# 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

$Start \rightarrow Expr$

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

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

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

priority ↓

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

| Sentential Form | Input |
|---|---|
| *Expr* | ↑ **name** + **name** × **name** |
| *Expr* + *Term* | ↑ **name** + **name** × **name** |
| *Term* + *Term* | ↑ **name** + **name** × **name** |
| *Factor* + *Term* | ↑ **name** + **name** × **name** |
| **name** + *Term* | ↑ **name** + **name** × **name** |
| **name** + *Term* | **name** ↑ + **name** × **name** |
| **name** + *Term* | **name** + ↑ **name** × **name** |
| **name** + *Term* × *Factor* | **name** + ↑ **name** × **name** |
| **name** + *Factor* × *Factor* | **name** + ↑ **name** × **name** |
| **name** + **name** × *Factor* | **name** + ↑ **name** × **name** |
| **name** + **name** × *Factor* | **name** + **name** ↑ × **name** |
| **name** + **name** × *Factor* | **name** + **name** × ↑ **name** |
| **name** + **name** × **name** | **name** + **name** × ↑ **name** |
| **name** + **name** × **name** | **name** + **name** × **name** ↑ |

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

| Sentential Form | Input |
|---|---|
| *Expr* | ↑ **name** + **name** × **name** |
| *Expr* + *Term* | ↑ **name** + **name** × **name** |
| *Term* + *Term* | ↑ **name** + **name** × **name** |
| *Factor* + *Term* | ↑ **name** + **name** × **name** |
| **name** + *Term* | ↑ **name** + **name** × **name** |

> The current input terminal being scanned is called the lookahead symbol

| | |
|---|---|
| **name** + *Factor* × *Factor* | **name** + ↑ **name** × **name** |
| **name** + **name** × *Factor* | **name** + ↑ **name** × **name** |
| **name** + **name** × *Factor* | **name** + **name** ↑ × **name** |
| **name** + **name** × *Factor* | **name** + **name** × ↑ **name** |
| **name** + **name** × **name** | **name** + **name** × ↑ **name** |
| **name** + **name** × **name** | **name** + **name** × **name** ↑ |

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

$$Start \overset{lm}{\Rightarrow} Expr \overset{lm}{\Rightarrow} Expr \overset{lm}{\Rightarrow} Expr \overset{lm}{\Rightarrow} Expr \overset{lm}{\Rightarrow} Expr \overset{lm}{\Rightarrow} Expr$$

