

---

# CMPSC 160

## Translation of Programming Languages

Lecture 4: Top-down parsing and LL(1)  
parsing

---

# Recap on Last Lecture

---

- RE vs. CFG
    - Expressiveness and performance
  - Key concepts
    - Parse tree, Derivation, Leftmost/rightmost derivation, Sentential form
  - **Ambiguity**: if there is multiple parse tree for one input, then there is ambiguity in our CFG grammar.
    - Break the tie between different choices: else-dangling example
    - Clarify Precedence and Associativity
  - **Existence**: If there is a parse tree, there must be a leftmost and a rightmost derivation.
  - But we have not yet talked about how to find/build such a parse tree efficiently? Two general parsing techniques:
    - **Top-down parses + Predicative LL parsers**
    - Bottle-up parses
-

# Top-down Parsing Algorithm

---

Construct the root node of the parse tree, label it with the start symbol, and set the current-node to root node

Repeat until all the input is consumed (i.e., until the frontier of the parse tree matches the input string)

- 1 If the label of the current node is a non-terminal node A, **select a (random) production** with A on its lhs and, for each symbol on its rhs, construct the appropriate child (**not terminating**)
- 2 If the current node is a terminal symbol:
  - If it matches the input string, consume it (advance the input pointer)
  - If it does not match the input string, **backtrack<sup>1</sup>**
- 3 Set the current node to the next node in the frontier of the parse tree
  - If there is no node left in the frontier of the parse tree and input is not consumed, then **backtrack<sup>2</sup>**

Performance issue: two sources of backtracking + one concern of just not terminating

---

# Example

---

Let's use the expression grammar with correct precedence and associativity as an example

1	$S$	$\rightarrow$	$Expr$
2	$Expr$	$\rightarrow$	$Expr + Term$
3			$Expr - Term$
4			$Term$
5	$Term$	$\rightarrow$	$Term * Factor$
6			$Term / Factor$
7			$Factor$
8	$Factor$	$\rightarrow$	num
9			id

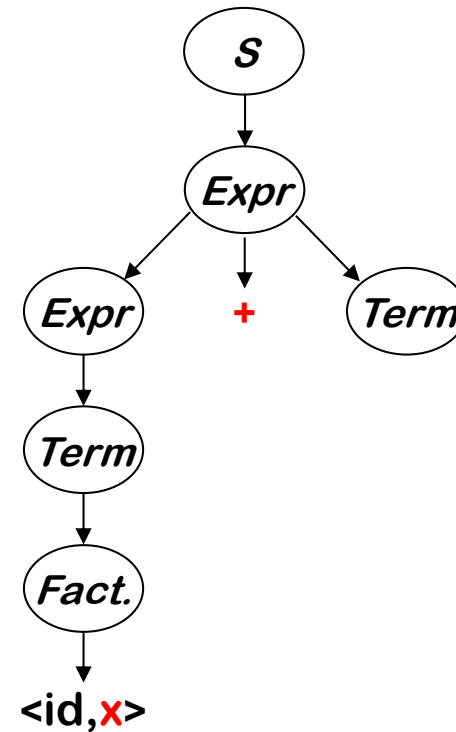
And the input:  $x - 2 * y$

---

# Example: backtrack<sup>1</sup>

Let's try  $x - 2 * y$  :

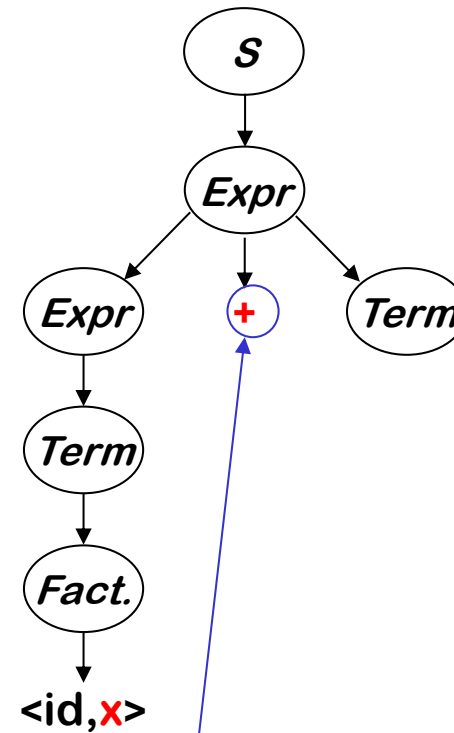
Rule	Sentential Form	Input
-	$S$	$\uparrow x - 2 * y$
1	$Expr$	$\uparrow x - 2 * y$
2	$Expr + Term$	$\uparrow x - 2 * y$
4	$Term + Term$	$\uparrow x - 2 * y$
7	$Factor + Term$	$\uparrow x - 2 * y$
9	$\langle id, x \rangle + Term$	$\uparrow x - 2 * y$
	$\langle id, x \rangle + Term$	$x \uparrow - 2 * y$



