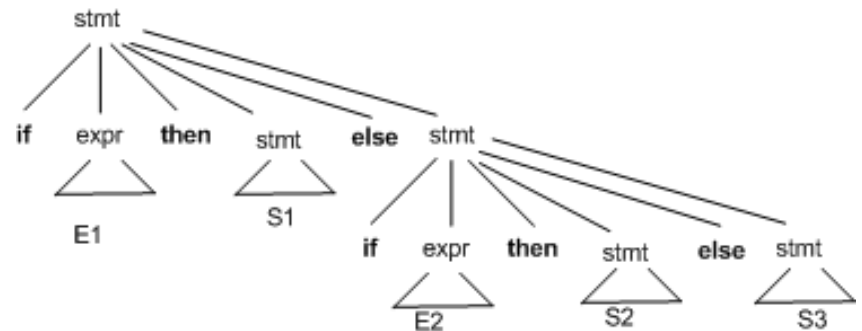
$R \rightarrow R+R / RR / R^* / a / b / c$

$a+(bc)$



# Elimination of ambiguity

stmt  $\rightarrow$  **if** expr **then** stmt  
| **if** expr **then** stmt **else** stmt  
| **other**





# Elimination of ambiguity (cont.)

- Idea:

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

```
stmt  →  matched_stmt
      |  open_stmt

matched_stmt → If expr then matched_stmt else matched_stmt
            | other

open_stmt  → If expr then stmt
            | If expr then matched_stmt else open_stmt
```

# REMOVING AMBIGUITY

## Take Original Grammar:

```
stmt → if expr then stmt  
      | if expr then stmt else stmt  
      | other (any other statement)
```

Rule: Match each **else** with the closest previous unmatched **then**.

## Revise to remove ambiguity:

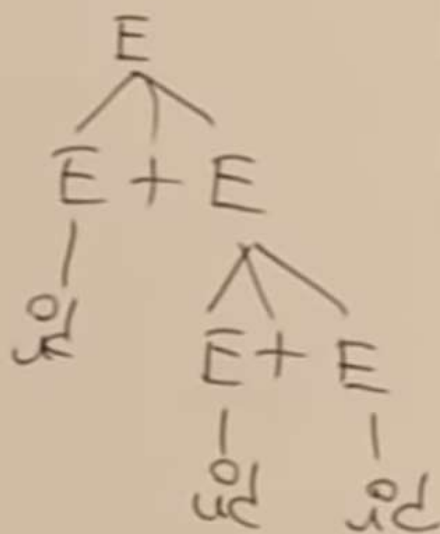
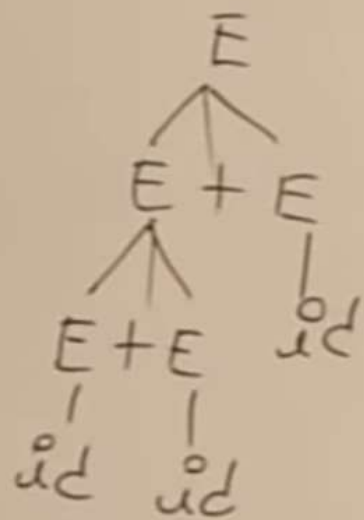
```
stmt → matched_stmt | unmatched_stmt  
matched_stmt → if expr then matched_stmt else matched_stmt /  
other  
unmatched_stmt → if expr then stmt  
                  | if expr then matched_stmt else unmatched_stmt
```

$$E \rightarrow E + E$$

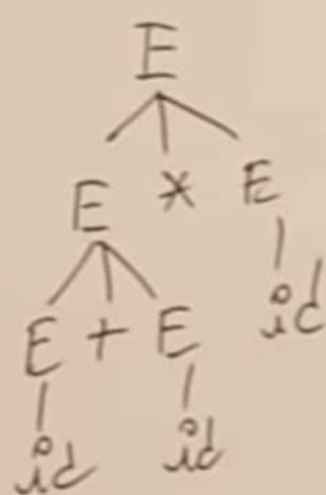
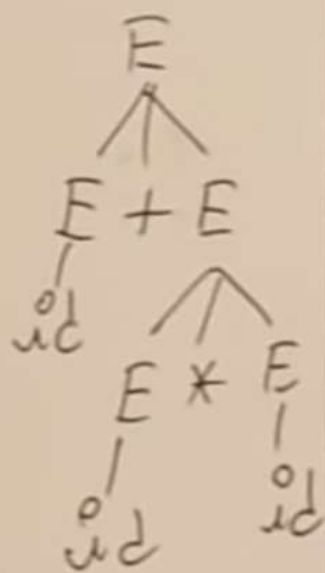
$$/ E * E$$

