

More Elimination of Left-Recursion

- In general

$$S \rightarrow S \alpha_1 \mid \dots \mid S \alpha_n \mid \beta_1 \mid \dots \mid \beta_m$$

- All strings derived from S start with one of β_1, \dots, β_m and continue with several instances of $\alpha_1, \dots, \alpha_n$

- Rewrite as

$$\begin{aligned} S &\rightarrow \beta_1 S' \mid \dots \mid \beta_m S' \\ S' &\rightarrow \alpha_1 S' \mid \dots \mid \alpha_n S' \mid \varepsilon \end{aligned}$$

Contd., Removal of left recursion due to many productions ...

- Left recursion may also be introduced by two or more grammar rules.

- For example: $S \rightarrow Aa \mid b$
 $A \rightarrow Ac \mid Sd \mid \epsilon$

$$\begin{aligned} S &\rightarrow Aa \mid b \\ \cancel{A} &\rightarrow \cancel{Ac} \mid \cancel{Sd} \mid \epsilon \\ A &\rightarrow \cancel{Ac} \mid \cancel{Aad} \mid \cancel{bd} \mid \epsilon \\ &\quad \quad \quad \alpha_1 \quad \quad \quad \alpha_2 \quad \quad \quad \beta \\ A &\rightarrow bdA' \mid A' \\ A' &\rightarrow CA' \mid adA' \end{aligned}$$

- After the first step (substitute S by its RHS in the rules) the grammar becomes

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Aad \mid bd \mid \epsilon$$

Removal of left recursion due to many productions ...

- Split productions into α and β
- General rule: If we have $A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid \beta_1 \mid \beta_2 \mid \dots$

Then rewrite as:

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

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Aad \mid bd \mid \epsilon$$

- **Left-recursive parts (α):**
 - $Ac \rightarrow \alpha_1 = c$
 - $Aad \rightarrow \alpha_2 = ad$
- **Non-left-recursive parts (β):**
 - $(\beta_1) bd$
 - $(\beta_2) \epsilon$

Removal of left recursion due to many productions ...

- Rewrite A without left recursion

$$A \rightarrow bd A' \mid A'$$

$$A' \rightarrow c A' \mid ad A' \mid \varepsilon$$

- After the second step (removal of left recursion) the grammar becomes

$$S \rightarrow Aa \mid b$$

$$A \rightarrow bdA' \mid A'$$

$$A' \rightarrow cA' \mid adA' \mid \epsilon$$

Left Factoring

- In top-down parsing when it is not clear which production to choose for expansion of a symbol

Defer (accept) the decision till we have seen enough input

- In general, if $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$
- Defer (Accept) decision by expanding $A \rightarrow \alpha A'$
- we can then expand $A' \rightarrow \beta_1 \text{ or } \beta_2$

Therefore $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$ $\rightarrow\rightarrow$ transforms to $\rightarrow\rightarrow$

$$\begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta_1 \mid \beta_2 \end{array}$$

Here:

- α = the common prefix
- A' = a new non-terminal that handles the “rest of the options.”

Dangling else problem again

- Dangling else problem can be handled by left factoring

stmt \rightarrow if expr then stmt else stmt
 | if expr then stmt

can be transformed to

stmt \rightarrow if expr then stmt S'
S' \rightarrow else stmt | ϵ

Left Factoring:

1. The common prefix α is if (Expr) Stmt.
2. The varying suffixes β_1 is else Stmt and β_2 is ϵ (the empty string, meaning "nothing").

Rule : $A \rightarrow \alpha A'$
 $A' \rightarrow \beta_1 \mid \beta_2$

Summary of Recursive Descent

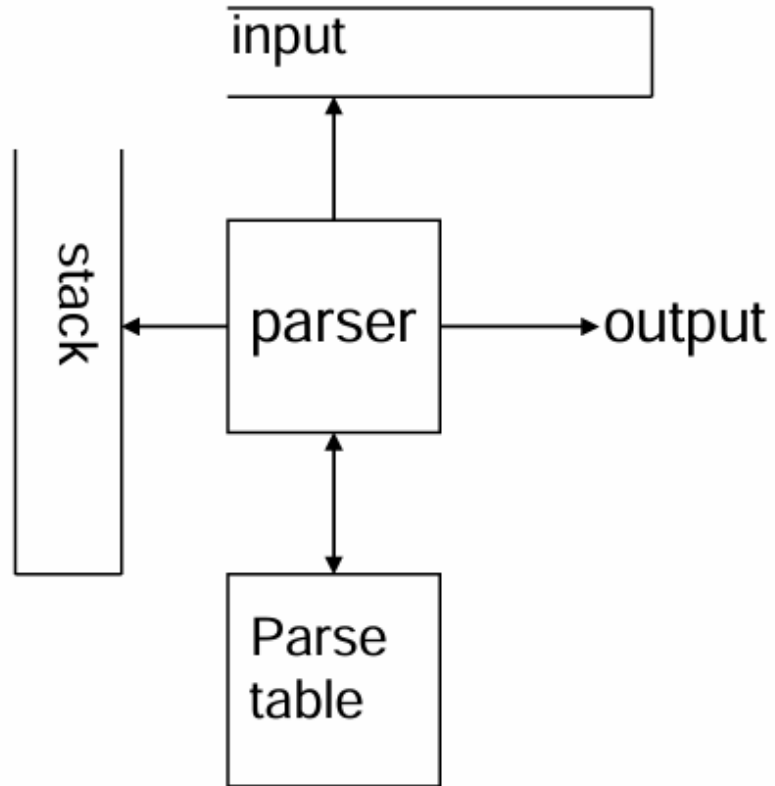
- Simple and general parsing strategy Left-recursion must be eliminated first ... but that can be done automatically
- Unpopular because of backtracking
- Thought to be too inefficient
- In practice, backtracking is eliminated by restricting the grammar

