

# CS 335: Top-Down Parsing

**Swarnendu Biswas**

Department of Computer Science and Engineering,  
Indian Institute of Technology Kanpur

Sem 2023-24-II



# Example Expression Grammar

$Start \rightarrow Expr$

$Expr \rightarrow Expr + Term \mid Expr - Term \mid Term$

$Term \rightarrow Term \times Factor \mid Term \div Factor \mid Factor$

$Factor \rightarrow (Expr) \mid \mathbf{num} \mid \mathbf{name}$

↓  
priority

# Derivation of **name + name × name** with Oracular Knowledge

Sentential Form	Input
<i>Expr</i>	↑ <b>name + name × name</b>
<i>Expr + Term</i>	↑ <b>name + name × name</b>
<i>Term + Term</i>	↑ <b>name + name × name</b>
<i>Factor + Term</i>	↑ <b>name + name × name</b>
<b>name + Term</b>	↑ <b>name + name × name</b>
<b>name + Term</b>	<b>name</b> ↑ <b>+ name × name</b>
<b>name + Term</b>	<b>name</b> + ↑ <b>name × name</b>
<b>name + Term × Factor</b>	<b>name</b> + ↑ <b>name × name</b>
<b>name + Factor × Factor</b>	<b>name</b> + ↑ <b>name × name</b>
<b>name + name × Factor</b>	<b>name</b> + ↑ <b>name × name</b>
<b>name + name × Factor</b>	<b>name</b> + <b>name</b> ↑ <b>× name</b>
<b>name + name × Factor</b>	<b>name</b> + <b>name</b> × ↑ <b>name</b>
<b>name + name × name</b>	<b>name</b> + <b>name</b> × ↑ <b>name</b>
<b>name + name × name</b>	<b>name</b> + <b>name</b> × <b>name</b> ↑

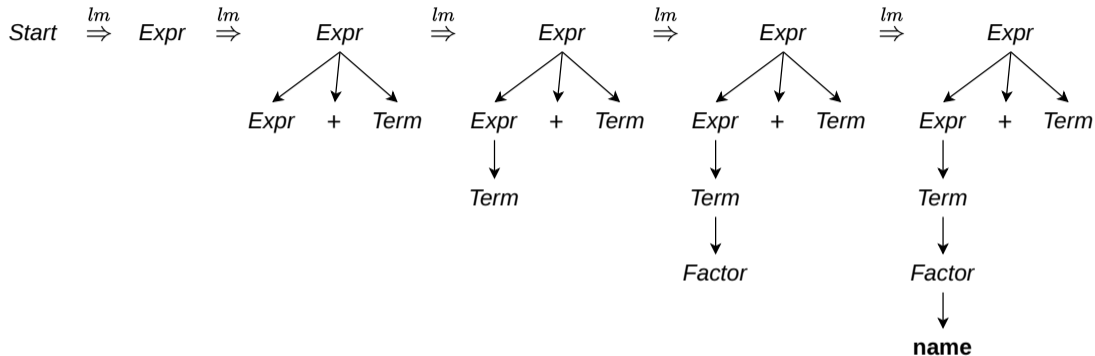
# Derivation of **name + name × name** with Oracular Knowledge

Sentential Form	Input
<i>Expr</i>	↑ <b>name + name × name</b>
<i>Expr + Term</i>	↑ <b>name + name × name</b>
<i>Term + Term</i>	↑ <b>name + name × name</b>
<i>Factor + Term</i>	↑ <b>name + name × name</b>
<b>name + Term</b>	↑ <b>name + name × name</b>

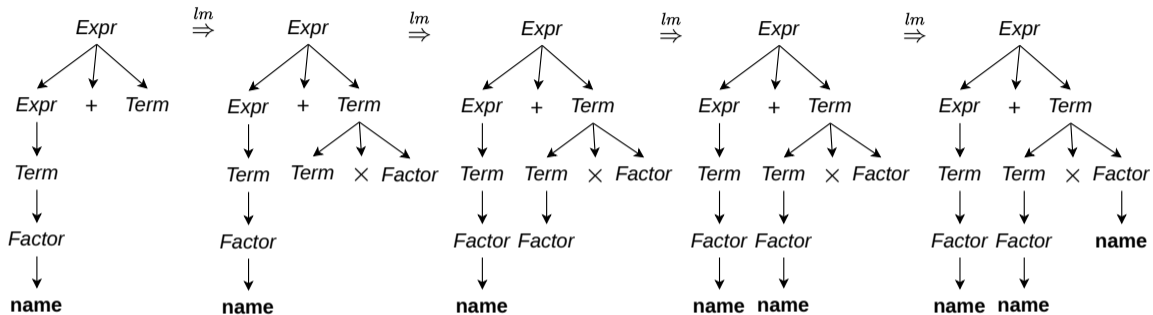
The current input terminal being scanned is called the lookahead symbol

<i>name + Factor × Factor</i>	<b>name +</b>   <b>name × name</b>
<b>name + name × Factor</b>	<b>name +</b> ↑ <b>name × name</b>
<i>name + name × Factor</i>	<b>name + name</b> ↑ <b>× name</b>
<i>name + name × Factor</i>	<b>name + name ×</b> ↑ <b>name</b>
<b>name + name × name</b>	<b>name + name ×</b> ↑ <b>name</b>
<b>name + name × name</b>	<b>name + name × name</b> ↑