# Example: backtrack<sup>1</sup>

Let's try  $x-2*y$  :

Rule	Sentential Form	Input
-	$S$	$\uparrow x-2*y$
1	$Expr$	$\uparrow x-2*y$
2	$Expr + Term$	$\uparrow x-2*y$
4	$Term + Term$	$\uparrow x-2*y$
7	$Factor + Term$	$\uparrow x-2*y$
9	$\langle id, x \rangle + Term$	$\uparrow x-2*y$
	$\langle id, x \rangle + Term$	$x \uparrow -2*y$



Note that “-” doesn't match “+”

The parser must backtrack to here

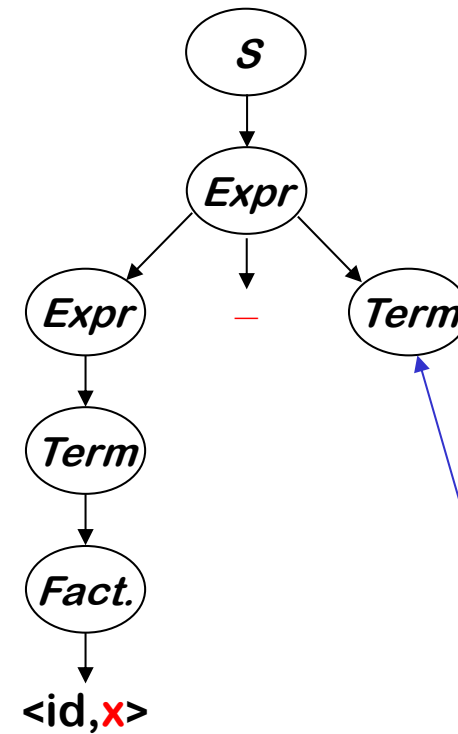
# Example

Continuing with  $x - 2 * y$  :

Rule	Sentential Form	Input
-	$S$	$\uparrow x - 2 * y$
1	$Expr$	$\uparrow x - 2 * y$
3	$Expr - Term$	$\uparrow x - 2 * y$
4	$Term - Term$	$\uparrow x - 2 * y$
7	$Factor - Term$	$\uparrow x - 2 * y$
9	$\langle id, x \rangle - Term$	$\uparrow x - 2 * y$
-	$\langle id, x \rangle - Term$	$x \uparrow - 2 * y$
-	$\langle id, x \rangle - Term$	$x - \uparrow 2 * y$

This time “-” and “-” matched

We can advance past “-” to look at “2”



Now we need to extend  $Term$ , the last  $NT$  in the fringe of the parse tree

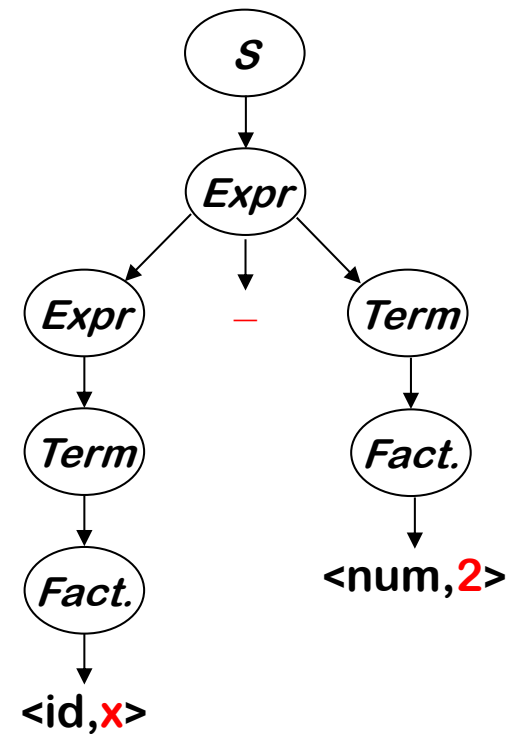
# Example: backtrack<sup>2</sup>

Trying to match the “2” in  $x - 2 * y$  :

Rule	Sentential Form	Input
–	$\langle \text{id}, x \rangle - \text{Term}$	$x - \uparrow 2 * y$
7	$\langle \text{id}, x \rangle - \text{Factor}$	$x - \uparrow 2 * y$
9	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle$	$x - \uparrow 2 * y$
–	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle$	$x - 2 \uparrow * y$

Where are we?

