

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 indeterminism at this point:

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

- Which non-terminal A_i choose to expand?
- If we have a rule of the form $A \rightarrow \gamma_1 | \dots | \gamma_k$, which γ_j —if any— will be used?
 \Rightarrow bear in mind the *lookahead* token.

at most one γ_j may make sense

(Counter)example: backtracking

$S \rightarrow aAc$

| c

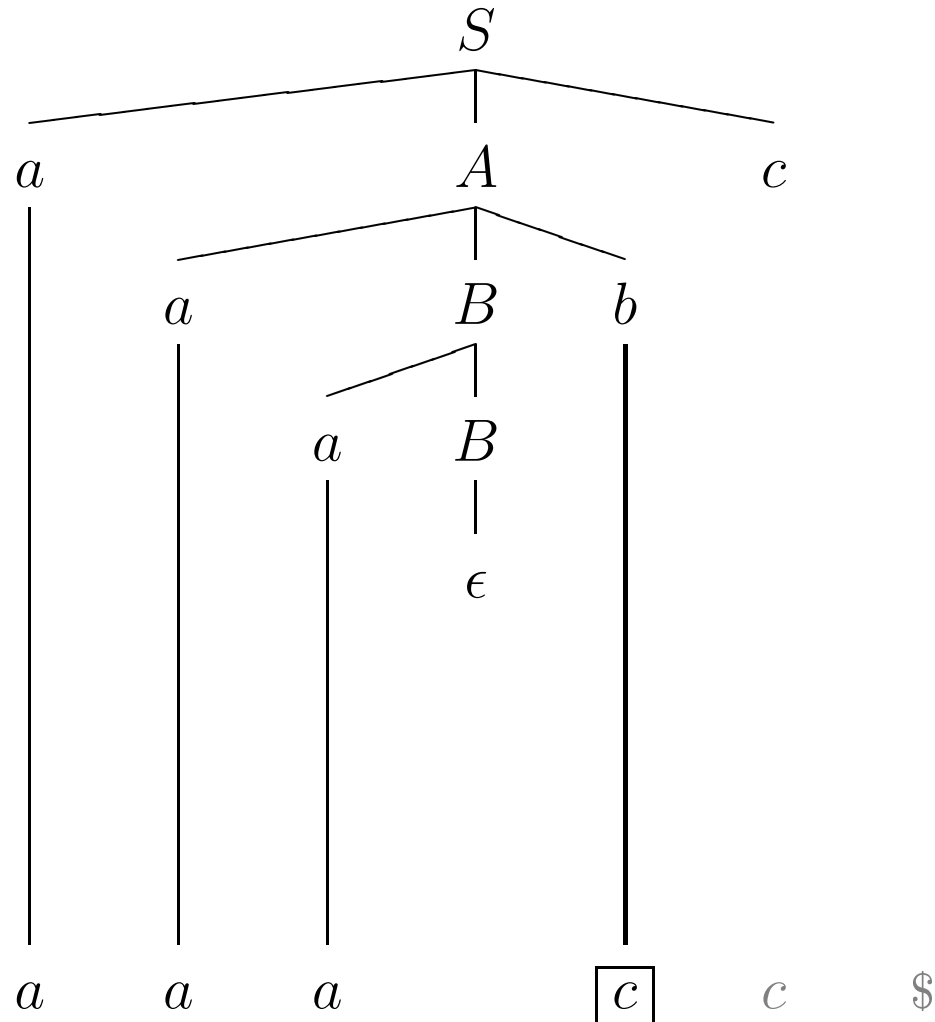
$A \rightarrow aBb$

| Bc

$B \rightarrow aB$

| ϵ

$w = a a a c c$



(Counter)example: backtracking

$S \rightarrow aAc$

| c

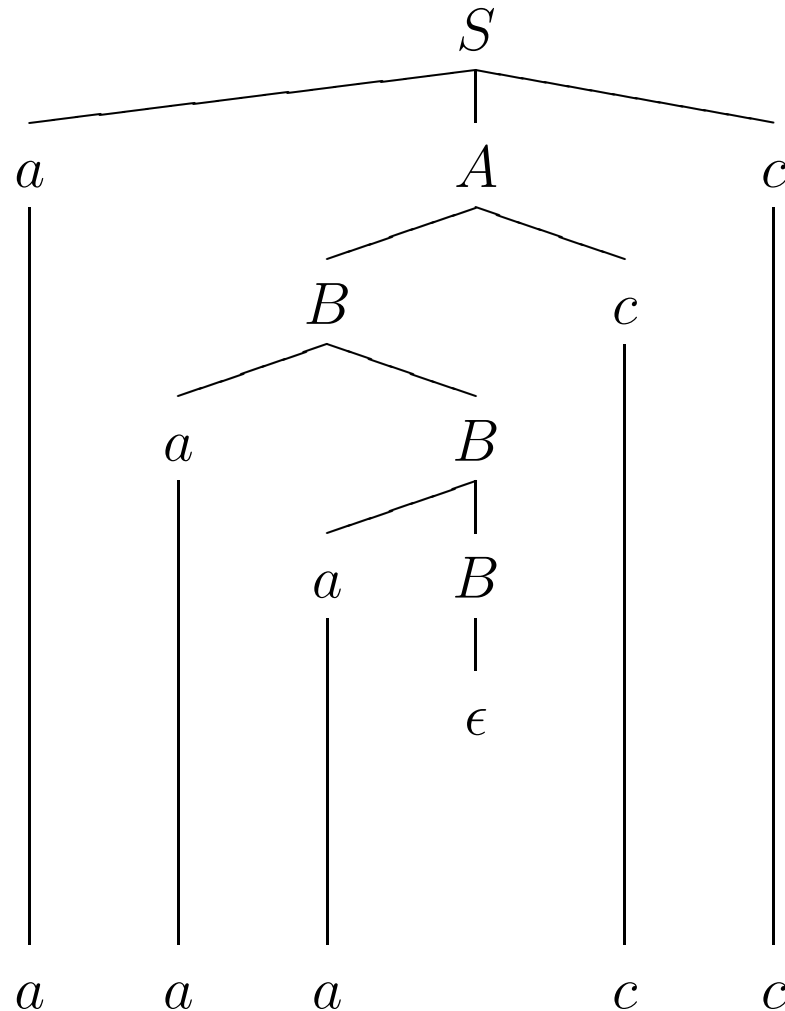
$A \rightarrow aBb$

| Bc

$B \rightarrow aB$

| ϵ

$w = a a a c c$



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, \dots , are non-terminal symbols
- a, b, c, \dots , 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 (ϵ)
- α, β, γ are words formed by terminal and non-terminal symbols, possibly empty (ϵ)
- $A \rightarrow \alpha$ is a production rule.
- $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$ is the group of rules for the non-terminal symbol A

Notation. Backus-Naur Form (BNF)

• $\gamma_1 A \gamma_2 \Rightarrow \gamma_1 \alpha \gamma_2$ is a derivation step with the rule $A \rightarrow \alpha$

• \Rightarrow^* is the reflexive transitive closure of \Rightarrow

• Leftmost-derivation (\Rightarrow_{ld}^*) if every step:

$$w A \beta \Rightarrow_{A \rightarrow \alpha} w \alpha \beta$$

• Rightmost-derivation (\Rightarrow_{rd}^*) if every step:

$$\beta A w \Rightarrow_{A \rightarrow \alpha} \beta \alpha w$$

Extended Backus-Naur Form (EBNF)

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

- $\alpha_1 \cdots \alpha_n$: concatenation, or ϵ
- $\alpha_1 | \dots | \alpha_n$: alternatives
- α_1^* : 0 or more times α_1
- α_1^+ : 1 or more times α_1 ($\alpha_1 \alpha_1^*$)
- $\alpha_1?$: 0 or 1 times α_1 ($\alpha_1 | \epsilon$)
- (α_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.