# Derivation of **name + name × name** with Oracular Knowledge



# Derivation of **name + name × name** with Oracular Knowledge



# Top-Down Parsing

## High-level idea in top-down parsing

- (i) Start with the root (i.e., start symbol) of the parse tree
- (ii) Grow the tree downwards by expanding the production at the lower levels of the tree
  - ▶ **Select a nonterminal** and extend it by adding children corresponding to the right side of **some production** for the nonterminal
- (iii) Repeat till the lower fringe consists only of terminals and the input is consumed

- Top-down parsing finds a **leftmost derivation** for an input string
- Expands the parse tree with a **preorder depth-first** traversal

# Top-Down Parsing

## High-level idea in top-down parsing

- (i) Start with the root (i.e., start symbol) of the parse tree
- (ii) Grow the tree downwards by expanding the production at the lower levels of the tree
  - ▶ **Select a nonterminal** and extend it by adding children corresponding to the right side of **some production** for the nonterminal
- (iii) Repeat till the lower fringe consists only of terminals and the input is consumed

## Mismatch in the lower fringe and the remaining input stream implies

- (i) Wrong choice of productions while expanding nonterminals, selection of a production may involve trial-and-error
- (ii) Input character stream is not part of the language



# Top-Down Parsing Algorithm

```
root = node for the Start symbol
curr = root
push(null) // Stack

word = getNextWord()
while (true)
  if curr ∈ Nonterminal
    pick next rule  $A \rightarrow \beta_1\beta_2\dots\beta_n$  to expand curr
    create nodes for  $\beta_1, \beta_2, \dots, \beta_n$  as children of curr
    push( $\beta_n\beta_{n-1}\dots\beta_1$ ) // reverse order
    curr =  $\beta_1$ 
  if curr == word
    word = getNextWord()
    curr = pop() // Consumed
  if word == EOF and curr == null
    accept input
  else
    backtrack
```

# Derivation of $\text{name} + \text{name} \times \text{name}$

Rule #	Production
0	$\text{Start} \rightarrow \text{Expr}$
1	$\text{Expr} \rightarrow \text{Expr} + \text{Term}$
2	$\text{Expr} \rightarrow \text{Expr} - \text{Term}$
3	$\text{Expr} \rightarrow \text{Term}$
4	$\text{Term} \rightarrow \text{Term} \times \text{Factor}$
5	$\text{Term} \rightarrow \text{Term} \div \text{Factor}$
6	$\text{Term} \rightarrow \text{Factor}$
7	$\text{Factor} \rightarrow (\text{Expr})$
8	$\text{Factor} \rightarrow \text{num}$
9	$\text{Factor} \rightarrow \text{name}$

Rule #	Sentential Form	Input
	$\text{Expr}$	$\uparrow \text{name} + \text{name} \times \text{name}$
1	$\text{Expr} + \text{Term}$	$\uparrow \text{name} + \text{name} \times \text{name}$
3	$\text{Term} + \text{Term}$	$\uparrow \text{name} + \text{name} \times \text{name}$
6	$\text{Factor} + \text{Term}$	$\uparrow \text{name} + \text{name} \times \text{name}$
9	$\text{name} + \text{Term}$	$\uparrow \text{name} + \text{name} \times \text{name}$
	$\text{name} + \text{Term}$	$\text{name} \uparrow + \text{name} \times \text{name}$
	$\text{name} + \text{Term}$	$\text{name} + \uparrow \text{name} \times \text{name}$
4	$\text{name} + \text{Term} \times \text{Factor}$	$\text{name} + \uparrow \text{name} \times \text{name}$
4	$\text{name} + \text{Factor} \times \text{Factor}$	$\text{name} + \uparrow \text{name} \times \text{name}$
9	$\text{name} + \text{name} \times \text{Factor}$	$\text{name} + \uparrow \text{name} \times \text{name}$
	$\text{name} + \text{name} \times \text{Factor}$	$\text{name} + \text{name} \uparrow \times \text{name}$
	$\text{name} + \text{name} \times \text{Factor}$	$\text{name} + \text{name} \times \uparrow \text{name}$
9	$\text{name} + \text{name} \times \text{name}$	$\text{name} + \text{name} \times \uparrow \text{name}$
	$\text{name} + \text{name} \times \text{name}$	$\text{name} + \text{name} \times \text{name} \uparrow$

# Derivation of $\text{name} + \text{name} \times \text{name}$

Rule #	Production
0	$\text{Start} \rightarrow \text{Expr}$
1	$\text{Expr} \rightarrow \text{Expr} + \text{Term}$
2	$\text{Expr} \rightarrow \text{Expr} - \text{Term}$
3	$\text{Expr} \rightarrow \text{Term}$
4	$\text{Term} \rightarrow \text{Term} \times \text{Factor}$
5	$\text{Term} \rightarrow \text{Term} / \text{Factor}$
6	$\text{Term} \rightarrow \text{num}$
7	$\text{Factor} \rightarrow \text{id}$
8	$\text{Factor} \rightarrow \text{num}$
9	$\text{Factor} \rightarrow \text{name}$