- num matches “2”
  - We have more input, but no *NTs* left to expand
  - The expansion terminated too soon
- ⇒ Need to backtrack

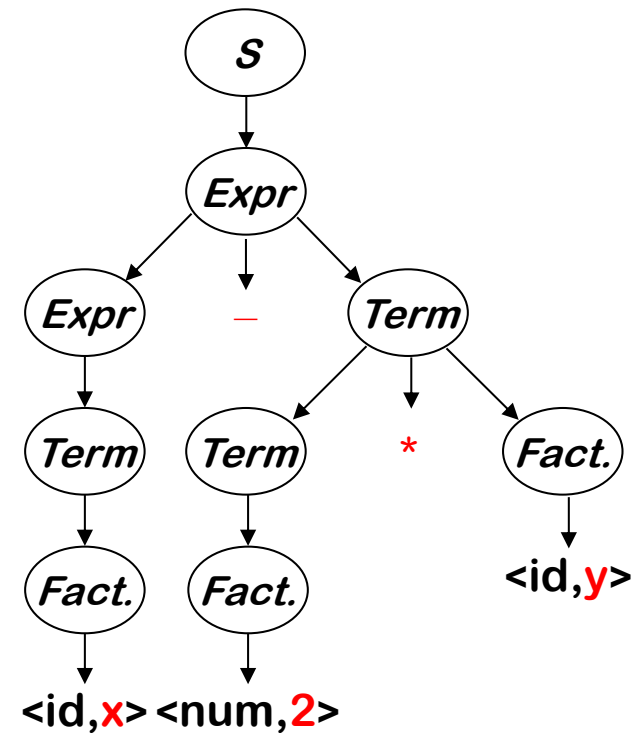




# Example: backtrack<sup>2</sup>

Trying again with “2” in  $x - 2 * y$  :

Rule	Sentential Form	Input
-	$\langle \text{id}, x \rangle - \text{Term}$	$x - \uparrow 2 * y$
5	$\langle \text{id}, x \rangle - \text{Term} * \text{Factor}$	$x - \uparrow 2 * y$
7	$\langle \text{id}, x \rangle - \text{Factor} * \text{Factor}$	$x - \uparrow 2 * y$
8	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle * \text{Factor}$	$x - \uparrow 2 * y$
-	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle * \text{Factor}$	$x - 2 \uparrow * y$
-	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle * \text{Factor}$	$x - 2 * \uparrow y$
9	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle$	$x - 2 * \uparrow y$
-	$\langle \text{id}, x \rangle - \langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle$	$x - 2 * y \uparrow$



This time, we matched and consumed all the input  
 $\Rightarrow$  Success!

# Another possible parsing

Other choices for expansion are possible

<i>Rule</i>	<i>Sentential Form</i>	<i>Input</i>
—	$S$	$\uparrow x - 2 * y$
1	$Expr$	$\uparrow x - 2 * y$
2	$Expr + Term$	$\uparrow x - 2 * y$
2	$Expr + Term + Term$	$\uparrow x - 2 * y$
2	$Expr + Term + Term + Term$	$\uparrow x - 2 * y$
2	$Expr + Term + Term + \dots + Term$	$\uparrow x - 2 * y$

consuming no input !

This does not terminate

- Wrong choice of expansion leads to non-termination, the parser will not backtrack since it does not get to a point where it can backtrack
- Non-termination is a bad property for a parser to have
- Parser must make the right choice

# Left Recursion

---

Top-down parsers cannot handle **left-recursive grammars**

Our expression grammar is left-recursive

- This can lead to non-termination in a top-down parser
- For a top-down parser, **any recursion must be right recursion**
- We would like to convert the **left recursion to right recursion**  
(without changing the language that is defined by the grammar)

A grammar is **left recursive** if there exists a non-terminal  $A$  such that there exists a derivation  $A \Rightarrow^+ A\alpha$ , for some string  $\alpha \in (NT \cup T)^+$

