# Syntactic Analysis (Parsing)

José Miguel Rivero

`rivero@cs.upc.edu`

Barcelona School of Informatics  (FIB)

Universitat Politècnica de Catalunya   BarcelonaTech   (UPC)

# Summary

- Linear Parsing Algorithms. (Counter)example
- Notations: BNF and Extended BNF
- Applying a Rule $A \rightarrow \alpha_i$
- Nullable, First and Follow
- Methods of Linear Parsing
  - Top-down Parsers
    - Grammar Restrictions in Top-down Parsing
      - Elimination of Left Recursion
      - Left Factoring
    - Types of Top-down Parsers
      - Table-driven Top-down Parser
      - Predictive Recursive Top-down Parser
  - Bottom-up Parsers

# Linear Parsing Algorithms

- The list of tokens $w$ will be visited only once, usually in a *left-to-right* traversal.
  The current token receives the name of *lookahead*

- Compute the derivation $S \Rightarrow^* w$ without *backtracking*.
  Sources of <u>indeterminism</u> at this point:

$$S \Rightarrow^* w_0 A_1 w_1 \ldots A_n w_n \Rightarrow$$

- Which non-terminal $A_i$ choose to expand?

- If we have a rule of the form $A \rightarrow \gamma_1 | \ldots | \gamma_k$,
  which $\gamma_j$ —if any— will be used?
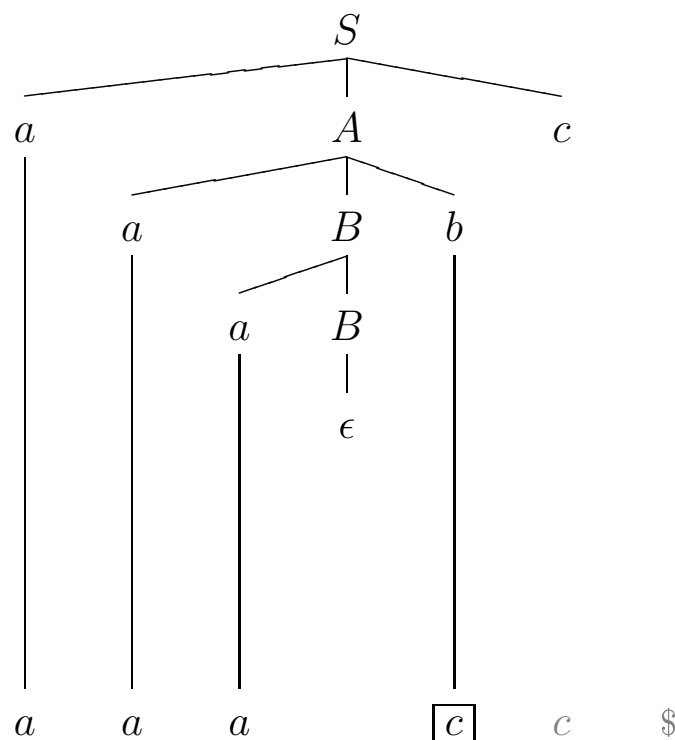  $\Rightarrow$ bear in mind the *lookahead* token.

  $\boxed{\text{at most one } \gamma_j \text{ may make sense}}$

# (Counter)example: backtracking

$$
\begin{aligned}
S &\rightarrow && aAc \\
  &\phantom{\rightarrow} && |\; c \\
A &\rightarrow && aBb \\
  &\phantom{\rightarrow} && |\; Bc \\
B &\rightarrow && aB \\
  &\phantom{\rightarrow} && |\; \epsilon
\end{aligned}
$$

$w = a\,a\,a\,c\,c$

# (Counter)example: backtracking

$$S \rightarrow aAc$$
$$\mid c$$
$$A \rightarrow aBb$$
$$\mid Bc$$
$$B \rightarrow aB$$
$$\mid \epsilon$$

$$w = a\,a\,a\,c\,c$$

# Exercises

- Find another simpler grammar that cannot be parsed with a linear algorithm.
  Give an input that demonstrate it.
- Characterize some conflicting grammars.