Rule #	Sentential Form	Input
	$\text{Expr}$	$\uparrow \text{name} + \text{name} \times \text{name}$
1	$\text{Expr} + \text{Term}$	$\uparrow \text{name} + \text{name} \times \text{name}$
3	$\text{Term} + \text{Term}$	$\uparrow \text{name} + \text{name} \times \text{name}$
6	$\text{Factor} + \text{Term}$	$\uparrow \text{name} + \text{name} \times \text{name}$
9	$\text{name} + \text{Term}$	$\uparrow \text{name} + \text{name} \times \text{name}$
	$\text{name} + \text{Term} \times \text{Factor}$	$\uparrow \text{name} + \text{name} \times \text{name}$
	$\text{name} + \text{Term} \times \text{Factor}$	$\uparrow \text{name} \times \text{name}$
	$\text{name} + \text{Term} \times \text{Factor}$	$\uparrow \text{name} \times \text{name}$
4	$\text{name} + \text{Factor} \times \text{Factor}$	$\text{name} + \uparrow \text{name} \times \text{name}$
9	$\text{name} + \text{name} \times \text{Factor}$	$\text{name} + \uparrow \text{name} \times \text{name}$
	$\text{name} + \text{name} \times \text{Factor}$	$\text{name} + \text{name} \uparrow \times \text{name}$
	$\text{name} + \text{name} \times \text{Factor}$	$\text{name} + \text{name} \times \uparrow \text{name}$
9	$\text{name} + \text{name} \times \text{name}$	$\text{name} + \text{name} \times \uparrow \text{name}$
	$\text{name} + \text{name} \times \text{name}$	$\text{name} + \text{name} \times \text{name} \uparrow$

How does a top-down parser choose which rule to apply?

# Deterministically Selecting a Production in Expression Grammar

Rule #	Production
0	$Start \rightarrow Expr$
1	$Expr \rightarrow Expr + Term$
2	$Expr \rightarrow Expr - Term$
3	$Expr \rightarrow Term$
4	$Term \rightarrow Term \times Factor$
5	$Term \rightarrow Term \div Factor$
6	$Term \rightarrow Factor$
7	$Factor \rightarrow (Expr)$
8	$Factor \rightarrow \mathbf{num}$
9	$Factor \rightarrow \mathbf{name}$

Rule #	Sentential Form	Input
	$Expr$	$\uparrow \mathbf{name + name \times name}$
1	$Expr + Term$	$\uparrow \mathbf{name + name \times name}$
1	$Expr + Term + Term$	$\uparrow \mathbf{name + name \times name}$
1	$Expr + Term + Term + \dots$	$\uparrow \mathbf{name + name \times name}$
1	$\dots$	$\uparrow \mathbf{name + name \times name}$
1	$\dots$	$\uparrow \mathbf{name + name \times name}$

# Deterministically Selecting a Production in Expression Grammar

Rule #	Production
0	$Start \rightarrow Expr$
1	$Expr \rightarrow Expr + Term$
2	$Expr \rightarrow Expr - Term$
3	$Expr \rightarrow Term$
4	$Term \rightarrow Term \times Factor$
5	$Term \rightarrow Term / Factor$
6	$Term \rightarrow (Expr)$
7	$Factor \rightarrow \text{num}$
8	$Factor \rightarrow \text{num}$
9	$Factor \rightarrow \text{name}$

Rule #	Sentential Form	Input
	$Expr$	$\uparrow \text{name} + \text{name} \times \text{name}$
1	$Expr + Term$	$\uparrow \text{name} + \text{name} \times \text{name}$
1	$Expr + Term + Term$	$\uparrow \text{name} + \text{name} \times \text{name}$
1	$Expr + Term + Term + \dots$	$\uparrow \text{name} + \text{name} \times \text{name}$
1		$\uparrow \text{name} + \text{name} \times \text{name}$
		$+ \text{name} \times \text{name}$

A top-down parser can loop indefinitely with left-recursive grammar

# Left Recursion

A grammar is left-recursive if it has a nonterminal  $A$  such that there is a derivation  $A \xRightarrow{+} A\alpha$  for some string  $\alpha$

**Direct** There is a production of the form  $A \rightarrow A\alpha$

**Indirect** The first symbol on the right-hand side of a rule can derive the symbol on the left

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Sd \mid \epsilon$$

We can often reformulate a grammar to avoid left recursion

# Remove Direct Left Recursion

Grammar with left recursion

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

Grammar without left recursion

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

Example

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

$\Rightarrow$

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

# Non-Left-Recursive Expression Grammar

## Expression Grammar with Recursion

Rule #	Production
0	$Start \rightarrow Expr$
1	$Expr \rightarrow Expr + Term$
2	$Expr \rightarrow Expr - Term$
3	$Expr \rightarrow Term$
4	$Term \rightarrow Term \times Factor$
5	$Term \rightarrow Term \div Factor$
6	$Term \rightarrow Factor$
7	$Factor \rightarrow (Expr)$
8	$Factor \rightarrow \mathbf{num}$
9	$Factor \rightarrow \mathbf{name}$

## Expression Grammar without Recursion

Rule #	Production
0	$Start \rightarrow Expr$
1	$Start \rightarrow Term Expr'$
2	$Expr' \rightarrow +Term Expr'$
3	$Expr' \rightarrow -Term Expr'$
4	$Expr' \rightarrow \epsilon$
5	$Term \rightarrow Factor Term'$
6	$Term \rightarrow \times Factor Term'$
7	$Term \rightarrow \div Factor Term'$
8	$Term' \rightarrow \epsilon$
9	$Factor \rightarrow (Expr)$
10	$Factor \rightarrow \mathbf{num}$
11	$Factor \rightarrow \mathbf{name}$