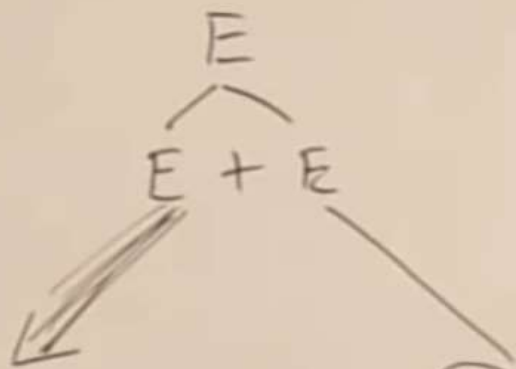
$$/ id$$

$$\boxed{id + id + id}$$



$$\boxed{id + id * id}$$

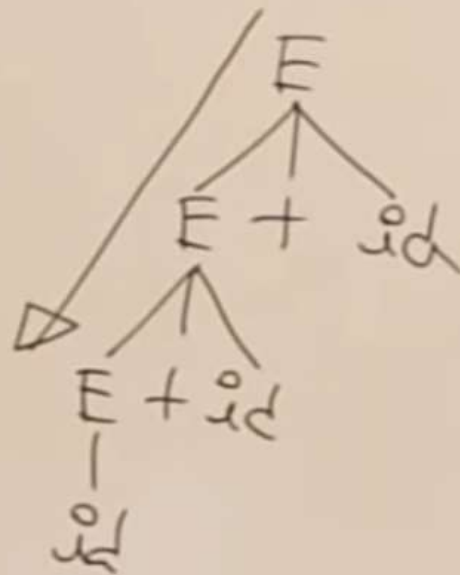
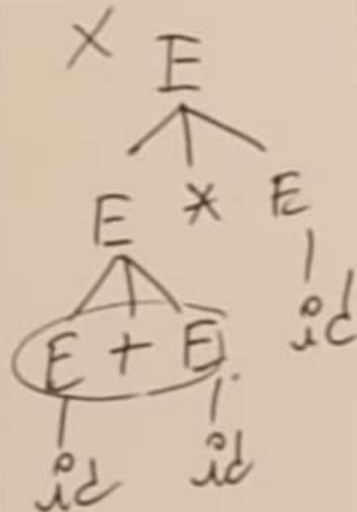
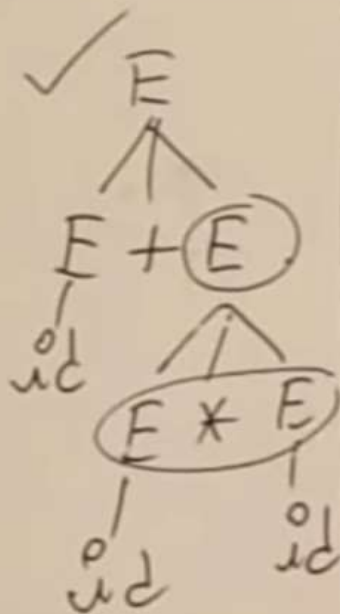




$E \rightarrow \textcircled{E} + id / id$

$id + id * id$

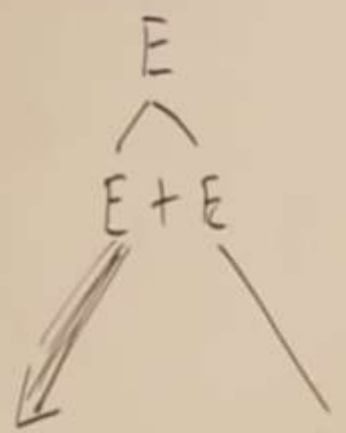
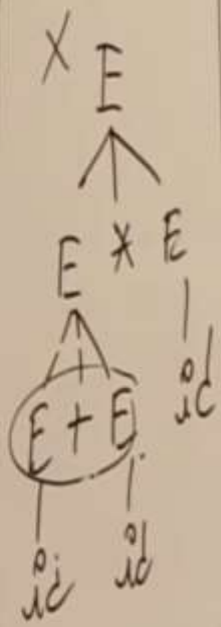
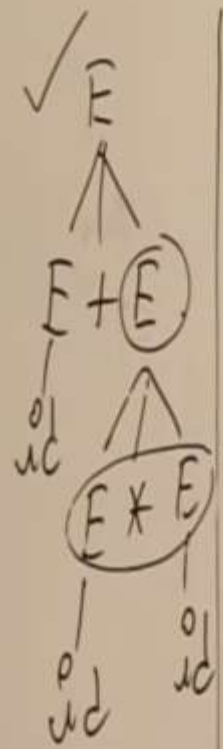
$id + id + id$





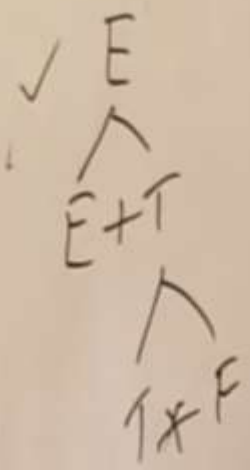
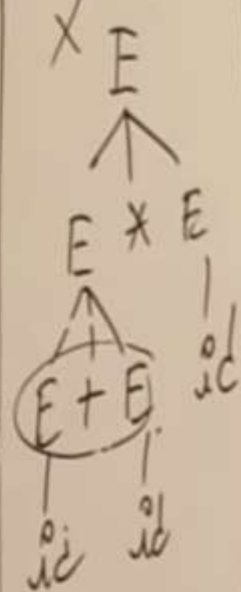
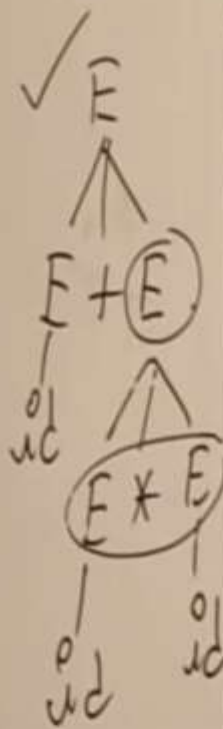
$id + id * id$

