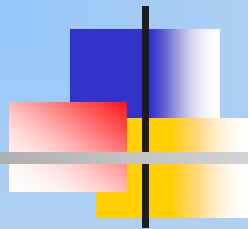
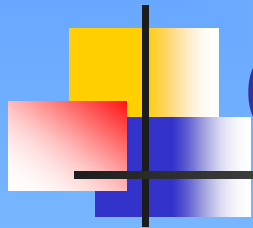


CSE302: Compiler Design



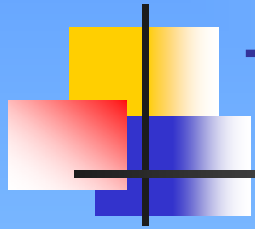
Instructor: Dr. Liang Cheng
Department of Computer Science and Engineering
P.C. Rossin College of Engineering & Applied Science
Lehigh University

March 13, 2007



Outline

- Recap
 - Top-down parsing (Section 4.4)
- Summary and homework



Top-Down Parsing

- Finding a **leftmost** derivation for an input string
 - At each step the key problem is determining the production to be applied for a nonterminal, say A
 - Recursive-descent parsing
 - Predictive parsing for LL(1) grammars



Recursive-Descent Parsing

- `void A() {`

- `record input pointer;`

- Choose an A -production, $A \rightarrow X_1 X_2 \dots X_n$

- for ($i=1$ to n) {

- if (X_i is a nonterminal) call $X_i()$;

- else if (X_i equals the current input a)

- advance the input to the next symbol;

- else /* an error occurred, **backtrack** */

- }

- }

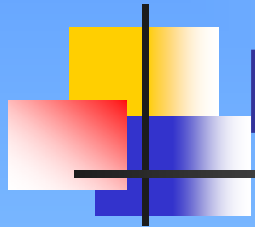
- Grammar

- $S \rightarrow cAd$

- $A \rightarrow ab \mid a$

- Input string

- cad



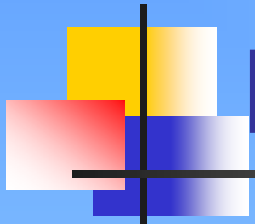
Predictive Parsers

- Recursive-descent parsers with one input symbol lookahead that requires no backtracking
 - No backtracking means deterministic in choosing a production
 - Can be constructed for a class of grammars called LL(1)
 - 1st L: scanning the input from left to right
 - 2nd L: producing a leftmost derivation



FIRST and FOLLOW Sets

- $\text{FIRST}(\alpha)$ is the set of **terminals** that begin strings derived from α
- FOLLOW sets for **ALL** nonterminals A
 - Place $\$$ in $\text{FOLLOW}(S)$, where S is the start symbol, and $\$$ is the input right endmarker
 - $\$$ is not a symbol of any grammar
 - If there is a production $A \rightarrow \alpha B \beta$, then everything in $\text{FIRST}(\beta)$ except ε is in $\text{FOLLOW}(B)$
 - If there is a production $A \rightarrow \alpha B$, or a production $A \rightarrow \alpha B \beta$, where $\text{FIRST}(\beta)$ contains ε (i.e. $\beta \Rightarrow \varepsilon$), then everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(B)$
 - Whatever followed A must follow B



LL(1) Grammars

- Whenever $A \rightarrow \alpha$ and $A \rightarrow \beta$ are two distinct A -productions of G , the following conditions hold
 - For no terminal a do both α and β derive strings beginning with a
 - At most one of α and β can derive the empty string
 - If $\beta \overset{*}{\Rightarrow} \varepsilon$, then α does not derive any string beginning with a terminal in $\text{FOLLOW}(A)$
 - If $\alpha \overset{*}{\Rightarrow} \varepsilon$, then β does not derive any string beginning with a terminal in $\text{FOLLOW}(A)$



Why Such Conditions?

- In top-down parsing
 - At each step the key problem is determining the production to be applied for a nonterminal, say A
- $S \xRightarrow{*} \gamma A \lambda$
 I_m
- $A \rightarrow \alpha$ and $A \rightarrow \beta$
 - $FIRST(\alpha)$ and $FIRST(\beta)$ should be disjoint sets
 - If a is in $FIRST(\beta)$ then choose $A \rightarrow \beta$
 - $S \xRightarrow{*} \gamma A \lambda \Rightarrow \gamma \beta \lambda$
 - Except the following possibility in applying $A \rightarrow \alpha$
 - $S \xRightarrow{*} \gamma A \lambda \Rightarrow \gamma \alpha \lambda \xRightarrow{*} \gamma \lambda$
Thus we should eliminate such possibility for deterministic decision of choosing a A -production
 - If ϵ is in $First(\alpha)$, then $FOLLOW(A)$ and $FIRST(\beta)$ should be disjoint sets