# Eliminating Indirect Left Recursion

- **Input:** Grammar  $G$  with no cycles or  $\epsilon$ -productions
- **Algorithm:**

```
Arrange nonterminals in some order  $A_1, A_2, \dots, A_n$ 
for  $i \leftarrow 1 \dots n$ 
  for  $j \leftarrow 1 \dots i-1$ 
    if  $\exists$  a production  $A_i \rightarrow A_j \gamma$ 
      Replace  $A_i \rightarrow A_j \gamma$  with one or more productions that expand  $A_j$ 
  Eliminate the immediate left recursion among the  $A_i$  productions
```

Loop invariant at the start of the outer iteration  $i$

$\forall k < i$ , no production expanding  $A_k$  has  $A_i$  in its body (i.e., right-hand side) for all  $l < k$

The algorithm establishes a topological ordering on nonterminals

# Eliminating Indirect Left Recursion

- **Input:** Grammar  $G$  with no cycles or  $\epsilon$ -productions
- **Algorithm:**

```
Arrange nonterminals in some order  $A_1, A_2, \dots, A_n$ 
for  $i \leftarrow 1 \dots n$ 
  for  $j \leftarrow 1 \dots i-1$ 
    if  $\exists$  a production  $A_i \rightarrow A_j \gamma$ 
      Replace  $A_i \rightarrow A_j \gamma$  with one or more productions that expand  $A_j$ 
  Eliminate the immediate left recursion among the  $A_i$  productions
```

$$S \rightarrow Aa \mid b$$
$$\Rightarrow$$
$$S \rightarrow Aa \mid b$$
$$A \rightarrow Ac \mid Sd \mid \epsilon$$
$$A \rightarrow bdA' \mid A'$$
$$A' \rightarrow cA' \mid adA' \mid \epsilon$$

# Implementing Backtracking

- A top-down parser may need to undo its actions after it detects a mismatch between the parse tree's leaves and the input
  - ▶ Implies a possible expansion with a wrong production
- Steps in backtracking
  - ▶ Set *curr* to parent and delete the children
  - ▶ Expand the node *curr* with **untried rules** if any
    - ▶ Create child nodes for each symbol in the right hand of the production
    - ▶ Push those symbols onto the stack in reverse order
    - ▶ Set *curr* to the first child node
  - ▶ **Move** *curr* **up the tree** if there are no untried rules
  - ▶ Report a syntax error when there are no more moves

# Backtracking is Expensive

- (i) Parser expands a nonterminal with the wrong rule
- (ii) Mismatch between the lower fringe of the parse tree and the input is detected
- (iii) Parser undoes the last few actions
- (iv) Parser tries other productions (if any)

A large subset of CFGs can be parsed without backtracking

The grammar may require transformations

# Avoid Backtracking

- Parser is to select the next rule
  - ▶ Compare the `curr` symbol and the next input symbol called the lookahead
  - ▶ Use the lookahead to disambiguate the possible production rules
- Intuition
  - ▶ Each alternative for the leftmost nonterminal leads to a **distinct terminal** symbol
  - ▶ Which rules to choose becomes obvious by comparing the next word in the input stream

## Definition

Backtrack-free grammar (also called predictive grammar) is a CFG for which a leftmost, top-down parser can always predict the correct rule with a one-word lookahead

# FIRST Set

## Definition

Given a string  $\gamma$  of terminal and nonterminal symbols,  $\text{FIRST}(\gamma)$  is the set of all terminal symbols that can begin any string derived from  $\gamma$

- We also need to keep track of which symbols can produce the empty string
- $\text{FIRST} : (NT \cup T \cup \{\epsilon, \text{EOF}\}) \rightarrow (T \cup \{\epsilon, \text{EOF}\})$

### ● Steps to compute FIRST set

1. If  $X$  is a terminal, then  $\text{FIRST}(X) = \{X\}$
2. If  $X \rightarrow \epsilon$  is a production, then  $\epsilon \in \text{FIRST}(X)$
3. If  $X$  is a nonterminal and  $X \rightarrow Y_1 Y_2 \dots Y_k$  is a production, then
  - (i)  $\text{FIRST}(X) = \text{FIRST}(Y_1)$  provided  $Y_1 \not\rightarrow \epsilon$
  - (ii) If for some  $i \leq k$  and  $1 \leq j < i$ ,  $a \in \text{FIRST}(Y_i)$ , and  $\forall j, \epsilon \in \text{FIRST}(Y_j)$ , then  $a \in \text{FIRST}(X)$
  - (iii) If  $\epsilon \in \text{FIRST}(Y_1, \dots, Y_k)$ , then  $\epsilon \in \text{FIRST}(X)$

### ● Generalization of FIRST relation to string of symbols

$$\text{FIRST}(X\gamma) = \text{FIRST}(X) \quad \text{if } X \not\rightarrow \epsilon$$

$$\text{FIRST}(X\gamma) = \text{FIRST}(X) \cup \text{FIRST}(\gamma) \quad \text{if } X \rightarrow \epsilon$$

# Example of FIRST Set Computation

## Grammar

$Start \rightarrow Expr$

$Expr \rightarrow Term Expr'$

$Expr' \rightarrow + Term Expr' \mid - Term Expr' \mid \epsilon$

$Term \rightarrow Factor Term'$

$Term' \rightarrow \times Factor Term' \mid \div Factor Term' \mid \epsilon$

$Factor \rightarrow (Expr) \mid \mathbf{num} \mid \mathbf{name}$

## FIRST Sets

$FIRST(Start) = \{\mathbf{name}, \mathbf{num}, ()\}$

$FIRST(Expr) = \{\mathbf{name}, \mathbf{num}, ()\}$

$FIRST(Expr') = \{+, -, \epsilon\}$

$FIRST(Term) = \{\mathbf{name}, \mathbf{num}, ()\}$

$FIRST(Term') = \{\times, \div, \epsilon\}$

$FIRST(Factor) = \{\mathbf{name}, \mathbf{num}, ()\}$

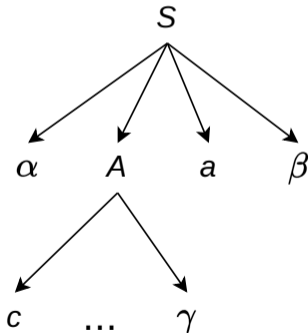
How does a parser decide when to apply the  $\epsilon$ -production?