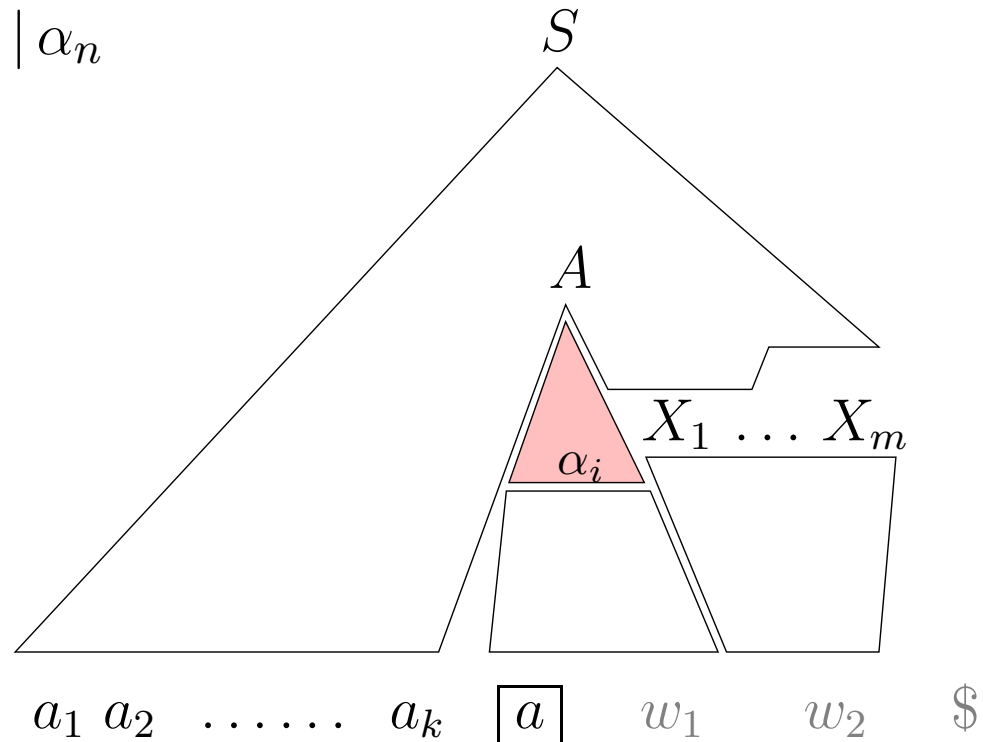
Rules take the form $id : reg_exp ;$

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 $A \rightarrow \alpha_i$

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

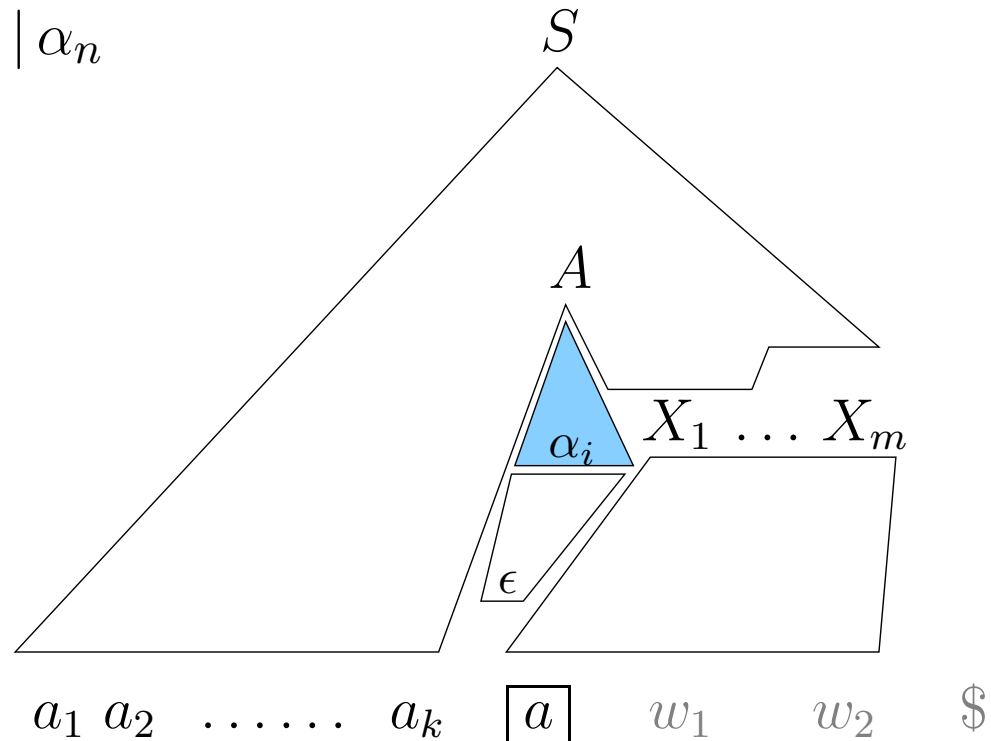


$A \Rightarrow \alpha_i \Rightarrow^* a w_1$
 $a \in \text{first}(\alpha_i)$

$X_1 \dots X_m \Rightarrow^* w_2$

Applying a Rule $A \rightarrow \alpha_i$

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