# Derivation of **name** + **name** $\times$ **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 → β₁β₂...βₙ to expand curr
    create nodes for β₁,β₂,...βₙ as children of curr
    push(βₙβₙ₋₁...β₁) // reverse order
    curr = β₁
  if curr == word
    word = getNextWord()
    curr = pop() // Consumed
  if word == EOF and curr == null
    accept input
  else
    backtrack
```

# Derivation of **name** + **name** × **name**

| Rule # | Production |
|--------|------------|
| 0 | *Start → Expr* |
| 1 | *Expr → Expr + Term* |
| 2 | *Expr → Expr − Term* |
| 3 | *Expr → Term* |
| 4 | *Term → Term × Factor* |
| 5 | *Term → Term ÷ Factor* |
| 6 | *Term → Factor* |
| 7 | *Factor → (Expr)* |
| 8 | *Factor → **num*** |
| 9 | *Factor → **name*** |

| Rule # | Sentential Form | Input |
|--------|-----------------|-------|
|   | *Expr* | ↑ **name** + **name** × **name** |
| 1 | *Expr + Term* | ↑ **name** + **name** × **name** |
| 3 | *Term + Term* | ↑ **name** + **name** × **name** |
| 6 | *Factor + Term* | ↑ **name** + **name** × **name** |
| 9 | **name** + *Term* | ↑ **name** + **name** × **name** |
|   | **name** + *Term* | **name** ↑ + **name** × **name** |
|   | **name** + *Term* | **name** + ↑ **name** × **name** |
| 4 | **name** + *Term × Factor* | **name** + ↑ **name** × **name** |
| 4 | **name** + *Factor × Factor* | **name** + ↑ **name** × **name** |
| 9 | **name** + **name** × *Factor* | **name** + ↑ **name** × **name** |
|   | **name** + **name** × *Factor* | **name** + **name** ↑ × **name** |
|   | **name** + **name** × *Factor* | **name** + **name** × ↑ **name** |
| 9 | **name** + **name** × **name** | **name** + **name** × ↑ **name** |
|   | **name** + **name** × **name** | **name** + **name** × **name** ↑ |

# Derivation of **name** + **name** × **name**

| Rule # | Production |
|--------|-----------|
| 0 | *Start → Expr* |
| 1 | *Expr → Expr + Term* |
| 2 | *Expr → Expr – Term* |
| 3 | *Expr → Term* |
| 4 | *Term → Term × Factor* |
| 5 | *Term* |
| 6 | *Term* |
| 7 | *Fac* |
| 8 | *Factor →* **num** |
| 9 | *Factor →* **name** |

| Rule # | Sentential Form | Input |
|--------|-----------------|-------|
|  | *Expr* | ↑ **name** + **name** × **name** |
| 1 | *Expr + Term* | ↑ **name** + **name** × **name** |
| 3 | *Term + Term* | ↑ **name** + **name** × **name** |
| 6 | *Factor + Term* | ↑ **name** + **name** × **name** |
| 9 | **name** + *Term* | ↑ **name** + **name** × **name** |
|  |  | + **name** × **name** |
|  |  | ↑ **name** × **name** |
|  |  | ↑ **name** × **name** |
| 4 | **name** + *Factor* × *Factor* | **name** + ↑ **name** × **name** |
| 9 | **name** + **name** × *Factor* | **name** + ↑ **name** × **name** |
|  | **name** + **name** × *Factor* | **name** + **name** ↑ × **name** |
|  | **name** + **name** × *Factor* | **name** + **name** × ↑ **name** |
| 9 | **name** + **name** × **name** | **name** + **name** × ↑ **name** |
|  | **name** + **name** × **name** | **name** + **name** × **name** ↑ |

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

# Deterministically Selecting a Production in Expression Grammar

| Rule # | Production |
|--------|------------|
| 0 | *Start → Expr* |
| 1 | *Expr → Expr + Term* |
| 2 | *Expr → Expr − Term* |
| 3 | *Expr → Term* |
| 4 | *Term → Term × Factor* |
| 5 | *Term → Term ÷ Factor* |
| 6 | *Term → Factor* |
| 7 | *Factor → (Expr)* |
| 8 | *Factor →* **num** |
| 9 | *Factor →* **name** |

| Rule # | Sentential Form | Input |
|--------|-----------------|-------|
| | *Expr* | ↑ **name** + **name** × **name** |
| 1 | *Expr + Term* | ↑ **name** + **name** × **name** |
| 1 | *Expr + Term + Term* | ↑ **name** + **name** × **name** |
| 1 | *Expr + Term + Term + . . .* | ↑ **name** + **name** × **name** |
| 1 | . . . | ↑ **name** + **name** × **name** |
| 1 | . . . | ↑ **name** + **name** × **name** |

# Deterministically Selecting a Production in Expression Grammar

| Rule # | Production |
|--------|------------|
| 0 | *Start → Expr* |
| 1 | *Expr → Expr + Term* |
| 2 | *Expr → Expr − Term* |
| 3 | *Expr → Term* |
| 4 | *Term → Term × Factor* |
| 5 | *Term* |
| 6 | *Term* |
| 7 | *Fac* |
| 8 | *Factor → num* |
| 9 | *Factor → name* |

| Rule # | Sentential Form | Input |
|--------|-----------------|-------|
| | *Expr* | ↑ **name** + **name** × **name** |
| 1 | *Expr + Term* | ↑ **name** + **name** × **name** |
| 1 | *Expr + Term + Term* | ↑ **name** + **name** × **name** |
| 1 | *Expr + Term + Term + . . .* | ↑ **name** + **name** × **name** |
| 1 | | ↑ **name** + **name** × **name** |
| | | + **name** × **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

$$A \rightarrow A\alpha_1 \,|\, A\alpha_2 \,|\, \ldots \,|\, A\alpha_m \,|\, \beta_1 \,|\, \ldots \,|\, \beta_n$$

### Grammar without left recursion

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

### Example

$$E \rightarrow E + T \,|\, T$$
$$T \rightarrow T * F \,|\, F$$
$$F \rightarrow (E) \,|\, \textbf{id}$$

$\Longrightarrow$

$$E \rightarrow TE'$$
$$E' \rightarrow +TE'$$
$$T \rightarrow FT'$$
$$T' \rightarrow *FT'$$
$$F \rightarrow (E) \,|\, \textbf{id}$$

# Non-Left-Recursive Expression Grammar

## Expression Grammar with Recursion

| Rule # | Production |
|--------|-----------|
| 0 | *Start* → *Expr* |
| 1 | *Expr* → *Expr* + *Term* |
| 2 | *Expr* → *Expr* − *Term* |
| 3 | *Expr* → *Term* |
| 4 | *Term* → *Term* × *Factor* |
| 5 | *Term* → *Term* ÷ *Factor* |
| 6 | *Term* → *Factor* |
| 7 | *Factor* → (*Expr*) |
| 8 | *Factor* → **num** |
| 9 | *Factor* → **name** |

## Expression Grammar without Recursion

| Rule # | Production |
|--------|-----------|
| 0 | *Start* → *Expr* |
| 1 | *Start* → *Term* *Expr*$'$ |
| 2 | *Expr*$'$ → +*Term* *Expr*$'$ |
| 3 | *Expr*$'$ → −*Term* *Expr*$'$ |
| 4 | *Expr*$'$ → $\epsilon$ |
| 5 | *Term* → *Factor* *Term*$'$ |
| 6 | *Term* → ×*Factor* *Term*$'$ |
| 7 | *Term* → ÷*Factor* *Term*$'$ |
| 8 | *Term*$'$ → $\epsilon$ |
| 9 | *Factor* → (*Expr*) |
| 10 | *Factor* → **num** |
| 11 | *Factor* → **name** |

# Eliminating Indirect Left Recursion

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