# FOLLOW Set

## Definition

FOLLOW ( $X$ ) is the set of terminals that can immediately follow  $X$

- That is,  $t \in \text{FOLLOW}(X)$  if there is any derivation containing  $Xt$



Terminal  $c$  is in FIRST ( $A$ ) and  $a$  is in FOLLOW ( $A$ )



# Steps to Compute FOLLOW Set

- (i) Place \$ in FOLLOW ( $S$ ) where  $S$  is the start symbol and the \$ is the end marker
- (ii) If there is a production  $A \rightarrow \alpha B \beta$ , then everything in FIRST ( $\beta$ ) except  $\epsilon$  is in FOLLOW ( $B$ )
- (iii) If there is a production  $A \rightarrow \alpha B$ , or a production  $A \rightarrow \alpha B \beta$  where FIRST ( $\beta$ ) contains  $\epsilon$ , then everything in FOLLOW ( $A$ ) is in FOLLOW ( $B$ )

# Example of FOLLOW Set Computation

## Grammar

$Start \rightarrow Expr$

$Expr \rightarrow Term Expr'$

$Expr' \rightarrow + Term Expr' \mid - Term Expr' \mid \epsilon$

$Term \rightarrow Factor Term'$

$Term' \rightarrow \times Factor Term' \mid \div Factor Term' \mid \epsilon$

$Factor \rightarrow (Expr) \mid \mathbf{num} \mid \mathbf{name}$

## FOLLOW Sets

$FOLLOW(Start) = \{\$\}$

$FOLLOW(Expr) = \{\$, )\}$

$FOLLOW(Expr') = \{\$, )\}$

$FOLLOW(Term) = \{\$, +, -, )\}$

$FOLLOW(Term') = \{\$, +, -, )\}$

$FOLLOW(Factor) = \{\$, +, -, \times, \div, )\}$

# Conditions for Backtrack-Free Grammar

- Consider a production  $A \rightarrow \beta$

$$\text{FIRST}^+(A \rightarrow \beta) = \begin{cases} \text{FIRST}(\beta) & \text{if } \epsilon \notin \text{FIRST}(\beta) \\ \text{FIRST}(\beta) \cup \text{FOLLOW}(A) & \text{otherwise} \end{cases}$$

- For any nonterminal  $A$  where  $A \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$ , a **backtrack-free grammar** has the property

$$\text{FIRST}^+(A \rightarrow \beta_i) \cap \text{FIRST}^+(A \rightarrow \beta_j) = \phi, \quad \forall 1 \leq i, j \leq n, i \neq j$$

Expression grammar on the previous slide is backtrack-free

# Not All Grammars are Backtrack-Free

$Start \rightarrow Expr$

$Expr \rightarrow Term Expr'$

$Expr' \rightarrow + Term Expr' \mid - Term Expr' \mid \epsilon$

$Term \rightarrow Factor Term'$

$Term' \rightarrow \times Factor Term' \mid \div Factor Term' \mid \epsilon$

$Factor \rightarrow (Expr) \mid \mathbf{num} \mid \mathbf{name}$

$Factor \rightarrow \mathbf{name} \mid \mathbf{name}[Arglist] \mid \mathbf{name} (Arglist)$

$Arglist \rightarrow Expr MoreArgs$

$MoreArgs \rightarrow , Expr MoreArgs \mid \epsilon$

# Not All Grammars are Backtrack-Free

$Start \rightarrow Expr$

$Expr \rightarrow Term Expr'$

$Expr' \rightarrow + Term Expr' \mid - Term Expr' \mid \epsilon$

$Term \rightarrow Factor Term'$

$Term' \rightarrow \times Factor Term' \mid \div Factor Term' \mid \epsilon$

$Factor \rightarrow (Expr) \mid \mathbf{num} \mid \mathbf{name}$

$Factor \rightarrow \mathbf{name} \mid \mathbf{name}[Arglist] \mid \mathbf{name} (Arglist)$

$Arglist \rightarrow Expr MoreArgs$

$MoreArgs \rightarrow , Expr MoreArgs \mid \epsilon$

Given a finite lookahead, we can always devise a non-backtrack-free grammar such that the lookahead is insufficient

# Left Factoring

## Definition

Left factoring is the process of extracting and isolating common prefixes in a set of productions

- **Algorithm:**

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma_1 \mid \dots \mid \gamma_j$$



$$\begin{aligned} A &\rightarrow \alpha B \mid \gamma_1 \mid \gamma_2 \dots \mid \gamma_j \\ B &\rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n \end{aligned}$$

# Summarizing Top-down Parsing

- Efficiency depends on the accuracy of selecting the correct production for expanding a nonterminal
  - ▶ Parser may not terminate in the worst-case
- A large subset of the context-free grammars can be parsed without backtracking