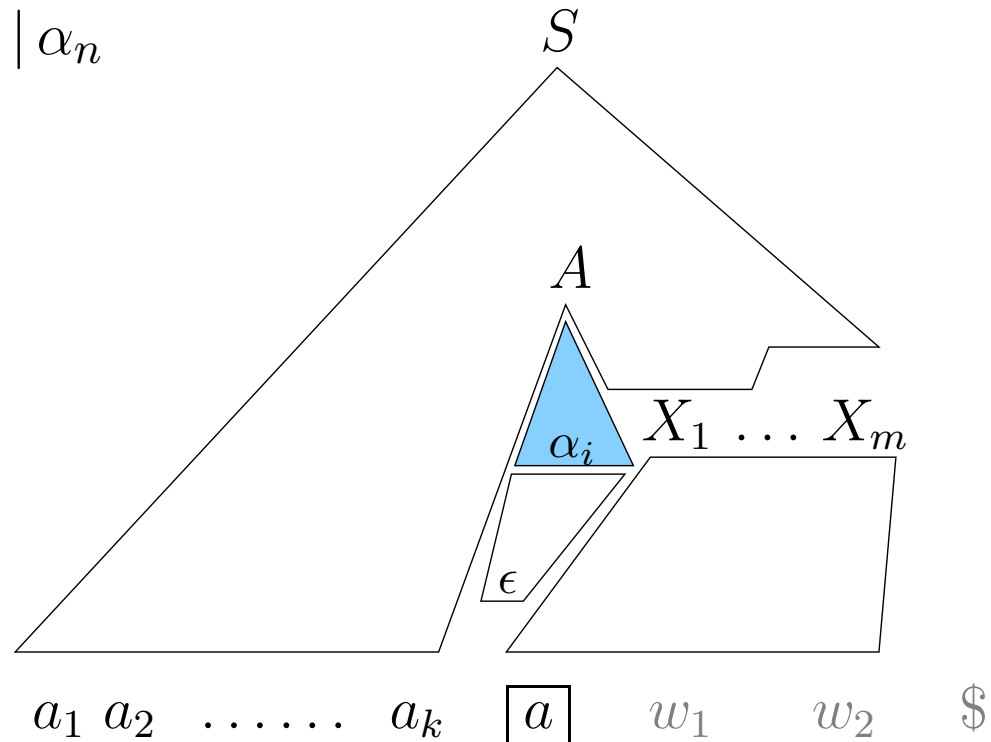


$A \Rightarrow \alpha_i \Rightarrow^* \epsilon$
nullable? (α_i)

$X_1 \dots X_m \Rightarrow^* a w_1 w_2$
 $S \$ \Rightarrow^* a_1 \dots a_k A a w_1 w_2 \$$

Applying a Rule $A \rightarrow \alpha_i$

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



$A \Rightarrow \alpha_i \Rightarrow^* \epsilon$
nullable?(α_i)

$S \$ \Rightarrow^* \gamma_1 A a \gamma_2$
 $a \in \text{follow}(A)$

Nullable, First, Follow

- Can α derive ϵ ?