Predictive Parsing For LL(1) Grammar

- The production $A \rightarrow \alpha$ is chosen if
 - The next input symbol a is in $\text{FIRST}(\alpha)$
 - The next input symbol a (or $\$$) is in $\text{FOLLOW}(A)$ and ϵ is in $\text{FIRST}(\alpha)$
 - The next symbol could be $\$$
- Thus we should construct a parsing table M where $M[A, a] = A \rightarrow \alpha$
 - In function A if the input is a , then call functions and/or match terminals of α



Constructing A Predictive Parsing Table M For **ANY** Grammar G

- For **each** production $A \rightarrow \alpha$
 - For each terminal a in $\text{FIRST}(A)$, add $A \rightarrow \alpha$ to $M[A, a]$
 - If ϵ is in $\text{FIRST}(\alpha)$, then for each terminal b in $\text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, b]$
 - If ϵ is in $\text{FIRST}(\alpha)$ and if $\$$ is in $\text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, \$]$
- If, after performing the above, there is no production at all in $M[A, a]$, then set $M[A, a]$ to **error**

Non-recursive Predictive Parsing

- A stack storing symbols, initialized with $\$S$
- An input buffer with an input pointer ip
- A parsing table M for grammar G

Point ip to the 1st input symbol

Set A to the top stack symbol

while($A \neq \$$) {

if (A is a) pop stack; advance ip

else if (A is a terminal) error();

else if ($M[A,a]$ is an error entry) error();

else if ($M[A,a] = A \rightarrow X_1 X_2 \dots X_k$) {

 output the production or other actions;

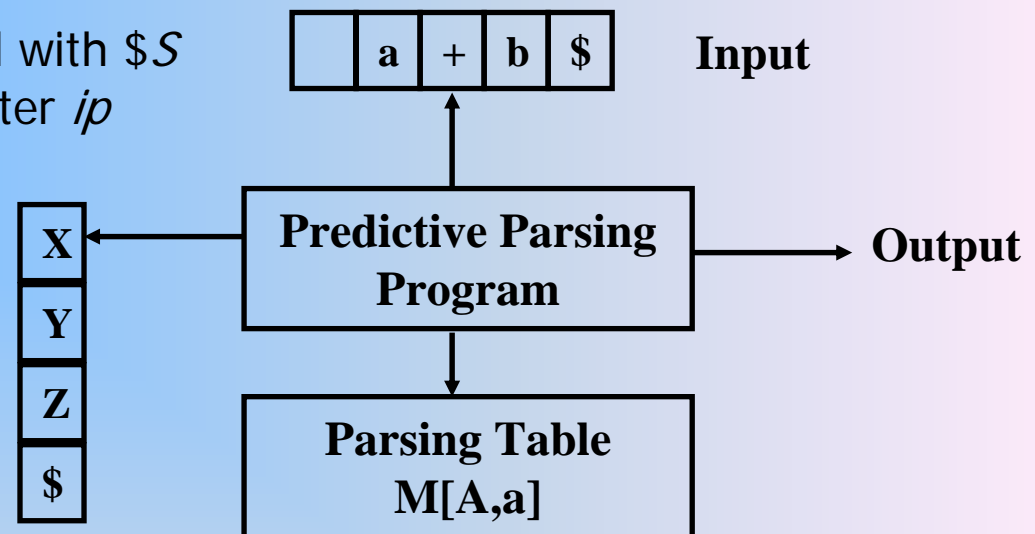
 pop the stack;

 push X_k, \dots, X_2, X_1 onto the stack with X_1 on top;

 }

 Set A to the top stack symbol;

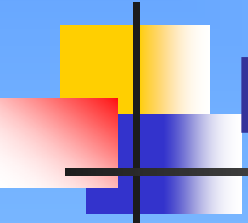
} Instructor: Dr. Liang Cheng





Error Recovery in Predictive Parsing

- A non-terminal on stack top and $M(A, a)$ empty
 - Idea: skip symbols on the input until a token in a selected set of synchronizing tokens is found.
 - The choice for a synchronizing set is important
 - Note that symbols in $FIRST(A)$ are by default in the synchronizing set because they have entries in the parsing table. In this case we can skip input and once we find a token in $FIRST(A)$ we resume parsing with A .
 - Define the synchronizing set of A to be $FOLLOW(A)$, then skip input until a token in $FOLLOW(A)$ appears and pop A from the stack. Resume parsing...
 - Productions that lead to ϵ if available might be used.
- A terminal on stack top and not match the input
 - Pop it and continue parsing (issuing an error message)
 - Or skip input until a matching terminal is reached