# Recursive-Descent Parsing



# Recursive-Descent Parsing

- Recursive-descent parsing is a form of top-down parsing that **may require** backtracking
  - ▶ Top-down approach is modeled by calls to functions, where there is one function for each nonterminal

```
void A() {  
    Choose an A-production  $A \rightarrow X_1X_2\dots X_k$   
    for  $i \leftarrow 1\dots k$   
        if  $X_i$  is a nonterminal  
            call function  $X_i$   
        else if  $X_i$  equals the current input symbol  $a$   
            advance the input to the next symbol  
        else  
            // error  
}
```

# Recursive-Descent Parsing with Backtracking

- Consider a grammar with two productions  $X \rightarrow \gamma_1$  and  $X \rightarrow \gamma_2$
- Suppose  $\text{FIRST}(\gamma_1) \cap \text{FIRST}(\gamma_2) \neq \phi$ 
  - ▶ Let us denote one of the common terminal symbols by  $a$
- The function for  $X$  will not know which production to use on the input token  $a$
- To support backtracking
  - ▶ All productions should be tried in some order
  - ▶ Failure for some production implies the parser needs to try the remaining productions
  - ▶ Report an error only when there are no other rules

# Predictive Parsing

## Definition

Predictive parsing is a special case of recursive-descent parsing that does not require backtracking

- Lookahead symbol unambiguously determines which production rule to use
- Advantage is that the algorithm is simple and the parser can be constructed by hand

$$\begin{aligned} \textit{stmt} &\rightarrow \mathbf{expr}; \\ &| \mathbf{if} (\textit{expr}) \textit{stmt} \\ &| \mathbf{for} (\textit{optexpr}; \textit{optexpr}; \textit{optexpr}) \textit{stmt} \\ &| \mathbf{other} \\ \textit{optexpr} &\rightarrow \mathbf{expr} \mid \epsilon \end{aligned}$$

# Pseudocode for a Predictive Parser

```
void stmt() {
    switch(lookahead) {
        case expr: { match(expr); match(';'); break; }
        case if: {
            match(if); match('('); match(expr); match(')'); stmt(); break;
        }
        case for: {
            match(for); match('('); optexpr(); match(';'); optexpr(); match(';');
            optexpr(); match(')'); stmt(); break;
        }
        case other: { match(other); break; }
        default: { print("syntax error"); }
    }
}
```

# Non-Recursive Predictive Parsing

# LL(k) Grammars

## Definition

A CFG  $G = (T, NT, S, P)$  is LL(1) if and only if for every nonterminal  $A \in NT$  where  $A \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$  such that  $\beta_i \in \Sigma^*$ , we have

$$\text{FIRST}^+(A \rightarrow \beta_i) \cap \text{FIRST}^+(A \rightarrow \beta_j) = \phi, \quad \forall 1 \leq i, j \leq n, i \neq j$$

- First L stands for left-to-right scan, second L stands for leftmost derivation, and  $k$  represents the number of lookahead tokens
- LL(k) grammars are the class of CFGs for which no backtracking is required
  - ▶ Predictive parsers accept LL(k) grammars
- Every LL(1) grammar is a LL(2) grammar
- Many programming language constructs are LL(1)

# Definition of LL(k) Grammar

- For a given word  $w \in T^*$  and non-negative integer  $k$ ,
  - ▶  $w/k$  is  $w$  if  $|w| \leq k$ , or
  - ▶  $w/k$  is a string consisting of the first  $k$  symbols of  $w$  if  $|w| > k$ .
- A CFG  $G = (T, NT, S, P)$  is LL(k) for some positive integer  $k$  if and only if given
  - (i) a word  $w \in T^*$  such that  $|w| \leq k$ ,
  - (ii) a nonterminal  $A \in NT$ , and
  - (iii) a word  $w_1 \in T^*$ ,there is at most one production  $p \in P$  such that for some  $w_2, w_3 \in T^*$ ,
  1.  $S \Rightarrow w_1Aw_3$ ,
  2.  $A \xRightarrow{+} w_2$  by first applying production  $p$ ,
  3.  $w_2w_3/k = w$ .

# Definition of LL(k) Grammar

- For a given word  $w \in T^*$  and non-negative integer  $k$ ,
  - ▶  $w/k$  is  $w$  if  $|w| \leq k$ , or
  - ▶  $w/k$  is a string consisting of the first  $k$  symbols of  $w$  if  $|w| > k$ .
- A CFG  $G = (T, NT, S, P)$  is LL(k) for some positive integer  $k$  if and only if given

Stated informally in terms of parsing, an LL(k) grammar is a CFG such that for any word in its language, each production in its derivation can be identified with certainty by inspecting the word from its beginning (left end) to the  $k^{\text{th}}$  symbol beyond the beginning of the production.

2.  $A \xRightarrow{+} w_2$  by first applying production  $p$ ,
3.  $w_2 w_3 / k = w$ .