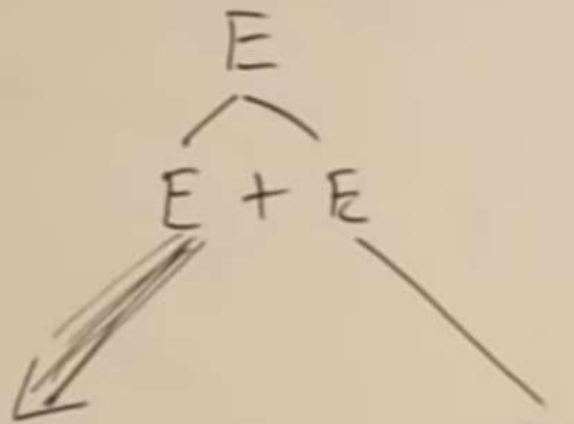
$E \rightarrow E + T$   
 $T \rightarrow T * F$



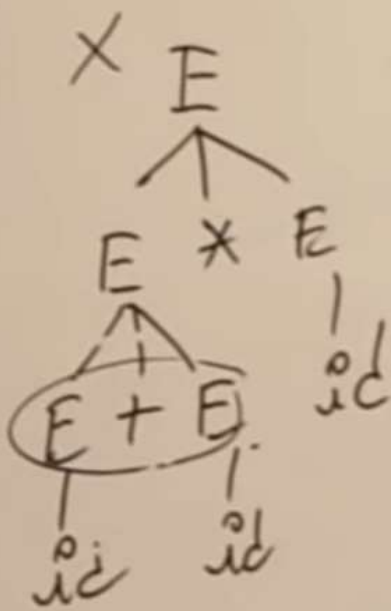
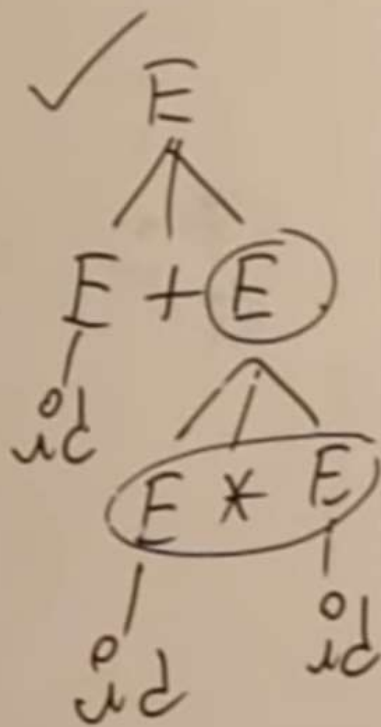
$id + id * id$

$E \rightarrow E + T / T$   
 $T \rightarrow T * F / F$   
 $F \rightarrow id$



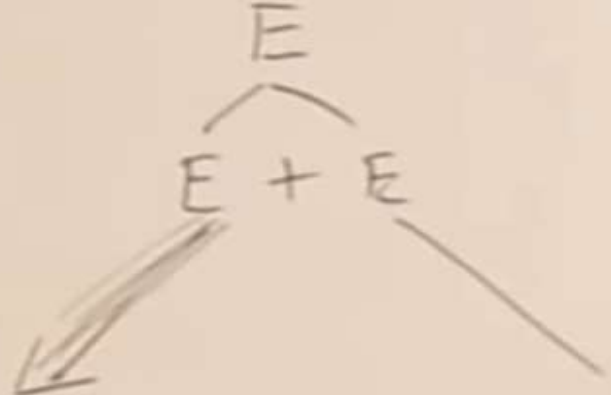


$id + id * id$

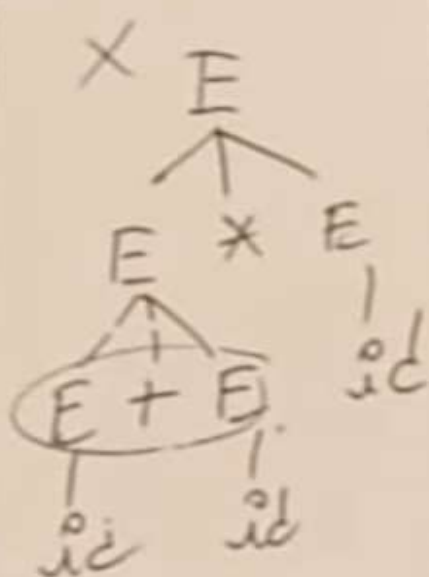
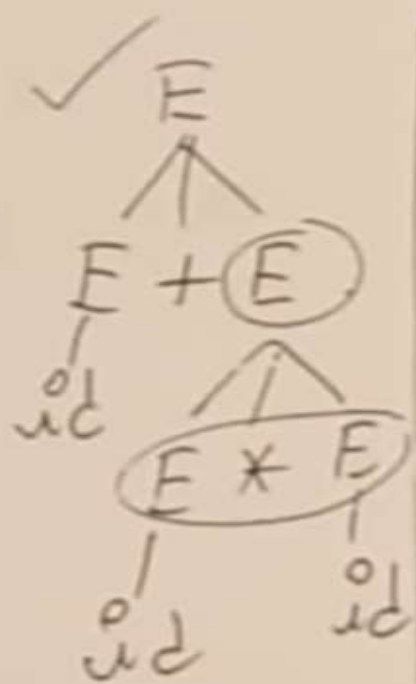


$E \rightarrow E + T / T$   
 $T \rightarrow T * F / F$   
 $F \rightarrow id$

$2 \uparrow 3 \uparrow 2$   
 $2^3 2$



$id + id * id$



$E \rightarrow E + T / T$

$T \rightarrow T * F / F$

$F \rightarrow G \uparrow F / G$

$G \rightarrow id$

$(+), (*), (\uparrow)$

$$R \rightarrow R + R$$

$$/ R R$$

$$/ R *$$

$$/ a$$

$$/ b$$

$$/ c$$

$$E \rightarrow E + T / T$$

$$T \rightarrow T F / F$$

$$F \rightarrow F * / a / b / c.$$



$bExp \rightarrow bExp \ \& \ bExp$

/  $bExp$  and  $bExp$

/ not  $bExp$

/ True

/ False.