$\Rightarrow^+$  means a bunch of (1 or more) productions applied in series

$\alpha \in (NT \cup T)^+$  means that  $\alpha$  is a non-empty sequence of nonterminal and terminal symbols

---

# Eliminating Immediate Left Recursion

To remove left recursion, we can transform the grammar

Consider a grammar fragment of the form

$$\begin{array}{l} A \rightarrow A \alpha \\ \quad | \beta \end{array}$$

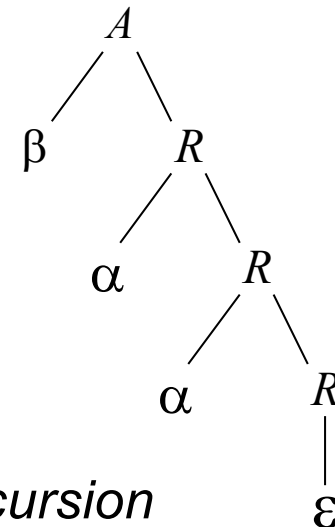
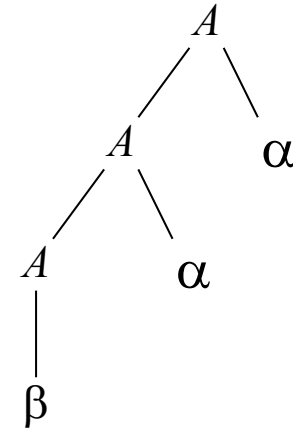
where  $\alpha$  or  $\beta$  are strings of terminal and non-terminal symbols  
and neither  $\alpha$  nor  $\beta$  start with  $A$

We can rewrite this as

$$\begin{array}{l} A \rightarrow \beta R \\ R \rightarrow \alpha R \\ \quad | \varepsilon \end{array}$$


where  $R$  is a new non-terminal

*This accepts the same language, but uses only right recursion*



# Question Time 😊

---

$A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m | \beta_1 | \beta_2 | \dots | \beta_n$   ?

Given

$A \rightarrow A\alpha$   
 $| \beta$



$A \rightarrow \beta R$   
 $R \rightarrow \alpha R$   
 $| \varepsilon$

---

# Eliminating Immediate Left Recursion

---

- ❖ The general form for left recursion is

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

is can be replaced by

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

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

Is the converted CFG much more complicated?

---

# Example

The expression grammar contains two cases of left recursion

<i>Expr</i>	→	<i>Expr + Term</i>	<i>Term</i>	→	<i>Term * Factor</i>
		<i>Expr - Term</i>			<i>Term / Factor</i>
		<i>Term</i>			<i>Factor</i>

Applying the transformation yields

<i>Expr</i>	→	<i>Term Expr'</i>	<i>Term</i>	→	<i>Factor Term'</i>
<i>Expr'</i>	→	<i>+ Term Expr'</i>	<i>Term'</i>	→	<i>* Factor Term'</i>
		<i>- Term Expr'</i>			<i>/ Factor Term'</i>
		$\epsilon$			$\epsilon$

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



$A \rightarrow \beta_1 A' | \beta_2 A' | \dots | \beta_n A'$   
 $A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_m A' | \epsilon$

# Left-Recursive and Right-Recursive Grammar

1	<i>S</i>	→	<i>Expr</i>
2	<i>Expr</i>	→	<i>Expr + Term</i>
3			<i>Expr - Term</i>
4			<i>Term</i>
5	<i>Term</i>	→	<i>Term * Factor</i>
6			<i>Term / Factor</i>
7			<i>Factor</i>
8	<i>Factor</i>	→	num
9			id

1	<i>S</i>	→	<i>Expr</i>
2	<i>Expr</i>	→	<i>Term Expr'</i>
3	<i>Expr'</i>	→	<i>+ Term Expr'</i>
4			<i>- Term Expr'</i>
5			$\epsilon$
6	<i>Term</i>	→	<i>Factor Term'</i>
7	<i>Term'</i>	→	<i>* Factor Term'</i>
8			<i>/ Factor Term'</i>
9			$\epsilon$
10	<i>Factor</i>	→	num
11			id