Implementing Error Recovery in Predictive Parsing

- General approach: modify the empty cells in the parsing table
 - if $M[A, a] = \{\text{empty}\}$ and a belongs to $\text{Follow}(A)$ then we set $M[A, a] = \text{"synch"}$
 - Error-recovery strategy for $X = \text{top-of-the-stack}$ and $a = \text{current-input}$,
 - If X is a nonterminal and $M[X, a] = \{\text{empty}\}$ then skip over a in the input
 - If X is a nonterminal and $M[X, a] = \{\text{synch}\}$ then pop X off of stack
 - If X is a terminal and $X \neq a$ then pop the terminal (essentially inserting it)

Revised Parsing Table

Non-terminal	INPUT SYMBOL					
	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	$F \rightarrow id$	_____	_____	$F \rightarrow (E)$	_____	_____

Skip input symbol.

From Follow sets.
Pop nonterminal on top of stack.

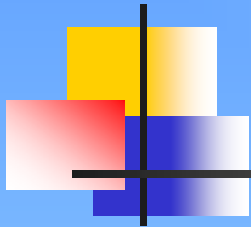
It's a "synch" action.

- $E \rightarrow TE'$
- $E' \rightarrow +TE' \mid \epsilon$
- $T \rightarrow FT'$
- $T' \rightarrow *FT' \mid \epsilon$
- $F \rightarrow (E) \mid id$



Example Error Messages

- Every non-terminal symbolizes an abstract language construct.
 - E means an expression
 - top-of-stack is E, input is +
"Error at location i, expressions cannot start with a '+' or
"error at location i, invalid expression"
 - T means a summation term
 - Top-of-stack is T, input is *
"error at location i, invalid operand."
 - When the top-of-the stack is a terminal that does not match...
 - Top-of-stack is id and the input is +
"error at location i: identifier expected"
 - Top-of-stack is) and the input is terminal other than)
"error at location i: left parenthesis at location m has no
closing right parenthesis"



Implement Top-down Parsers

- Stack: write ADT to manipulate its contents
- Input stream: responsibility of lexical analyzer
- Key issue: how is parsing table implemented? Assign unique IDs

Non-terminal	INPUT SYMBOL					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$	synch	synch
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	<u>$T \rightarrow FT'$</u>	synch		$T \rightarrow FT'$	synch	synch
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	<u>$F \rightarrow id$</u>	synch	synch	$F \rightarrow (E)$	synch	synch

All rules have unique IDs

Ditto for synch actions

Also for blanks which handle errors

Implement Top-down Parsers

Non-terminal	INPUT SYMBOL					
	id	+	*	()	\$
E	1	18	19	1	9	10
E'	20	2	21	22	3	3
T	4	11	23	4	12	13
T'	24	6	5	25	6	6
F	8	14	15	7	16	17

1 $E \rightarrow TE'$

2 $E' \rightarrow +TE'$

3 $E' \rightarrow \epsilon$

4 $T \rightarrow FT'$

5 $T' \rightarrow *FT'$

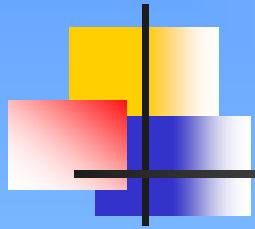
6 $T' \rightarrow \epsilon$

7 $F \rightarrow (E)$

8 $F \rightarrow id$

9 – 17 :
Sync
Actions

18 – 25 :
Error
Handlers



Implement Top-down Parsers

- Each # (or set of #s) corresponds to a procedure that:
 - Uses stack ADT
 - Gets tokens
 - Prints error messages
 - Prints diagnostic messages
 - Handles errors



Implement Top-down Parsers

```
state = M[ top(s), current_token ]
```

```
switch (state)
```

```
{
```

```
  case 1: proc_E_TE'();  
          break ;
```

```
  ...
```

```
  case 8: proc_F_id();  
          break ;
```

```
  case 9: proc_sync_9();  
          break ;
```

```
  ...
```

```
  case 17: proc_sync_17();  
           break ;
```

```
  case 18:
```

```
  ...
```

```
  case 25:
```

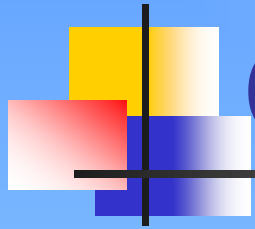
```
}
```

Combine → put
in another switch

Some sync
actions may be
same

Procs to handle errors

Some error
handlers may be
similar



Outline

- Recap
- Top-down parsing (Section 4.4)
- **Summary and homework**