$E \rightarrow E \ \& \ F / F$

$F \rightarrow F \text{ and } G / G$

$G \rightarrow \text{Not } G / \text{True} / \text{False.}$

$$\begin{aligned}
 A &\rightarrow A \$ B / B \\
 B &\rightarrow B \# C / C \\
 C &\rightarrow C @ D / D \\
 D &\rightarrow d
 \end{aligned}$$

$$\begin{aligned}
 \$ &\succ \$ \\
 \# &\succ \# \\
 @ &\succ @ \\
 \$ &\prec \# \prec @
 \end{aligned}$$

$$\begin{aligned}
 E &\rightarrow E * F \\
 &\quad / F + E \\
 &\quad / F \\
 &\quad / F \\
 &\quad / d \\
 F &\rightarrow F - F \\
 &\quad / d
 \end{aligned}$$

~~$F \rightarrow F - F$~~

$$* = +$$

$$\begin{aligned}
 \$ &\succ \$ \\
 \# &\succ \# \\
 @ &\succ @ \\
 \$ &\prec \# \prec @
 \end{aligned}$$

$$\begin{aligned}
 * &\succ * \\
 + &\prec +
 \end{aligned}$$

# Elimination of left recursion

- A grammar becomes left-recursive if it has any non-terminal 'A' whose derivation contains 'A' itself as the left-most symbol.
- Left-recursive grammar is considered to be a problematic situation for top-down parsers. Top-down parsers start parsing from the Start symbol, which in itself is non-terminal.
- So, when the parser encounters the same non-terminal in its derivation, it becomes hard for it to judge when to stop parsing the left non-terminal and it goes into an infinite loop.
- 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 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
    - }

# RESOLVING DIFFICULTIES : LEFT RECURSION

A left recursive grammar has rules that support the derivation :  $A \Rightarrow^+ A\alpha$  , for some  $\alpha$ .

Top-Down parsing can't reconcile this type of grammar, since it could consistently make choice which wouldn't allow termination.

$$A \Rightarrow A\alpha \Rightarrow A\alpha\alpha \Rightarrow A\alpha\alpha\alpha \dots \text{etc.} \quad A \rightarrow A\alpha \mid \beta$$

**Take left recursive grammar:**

$$A \rightarrow A\alpha \mid \beta$$

**To the following:**

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \epsilon$$

# WHY IS LEFT RECURSION A PROBLEM ?

**Consider:**

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

**id**

**Derive :  $id + id + id$**

$E \Rightarrow E + T \Rightarrow$

**How can left recursion be removed ?**

$E \rightarrow E + T \mid T$

**What does this generate?**

$E \Rightarrow E + T \Rightarrow T + T$

$E \Rightarrow E + T \Rightarrow E + T + T \Rightarrow T + T + T$

...

**How does this build strings ?**

**What does each string have to start with ?**

# RESOLVING DIFFICULTIES : LEFT RECURSION (2)

**Informal Discussion:**

Take all productions for A and order as:

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

Where no  $\beta_i$  begins with A.

Now apply concepts of previous slide: A

$$\rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \in$$

**For our example:**

$$\begin{array}{lcl}
 E \rightarrow E + T \mid T & \longrightarrow & \left\{ \begin{array}{l} E \rightarrow TE' \\ E' \rightarrow + TE' \mid \in \end{array} \right. \\
 T \rightarrow T * F & \longrightarrow & \left\{ \begin{array}{l} T \rightarrow FT' \\ T' \rightarrow * FT' \mid \in \end{array} \right. \\
 F \rightarrow ( E & \longrightarrow & F \rightarrow ( E ) \mid id
 \end{array}$$

# RESOLVING DIFFICULTIES : LEFT

**RECURSION** (2)  
Problem: If left recursion is two-or-more levels deep,  
this isn't enough

$$\left. \begin{array}{l} S \rightarrow Aa \mid b \\ A \rightarrow Ac \mid Sd \mid \epsilon \end{array} \right\} \quad S \Rightarrow Aa \Rightarrow Sda$$

## Algorithm:

*Input:* Grammar G with ordered Non-Terminals  $A_1, \dots, A_n$

*Output:* An equivalent grammar with no left recursion

1. Arrange the non-terminals in some order  $A_1 = \text{start NT}, A_2, \dots, A_n$
2. for  $i := 1$  to  $n$  do begin  
    for  $j := 1$  to  $i - 1$  do begin  
        replace each production of the form  $A_i \rightarrow A_j \gamma$   
        by the productions  $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$   
        where  $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  are all current  $A_j$  productions ;  
    end  
    eliminate the immediate left recursion among  $A_i$  productions



# USING THE ALGORITHM

Apply the algorithm to:  $A_1 \rightarrow A_2 a \mid b \mid \epsilon$

$A_2 \rightarrow A_2 c \mid A_1 d$

$i = 1$

For  $A_1$  there is no left recursion

$i = 2$

for  $j=1$  to 1 do

Take productions:  $A_2 \rightarrow A_1 \gamma$  and replace  
with

$A_2 \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$  where  $A_1 \rightarrow \delta_1 \mid \delta_2$

$\mid \dots \mid \delta_k$  are  $A_1$

productions  
What's left: In our case  $A_1 \rightarrow A_2 \mid b \mid \epsilon$  becomes  $A_2 \rightarrow A_2 ad \mid bd \mid d$