```
Arrange nonterminals in some order A₁, A₂,...Aₙ
for i ← 1...n
   for j ← 1...i − 1
      if ∃ a production Aᵢ → Aⱼγ
         Replace Aᵢ → Aⱼγ with one or more productions that expand Aⱼ
   Eliminate the immediate left recursion among the Aᵢ productions
```

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

$\forall k < i$, no production expanding $A_k$ has $A_l$ 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 A₁, A₂, ... Aₙ
for i ← 1...n
  for j ← 1...i-1
    if ∃ a production Aᵢ → Aⱼγ
       Replace Aᵢ → Aⱼγ with one or more productions that expand Aⱼ
  Eliminate the immediate left recursion among the Aᵢ productions
```

$$S \rightarrow Aa \mid b$$
$$A \rightarrow Ac \mid Sd \mid \epsilon$$

$$\Longrightarrow$$

$$S \rightarrow Aa \mid b$$
$$A \rightarrow bdA^{'} \mid A^{'}$$
$$A^{'} \rightarrow cA^{'} \mid adA^{'} \mid \epsilon$$

# 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, 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
- 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 FIRST $(X) = \{X\}$
  2. If $X \rightarrow \epsilon$ is a production, then $\epsilon \in$ FIRST $(X)$
  3. If $X$ is a nonterminal and $X \rightarrow Y_1 Y_2 \ldots Y_k$ is a production
     (i) Everything in FIRST $(Y_1)$ is in FIRST $(X)$
     (ii) If for some $i$, $a \in$ FIRST $(Y_i)$ and $\forall i \leq j < i, \epsilon \in$ FIRST $(Y_j)$, then $a \in$ FIRST $(X)$
     (iii) If $\epsilon \in$ FIRST $(Y_1, \ldots Y_k)$, then $\epsilon \in$ FIRST $(X)$
- Generalize FIRST relation to string of symbols

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

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

# Example of FIRST Set Computation

## Grammar

$$Start \rightarrow Expr$$
$$Expr \rightarrow Term\ Expr'$$
$$Expr' \rightarrow +\ Term\ Expr'\ |\ -\ Term\ Expr'\ |\ \epsilon$$
$$Term \rightarrow Factor\ Term'$$
$$Term' \rightarrow \times\ Factor\ Term'\ |\ \div\ Factor\ Term'\ |\ \epsilon$$
$$Factor \rightarrow (Expr)\ |\ \mathbf{num}\ |\ \mathbf{name}$$