# Notation. Backus-Naur Form (BNF)

- $S$, $S'$ are the initial symbol of the grammar
- $A$, $B$, $C$, $\ldots$, are non-terminal symbols
- $a$, $b$, $c$, $\ldots$, and $\$$ are terminal symbols (tokens)
- $X$, $Y$, $Z$ are terminal or non-terminal symbols
- $u$, $v$, $w$ are words formed by terminal symbols, possibly empty ($\epsilon$)
- $\alpha$, $\beta$, $\gamma$ are words formed by terminal and non-terminal symbols, possibly empty ($\epsilon$)
- $A \to \alpha$ is a production rule.
- $A \to \alpha_1 \mid \ldots \mid \alpha_n$ is the group of rules for the non-terminal symbol $A$

FIB

# Notation. Backus-Naur Form (BNF)

- $\gamma_1 A \gamma_2 \Rightarrow \gamma_1 \alpha \gamma_2$ is a derivation step with the rule $A \to \alpha$
- $\Rightarrow^*$ is the reflexive transitive closure of $\Rightarrow$
- Leftmost-derivation ($\Rightarrow^*_{ld}$) if every step:
$$w A \beta \Rightarrow_{A \to \alpha} w \alpha \beta$$
- Rightmost-derivation ($\Rightarrow^*_{rd}$) if every step:
$$\beta A w \Rightarrow_{A \to \alpha} \beta \alpha w$$

FIB

# Extended Backus-Naur Form (EBNF)

Rules $A \to \alpha$ where $\alpha$ can take the form of a regular expression:

- $\alpha_1 \cdots \alpha_n$ : concatenation, or $\epsilon$
- $\alpha_1 | \ldots | \alpha_n$ : alternatives
- $\alpha_1^*$ : 0 or more times $\alpha_1$
- $\alpha_1^+$ : 1 or more times $\alpha_1$ ($\alpha_1 \alpha_1^*$)
- $\alpha_1$? : 0 or 1 times $\alpha_1$ ($\alpha_1 | \epsilon$)
- $(\alpha_1)$ : parenthesis for breaking the standard precedence between operators
$$\{\,|\,\} \prec_p \{\,\cdot\,\} \prec_p \{\,^*,\,^+,\,?\,\}$$
- a non-terminal symbol $A$, or a terminal symbol $a$

# Extended Backus-Naur Form (EBNF)

In *ANTLR*, identifiers in capital letters denote terminals symbols, and identifiers beginning in lower case denote non-terminals.
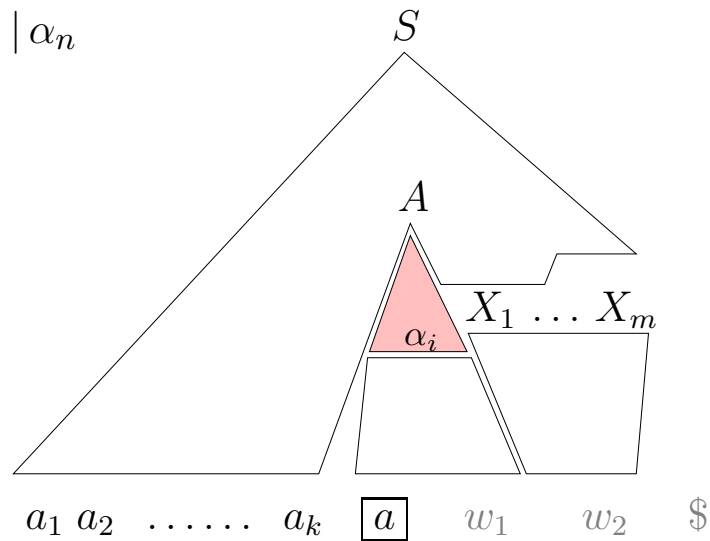
Rules take the form $\quad id : reg\_exp\,;$

Example:

```
expr : term ((PLUS|MINUS) term)* ;

term : factor ((TIMES|QUOTIENT) factor)* ;

factor : IDENT (LEFT_BRK expr RIGHT_BRK)?
       | NUM
       | LEFT_PAR expr RIGHT_PAR
       ;
```

# Applying a Rule $A \to \alpha_i$
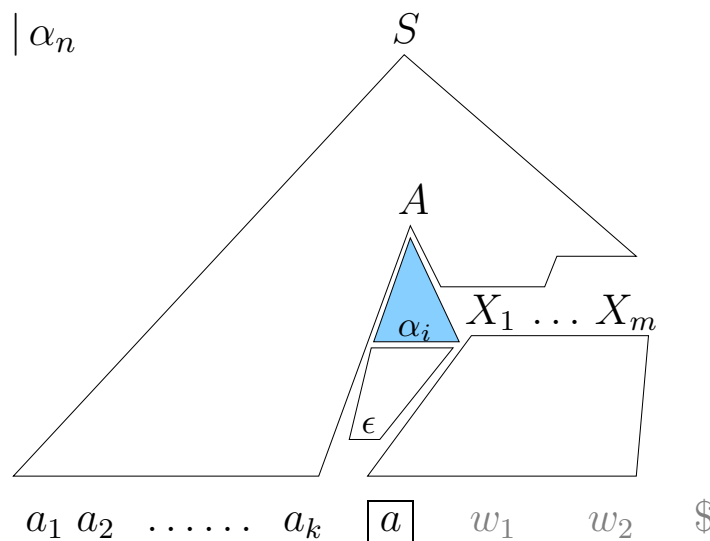
$A \;\to\; \alpha_1 \mid \ldots \mid \alpha_n$