**Are we done ?**

$A_2 \rightarrow A_2 c \mid A_2 ad \mid bd \mid d$

## USING THE ALGORITHM (2)

**No ! We must still remove  $A_2$  left recursion !**

$$A_1 \rightarrow A_2 a \mid b \mid \epsilon$$

$$A_2 \rightarrow A_2 c \mid A_2 ad \mid bd \mid d$$

**Recall:**

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$$

$$A_1 \rightarrow A_2 a \mid b \mid \epsilon$$

$$A_2 \rightarrow bdA_2' \mid dA_2'$$

$$A_2' \rightarrow c A_2' \mid adA_2' \mid \epsilon$$

**Apply to above case. What do you get ?**

# REMOVING DIFFICULTIES : $\epsilon$ -MOVES

**Transformation:** In order to remove  $A \rightarrow \epsilon$  find all rules of the form  $B \rightarrow uAv$  and *add* the rule  $B \rightarrow uv$  to the grammar  $G$ . Why does

this work ?

**Examples:**

$$E \rightarrow TE'$$

$$T_{\epsilon} \rightarrow E' \quad F \rightarrow T, + TE' \mid$$

$$F_{\epsilon} T \rightarrow ' \rightarrow (E^*) FT \mid i'd \mid$$

$$A_1 \rightarrow A_2 a \mid b$$

$$A_2 \rightarrow bd A_2' \mid A_2'$$

$$A_2' \rightarrow c A_2' \mid bd A_2' \mid \epsilon$$

**A is Grammar  $\epsilon$ -free if:**

1. It has no  $\epsilon$ -production **or**
2. There is exactly one  $\epsilon$ -production  
 $S \rightarrow \epsilon$  and then the start symbol  $S$  does not appear on the right side of any production.

# REMOVING DIFFICULTIES : CYCLES


How would cycles be removed ?

Make sure every production is adding some **terminal(s)** (except a single  $\epsilon$ -production in the start NT)...

e.g.

$S \rightarrow SS \mid (S) \mid \epsilon$

Has a cycle:  $S \Rightarrow SS \Rightarrow S$

  
 $S \rightarrow \epsilon$

Transform to:

$S \rightarrow S(S) \mid (S) \mid \epsilon$

$A()$   
 $\{$   
 $A()$   
 $\}$

Re

LR

RR

$\underline{A} \rightarrow \underline{A}\alpha / \beta$

$A \rightarrow \underline{\alpha} \underline{A} / \beta$

$A()$

$\alpha$

$A()$

$A$   
 $\swarrow \searrow$   
 $\alpha \quad A$   
 $\swarrow \searrow$   
 $\alpha \quad \alpha$   
 $\downarrow$   
 $\beta$

$A$   
 $\downarrow$   
 $\beta$

$\beta\alpha^*$

$\alpha^* \beta$

$\}$   
 $A$   
 $\swarrow \searrow$   
 $\alpha \quad \alpha$   
 $\swarrow \searrow$   
 $\beta \quad \beta$

