

CSE401: Parsing (C)

David Notkin
Autumn 2000

UW CSE401 AQ 2000 • D. Notkin • All rights reserved • Parsing B • Slide 2

Objectives: today

- Recap and clarify PREDICT table
- Describe computation of FIRST and FOLLOW
 - And the relationship to PREDICT
- Recursive descent parsing
 - High-level issues and
 - (time-permitting) a walk through the PL/0 parser

UW CSE401 AQ 2000 • D. Notkin • All rights reserved • Parsing B • Slide 3

Next week

- Monday: no class
- Wednesday: Mark Seigle
- Friday: Matthai Philipose

UW CSE401 AQ 2000 • D. Notkin • All rights reserved • Parsing B • Slide 4

Example

```

1 Stmt ::= 1 if expr then Stmt else Stmt |
2 while Expr do Stmt |
3 begin Stmts end
4 Stmts ::= 4 Stmt ; Stmts | 5 ε
6 Expr ::= 6 id
    
```

	if	then	else	while	do	begin	end	id	;
Stmt	1			2		3			
Stmts	4			4		4	5		
Expr								6	

UW CSE401 AQ 2000 • D. Notkin • All rights reserved • Parsing B • Slide 5

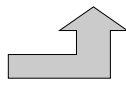
```

if id then
  begin end
else
  while id do
    begin end
  end
    
```

```

> Stmt
1. if Expr then Stmt else Stmt end
2. if id then Stmt else Stmt end
3. if id then begin Stmts end
   else Stmt end
4. if id then begin end
   else while Expr do Stmt end
5. if id then begin end
   else while id do
     begin Stmts end end
6. if id then begin end
   else while id do
     begin end end
    
```

Production applied



UW CSE401 AQ 2000 • D. Notkin • All rights reserved • Parsing B • Slide 6

Compute PREDICT

In terms of FIRST and FOLLOW

- FIRST(α)
 - ✓ α is any string of tokens
 - ✓ FIRST(α) is the set of terminals that being strings derived from α
 - ✓ If α can derive ϵ , then include ϵ in FIRST
- FOLLOW(A)
 - ✓ A is a non-terminal
 - ✓ FOLLOW(A) is the set of terminals that can appear immediately after A in some derivation in the language
 - ✓ If A can be the end of such a derivation, then include \$ in FOLLOW
 - ✓ \$ represents "end of input"

Computing FIRST(X) for all grammar symbols X

- If X is a terminal, then FIRST(X) is {X}
- If $X ::= \epsilon$, then add ϵ to FIRST(X)
- If X is a non-terminal and $X ::= Y_1 Y_2 \dots Y_k$
 - FIRST(X) is initialized to FIRST(Y_1)
 - If $\epsilon \in \text{FIRST}(Y_1)$ then union in FIRST(Y_2)
 - If $\epsilon \in \text{FIRST}(Y_1)$ and $\epsilon \in \text{FIRST}(Y_2)$ then union in FIRST(Y_3)
 - Etc.
 - If $\epsilon \in \text{FIRST}(Y_i)$ for all i, then add ϵ to FIRST(X)

Computing FIRST($X_1 X_2 \dots X_n$)

- Start with all the non- ϵ symbols in FIRST(X_1)
- If $\epsilon \in \text{FIRST}(X_1)$, then union in all the non- ϵ symbols in FIRST(X_2)
- If $\epsilon \in \text{FIRST}(X_1)$ and $\epsilon \in \text{FIRST}(X_2)$, then union in all the non- ϵ symbols in FIRST(X_3)
- Etc.
- If $\epsilon \in \text{FIRST}(X_i)$ for all i, then add ϵ to FIRST($X_1 X_2 \dots X_n$)

Computing FOLLOW(A) for all non-terminals A

- Place \$ in FOLLOW(S), where S is the start symbol and \$ is the end-of-input marker
- If there is a production $A ::= \alpha B \beta$ then everything in FIRST(β) except for ϵ is placed in FOLLOW(B)
- If there is a production $A ::= \alpha B$ (or $A ::= \alpha B \beta$ where $\epsilon \in \text{FIRST}(\beta)$) then everything in FOLLOW(A) is placed in FOLLOW(B)
- Apply until a fixpoint is reached

	FIRST	FOLLOW
$S ::= \text{if } E \text{ then } S \text{ else } S$		
while E do S		
begin Ss end		
$Ss ::= S ; Ss$		
ϵ		
$E ::= \text{id}$		

Another example

$E ::= T \{ + T \}$
 $T ::= F \{ * F \}$
 $F ::= - F \mid \text{id} \mid (E)$

Sugared

$E ::= \mathbf{1} T E'$
 $E' ::= \mathbf{2} + T E' \mid \mathbf{3} \epsilon$
 $T ::= \mathbf{4} F T'$
 $T' ::= \mathbf{5} * F T' \mid \mathbf{6} \epsilon$
 $F ::= \mathbf{7} - F \mid \mathbf{8} \text{id} \mid \mathbf{9} (E)$

Unsugared

Construct predictive parsing table

- Input is a grammar G
- Output is a table M indexed by non-terminals (rows) and terminals (columns)
- ```

 foreach production $A ::= \alpha$ in G do
 foreach terminal a in FIRST(α)
 add the production to M[A,a];
 if $\epsilon \in \text{FIRST}(\alpha)$
 add the production to M[A,b] for
 each terminal b in FOLLOW(A)
 if $\epsilon \in \text{FIRST}(\alpha)$ and $\$ \in \text{FOLLOW}(A)$
 add the production to M[A,$]

```

 Every remaining undefined entry is an error

## PREDICT

|    | id | + | - | * | / | ( | ) |
|----|----|---|---|---|---|---|---|
| E  |    |   |   |   |   |   |   |
| E' |    |   |   |   |   |   |   |
| T  |    |   |   |   |   |   |   |
| T' |    |   |   |   |   |   |   |
| F  |    |   |   |   |   |   |   |

## PREDICT and LL(1)

- If the PREDICT table has at most one entry in each cell, then the grammar is LL(1)
  - And there is only one choice (it's predictive) , making it fast to parse and easy to implement
- Multiple entries in a cell
  - Arise with left recursion, ambiguity, common prefixes, etc.
  - Can patch by hand
  - Or use more powerful parsing techniques

## Recursive descent parsers

- Write procedure for each non-terminal
- Each procedure selects the correct right-hand side by peeking at the input tokens
- Then the r.h.s. is consumed
  - If it's a terminal symbol, verify it is next and then advance through the token stream
  - If it's a non-terminal, call corresponding procedure
- Construct and return AST representing the r.h.s.

## It's demo time...

- Let's look at some of the PL/0 code to see how the recursive descent parsing works in practice