$S$

$A$

$\alpha_i$

$X_1 \ldots X_m$

$a_1 \; a_2 \quad \ldots \ldots \quad a_k \quad \boxed{a} \quad w_1 \qquad w_2 \qquad \$$

$A \Rightarrow \alpha_i \Rightarrow^* a \, w_1 \qquad\qquad X_1 \ldots X_m \Rightarrow^* w_2$
$a \in first(\alpha_i)$

FIB

# Applying a Rule $A \to \alpha_i$

$A \;\to\; \alpha_1 \mid \ldots \mid \alpha_n$

$S$

$A$

$\alpha_i$

$X_1 \ldots X_m$

$\epsilon$

$a_1 \; a_2 \quad \ldots \ldots \quad a_k \quad \boxed{a} \quad w_1 \qquad w_2 \qquad \$$

$A \Rightarrow \alpha_i \Rightarrow^* \epsilon \qquad\qquad X_1 \ldots X_m \Rightarrow^* a \, w_1 \, w_2$
$nullable?(\alpha_i) \qquad\qquad S \,\$ \;\Rightarrow^* a_1 \ldots a_k \, A \, a \, w_1 \, w_2 \, \$$

FIB

# Applying a Rule $A \to \alpha_i$

$$A \to \alpha_1 \mid \ldots \mid \alpha_n$$



$$A \Rightarrow \alpha_i \Rightarrow^* \epsilon \qquad\qquad S \, \$ \Rightarrow^* \gamma_1 \, A \, a \, \gamma_2$$
$$nullable?(\alpha_i) \qquad\qquad\qquad a \in follow(A)$$

# Nullable, First, Follow

- Can $\alpha$ derive $\epsilon$ ?

$$nullable?(\alpha) \text{ iff } \alpha \Rightarrow^* \epsilon$$

- Set of terminals $a$ that can be the first symbol of words derived from $\alpha$

$$first(\alpha) \;=\; \{\, a \mid \alpha \Rightarrow^* a\,w \,\}$$

- Set of terminals $a$ that can follow $A$ in a derivation from $S\,\$$

$$follow(A) \;=\; \{\, a \mid S\,\$ \Rightarrow^* \gamma_1 A\, a\, \gamma_2 \,\}$$

# To Sum Up . . .

$$S \rightarrow \quad aAc$$
$$\mid \quad c$$
$$A \rightarrow \quad aBb$$
$$\mid \quad Bc$$
$$B \rightarrow \quad aB$$
$$\mid \quad \epsilon$$

$$first(aB) \quad = \quad \{\, a \,\}$$
$$first(\epsilon) \quad = \quad \emptyset$$
$$first(B) \quad = \quad first(aB) \cup first(\epsilon)$$
$$= \quad \{\, a \,\}$$
$$first(aBb) \quad = \quad \{\, a \,\}$$
$$first(Bc) \quad = \quad \{\, a,\, c \,\}$$
$$first(A) \quad = \quad first(aBb) \cup first(Bc)$$
$$= \quad \{\, a,\, c \,\}$$
$$first(aAc) \quad = \quad \{\, a \,\}$$
$$first(c) \quad = \quad \{\, c \,\}$$
$$first(S) \quad = \quad first(aAc) \cup first(c)$$
$$= \quad \{\, a,\, c \,\}$$

FIB

# To Sum Up . . .

$$S \rightarrow \quad aAc$$
$$\mid \quad c$$
$$A \rightarrow \quad aBb$$
$$\mid \quad Bc$$
$$B \rightarrow \quad aB$$
$$\mid \quad \epsilon$$

$$first(aB) \quad = \quad \{\, a \,\}$$
$$first(\epsilon) \quad = \quad \emptyset$$
$$nullable?(\epsilon) \quad = \quad true$$
$$follow(B) \quad = \quad \{\, b,\, c \,\}$$

FIB

# $nullable?(\alpha)$ (BNF)

$$nullable?(\alpha) \text{ iff } \alpha \Rightarrow^* \epsilon$$

- $nullable?(a) = false$

- $nullable?(A) = true \text{ if } \exists A \rightarrow \alpha \in G \wedge nullable(\alpha)$
  $= false \text{ otherwise}$