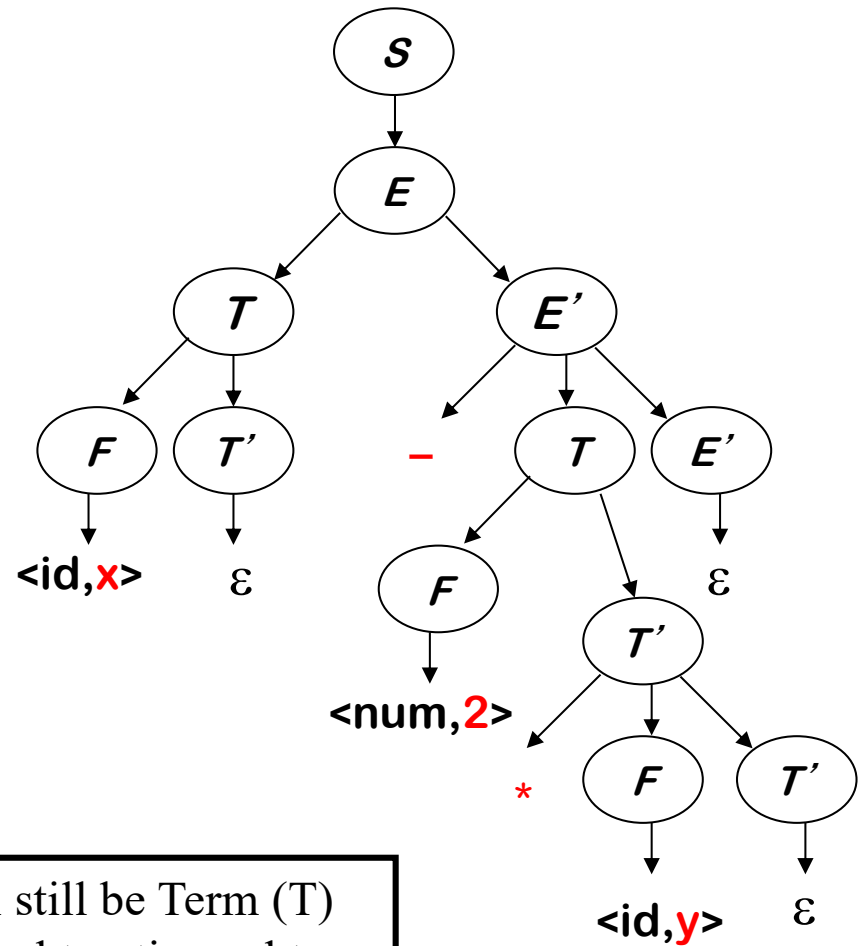
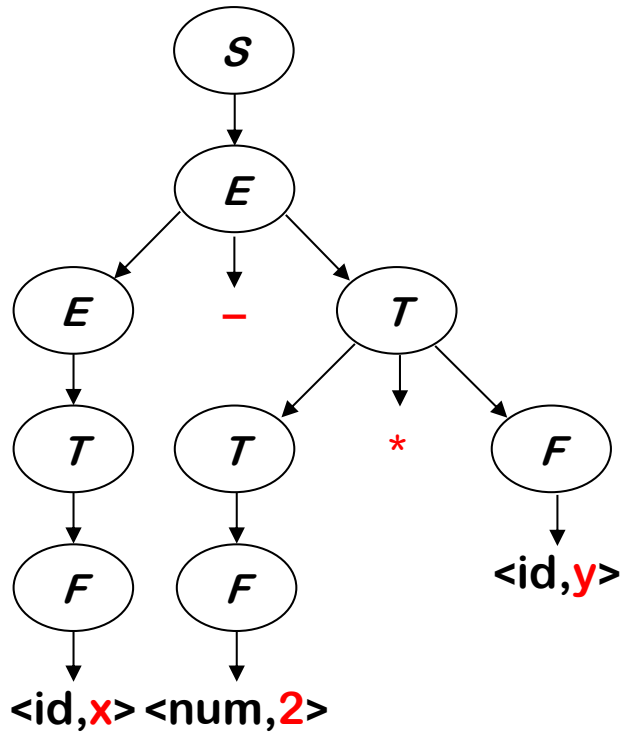
This grammar is correct, if somewhat non-intuitive.

A top-down parser will terminate using it.

Q: Will the transformation change the associativity and precedence?



# Preserves Precedence



Key: the root node for a multiplication subtree will still be Term (T) must be a child node of some E, a root node for a subtraction subtree.

# Eliminating Indirect Left Recursion

Example:

A  $\rightarrow$  Cx  
B  $\rightarrow$  Cy  
C  $\rightarrow$  A | B | z

Rule: first establish some kind of order for non-terminals, and then find all paths where indirect recursion happens.