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

- Set of terminals a that can be the first symbol of words derived from α

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

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

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

To Sum Up ...

$S \rightarrow aAc$

| c

$A \rightarrow aBb$

| Bc

$B \rightarrow aB$

| ϵ

$first(aB) = \{a\}$

$first(\epsilon) = \emptyset$

$first(B) = first(aB) \cup first(\epsilon)$

$= \{a\}$

$first(aBb) = \{a\}$

$first(Bc) = \{a, c\}$

$first(A) = first(aBb) \cup first(Bc)$

$= \{a, c\}$

$first(aAc) = \{a\}$

$first(c) = \{c\}$

$first(S) = first(aAc) \cup first(c)$

$= \{a, c\}$

To Sum Up ...

$S \rightarrow aAc$
| c

$A \rightarrow aBb$
| Bc

$B \rightarrow aB$
| ϵ

$first(aB) = \{a\}$

$first(\epsilon) = \emptyset$

$nullable?(\epsilon) = true$

$follow(B) = \{b, c\}$

nullable?(α) (BNF)

nullable?(α) iff $\alpha \Rightarrow^* \epsilon$

• *nullable?*(a) = *false*

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

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

nullable?(α) (Extended BNF)

nullable?(α) iff $\alpha \Rightarrow^* \epsilon$

- *nullable?*($\alpha_1 \alpha_2$) = *true* if *nullable?*(α_1) \wedge *nullable?*(α_2)
= *false* otherwise
- *nullable?*($\alpha_1 \mid \alpha_2$) = *true* if *nullable?*(α_1) \vee *nullable?*(α_2)
= *false* otherwise
- *nullable?*(α_1^*) = *true*
- *nullable?*(α_1^+) = *true* if *nullable?*(α_1)
= *false* otherwise
- *nullable?*($\alpha_1?$) = *true*

first(α) (BNF)

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

- $\textit{first}(a) = \{a\}$
- $\textit{first}(A) \supseteq \textit{first}(\alpha_j)$ if $\exists A \rightarrow \alpha_j \in G$
- $\textit{first}(X_1 \dots X_n) \supseteq \textit{first}(X_i)$ if *nullable?*($X_1 \dots X_{i-1}$)

$first(\alpha)$ (Extended BNF)

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

- $first(\alpha_1 \alpha_2) = first(\alpha_1)$ if $\neg nullable?(\alpha_1)$
 $first(\alpha_1 \alpha_2) = first(\alpha_1) \cup first(\alpha_2)$ if $nullable?(\alpha_1)$
- $first(\alpha_1 \mid \alpha_2) = first(\alpha_1) \cup first(\alpha_2)$
- $first(\alpha_1^*) = first(\alpha_1)$
- $first(\alpha_1^+) = first(\alpha_1)$
- $first(\alpha_1?) = first(\alpha_1)$

follow(A) (BNF)

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

- $\text{follow}(S) \supseteq \{ \$ \}$
- $\text{follow}(B) \supseteq \text{first}(\beta)$ if $\exists A \rightarrow \alpha B \beta \in G$
- $\text{follow}(B) \supseteq \text{follow}(A)$ if $\exists A \rightarrow \alpha B \in G$
- $\text{follow}(B) \supseteq \text{follow}(A)$ if $\exists A \rightarrow \alpha B \beta \in G \wedge \text{nullable?}(\beta)$

$follow(\alpha)$ (Extended BNF)

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

- $follow(\alpha) \supseteq follow(A)$ if $\exists A \rightarrow \alpha \in G$
- $follow(\alpha_1) \supseteq first(\alpha_2)$ 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 α_1^* occurs in G
- $follow(\alpha_1) \supseteq first(\alpha_1) \cup follow(\alpha_1^+)$ if α_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:

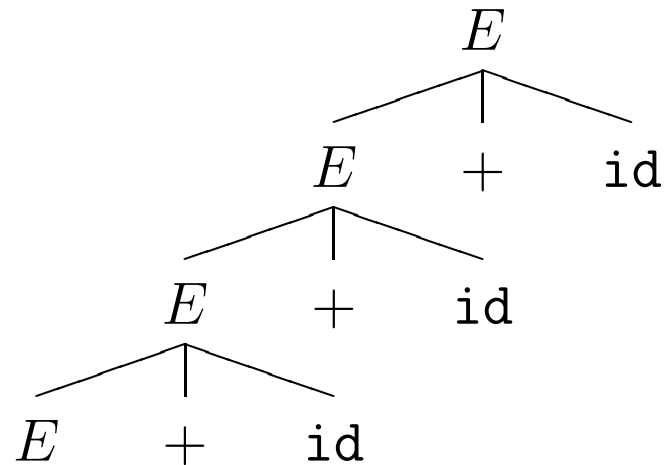
$$\begin{array}{l} E \rightarrow E + \text{id} \\ \quad | \text{id} \end{array}$$
$$\text{first}(E + \text{id}) = \{\text{id}\}$$
$$\text{first}(\text{id}) = \{\text{id}\}$$
$$\text{first}(E + \text{id}) \cap \text{first}(\text{id}) \neq \emptyset$$
$$w = \text{id} + \text{id} + \text{id}$$