- $nullable?(X_1 \ldots X_n) = true \text{ if } \forall i : 1 \leq i \leq n : nullable(X_i)$
  $= false \text{ otherwise}$

# $nullable?(\alpha)$ (Extended BNF)

$$nullable?(\alpha) \text{ iff } \alpha \Rightarrow^* \epsilon$$

- $nullable?(\alpha_1\, \alpha_2) = true \text{ if } nullable?(\alpha_1) \wedge nullable?(\alpha_2)$
  $= false \text{ otherwise}$

- $nullable?(\alpha_1 \,|\, \alpha_2) = true \text{ if } nullable?(\alpha_1) \vee nullable?(\alpha_2)$
  $= false \text{ otherwise}$

- $nullable?(\alpha_1^*) = true$

- $nullable?(\alpha_1^+) = true \text{ if } nullable?(\alpha_1)$
  $= false \text{ otherwise}$

- $nullable?(\alpha_1?) = true$

# $first(\alpha)$ **(BNF)**

$$first(\alpha) \;=\; \{\, a \mid \alpha \Rightarrow^* a\,w \,\}$$

- $first(a) \qquad\quad = \{a\}$

- $first(A) \qquad\quad \supseteq first(\alpha_j) \;\;$ **if** $\exists\, A \to \alpha_j \in G$

- $first(X_1 \ldots X_n) \supseteq first(X_i) \;\;$ **if** $nullable?(X_1 \ldots X_{i-1})$

# $first(\alpha)$ **(Extended BNF)**

$$first(\alpha) \;=\; \{\, a \mid \alpha \Rightarrow^* a\,w \,\}$$

- $first(\alpha_1\,\alpha_2) \;= first(\alpha_1) \qquad\qquad\quad$ **if** $\neg\,nullable?(\alpha_1)$
  $first(\alpha_1\,\alpha_2) \;= first(\alpha_1) \cup first(\alpha_2) \;\;$ **if** $nullable?(\alpha_1)$

- $first(\alpha_1 \mid \alpha_2) = first(\alpha_1) \cup first(\alpha_2)$

- $first(\alpha_1^*) \qquad = first(\alpha_1)$

- $first(\alpha_1^+) \qquad = first(\alpha_1)$

- $first(\alpha_1?) \qquad = first(\alpha_1)$

# $follow(A)$ (BNF)

$$follow(A) \;=\; \{\, a \mid S\,\$ \Rightarrow^* \gamma_1 A\, a\, \gamma_2 \,\}$$

- $follow(S) \supseteq \{\,\$\,\}$

- $follow(B) \supseteq first(\beta)$    **if** $\exists\, A \to \alpha B \beta \in G$

- $follow(B) \supseteq follow(A)$ **if** $\exists\, A \to \alpha B \in G$
  $follow(B) \supseteq follow(A)$ **if** $\exists\, A \to \alpha B \beta \in G \;\wedge$
  $$nullable?(\beta)$$

# $follow(\alpha)$ (Extended BNF)

**If** $\alpha$ occurs in $G$: $follow(\alpha) \;=\; \{\, a \mid S\,\$ \Rightarrow^* \gamma_1 \alpha\, a\, \gamma_2 \,\}$

- $follow(\alpha) \;\supseteq\; follow(A)$              **if** $\exists\, A \to \alpha \in G$
- $follow(\alpha_1) \supseteq first(\alpha_2)$            **if** $\alpha_1\,\alpha_2$ occurs in $G$
- $follow(\alpha_2) \supseteq follow(\alpha_1\,\alpha_2)$     **if** $\alpha_1\,\alpha_2$ occurs in $G$
- $follow(\alpha_1) \supseteq follow(\alpha_1\,\alpha_2)$     **if** $\alpha_1\,\alpha_2$ occurs in $G$
  $$\wedge\; nullable?(\alpha_2)$$

- $follow(\alpha_1) \supseteq follow(\alpha_1\,|\,\alpha_2)$    **if** $\alpha_1\,|\,\alpha_2$ occurs in $G$
- $follow(\alpha_2) \supseteq follow(\alpha_1\,|\,\alpha_2)$    **if** $\alpha_1\,|\,\alpha_2$ occurs in $G$
- $follow(\alpha_1) \supseteq first(\alpha_1) \cup follow(\alpha_1^*)$ **if** $\alpha_1^*$ occurs in $G$
- $follow(\alpha_1) \supseteq first(\alpha_1) \cup follow(\alpha_1^+)$ **if** $\alpha_1^+$ occurs in $G$
- $follow(\alpha_1) \supseteq follow(\alpha_1?)$           **if** $\alpha_1?$ occurs in $G$

# Methods of Linear Parsing

The list of tokens will be traversed *left-to-right*. Decisions to proceed take into account one token of lookahead.