Order: C < B < A

NT can only have NT with lower order on the RHS of its production rule

C  $\rightarrow$  Cx | Cy | z

C  $\rightarrow$  zC'  
C'  $\rightarrow$  xC' | yC' |  $\epsilon$

A  $\rightarrow$  Cx  
B  $\rightarrow$  Cy  
C  $\rightarrow$  zC'  
C'  $\rightarrow$  xC' | yC' |  $\epsilon$

Done!!!

# Eliminating Left Recursion

---

The previous transformation eliminates immediate left recursion  
What about more general, indirect left recursion?

The general algorithm (Algorithm 4.1 in the Textbook):

*Arrange the NTs into some order  $A_1, A_2, \dots, A_n$*

*for  $i \leftarrow 1$  to  $n$*

*for  $j \leftarrow 1$  to  $i-1$*

*replace each production  $A_i \rightarrow A_j \gamma$  with*

*$A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma,$*

*where  $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  are all the current productions for  $A_j$*

*eliminate any immediate left recursion on  $A_i$  using the direct transformation*

---

# Efficiency with Backtracking?

---

If it picks the wrong production, a top-down parser may backtrack  
Alternative is to look ahead in input & use context to pick correctly.

Solution: **Predictive Parsing**

## Basic idea

The main idea is to look ahead at the next few tokens and use that token to pick the production that you should apply

$X \rightarrow + X$	Here we can use the + and – to decide which rule to apply
$  - Y$	

What is the potential problem here?

$A \rightarrow \alpha B \mid \alpha C$

---

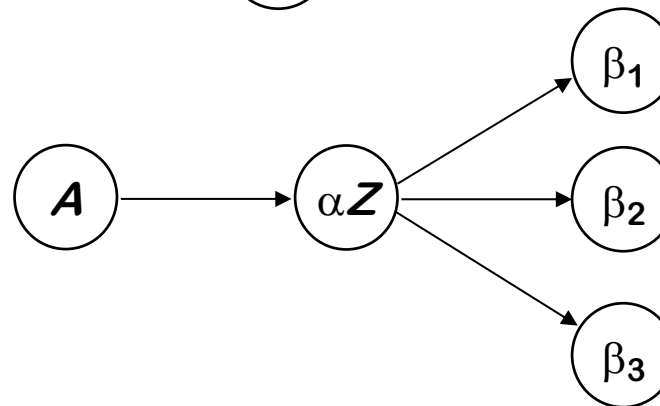
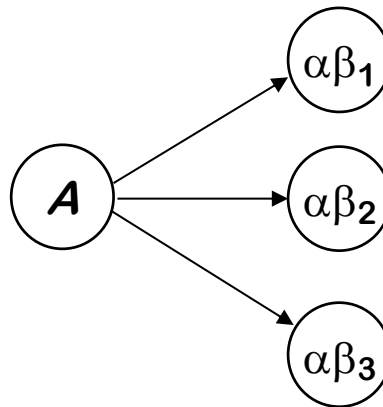
# Left Factoring

A graphical explanation for the left-factoring

$$\begin{array}{l} \mathbf{A} \rightarrow \alpha\beta_1 \\ | \alpha\beta_2 \\ | \alpha\beta_n \end{array}$$

becomes ...

$$\begin{array}{l} \mathbf{A} \rightarrow \alpha \mathbf{Z} \\ \mathbf{Z} \rightarrow \beta_1 \\ | \beta_2 \\ | \beta_n \end{array}$$



# Left Factoring

- We already learned one transformation: Removing left-recursion
- There is another transformation called left-factoring

Left-Factoring Algorithm:

$\forall A \in NT,$

*find the longest prefix  $\alpha$  that occurs in two or more right-hand sides of  $A$*

*if  $\alpha \neq \varepsilon$  then replace all of the  $A$  productions,*

$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_k,$

*with*

$A \rightarrow \alpha Z \mid \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_k$

$Z \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

*where  $Z$  is a new element of  $NT$*

*Repeat until no common prefixes remain*

# More on Predictive Parsing

---

$A \rightarrow \alpha_1$

$A \rightarrow \alpha_2$

$A \rightarrow \alpha_3$

- Suppose all three rules have some NT on the RHS?
- Which rule to select if the next token is “x”? No idea.