Elimination of Left Recursion

Grammar G:

$$E \rightarrow E + id$$
$$E \rightarrow id$$

$w = id + id + id$



id + id + id

Elimination of Left Recursion (BNF)

$$\begin{array}{l} A \rightarrow A \alpha_1 \\ | \dots \\ | A \alpha_n \\ | \beta_1 \\ | \dots \\ | \beta_m \end{array}$$

Transform into right recursion:

$$\begin{array}{l} A \rightarrow \beta_1 A' \\ | \dots \\ | \beta_m A' \end{array}$$

$$\begin{array}{l} A' \rightarrow \alpha_1 A' \\ | \dots \\ | \alpha_n A' \\ | \epsilon \end{array}$$

Elimination of Left Recursion (EBNF)

$$\begin{array}{l} A \rightarrow A \alpha_1 \\ | \dots \\ | A \alpha_n \\ | \beta_1 \\ | \dots \\ | \beta_m \end{array}$$

Extended BNF: use regular expressions

$$A \rightarrow (\beta_1 | \dots | \beta_m) (\alpha_1 | \dots | \alpha_n)^*$$

$$A \rightarrow B (A')^*$$

$$B \rightarrow \beta_1 | \dots | \beta_m$$

$$A' \rightarrow \alpha_1 | \dots | \alpha_n$$

Exercises

$$\begin{array}{l} LI \rightarrow LI I \\ \quad | I \end{array}$$

$$\begin{array}{l} E \rightarrow E + T \\ \quad | T \end{array}$$

$$\begin{array}{l} T \rightarrow T * F \\ \quad | F \end{array}$$

$$\begin{array}{l} F \rightarrow (E) \\ \quad | id \end{array}$$

Indirect left recursion:

$$A \rightarrow B d$$

$$B \rightarrow C e$$

$$C \rightarrow A f$$

$$\quad | g$$

Left Factoring

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

Left Factoring (BNF)

$$\begin{array}{l} A \rightarrow \beta \alpha_1 \\ | \dots \\ | \beta \alpha_n \\ | \gamma_1 \\ | \dots \\ | \gamma_m \end{array}$$

$$\begin{array}{l} A \rightarrow \beta A' \\ | \gamma_1 \\ | \dots \\ | \gamma_m \\ \\ A' \rightarrow \alpha_1 \\ | \dots \\ | \alpha_n \end{array}$$

Left Factoring (EBNF)

$$\begin{array}{l} A \rightarrow \beta \alpha_1 \\ | \dots \\ | \beta \alpha_n \\ | \gamma_1 \\ | \dots \\ | \gamma_m \end{array}$$

Extended BNF:

$$A \rightarrow \beta (\alpha_1 | \dots | \alpha_n) | \gamma_1 | \dots | \gamma_m$$
$$A \rightarrow \beta A' | \gamma_1 | \dots | \gamma_m$$
$$A' \rightarrow \alpha_1 | \dots | \alpha_n$$

Exercises

$$\begin{aligned} P &\rightarrow \text{if } C \text{ then } P \text{ endif} \\ &\quad | \text{if } C \text{ then } P \text{ else } P \text{ endif} \\ &\quad | p \\ C &\rightarrow c \end{aligned}$$
$$\begin{aligned} I &\rightarrow LE \text{ ':=' } E \\ &\quad | \text{write '(' } E \text{ ')'} \\ &\quad | \text{id '(' } E \text{ ')'} \\ E &\rightarrow \text{id} \\ &\quad | \text{num} \\ LE &\rightarrow \text{id} \end{aligned}$$

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	...	a	...	a_n	\$
A_1						
⋮						
A			$A \rightarrow \alpha_k$			
⋮						
A_m						

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
- Recursive predictive parsers
 - Parsing algorithm is formed by a set of mutually recursive functions
 - Each rule $A \rightarrow \alpha$ generates the code of its function

```
void A(void) {  
    // Code generated from  $\alpha$   
}
```
 - Gencode describes how to translate a rule to the associated function

Table-driven Top-down Parser

M	a_1	...	a	...	a_n	$\$$
A_1						
\vdots						
A			$A \rightarrow \alpha_k$			
\vdots						
A_m						

$A \rightarrow \alpha_1$
 \vdots
 α_k
 \vdots
 α_o

$A \rightarrow \alpha_k \in M[A, a]$ 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 $A \rightarrow \alpha \in G$ do

add $A \rightarrow \alpha$ 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 :

$$E \rightarrow E + T$$

$$| T$$

$$T \rightarrow T * F$$

$$| F$$

$$F \rightarrow \text{id}$$

$$| (E)$$

Grammar G_2 :

$$E \rightarrow T E'$$

$$E' \rightarrow + T E'$$

$$| \epsilon$$

$$T \rightarrow F T'$$

$$T' \rightarrow * F T'$$

$$| \epsilon$$

$$F \rightarrow \text{id}$$

$$| (E)$$

Exercises

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

Grammar G_3 :

$$\begin{array}{l} P \rightarrow \text{if } C \text{ then } P \\ \quad | \text{if } C \text{ then } P \text{ else } P \\ \quad | p \\ C \rightarrow c \end{array}$$