- Top-down parsers (LL(1))
  - Build the AST from the root to the leaves (top-down)
  - Follow a left-most derivation in forward direction
  - More intuitive: can be *manually* written
  - **Grammars may need preprocessing**
- Bottom-up parsers (LR(1))
  - Build the AST from the leaves to the root (bottom-up)
  - Follow a right-most derivation in *backward* direction
  - Less intuitive than top-down parsers
  - Slightly more powerful

# Elimination of Left Recursion

Grammar G:

$$E \; \rightarrow \; E \; + \; \text{id}$$
$$\mid \; \text{id}$$

$$first(E \; + \; \text{id}) \; = \; \{\text{id}\}$$
$$first(\text{id}) \; = \; \{\text{id}\}$$
$$first(E \; + \; \text{id}) \; \cap \; first(\text{id}) \; \neq \; \emptyset$$

$$w = \text{id} + \text{id} + \text{id}$$

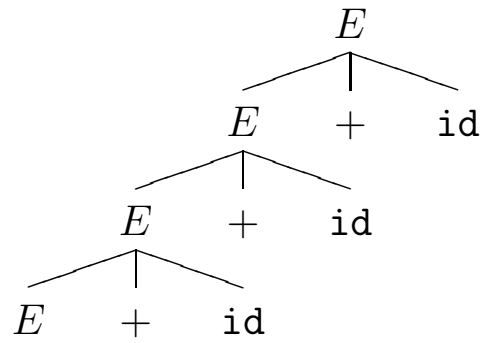# Elimination of Left Recursion

Grammar G:

$$E \rightarrow E + \text{id}$$
$$\mid \text{id}$$

$w = \text{id} + \text{id} + \text{id}$

# Elimination of Left Recursion (BNF)

$$A \rightarrow A\,\alpha_1$$
$$\mid \dots$$
$$\mid A\,\alpha_n$$
$$\mid \beta_1$$
$$\mid \dots$$
$$\mid \beta_m$$

Transform into right recursion:

$$A \rightarrow \beta_1\,A'$$
$$\mid \dots$$
$$\mid \beta_m\,A'$$

$$A' \rightarrow \alpha_1\,A'$$
$$\mid \dots$$
$$\mid \alpha_n\,A'$$
$$\mid \epsilon$$

# Elimination of Left Recursion (EBNF)

$$A \;\rightarrow\; A\,\alpha_1$$
$$|\;\ldots$$
$$|\;A\,\alpha_n$$
$$|\;\beta_1$$
$$|\;\ldots$$
$$|\;\beta_m$$

Extended BNF: use regular expressions

$$A \;\rightarrow\; (\,\beta_1 \,|\, \ldots \,|\, \beta_m\,)\,(\,\alpha_1 \,|\, \ldots \,|\, \alpha_n\,)^*$$

$$A \;\rightarrow\; B\,(A')^*$$
$$B \;\rightarrow\; \beta_1 \,|\, \ldots \,|\, \beta_m$$
$$A' \;\rightarrow\; \alpha_1 \,|\, \ldots \,|\, \alpha_n$$

# Exercises

$$LI \;\rightarrow\; LI\,I$$
$$|\;I$$

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

Indirect left recursion:

$$A \;\rightarrow\; B\,d$$
$$B \;\rightarrow\; C\,e$$
$$C \;\rightarrow\; A\,f$$
$$|\;g$$

# Left Factoring

$$E \rightarrow T + E$$
$$| \; T$$

$$first(T + E) = first(T) = \{\, \texttt{id}, (\, \}$$
$$first(T) = \{\, \texttt{id}, (\, \}$$
$$first(T + E) \cap first(T) \neq \emptyset$$

$$T \rightarrow \texttt{id}$$
$$| \; (E)$$

$$first(\texttt{id}) = \{\, \texttt{id} \,\}$$
$$first((E)) = \{\, (\, \}$$
$$first(\texttt{id}) \cap first((E)) = \emptyset$$