$\}$

$A(C)$   
 $\{ A(C) \}$   
 $\alpha$

LR

Re

RR

$\{ A \rightarrow A\alpha / \beta \}$

$A \rightarrow \alpha A / \beta$

$A(C)$

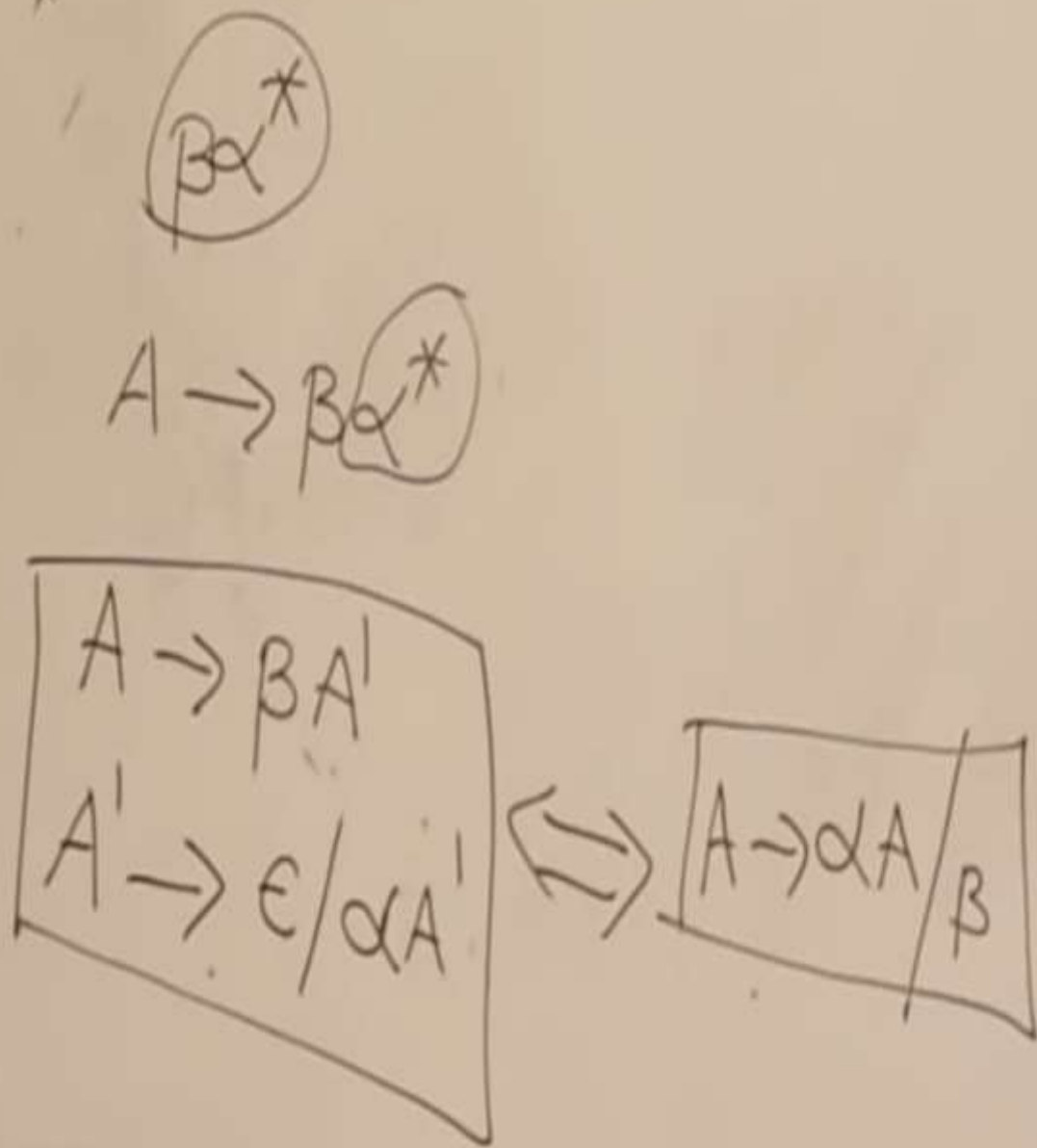
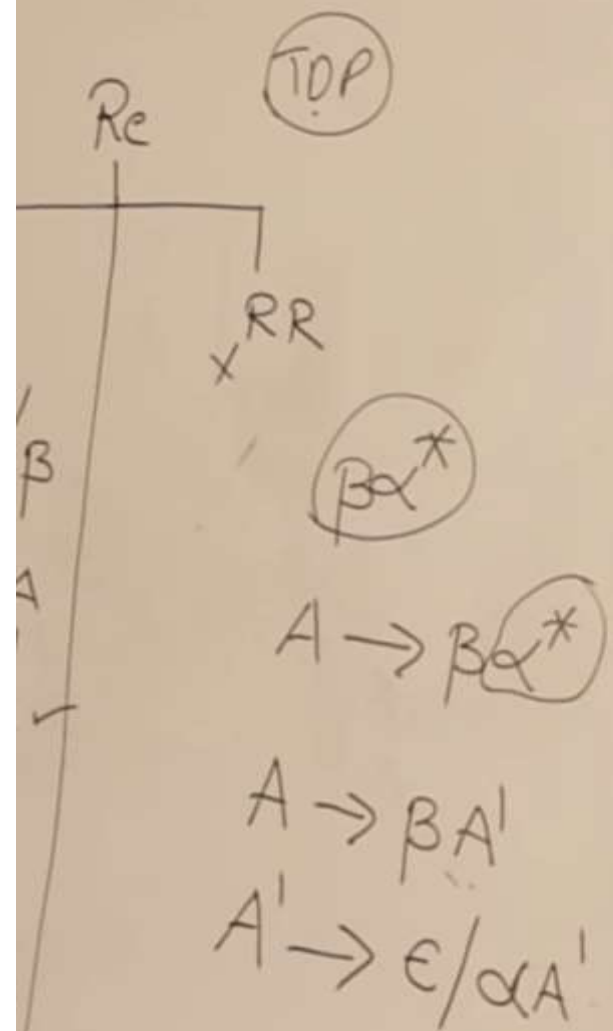
$\{ \alpha A(C) \}$

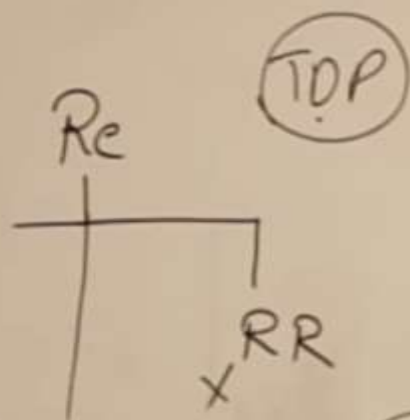
$A$   
 $\alpha$   
 $A$   
 $\alpha$   
 $\beta$

$A$   
 $\beta$

$\{ A \rightarrow A\alpha / \beta \}$   
 $A$   
 $\alpha$   
 $A$   
 $\alpha$   
 $\beta$   
 $\beta\alpha^*$

$\alpha^* \beta$





$$\boxed{\begin{array}{l} \bar{E} \rightarrow \bar{E} + \bar{T} / \bar{T} \\ \bar{A} \quad \bar{A} \quad \alpha \quad \beta \end{array}}$$

$$\boxed{\beta \alpha^*}$$

$$A \rightarrow \beta \alpha^*$$

$$\boxed{\begin{array}{l} E \rightarrow TE' \\ E' \rightarrow \epsilon / +TE' \end{array}}$$

$$\boxed{\begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \epsilon / \alpha A' \end{array}}$$

$$\Leftrightarrow \boxed{A \rightarrow A\alpha / \beta}$$



TOP

xRR

$\beta\alpha^*$

$A \rightarrow \beta\alpha^*$

$\frac{S}{A} \rightarrow \frac{S\alpha}{A} \frac{S\beta}{\alpha} \frac{S}{\beta}$

$S \rightarrow \alpha S'$

$S' \rightarrow \epsilon / \alpha S S'$

$A \rightarrow \beta A'$   
 $A' \rightarrow \epsilon / \alpha A'$

$\Leftrightarrow$

$A \rightarrow A\alpha / \beta$

$$S \rightarrow (L) / x$$

$$\underbrace{L \rightarrow \underbrace{L_1 S}_{\alpha} / \underbrace{S}_{\beta}}_{A}$$

$$A \rightarrow A \alpha / \beta$$

$\Downarrow$

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' / \epsilon$$

$L \rightarrow S L'$
$L' \rightarrow , S' / \epsilon$

# Left factoring

- If more than one grammar production rules has a common prefix string, then the top-down parser cannot make a choice as to which of the production it should take to parse the string in hand
- Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive or top-down parsing.
- Then it cannot determine which production to follow to parse the string as both productions are starting from the same terminal (or non-terminal). To remove this confusion, we use a technique called left factoring.
- Left factoring transforms the grammar to make it useful for top-down parsers. In this technique, we make one production for each common prefixes and the rest of the derivation is added by new productions.
- 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.)