❖ First Set for a Non-Terminal (NT)  $\alpha$  :

$x \in \text{FIRST}(\alpha)$       iff 1)  $\alpha \Rightarrow^* \mathbf{x} \gamma$ , for some  $\gamma \in (NT \cup T)^*$  and  $x \in T$   
2)  $\alpha \Rightarrow^* \varepsilon$  and  $x = \varepsilon$

ALL tokens that can be at the beginning of a string that can be derived from  $\alpha$

❖ First set for a Terminal (T)  $x$

$x = \text{FIRST}(x)$

It seems that if all first sets are distinct, we are done!

$\text{Predict}(A \rightarrow \alpha_1) = \text{First}(\alpha_1)$

$\text{Predict}(A \rightarrow \alpha_2) = \text{First}(\alpha_2)$

$\text{Predict}(A \rightarrow \alpha_3) = \text{First}(\alpha_3)$

---

# Question Time ☺

---

$S \rightarrow AB$

$A \rightarrow x | y$

$B \rightarrow 0 | 1$

$\text{FIRST}(S) = \{ \}$

Definition: ALL terminals that can be at the beginning of a string that can be derived from S.

$S \rightarrow AB$

$A \rightarrow x | y | \varepsilon$

$B \rightarrow 0 | 1$

$\text{FIRST}(S) = \{ \}$

$S \rightarrow AB$

$A \rightarrow x | y | \varepsilon$

$B \rightarrow 0 | 1 | \varepsilon$

$\text{FIRST}(S) = \{ \}$

---



# Examples: Compute FIRST(S)

---

$S \rightarrow AB$   
 $A \rightarrow x \mid y$   
 $B \rightarrow 0 \mid 1$

$\text{FIRST}(S) = \{ x, y \}$

$S \rightarrow AB$   
 $A \rightarrow x \mid y \mid \varepsilon$   
 $B \rightarrow 0 \mid 1$

$\text{FIRST}(S) = \{ x, y, 0, 1 \}$

$S \rightarrow AB$   
 $A \rightarrow x \mid y \mid \varepsilon$   
 $B \rightarrow 0 \mid 1 \mid \varepsilon$

$\text{FIRST}(S) = \{ x, y, 0, 1, \varepsilon \}$

---

# How to Compute FIRST Sets

---

To construct  $\text{FIRST}(X)$  for a grammar symbol  $X$ , apply the following rules until no more symbols can be added to  $\text{FIRST}(X)$

1. If  $X$  is a terminal, then  $\text{FIRST}(X) = X$
  2. If  $X$  is a non-terminal and  $X \rightarrow \varepsilon$  is a production, then put  $\varepsilon$  in  $\text{FIRST}(X)$
  3. If  $X$  is a non-terminal and  $X \rightarrow t$  is a production ( $t$  is a terminal), then put  $t$  in  $\text{FIRST}(X)$
  4. If  $X$  is a non-terminal and  $X \rightarrow Y_1 | Y_2 \dots | Y_k$  is a production, then let  $\text{FIRST}(X) = \text{FIRST}(Y_1) \cup \text{FIRST}(Y_2) \dots \cup \text{FIRST}(Y_k)$
  5. If  $X$  is a non-terminal and  $X \rightarrow Y_1 Y_2 \dots Y_k$  is a production, then
    - put every symbol in  $\text{FIRST}(Y_1)$  other than  $\varepsilon$  to  $\text{FIRST}(X)$
    - if  $\varepsilon$  is in  $\text{FIRST}(Y_j)$  for all  $1 \leq j < k$ , put every symbol in  $\text{FIRST}(Y_k)$  other than  $\varepsilon$  to  $\text{FIRST}(X)$
    - put  $\varepsilon$  in  $\text{FIRST}(X)$  if  $\varepsilon$  is in  $\text{FIRST}(Y_i)$  for all  $1 \leq i \leq k$
-

# We still have those pesky epsilons ...

---

$T \rightarrow S z$

$S \rightarrow AB$

$\text{FIRST}(S) = \{ x, y, 0, 1, \varepsilon \}$

$A \rightarrow x \mid y \mid \varepsilon$

$B \rightarrow 0 \mid 1 \mid \varepsilon$

- Suppose we current node to process is  $S$ , and the next token is  $z$ . Are we having an error?
  - We shall also examine the set of characters that can *follow* the current non-terminal If we have  $\varepsilon$  in our FIRST sets
  - This is what the FOLLOW set defines.
    - FOLLOW( $A$ ) for a non-terminal symbol  $A$ .
      - The set of terminal symbols that can appear immediately to the right of  $A$  in some sentential form.
-