# Left Factoring (BNF)

$$A \rightarrow \beta \, \alpha_1$$
$$| \; \ldots$$
$$| \; \beta \, \alpha_n$$
$$| \; \gamma_1$$
$$| \; \ldots$$
$$| \; \gamma_m$$

$$A \rightarrow \beta \, A'$$
$$| \; \gamma_1$$
$$| \; \ldots$$
$$| \; \gamma_m$$

$$A' \rightarrow \alpha_1$$
$$| \; \ldots$$
$$| \; \alpha_n$$

# Left Factoring (EBNF)

$$A \;\rightarrow\; \beta\,\alpha_1$$
$$\mid\; \ldots$$
$$\mid\; \beta\,\alpha_n$$
$$\mid\; \gamma_1$$
$$\mid\; \ldots$$
$$\mid\; \gamma_m$$

Extended BNF:

$$A \;\rightarrow\; \beta\,(\,\alpha_1 \mid \ldots \mid \alpha_n\,) \mid \gamma_1 \mid \ldots \mid \gamma_m$$

$$A \;\rightarrow\; \beta\,A' \mid \gamma_1 \mid \ldots \mid \gamma_m$$
$$A' \;\rightarrow\; \alpha_1 \mid \ldots \mid \alpha_n$$

# Exercises

$$P \;\rightarrow\; \text{if } C \text{ then } P \text{ endif}$$
$$\mid\; \text{if } C \text{ then } P \text{ else } P \text{ endif}$$
$$\mid\; \text{p}$$
$$C \;\rightarrow\; \text{c}$$

$$I \;\rightarrow\; LE \text{ ':=' } E$$
$$\mid\; \text{write '(' } E \text{ ')'}$$
$$\mid\; \text{id '(' } E \text{ ')'}$$
$$E \;\rightarrow\; \text{id}$$
$$\mid\; \text{num}$$
$$LE \;\rightarrow\; \text{id}$$

# Types of Top-down Parsers

- Table Driven parsers (iterative)
  - Parsing algorithm is fixed, driven by a decision table
  - Table $M$ is built from the grammar $G$.
    Empty boxes correspond to syntax errors

| $M$ | $a_1$ | $\ldots$ | $a$ | $\ldots$ | $a_n$ | $\$$ |
|---|---|---|---|---|---|---|
| $A_1$ | | | | | | |
| $\vdots$ | | | | | | |
| $A$ | | | $A \to \alpha_k$ | | | |
| $\vdots$ | | | | | | |
| $A_m$ | | | | | | |

# Types of Top-down Parsers

- Recursive predictive parsers
  - Parsing algorithm is formed by a set of mutually recursive functions
  - Each rule $A \to \alpha$ generates the code of its function
    ```
    void A(void) {
        // Code generated from α
    }
    ```
  - `Gencode` describes how to translate a rule to the associated function

# Table-driven Top-down Parser

| $M$ | $a_1$ | $\ldots$ | $a$ | $\ldots$ | $a_n$ | $\$$ |
|---|---|---|---|---|---|---|
| $A_1$ | | | | | | |
| $\vdots$ | | | | | | |
| $A$ | | | $A \to \alpha_k$ | | | |
| $\vdots$ | | | | | | |
| $A_m$ | | | | | | |

$$A \to \begin{array}{l} \alpha_1 \\ \mid \quad \ldots \\ \mid \quad \alpha_k \\ \mid \quad \ldots \\ \mid \quad \alpha_o \end{array}$$

$A \to \alpha_k \in M[A, a]$ if :

- $a \in first(\alpha_k)$ , or
- $nullable?(\alpha_k)$ and $a \in follow(A)$

# Table-driven Top-down Parser

Algorithm to build the parser table $M[A, a]$

```
for all rule A → α ∈ G do
    add A → α to M[A, a] if:
```
- $a \in first(\alpha)$ or,
- $nullable?(\alpha)$ and $a \in follow(A)$

# Exercises

Complete the table $M$ for the following grammars, indicating the possible problems.

Grammar $G_1$ :

$$
\begin{aligned}
E \to\ & E + T \\
|\ & T \\
T \to\ & T * F \\
|\ & F \\
F \to\ & \texttt{id} \\
|\ & (\,E\,)
\end{aligned}
$$

Grammar $G_2$ :

$$
\begin{aligned}
E \to\ & T\,E' \\
E' \to\ & +\,T\,E' \\
|\ & \epsilon \\
T \to\ & F\,T' \\
T' \to\ & *\,F\,T' \\
|\ & \epsilon \\
F \to\ & \texttt{id} \\
|\ & (\,E\,)
\end{aligned}
$$

# Exercises

Complete the table $M$ for the following grammars, indicating the possible problems.