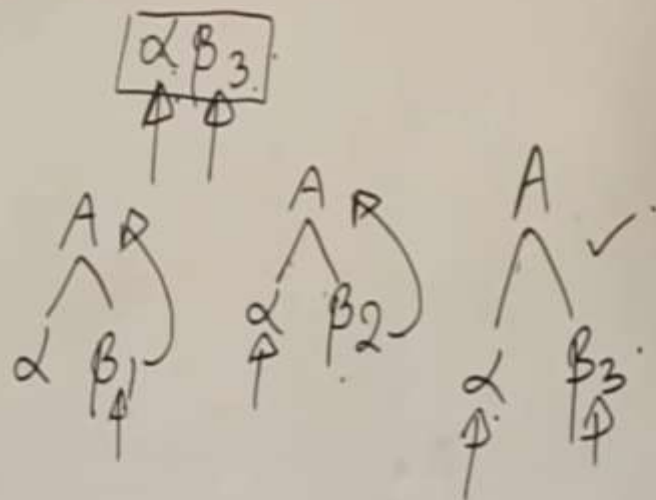
- Algorithm

- For each non-terminal  $A$ , find the longest prefix  $\alpha$  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
  - $A \rightarrow \alpha A' \mid \gamma$
  - $A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

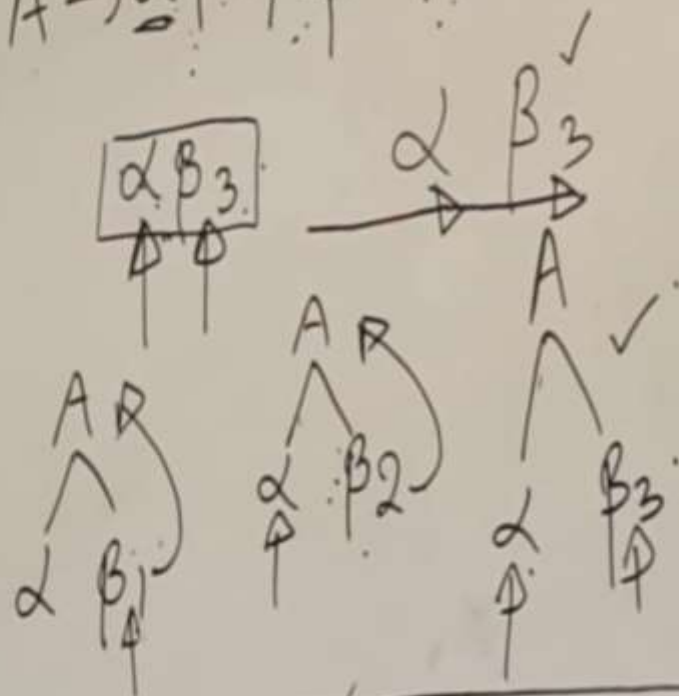
- Example:

- $S \rightarrow I E t S \mid i E t S e S \mid a$
- $E \rightarrow b$

$$A \rightarrow \check{\alpha\beta_1} / \alpha\beta_2 / \alpha\beta_3$$



$$A \rightarrow \check{\alpha\beta_1} / \alpha\beta_2 / \alpha\beta_3$$



$$A \rightarrow \check{\alpha A'}$$

$$A' \rightarrow \beta_1 / \beta_2 / \beta_3$$



$$S \rightarrow \underline{iEtS} \cdot$$

$$\quad \quad \quad / \underline{iEtSeS}$$

$$\quad \quad \quad / a \cdot$$

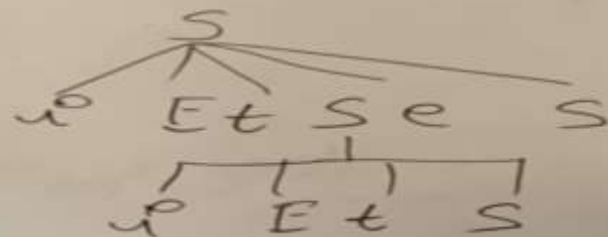
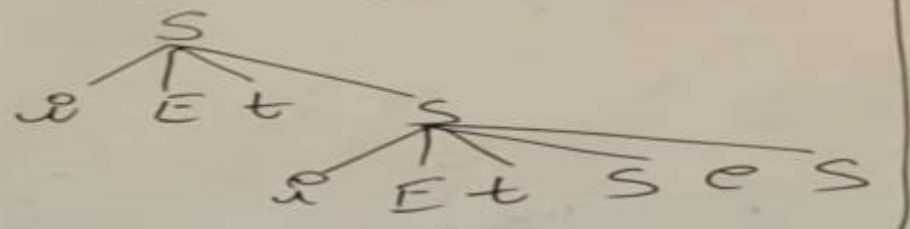
$$E \rightarrow b$$

$$S \rightarrow \underline{iEtSS'} / a \cdot$$

$$S' \rightarrow \epsilon / eS$$

$$E \rightarrow b \cdot$$

$\underline{aEt} \underline{aEt} \underline{SeS}$ .