# How to Compute FOLLOW Sets?

---

To construct  $\text{FOLLOW}(A)$  for a non-terminal symbol  $A$ , apply the following rules until no more symbols can be added to  $\text{FOLLOW}(A)$ :

1. Put  $\$$  in  $\text{FOLLOW}(S)$  ( $\$$  is the end-of-file symbol,  $S$  is the start symbol)
  2. If there is a production  $B \rightarrow \alpha A \beta$ ,
    - then put everything in  $\text{FIRST}(\beta)$  - except  $\epsilon$  - in  $\text{FOLLOW}(A)$
    - if  $\epsilon$  is in  $\text{FIRST}(\beta)$ , then put everything in  $\text{FOLLOW}(B)$  in  $\text{FOLLOW}(A)$
  3. If there is a production  $B \rightarrow \alpha A$ , then put everything in  $\text{FOLLOW}(B)$  in  $\text{FOLLOW}(A)$
-

# FIRST/FOLLOW Example

---

Example Input:  $x + y ( z + a ( b ) )$

```
Expression → Function
            | ( Expression )
            | Primary + Expression
            | Primary
Primary     → id
            | num
Function    → id ( ParamList )
ParamList  → Expression ParamList
            | ε
```

- $\text{First}(\alpha)$ : ALL tokens that can be at the beginning of a string that can be derived from  $\alpha$ .
- $\text{FOLLOW}(\alpha)$ : The set of terminal symbols that can appear immediately to the right of  $\alpha$  in some sentential form.

$\text{FIRST}(\text{Expression}) = \{ \}$   
 $\text{FIRST}(\text{Primary}) = \{ \}$   
 $\text{FIRST}(\text{Function}) = \{ \}$   
 $\text{FIRST}(\text{ParamList}) = \{ \}$

$\text{FOLLOW}(\text{Expression}) = \{ \}$   
 $\text{FOLLOW}(\text{Primary}) = \{ \}$   
 $\text{FOLLOW}(\text{Function}) = \{ \}$   
 $\text{FOLLOW}(\text{ParamList}) = \{ \}$

---

# FIRST/FOLLOW Example

Example Input:  $x + y ( z + a ( b ) )$

```
Expression → Function
            | ( Expression )
            | Primary + Expression
            | Primary
Primary     → id
            | num
Function    → id ( ParamList )
ParamList  → Expression ParamList
            | ε
```

- $\text{First}(\alpha)$ : ALL tokens that can be at the beginning of a string that can be derived from  $\alpha$ .
- $\text{FOLLOW}(\alpha)$ : The set of terminal symbols that can appear immediately to the right of  $\alpha$  in some sentential form.

$\text{FIRST}(\text{Expression}) = \{ (, \text{num}, \text{id} \}$   
 $\text{FIRST}(\text{Primary}) = \{ \text{num}, \text{id} \}$   
 $\text{FIRST}(\text{Function}) = \{ \text{id} \}$   
 $\text{FIRST}(\text{ParamList}) = \{ \text{id}, \text{num}, (, \varepsilon \}$

$\text{FOLLOW}(\text{Expression}) = \{ \$, (, ), \text{id}, \text{num} \}$   
 $\text{FOLLOW}(\text{Primary}) = \{ \$, (, ), +, \text{id}, \text{num} \}$   
 $\text{FOLLOW}(\text{Function}) = \{ \$, (, ), \text{id}, \text{num} \}$   
 $\text{FOLLOW}(\text{ParamList}) = \{ ) \}$

# LL(1) Grammars

---

Left-to-right scan of the input, Leftmost derivation, 1-token look-ahead

A grammar  $G$  is LL(1) if for each set of its productions

$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ :

$\text{FIRST}(\alpha_1), \text{FIRST}(\alpha_2), \dots, \text{FIRST}(\alpha_n)$ , are all pair-wise disjoint

If  $\alpha_j \Rightarrow^* \varepsilon$ , then  $\text{FIRST}(\alpha_j) \cap \text{FOLLOW}(A) = \emptyset$  for all  $1 \leq i \leq n, i \neq j$

- In other words, LL(1) grammars
    - during leftmost derivation, productions are uniquely predictable with a **one** token lookahead
-

# Where are we in the process?

---

