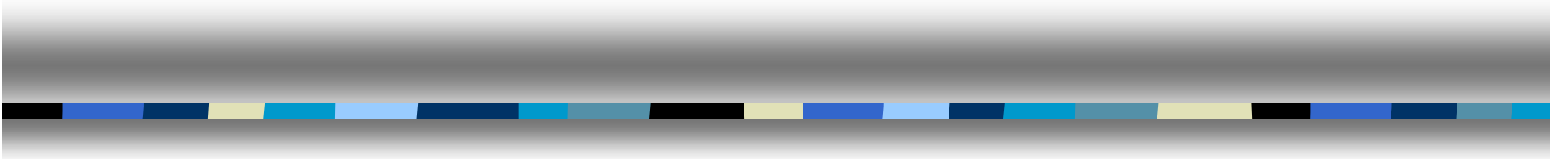


# Compiler Design





# Lecture-6

## Predictive Parsing



# Topics Covered

- n Introduction to Predictive Parsing
- n How we make Parsing Table.



# Predictive parsing

- n Recall the main idea of top-down parsing:
  - Start at the root, grow towards leaves
  - Pick a production and try to match input
  - May need to backtrack
- n Can we avoid the backtracking?
  - Given  $A \rightarrow \alpha \mid \beta$  the parser should be able to choose between  $\alpha$  and  $\beta$
- n How?
  - What if we do some "preprocessing" to answer the question: Given a non-terminal  $A$  and lookahead  $t$ , which (if any) production of  $A$  is guaranteed to start with a  $t$ ?



# Predictive parsing

- n If we have two productions:  $A \rightarrow \alpha \mid \beta$ , we want a distinct way of choosing the correct one.
- n Define:
  - for  $\alpha \in G$ ,  $x \in \text{FIRST}(\alpha)$  iff  $\alpha \Rightarrow^* x\gamma$
- n If  $\text{FIRST}(\alpha)$  and  $\text{FIRST}(\beta)$  contain no common symbols, we will know whether we should choose  $A \rightarrow \alpha$  or  $A \rightarrow \beta$  by looking at the lookahead symbol.



# Predictive parsing

- n Compute  $\text{FIRST}(X)$  as follows:
  - if  $X$  is a terminal, then  $\text{FIRST}(X) = \{X\}$
  - if  $X \rightarrow \varepsilon$  is a production, then add  $\varepsilon$  to  $\text{FIRST}(X)$
  - if  $X$  is a non-terminal and  $X \rightarrow Y_1 Y_2 \dots Y_n$  is a production, add  $\text{FIRST}(Y_i)$  to  $\text{FIRST}(X)$  if the preceding  $Y_j$ s contain  $\varepsilon$  in their  $\text{FIRST}$ s



# Predictive parsing

- n What if we have a "candidate" production  $A \rightarrow \alpha$  where  $\alpha = \varepsilon$  or  $\alpha \Rightarrow^* \varepsilon$ ?
- n We could expand if we knew that there is some sentential form where the current input symbol appears after  $A$ .
- n Define:
  - for  $A \in N$ ,  $x \in \text{FOLLOW}(A)$  iff  $\exists S \Rightarrow^* \alpha A x \beta$



# Predictive parsing

- n Compute FOLLOW as follows:
  - FOLLOW(S) contains EOF
  - For productions  $A \rightarrow \alpha B \beta$ , everything in FIRST( $\beta$ ) except  $\epsilon$  goes into FOLLOW(B)
  - For productions  $A \rightarrow \alpha B$  or  $A \rightarrow \alpha B \beta$  where FIRST( $\beta$ ) contains  $\epsilon$ , FOLLOW(B) contains everything that is in FOLLOW(A)





# Predictive parsing

n Armed with

- FIRST

- FOLLOW

n we can build a parser where **no** backtracking is required!



# Predictive parsing (w/table)

- n For each production  $A \rightarrow \alpha$  do:
  - For each terminal  $\mathbf{a} \in \text{FIRST}(\alpha)$  add  $A \rightarrow \alpha$  to entry  $M[A, \mathbf{a}]$
  - If  $\varepsilon \in \text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to entry  $M[A, \mathbf{b}]$  for each terminal  $\mathbf{b} \in \text{FOLLOW}(A)$ .
  - If  $\varepsilon \in \text{FIRST}(\alpha)$  and  $\text{EOF} \in \text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, \text{EOF}]$
- n Use table and stack to simulate recursion.



# Recursive Descent Parsing

## n Basic idea:

- Write a routine to recognize each lhs
- This produces a parser with mutually recursive routines.
- Good for hand-coded parsers.