## FIRST Sets

FIRST ($Start$) = {**name**, **num**, ( }

FIRST ($Expr$) = {**name**, **num**, ( }

FIRST $\left(Expr'\right)$ = {+, −, $\epsilon$}

FIRST ($Term$) = {**name**, **num**, ( }

FIRST $\left(Term'\right)$ = {×, ÷, $\epsilon$}

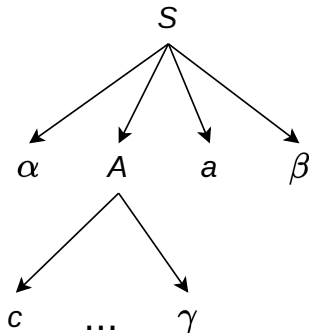FIRST ($Factor$) = {**name**, **num**, ( }

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, $t \in$ 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

$Start \rightarrow Expr$

$Expr \rightarrow Term\,Expr^{'}$

$Expr^{'} \rightarrow +\,Term\,Expr^{'} \mid -\,Term\,Expr^{'} \mid \epsilon$

$Term \rightarrow Factor\,Term^{'}$

$Term^{'} \rightarrow \times\,Factor\,Term^{'} \mid \div\,Factor\,Term^{'} \mid \epsilon$

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

## FOLLOW Sets

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

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

FOLLOW $\left(Expr^{'}\right) = \{\$, )\}$

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

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

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 $A \rightarrow \beta_1 | \beta_2 | \ldots | \beta_n$, a **backtrack-free grammar** has the property

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

Expression grammar on the previous slide is backtrack-free

# Not All Grammars are Backtrack-Free

$Start \rightarrow Expr$

$Expr \rightarrow Term\,Expr'$

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

$Term \rightarrow Factor\,Term'$

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

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

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

$Arglist \rightarrow Expr\,MoreArgs$

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

# Not All Grammars are Backtrack-Free

*Start* → *Expr*

*Expr* → *Term Expr'*

*Expr'* → + *Term Expr'* | − *Term Expr'* | $\epsilon$

*Term* → *Factor Term'*

*Term'* → × *Factor Term'* | ÷ *Factor Term'* | $\epsilon$

*Factor* → ( *Expr* ) | **num** | **name**

*Factor* → **name** | **name** [ *Arglist* ] | **name** ( *Arglist* )

*Arglist* → *Expr MoreArgs*

*MoreArgs* → , *Expr MoreArgs* | $\epsilon$

> 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**:

$$A \rightarrow \alpha\beta_1 \,|\, \alpha\beta_2 \,|\, \ldots \,|\alpha\beta_n \,|\, \gamma_1 \,|\, \ldots \,|\, \gamma_j$$

$$\Downarrow$$

$$A \rightarrow \alpha B \,|\, \gamma_1 \,|\, \gamma_2 \ldots \,|\, \gamma_j$$
$$B \rightarrow \beta_1 \,|\, \beta_2 \,|\, \ldots \,|\, \beta_n$$

# 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 → X₁X₂...Xₖ
  for i ← 1...k
    if Xᵢ is a nonterminal
      call function Xᵢ
    else if Xᵢ 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 $X \to \gamma_1$ and $X \to \gamma_2$
- Suppose FIRST $(\gamma_1) \cap$ 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

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

# 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"); }
  }
}
```

# LL(1) Grammars

## Definition
LL(*k*) grammars are the class of grammars for which no backtracking is required

- First L stands for left-to-right scan, second L stands for leftmost derivation
- There is one lookahead token in LL(1) and *k* lookahead tokens in LL(*k*)
- Predictive parsers accept LL(*k*) grammars
- Most programming language constructs are LL(1)
- Every LL(1) grammar is a LL(2) grammar

# 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 != $
  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 → Y₁Y₂...Yₖ
    // Expand with the production X → Y₁Y₂...Yₖ
    pop the stack
    // Simulate depth-first traversal
    push YₖYₖ₋₁...Y₁ onto the stack
  X ← top stack symbol
```