Grammar G_4 :

$$\begin{array}{l} P \rightarrow \text{if } C \text{ then } P P' \\ \quad | p \\ P' \rightarrow \epsilon \\ \quad | \text{else } P \\ C \rightarrow c \end{array}$$

Exercises

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

Grammar G_5 :

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

Table-driven Top-down Parser Algorithm

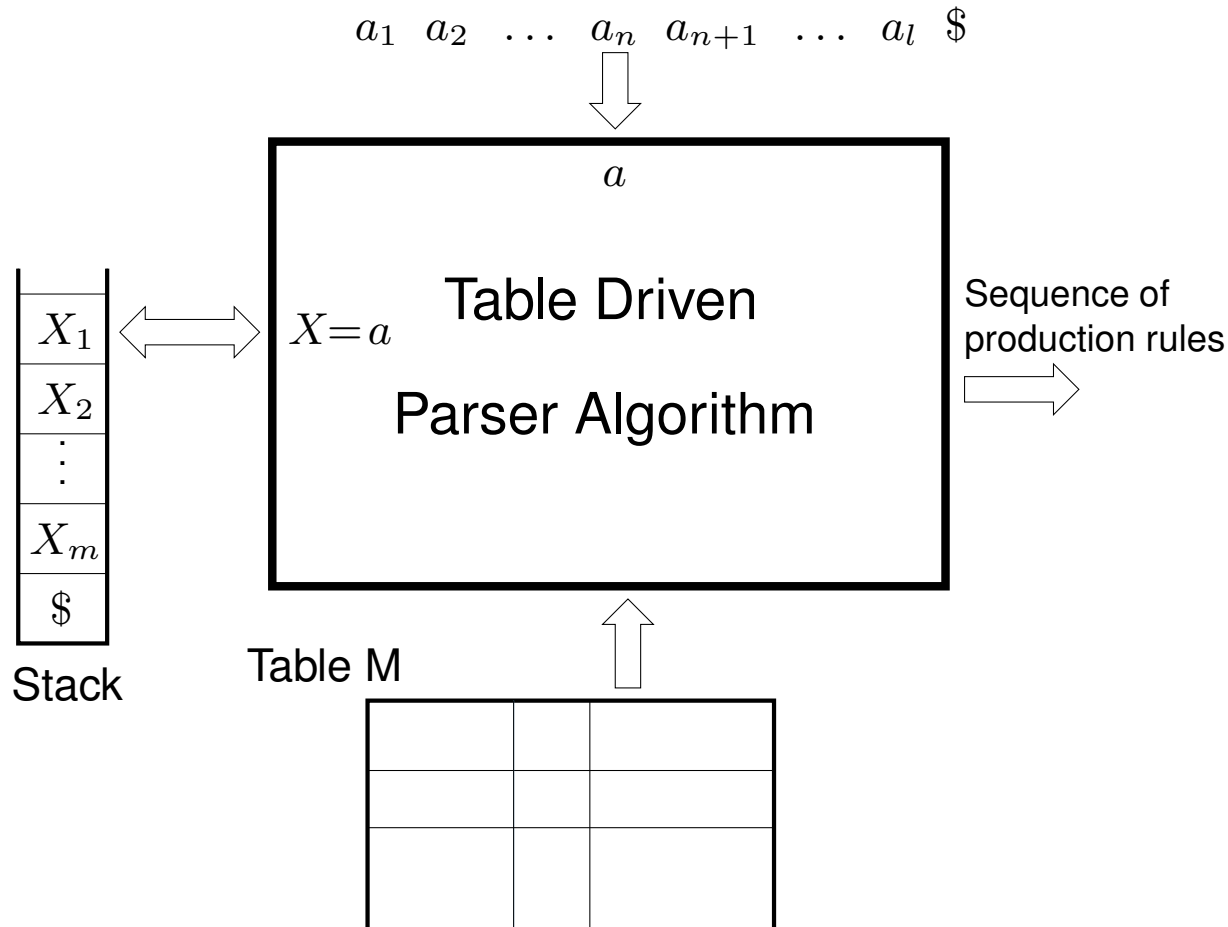


Table-driven Top-down Parser Algorithm

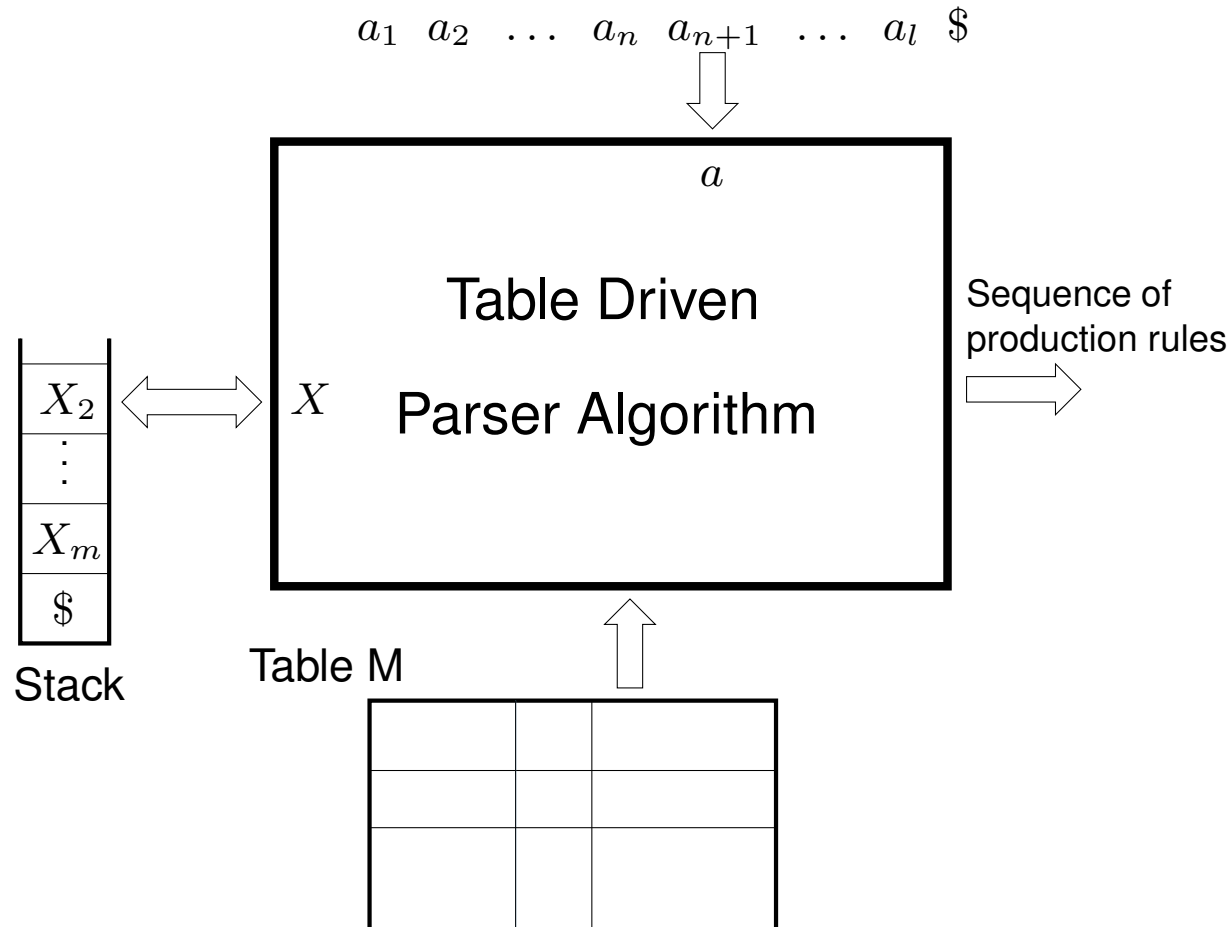


Table-driven Top-down Parser Algorithm

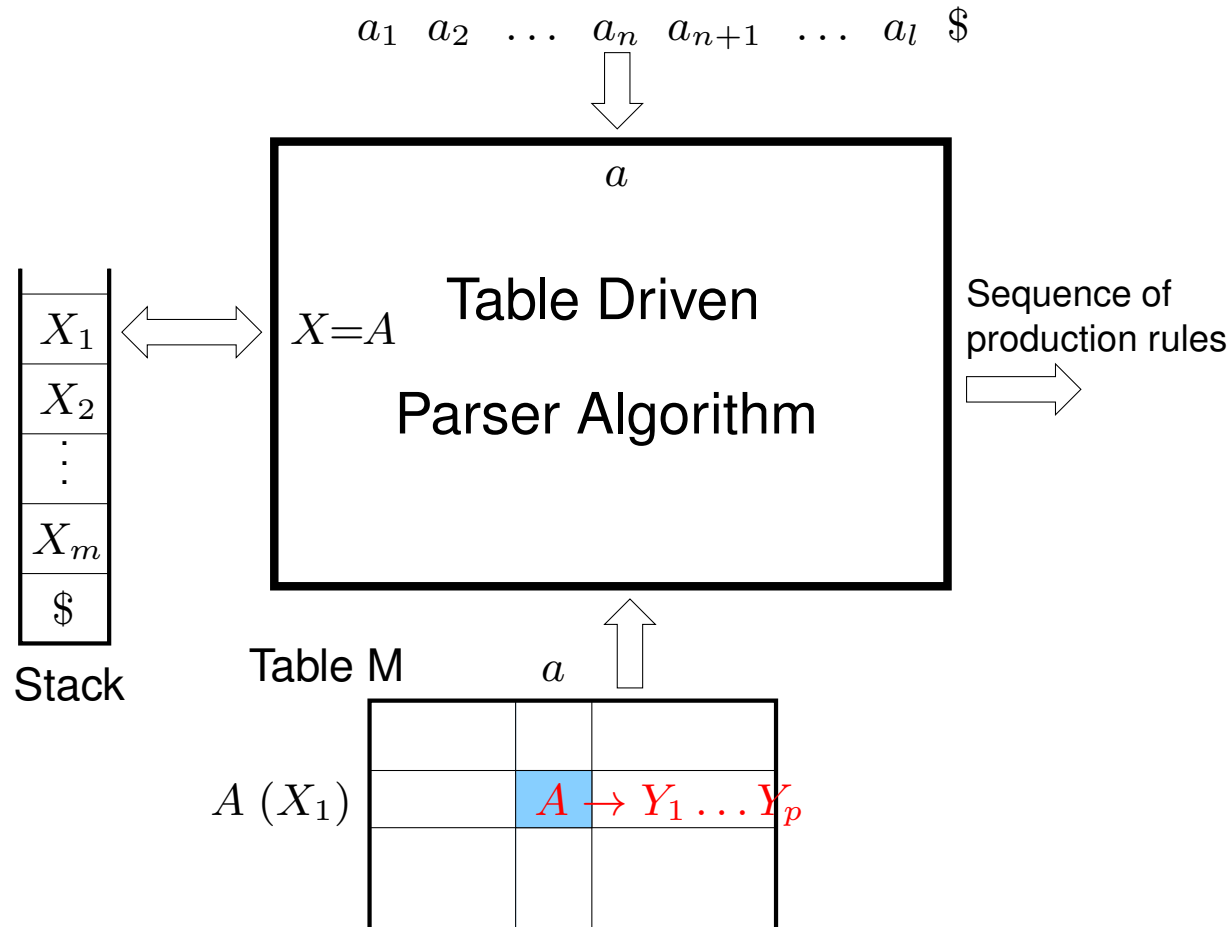


Table-driven Top-down Parser Algorithm

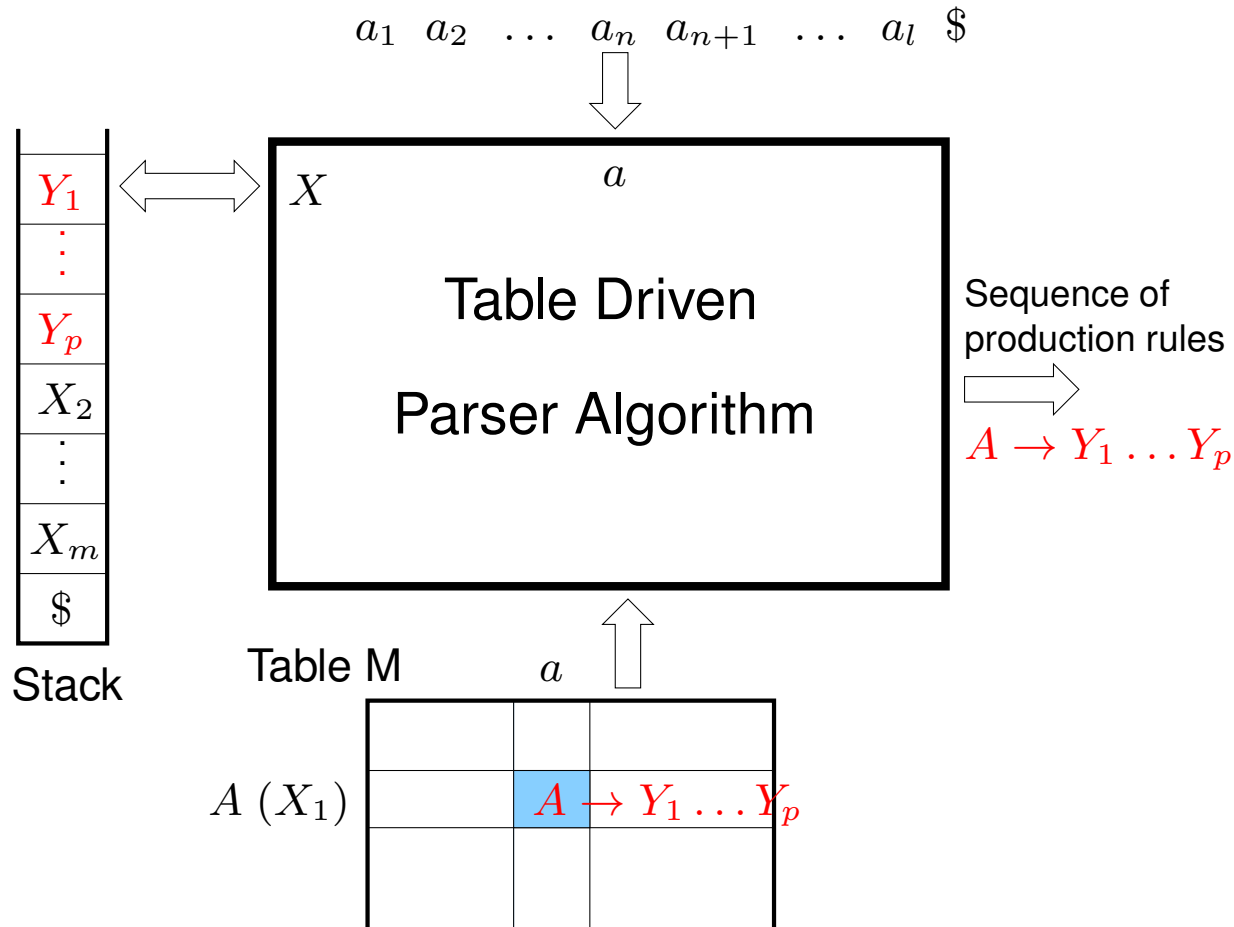


Table-driven Top-down Parser Algorithm

```
Stk := EmptyStack(); PushStack(Stk, $); PushStack(Stk, S);  
X := TopStack(Stk); a := FirstToken();  
while X ≠ $ 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 → Y1 ... Yp  
            emit production rule X → Y1 ... Yp  
            PopStack(Stk); for i := p downto 1 do PushStack(Stk, Yi);  
        X := TopStack(Stk);  
endwhile
```