

CSc 453

Compilers and Systems Software

8 : Top-Down Parsing III

Department of Computer Science
University of Arizona

collberg@gmail.com

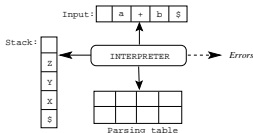
Copyright © 2009 Christian Collberg

Predictive Parsing

Predictive Parsing

Predictive Parsing . .

- Just as with lexical analysis, we can either hard-code a top-down parser, or build a generic table-driven interpreter. The latter is called a *Predictive Parser*.
- Instead of using recursion we store the current state of the parse on a stack:



- The predictive parser has
 - 1 an input stream (list of tokens followed by the *end-of-file-marker* \$),
 - 2 a stack with a sequence of grammar symbols, and an *end-of-stack-marker* \$,
 - 3 a parsing table $M[A, a] \rightarrow P$ mapping a non-terminal A and a terminal a to a grammar production P .
- Initially, the stack holds the start symbol, S .

At each step, the interpreter looks at the top stack element (X) and the current input symbol (a). There are three cases:

- 1 $X = a = \$ \Rightarrow$ success!
- 2 $M[X, a] = \text{error} \Rightarrow$ bail!
- 3 $X = a \neq \$ \Rightarrow$ match(). Move to next token and pop off X .
- 4 $M[X, a] = \{X \rightarrow UVW\} \Rightarrow$
 - 1 Pop X off the stack, then
 - 2 Push(W), Push(V), Push(U).

The next slide shows the algorithm in more detail.

```

a := first token
repeat
  X := top()
  if X is a terminal or $ then
    if X = a then
      pop()
      a := next token
    else error
  else
    if M[X,a] = X → Y1Y2...Yk then
      pop()
      push(Yk); ...; Push(Y1)
    else error
until X = $
    
```

- The Predictive parsing table for the grammar below:

	<u>id</u>	<u>±</u>	<u>*</u>	<u>(</u>	<u>)</u>	<u>\$</u>
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow \pm TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

$E \rightarrow TE'$ $E' \rightarrow \pm TE' \mid \epsilon$	$T \rightarrow FT'$ $T' \rightarrow *FT' \mid \epsilon$	$F \rightarrow (E) \mid \text{id}$
---	--	------------------------------------

Stack	Input	$M[A, a] =$
$\$E$	<u>id+id*id</u> \$	
$\$E'T$	<u>id+id*id</u> \$	$E \rightarrow TE'$
$\$E'T'F$	<u>id+id*id</u> \$	$T \rightarrow FT'$
$\$E'T'\text{id}$	<u>id+id*id</u> \$	$F \rightarrow \text{id}$
$\$E'T'$	<u>+id*id</u> \$	
$\$E'$	<u>+id*id</u> \$	$T' \rightarrow \epsilon$
$\$E'T\pm$	<u>+id*id</u> \$	$E' \rightarrow \pm TE'$
$\$E'T$	<u>id*id</u> \$	

Stack	Input	$M[A, a] =$
$\$E'T'F$	<u>id</u> *id\$	$T \rightarrow FT'$
$\$E'T'id$	<u>id</u> *id\$	$F \rightarrow id$
$\$E'T'$	*id\$	
$\$E'T'F*$	*id\$	$T' \rightarrow *FT'$
$\$E'T'F$	<u>id</u> \$	
$\$E'T'id$	<u>id</u> \$	$F \rightarrow id$
$\$E'T'$	\$	
$\$E'$	\$	$T' \rightarrow \epsilon$
\$	\$	$E' \rightarrow \epsilon$

Stack	Input	$M[A, a] =$
$\$E$	<u>id</u> ++\$	
$\$E'T$	<u>id</u> ++\$	$E \rightarrow TE'$
$\$E'T'F$	<u>id</u> ++\$	$T \rightarrow FT'$
$\$E'T'id$	<u>id</u> ++\$	$F \rightarrow id$
$\$E'T'$	++\$	
$\$E'$	++\$	$T' \rightarrow \epsilon$
$\$E'T_+$	++\$	$E' \rightarrow +TE'$
$\$E'T$	*\$	Error!

◀ ▶ 🔍 🔄

◀ ▶ 🔍 🔄

Building the Parse Table...

for each production $A \rightarrow \alpha$ do
 for each terminal a in $\text{FIRST}(\alpha)$ do
 $M[A, a] := \{A \rightarrow \alpha\}$

if ϵ is in $\text{FIRST}(\alpha)$ then
 for each terminal b in $\text{FOLLOW}(A)$ do
 $M[A, b] := \{A \rightarrow \alpha\}$

if ϵ is in $\text{FIRST}(\alpha)$ and $\$$ is in $\text{FOLLOW}(A)$ then
 $M[A, \$] := \{A \rightarrow \alpha\}$

for all undefined entries $M[A, a]$ do
 $M[A, a] := \text{error}$

$$E \rightarrow TE'$$

$$E' \rightarrow \pm TE' \mid \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

$$F \rightarrow \underline{_} E \underline{_} \mid id$$
 $\text{FIRST}(E) = \{\underline{_}, id\}$
 $\text{FIRST}(E') = \{\pm, \epsilon\}$
 $\text{FIRST}(T) = \{\underline{_}, id\}$
 $\text{FIRST}(T') = \{\pm, \epsilon\}$
 $\text{FIRST}(F) = \{\underline{_}, id\}$
 $\text{FOLLOW}(E) = \{\underline{_}, \$\}$
 $\text{FOLLOW}(E') = \{\underline{_}, \$\}$
 $\text{FOLLOW}(T) = \{\pm, \underline{_}, \$\}$
 $\text{FOLLOW}(T') = \{\pm, \underline{_}, \$\}$
 $\text{FOLLOW}(F) = \{\pm, \pm, \underline{_}, \$\}$

◀ ▶ 🔍 🔄

◀ ▶ 🔍 🔄

	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'						
T						
T'						
F						

- We start by looking at production $E \rightarrow TE'$.
- Since $\text{FIRST}(TE') = \text{FIRST}(T) = \{\langle, \text{id}\}$ we set $M[E, \langle] = M[E, \text{id}] = E \rightarrow TE'$.