## n Example:

- $A \rightarrow aB \mid b$  will correspond to

```
A() {  
    if (lookahead == 'a')  
        match('a');  
    B();  
    else if (lookahead == 'b')  
        match ('b');  
    else error();  
}
```



# Building a parser

n Original grammar:

$$\begin{aligned} E &\rightarrow E + E \\ E &\rightarrow E * E \\ E &\rightarrow (E) \\ E &\rightarrow \text{id} \end{aligned}$$

n This grammar is left-recursive, ambiguous and requires left-factoring. It needs to be modified before we build a predictive parser for it:

Remove ambiguity:

$$\begin{aligned} E &\rightarrow E + T \\ T &\rightarrow T * F \\ F &\rightarrow (E) \\ F &\rightarrow \text{id} \end{aligned}$$

Remove left recursion:

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \varepsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \varepsilon \\ F &\rightarrow (E) \\ F &\rightarrow \text{id} \end{aligned}$$



# Building a parser

The grammar:

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \varepsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \varepsilon \\ F &\rightarrow (E) \\ F &\rightarrow \mathbf{id} \end{aligned}$$
$$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{(, \mathbf{id}\}$$
$$\text{FIRST}(E') = \{+, \varepsilon\}$$
$$\text{FIRST}(T') = \{*, \varepsilon\}$$
$$\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{\$, \}$$
$$\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{+, \$, \}$$
$$\text{FOLLOW}(F) = \{*, +, \$, \}$$

Now, we can either build a table or design a recursive descend parser.

# Parsing table

	+	*	(	)	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			$F \rightarrow (E)$		<b><math>F \rightarrow id</math></b>	
+	match					
*		match				
(			match			
)				match		
id					match	
\$						accept

# Parsing table

Parse the input  $id*id$  using the parse table and a stack

Step	Stack	Input	Next Action
1	\$E	id*id\$	$E \rightarrow TE'$
2	\$E'T	id*id\$	$T \rightarrow FT'$
3	\$E'T'F	id*id\$	$F \rightarrow id$
4	\$E'T'id	id*id\$	match <b>id</b>
5	\$E'T'	*id\$	$T' \rightarrow *FT'$
6	\$T'F*	*id\$	match *
7	\$T'F	id\$	$F \rightarrow id$
8	\$T'id	id\$	match <b>id</b>
9	\$T'	\$	$T' \rightarrow \epsilon$
10	\$	\$	accept

# Recursive descend parser

```
parse() {  
    token = get_next_token();  
    if (E() and token == '$')  
    then return true  
    else return false  
}
```

```
E() {  
    if (T())  
    then return Eprime()  
    else return false  
}
```

```
Eprime() {  
    if (token == '+')  
    then token=get_next_token()  
    if (T())  
    then return Eprime()  
    else return false  
    else if (token==' ' or token=='$')  
    then return true  
    else return false  
}
```

The remaining procedures are similar.





# LL(1) parsing

- n Our parser scans the input **L**eft-to-right, generates a **L**eftmost derivation and uses **1** symbol of lookahead.
- n It is called an LL(1) parser.
- n If you can build a parsing table with no multiply-defined entries, then the grammar is LL(1).
- n Ambiguous grammars are never LL(1)
- n Non-ambiguous grammars are not necessarily LL(1)



# LL(1) parsing

- n For example, the following grammar will have two possible ways to expand  $S'$  when the lookahead is **else**.

```
S → if E then S S' | other
S' → else S | ε
E → id
```

- It may expand  $S' \rightarrow \mathbf{else} S$  or  $S' \rightarrow \varepsilon$
- We can resolve the ambiguity by instructing the parser to always pick  $S' \rightarrow \mathbf{else} S$ . This will match each **else** to the closest previous **then**.



# LL(1) parsing

- n Here's an example of a grammar that is NOT LL(k) for any k:

$$\begin{array}{l} S \rightarrow Ca \mid Cb \\ C \rightarrow cC \mid c \end{array}$$

- Why? Suppose the grammar was LL(k) for some k. Consider the input string  $c^{k+1}a$ . With only k lookaheads, the parser would not be able to decide whether to expand using  $S \rightarrow Ca$  or  $S \rightarrow Cb$
- Note that the grammar is actually regular: it generates strings of the form  $c^+(a|b)$



# Error detection in LL(1) parsing

- n An error is detected whenever an empty table slot is encountered.
- n We would like our parser to be able to recover from an error and continue parsing.
- n Phase-level recovery
  - We associate each empty slot with an error handling procedure.
- n Panic mode recovery
  - Modify the stack and/or the input string to try and reach state from which we can continue.



# Error recovery in LL(1) parsing

## n Panic mode recovery

### – Idea:

- Decide on a set of synchronizing tokens.
- When an error is found and there's a nonterminal at the top of the stack, discard input tokens until a synchronizing token is found.
- Synchronizing tokens are chosen so that the parser can recover quickly after one is found
  - e.g. a semicolon when parsing statements.
- If there is a terminal at the top of the stack, we could try popping it to see whether we can continue.
  - Assume that the input string is actually missing that terminal.



# Error recovery in LL(1) parsing

## n Panic mode recovery

- Possible synchronizing tokens for a nonterminal A
  - the tokens in  $FOLLOW(A)$ 
    - When one is found, pop A of the stack and try to continue
  - the tokens in  $FIRST(A)$ 
    - When one is found, match it and try to continue
  - tokens such as semicolons that terminate statements