# Construction of a LL(1) Parsing Table

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

```
for each production A → α in G
  for each terminal a in FIRST (α)
    add A → α to M[A,a]
  if ε ∈ FIRST (α)
    for each terminal b in FOLLOW (A)
      add A → α to M[A,b]
  if ε ∈ FIRST (α) and $ ∈ FOLLOW (A)
      add A → α to M[A,$]

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

# LL(1) Parsing Table

## Grammar

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

## FIRST Sets

FIRST $(E) = \{\textbf{id}, ( \}$

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

FIRST $(T) = \{\textbf{id}, ( \}$

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

FIRST $(F) = \{\textbf{id}, ( \}$

## FOLLOW Sets

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

FOLLOW $\left(E'\right) = \{\$, )\}$

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

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

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

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

# Working of a LL(1) Parser

| Stack | Input | Remark |
|---|---|---|
| $\$E$ | ↑ **id** + **id** ∗ **id**$\$$ | Expand $E \rightarrow TE'$ |
| $\$E'T$ | ↑ **id** + **id** ∗ **id**$\$$ | Expand $T \rightarrow FT'$ |
| $\$E'T'F$ | ↑ **id** + **id** ∗ **id**$\$$ | Expand $F \rightarrow$ **id** |
| $\$E'T'$**id** | ↑ **id** + **id** ∗ **id**$\$$ | Match **id** |
| $\$E'T'$ | ↑ + **id** ∗ **id**$\$$ | Expand $T \rightarrow \epsilon$ |
| $\$E'$ | ↑ + **id** ∗ **id**$\$$ | Expand $E' \rightarrow +TE'$ |
| $\$E'T+$ | ↑ + **id** ∗ **id**$\$$ | Match + |
| $\$E'T$ | ↑ **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**$\$$ | Expand $F \rightarrow$ **id** |
| $\$E'T'$**id** | ↑ **id**$\$$ | Match **id** |
| $\$E'T'$ | ↑ $\$$ | Expand $T' \rightarrow \epsilon$ |
| $\$E'$ | ↑ $\$$ | Expand $E' \rightarrow \epsilon$ |
| $\$$ | ↑ $\$$ | |

# More on LL(1) Parsing

- Grammars whose predictive parsing tables contain no duplicate entries are called 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

$$S \rightarrow iEtSS^{'} \mid a$$
$$S^{'} \rightarrow eS \mid \epsilon$$
$$E \rightarrow b$$

# LL(1) Parsing Table for an Ambiguous Grammar

$$S \rightarrow iEtSS' \mid a$$
$$S' \rightarrow eS \mid \epsilon$$
$$E \rightarrow b$$

| Nonterminal | a | b | e | i | t | \$ |
|---|---|---|---|---|---|---|
| $S$ | $S \rightarrow a$ | | | $S \rightarrow iEtSS'$ | | |
| $S'$ | | | $S' \rightarrow \epsilon$ $S' \rightarrow eS$ | | | $S' \rightarrow \epsilon$ |
| $E$ | | $E \rightarrow b$ | | | | |

# Detecting Errors in 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 error

## 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 Predictive Parsing

- Panic mode – skip over symbols until a token in a set of synchronizing (synch) tokens appear
  - Add all tokens in FOLLOW ($A$) to the synch set for $A$, parsing can continue if the parser sees an input symbol in FOLLOW ($A$)
  - Add symbols in FIRST ($A$) to the synch set for $A$, parsing can continue with the nonterminal $A$ that is at the top of the stack
  - Add keywords that can begin constructs
  - …

- Other error handling policies
  - Skip input if the table does not have an entry
  - Pop nonterminal if the table entry is synch

# Predictive Parsing Table with Synchronizing Tokens

## Grammar

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

## FOLLOW Sets

FOLLOW $(E)$ = FOLLOW $\left(E^{'}\right)$ = $\{\$, )\}$

FOLLOW $(T)$ = FOLLOW $\left(T^{'}\right)$ = $\{\$, +, )\}$

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 |
| $T^{'}$ | | $T^{'} \rightarrow \epsilon$ | $T^{'} \rightarrow *FT^{'}$ | | $T^{'} \rightarrow \epsilon$ | $T^{'} \rightarrow \epsilon$ |
| $F$ | $F \rightarrow \textbf{id}$ | synch | synch | $F \rightarrow (E)$ | synch | synch |

# Error Recovery Moves by 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, $M[F, +] =$ 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, 2nd edition, Pearson Education.

📕 K. Cooper and L. Torczon. Engineering a Compiler. Section 3.3, 2nd edition, Morgan Kaufmann.