Grammar $G_3$ :

$$
\begin{aligned}
P \to\ & \texttt{if } C \texttt{ then } P \\
|\ & \texttt{if } C \texttt{ then } P \texttt{ else } P \\
|\ & \texttt{p} \\
C \to\ & \texttt{c}
\end{aligned}
$$

Grammar $G_4$ :

$$
\begin{aligned}
P \to\ & \texttt{if } C \texttt{ then } P\ P' \\
|\ & \texttt{p} \\
P' \to\ & \epsilon \\
|\ & \texttt{else } P \\
C \to\ & \texttt{c}
\end{aligned}
$$

# Exercises

Complete the table $M$ for the following grammars, indicating the possible problems.

Grammar $G_5$ :

$$P \rightarrow \text{ if } C \text{ then } P \; P'$$
$$\mid \; \text{p}$$
$$P' \rightarrow \text{ endif}$$
$$\mid \; \text{else } P \text{ endif}$$
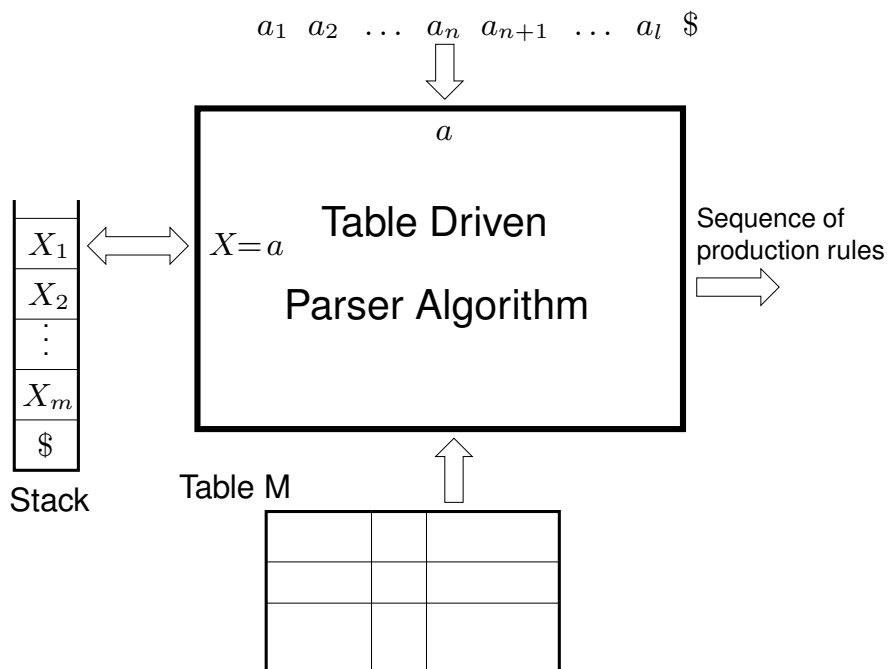$$C \rightarrow \text{ c}$$

# Table-driven Top-down Parser Algorithm

# Table-driven Top-down Parser Algorithm

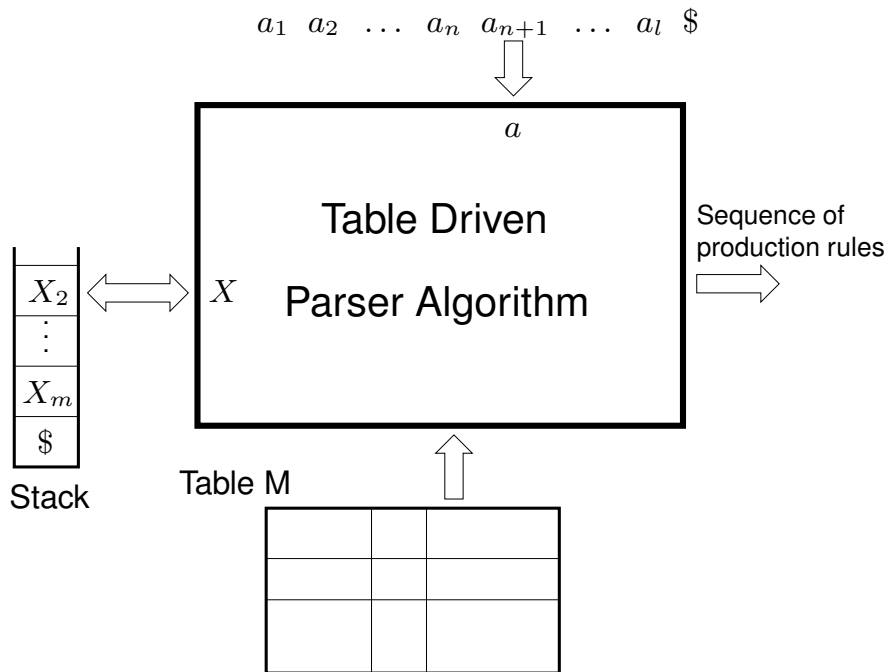$$a_1 \quad a_2 \quad \ldots \quad a_n \quad a_{n+1} \quad \ldots \quad a_l \quad \$$$