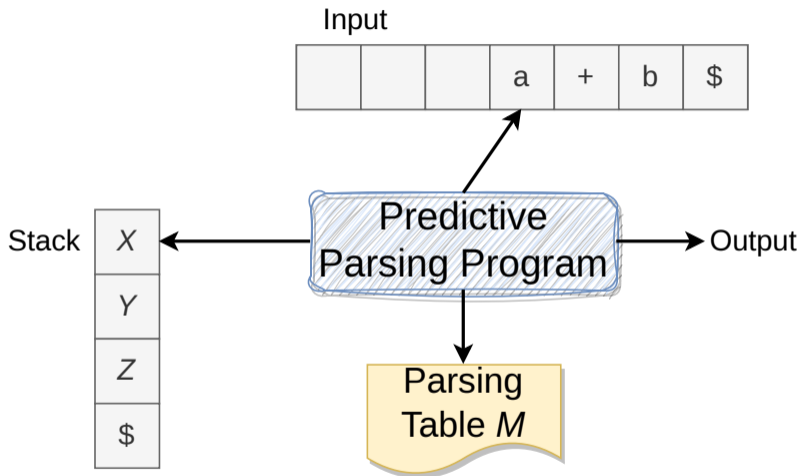
# Example LL(2) Parser

$$S \rightarrow AXQ \mid AYR$$

lookahead is the set of 2-sequence tokens that indicate which alternative will succeed

```
void S() {  
    if (lookahead(1) == A && lookahead(2) == X) {  
        match(A); match(X); match(Q);  
    } else if (lookahead(1) == A && lookahead(2) == Y) {  
        match(A); match(Y); match(R);  
    } else {  
        // Raise error  
    }  
}
```

# Nonrecursive Table-Driven LL(1) Parser



# LL(1) Parsing Algorithm

- **Input:** String  $w$  and parsing table  $M$  for grammar  $G$
- **Output:** A leftmost derivation of  $w$  if  $w \in L(G)$ ; otherwise, report an error
- **Algorithm:**

```
Let  $a$  be the first symbol in  $w$ 
Let  $X$  be the symbol at the top of the stack
while  $X \neq \$$ 
  if  $X == a$ 
    pop the stack and advance the input
  else if  $X$  is a terminal or  $M[X,a]$  is an error entry
    report error
  else if  $M[X,a] == X \rightarrow Y_1Y_2\dots Y_k$ 
    // Expand with the production  $X \rightarrow Y_1Y_2\dots Y_k$ 
    pop the stack
    // Simulate depth-first traversal
    push  $Y_kY_{k-1}\dots Y_1$  onto the stack
   $X \leftarrow$  top stack symbol
```

# Construction of a LL(1) Parsing Table

- **Input:** Grammar  $G$
- **Algorithm:**

```
for each production  $A \rightarrow \alpha$  in  $G$ 
  for each terminal  $a$  in  $FIRST(\alpha)$ 
    add  $A \rightarrow \alpha$  to  $M[A, a]$ 
  if  $\epsilon \in FIRST(\alpha)$ 
    for each terminal  $b$  in  $FOLLOW(A)$ 
      add  $A \rightarrow \alpha$  to  $M[A, b]$ 
  if  $\epsilon \in FIRST(\alpha)$  and  $\$ \in FOLLOW(A)$ 
    add  $A \rightarrow \alpha$  to  $M[A, \$]$ 

// No production in  $M[A, a]$  indicates error
```

# LL(1) Parsing Table

## Grammar

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

## FIRST Sets

$$\text{FIRST}(E) = \{\mathbf{id}, (\}$$

$$\text{FIRST}(E') = \{+, \epsilon\}$$

$$\text{FIRST}(T) = \{\mathbf{id}, (\}$$

$$\text{FIRST}(T') = \{*, \epsilon\}$$

$$\text{FIRST}(F) = \{\mathbf{id}, (\}$$

## FOLLOW Sets

$$\text{FOLLOW}(E) = \{\$, )\}$$

$$\text{FOLLOW}(E') = \{\$, )\}$$

$$\text{FOLLOW}(T) = \{\$, +, )\}$$

$$\text{FOLLOW}(T') = \{\$, +, )\}$$

$$\text{FOLLOW}(F) = \{\$, +, *, )\}$$

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

# Working of a LL(1) Parser

Stack	Input	Remark
$\$E$	$\uparrow \text{id} + \text{id} * \text{id}\$$	Expand $E \rightarrow TE'$
$\$E'T$	$\uparrow \text{id} + \text{id} * \text{id}\$$	Expand $T \rightarrow FT'$
$\$E'T'F$	$\uparrow \text{id} + \text{id} * \text{id}\$$	Expand $F \rightarrow \text{id}$
$\$E'T'\text{id}$	$\uparrow \text{id} + \text{id} * \text{id}\$$	Match <b>id</b>
$\$E'T'$	$\uparrow + \text{id} * \text{id}\$$	Expand $T \rightarrow \epsilon$
$\$E'$	$\uparrow + \text{id} * \text{id}\$$	Expand $E' \rightarrow +TE'$
$\$E'T+$	$\uparrow + \text{id} * \text{id}\$$	Match <b>+</b>
$\$E'T$	$\uparrow \text{id} * \text{id}\$$	Expand $T \rightarrow FT'$
$\$E'T'F$	$\uparrow \text{id} * \text{id}\$$	Expand $F \rightarrow \text{id}$
$\$E'T'\text{id}$	$\uparrow \text{id} * \text{id}\$$	Match <b>id</b>
$\$E'T'$	$\uparrow * \text{id}\$$	Expand $T' \rightarrow *FT'$
$\$E'T'F*$	$\uparrow * \text{id}\$$	Match <b>*</b>
$\$E'T'F$	$\uparrow \text{id}\$$	Expand $F \rightarrow \text{id}$
$\$E'T'\text{id}$	$\uparrow \text{id}\$$	Match <b>id</b>
$\$E'T'$	$\uparrow \$$	Expand $T' \rightarrow \epsilon$
$\$E'$	$\uparrow \$$	Expand $E' \rightarrow \epsilon$
$\$$	$\uparrow \$$	