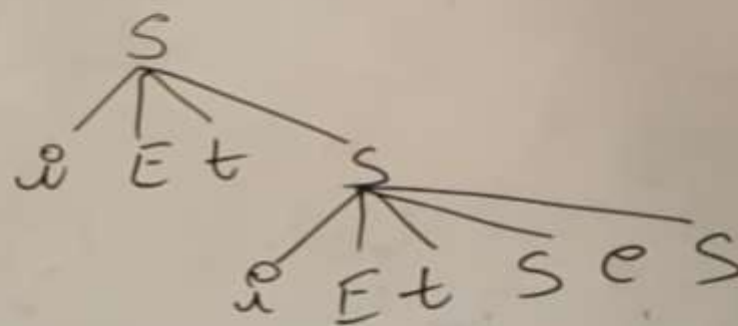
$S \rightarrow \underline{aEtS}$   
 $\quad \quad \underline{aEtSeS}$

$\quad \quad \quad /a$   
 $E \rightarrow b$

---

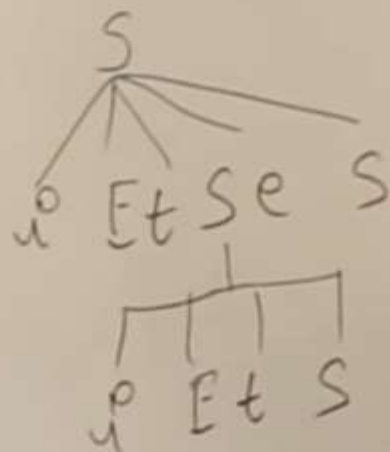
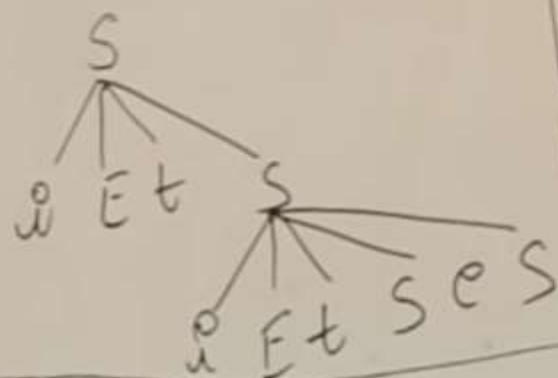
$S \rightarrow \underline{aEtSS'}/a$   
 $S' \rightarrow \epsilon / eS$   
 $E \rightarrow b$

$\underline{aEt} \underline{aEt} \underline{SeS}$ .

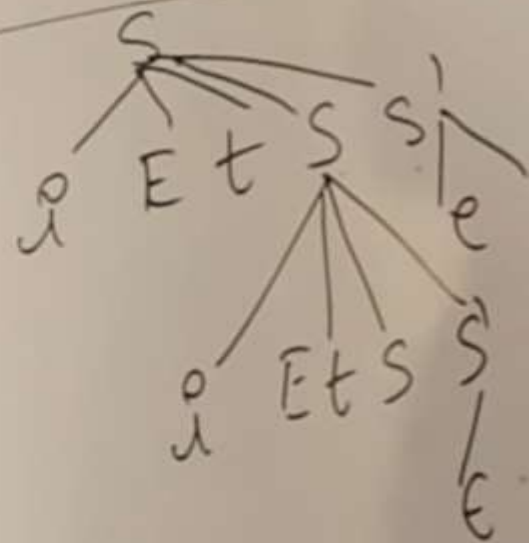
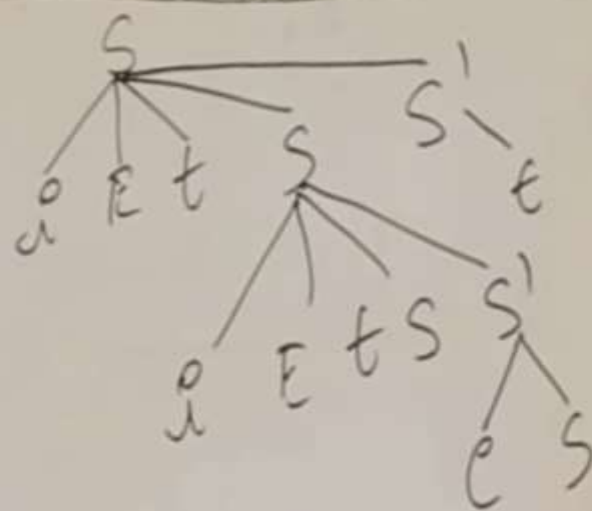


$S \rightarrow \underline{aEtS}$   
 $\quad \underline{aEtSeS}$   
 $\quad \underline{a}$   
 $E \rightarrow b$

$aEt aEtSeS$  ✓



$S \rightarrow \underline{aEtSS'S'}/a$   
 $S' \rightarrow \epsilon / eS$   
 $E \rightarrow b$





$S \rightarrow \underline{a}SSbS$   
/ $\underline{a}SaSb$   
/ $\underline{a}bb$   
/ $\underline{b}$

$S \rightarrow bSSaas$   
/ $bSSaSb$   
/ $bSb$   
/ $a$

# EXAMPLES OF LEFT FACTORING

1.  $S \rightarrow iEtS|iEtSES|a$                        $E \rightarrow b$

2.  $S \rightarrow aSSbS|aSaSb|abb|b$

3.  $S \rightarrow bSSaaS|bSSaSb|bSb|a$