$a$

Table Driven

Parser Algorithm

$X$

Sequence of production rules

$X_2$
$\vdots$
$X_m$
$\$$

**Stack**

**Table M**

# Table-driven Top-down Parser Algorithm

$$a_1 \quad a_2 \quad \ldots \quad a_n \quad a_{n+1} \quad \ldots \quad a_l \quad \$$$

$a$

Table Driven

Parser Algorithm

$X = A$

Sequence of production rules

$X_1$
$X_2$
$\vdots$
$X_m$
$\$$

**Stack**

**Table M**   $a$

$A \, (X_1)$     $A \rightarrow Y_1 \ldots Y_p$

# Table-driven Top-down Parser Algorithm



$$a_1 \quad a_2 \quad \ldots \quad a_n \quad a_{n+1} \quad \ldots \quad a_l \quad \$$$

Stack: $Y_1$, $\vdots$, $Y_p$, $X_2$, $\vdots$, $X_m$, $\$$

Table Driven Parser Algorithm

$X$ ... $a$

Sequence of production rules
$A \to Y_1 \ldots Y_p$

Table M $\quad a$

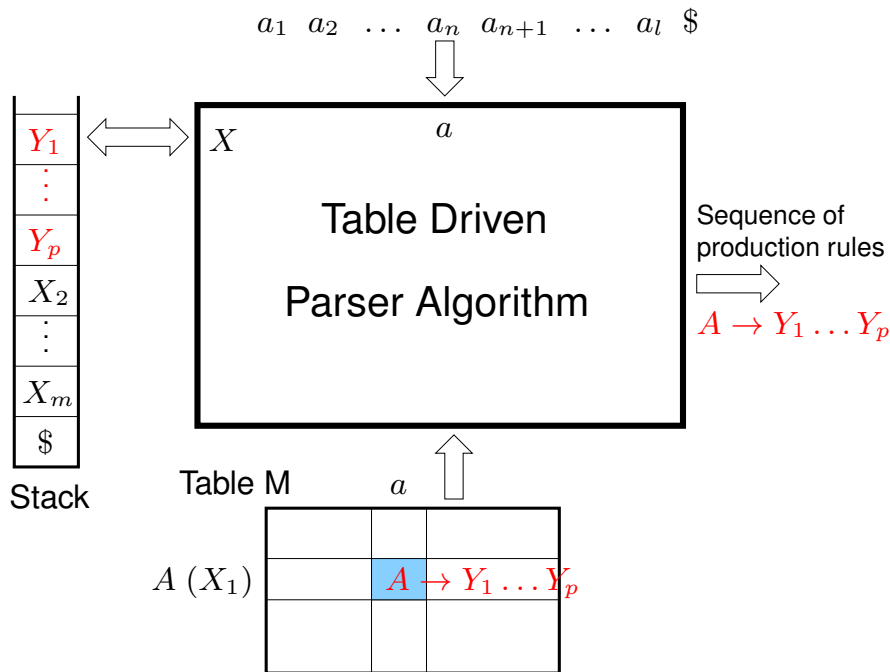$A\,(X_1)$ ... $A \to Y_1 \ldots Y_p$

# Table-driven Top-down Parser Algorithm

$Stk := EmptyStack(\,); \ PushStack(Stk, \$); \ PushStack(Stk, S);$

$X := TopStack(Stk); \ a := FirstToken(\,);$

`while` $X \neq \$$ `do`

    `if` $X$ *is terminal* `then`

        `if` $X = a$ `then`

            $PopStack(Stk); \ a := NextToken(\,);$

        `else`

            *throw syntax error*

    `else`  // $X$ *is non−terminal*

        `if` $M[X, a]$ *is empty* $(is\,error)$ `then`

            *throw syntax error*

        `else`  // $M[X, a] = X \to Y_1 \ldots Y_p$

            *emit production rule* $X \to Y_1 \ldots Y_p$

            $PopStack(Stk);$ `for` $i := p$ `downto` $1$ `do` $PushStack(Stk, Y_i);$

    $X := TopStack(Stk);$

`endwhile`