## More on LL(1) Parsing

- Grammars whose predictive parsing tables contain no duplicate entries are LL(1)
- No left-recursive or ambiguous grammar can be LL(1)
  - ▶ If grammar  $G$  is left-recursive or is ambiguous, then parsing table  $M$  will have at least one multiply-defined cell
- Some grammars cannot be transformed into LL(1)

The below grammar is ambiguous

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow eS \mid \epsilon$$

$$E \rightarrow b$$

# Limitations with LL(k) Parsing

LL(k) cannot see past **arbitrarily** long constructs from the left edge

$$S \rightarrow A+XQ \mid A+YR$$

Could left factor, but not always possible and natural

$$S \rightarrow A+(XQ \mid YR)$$

Programming language grammars may not be LL(k) (e.g., C function declaration vs definition)

$$\begin{aligned} \text{func} &\rightarrow \text{type ID '(' arg* ')' ';' } \\ &\rightarrow \text{type ID '(' arg* ')' '{' body '}' } \end{aligned}$$



# Using Ambiguous Grammars

# LL(1) Parsing Table for an Ambiguous Grammar

## Grammar

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow eS \mid \epsilon$$

$$E \rightarrow b$$

## FIRST Sets

$$\text{FIRST}(S) = \{i, a\}$$

$$\text{FIRST}(S') = \{e, \epsilon\}$$

$$\text{FIRST}(E) = \{b\}$$

## FOLLOW Sets

$$\text{FOLLOW}(S) = \{\$, e\}$$

$$\text{FOLLOW}(S') = \{\$, e\}$$

$$\text{FOLLOW}(E) = \{t\}$$

Nonterminal	<i>a</i>	<i>b</i>	<i>e</i>	<i>i</i>	<i>t</i>	$\$$
<i>S</i>	$S \rightarrow a$			$S \rightarrow iEtSS'$		
<i>S'</i>			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
<i>E</i>		$E \rightarrow b$				

# Detecting Errors in Table-Driven Predictive Parsing

## Error conditions

- (i) Terminal on top of the stack does not match the next input symbol
- (ii) Nonterminal  $A$  is on top of the stack,  $a$  is the next input symbol, and  $M[A, a]$  is empty

## Choices

- (i) Raise an error and quit parsing
- (ii) Print an error message, try to recover from the error, and continue with the compilation

# Error Recovery in Table-Driven Predictive Parsing

Assume  $A$  is the nonterminal at the top of the stack

Panic mode recovery skips over symbols until a token in a set of synchronizing (synch) tokens is found

- (i) Add all tokens in FOLLOW ( $A$ ) to the synch set for  $A$ 
  - ▶ Parsing can continue if the parser skips all input tokens until it sees an input symbol in FOLLOW ( $A$ )
- (ii) Add symbols in FIRST ( $A$ ) to the synch set for  $A$ 
  - ▶ Parsing can continue with  $A$  if the parser skips all input tokens until it sees an input symbol in FIRST ( $A$ )
- (iii) Add keywords that begin constructs
- (iv) Skip input if the table does not have an entry
- (v) ...

# Using FOLLOW Sets as Synchronizing Tokens

## Grammar

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

## FOLLOW Sets

$$\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{\$, \})$$

$$\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{\$, +, \})$$



$$\text{FOLLOW}(F) = \{\$, +, \times, \})$$

Nonterminal	id	+	*	(	)	\$	
$E$	$E \rightarrow TE'$				$E \rightarrow TE'$	synch	synch
$E'$		$E' \rightarrow +TE'$				$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
$T$	$T \rightarrow FT'$	synch			$T \rightarrow FT'$	synch	synch
$T'$		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$			$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
$F$	$F \rightarrow \mathbf{id}$	synch	synch		$F \rightarrow (E)$	synch	synch

# Error Recovery Moves by Table-Driven Predictive Parser

Stack	Input	Remark
$\$E$	$+id * +id\$$	Error, skip +
$\$E$	$id * +id\$$	Expand $E \rightarrow TE'$
$\$E' T$	$id * +id\$$	Expand $T \rightarrow FT'$
$\$E' T' F$	$id * +id\$$	Expand $F \rightarrow id$
$\$E' T' id$	$id * +id\$$	Match $id$
$\$E' T'$	$* +id\$$	Expand $T \rightarrow *FT'$
$\$E' T' F*$	$* +id\$$	Match *
$\$E' T' F$	$+id\$$	Error, $M[F, +] = \text{synch}$ , pop $F$
$\$E' T'$	$+id\$$	Expand $T \rightarrow \epsilon$
$\$E'$	$+id\$$	Expand $E' \rightarrow +TE'$
$\$E' T+$	$+id\$$	Match +
$\$E' T$	$id\$$	Expand $T \rightarrow FT'$
$\$E' T' F$	$id\$$	Expand $F \rightarrow id$
$\$E' T' id$	$id\$$	Match $id$
$\$E' T'$	$\$$	Expand $T' \rightarrow \epsilon$
$\$E'$	$\$$	Expand $E' \rightarrow \epsilon$
$\$$	$\$$	

# References

-  A. Aho et al. *Compilers: Principles, Techniques, and Tools*. Sections 2.4, 4.2–4.4, 2<sup>nd</sup> edition, Pearson Education.
-  K. Cooper and L. Torczon. *Engineering a Compiler*. Section 3.3, 2<sup>nd</sup> edition, Morgan Kaufmann.