Predictive Parsers

- Like recursive-descent but parser can “predict” which production to use
 - By looking at the next few tokens
 - No backtracking
- Predictive parsers accept **LL(k)** grammars
 - **L** means “left-to-right” scan of input
 - **L** means “leftmost derivation”
 - **k** means “predict based on k tokens of lookahead”
- In practice, LL(1) is used

Predictive parsing Contd.,

- Predictive parser can be implemented by maintaining an external stack



Parse table is a two dimensional array $M[X,a]$ where “X” is a non terminal and “a” is a terminal of the grammar

Parsing algorithm

- The parser considers 'X' the symbol on top of stack, and 'a' the current input symbol
- ~~These two symbols determine the action to be taken by the parser~~
- Assume that '\$' is a special token that is at the bottom of the stack and terminates the input string.

if $X = a = \$$ then halt

if $X = a \neq \$$ then pop(x) and ip++

if X is a non terminal

 then if $M[X,a] = \{X \rightarrow UVW\}$

 then begin pop(X); push(W,V,U)

 end

 else error

LL(1) Languages

- In recursive-descent, for each non-terminal and input token there may be a choice of production
- LL(1) means that for each non-terminal and token there is only one production
- Can be specified via 2D tables
 - One dimension for current non-terminal to expand
 - One dimension for next token
 - A table entry contains one production

Contd., with Left Factoring

- Formula of left factoring

$A \rightarrow \alpha\beta_1 | \alpha\beta_2$ can be written as $A \rightarrow \alpha A', \quad A' \rightarrow \beta_1 | \beta_2$

Here:

- α = the common prefix
- A' = a new non-terminal that handles the “rest of the options.”

Apply to Expression Grammar: $E \rightarrow T + E \mid T$

Here, both productions start with **common prefix T**.

So, $\alpha = T$, $\beta_1 = +E$, $\beta_2 = \epsilon$.

After factoring

$E \rightarrow TX,$
 $X \rightarrow +E | \epsilon$

Here X, Y are added as new Non-terminals

Apply to Term Grammar $T \rightarrow (E) \mid \text{int} \mid \text{int} * T$

Here, **int** is the common prefix.

So $\alpha = \text{int}$, and $\beta_1 = *T$, $\beta_2 = \epsilon$

$T \rightarrow (E) \mid \text{int}Y$

$Y \rightarrow *T | \epsilon$

final grammar became:

$E \rightarrow TX,$ $X \rightarrow +E \mid \epsilon$
 $T \rightarrow (E) \mid \text{int}Y,$ $Y \rightarrow *T \mid \epsilon$

Left-Factoring Example

- Recall the grammar

$$E \rightarrow T + E \mid T$$

$$T \rightarrow (E) \mid \text{int} \mid \text{int} * T$$

- Factor out common prefixes of productions

$$E \rightarrow T X$$

$$X \rightarrow + E \mid \varepsilon$$

$$T \rightarrow (E) \mid \text{int} Y$$

$$Y \rightarrow * T \mid \varepsilon$$

LL(1) Parsing Table Example

- Left-factored grammar

$$E \rightarrow TX$$

$$T \rightarrow (E) \mid \text{int } Y$$

$$X \rightarrow +E \mid \varepsilon$$

$$Y \rightarrow *T \mid \varepsilon$$

- The LL(1) parsing table:

	int	*	+	()	\$
E	TX			TX		
X			$+E$		ε	ε
T	$\text{int } Y$			(E)		
Y		$*T$	ε		ε	ε

Contd.,

$$E \rightarrow TX$$

$$T \rightarrow (E) \mid int Y$$

$$X \rightarrow +E \mid \varepsilon$$

$$Y \rightarrow *T \mid \varepsilon$$

FIRST sets:

- $FIRST(E) = FIRST(T) = \{ (, int \}$
- $FIRST(T) = \{ (, int \}$
- $FIRST(X) = \{ +, \varepsilon \}$
- $FIRST(Y) = \{ *, \varepsilon \}$

FOLLOW sets:

- $FOLLOW(E) = \{), \$ \}$
- $FOLLOW(T) = FIRST(X) \cup FOLLOW(E) = \{ +,), \$ \}$
- $FOLLOW(X) = FOLLOW(E) = \{), \$ \}$
- $FOLLOW(Y) = FOLLOW(T) = \{ +,), \$ \}$

	int	*	+	()	\$
E	TX			TX		
X			$+E$		ε	ε
T	$int Y$			(E)		
Y		$*T$	ε		ε	ε