	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow \pm TE'$				
T						
T'						
F						

- We next consider production $E' \rightarrow \pm TE'$.
- Since $\text{FIRST}(\pm TE') = \{\pm\}$ we set $M[E', \pm] = E' \rightarrow \pm TE'$.

	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow \pm TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T						
T'						
F						

- We next consider production $E' \rightarrow \epsilon$.
- Since $\text{FOLLOW}(E') = \{\rangle, \$\}$ we set $M[E', \rangle] = M[E', \$] = E' \rightarrow \epsilon$.

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

- We next consider production $T \rightarrow FT'$.
- Since $\text{FIRST}(FT') = \{\langle, \text{id}\}$ we set $M[T, \langle] = M[T, \text{id}] = T \rightarrow FT'$.

	<u>id</u>	<u>+</u>	<u>*</u>	<u>(</u>	<u>)</u>	<u>\$</u>
<u>E</u>	$E \rightarrow TE'$			$E \rightarrow TE'$		
<u>E'</u>		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
<u>T</u>	$T \rightarrow FT'$			$T \rightarrow FT'$		
<u>T'</u>			$T' \rightarrow *FT'$			
<u>F</u>						

- We next consider production $T' \rightarrow *FT'$.
- Since $\text{FIRST}(*FT') = \{*\}$ we set $M[T', *] = T' \rightarrow *FT'$.

	<u>id</u>	<u>+</u>	<u>*</u>	<u>(</u>	<u>)</u>	<u>\$</u>
<u>E</u>	$E \rightarrow TE'$			$E \rightarrow TE'$		
<u>E'</u>		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
<u>T</u>	$T \rightarrow FT'$			$T \rightarrow FT'$		
<u>T'</u>		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
<u>F</u>						

- We next consider production $T' \rightarrow \epsilon$.
- Since $\text{FOLLOW}(T') = \{+, _, \$\}$ we set $M[T', \pm] = M[T', _] = M[T', \$] = T' \rightarrow \epsilon$.

	<u>id</u>	<u>+</u>	<u>*</u>	<u>(</u>	<u>)</u>	<u>\$</u>
<u>E</u>	$E \rightarrow TE'$			$E \rightarrow TE'$		
<u>E'</u>		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
<u>T</u>	$T \rightarrow FT'$			$T \rightarrow FT'$		
<u>T'</u>		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
<u>F</u>	$F \rightarrow \underline{id}$			$F \rightarrow \underline{(E)}$		

- We next consider production $F \rightarrow \underline{(E)}$.
- Since $\text{FIRST}(\underline{(E)}) = \{(\}$ we set $M[F, \underline{(}] = F \rightarrow \underline{(E)}$.
- We finally consider production $F \rightarrow \underline{id}$.
- Since $\text{FIRST}(\underline{id}) = \{id\}$ we set $M[F, \underline{id}] = F \rightarrow \underline{id}$.

- Here's a simple implementation of predictive parsing in Java.
- Grammar symbols (both terminals and non-terminals) are represented by single characters: e represents E', f represents F', i represents id, "" represents epsilon, null represents error.
- The function `table(X,a)` looks up the relevant production.

```

class Expr {
    static String[] input = {"i","+","i","*","i","$"};
    static int current = 0;

    static boolean isTerminal(String S) {
        return S.equals("*") | S.equals("+") |
            S.equals("i") |
            S.equals("(") | S.equals(")");
    }

    static String[][] TABLE = {
        // id + * ( ) $
        {"Te", null, null, "Te", null, null}, // E
        {null, "+Te", null, null, "", ""}, // E'
        {"Ft", null, null, "Ft", null, null}, // T
        {null, "", "*Ft", null, "", ""}, // T'
        {"i", null, null, "(E)", null, null} // F
    };
}

```

```

static String table(String X, String a){
    int aIdx, XIdx=-1;
    switch (a.charAt(0)) {
        case 'i' : aIdx=0; break;
        case '+' : aIdx=1; break;
        case '*' : aIdx=2; break;
        case '(' : aIdx=3; break;
        case ')' : aIdx=4; break;
        case '$' : aIdx=5; break; }
    switch (X.charAt(0)) {
        case 'E' : XIdx=0; break;
        case 'E' : XIdx=1; break;
        case 'T' : XIdx=2; break;
        case 'T' : XIdx=3; break;
        case 'F' : XIdx=4; break; }
    return TABLE[XIdx][aIdx];
}

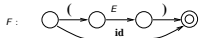
```

```

static boolean interpret(){
    push("$"); push("E");
    String a, X;
    while(true) {
        a = input[current];
        X = top();
        if (isTerminal(X) | X.equals("$")) {
            if (a.equals(X)) {pop(); current++;}
            else return false;
        } else {
            pop();
            String prod = table(X,a);
            if (prod==null) return false;
            for(int i=prod.length()-1;i>=0;i--)
                push(prod.charAt(i));
        }
        if (X.equals("$")) return true;
    }
}
}

```

Error Recovery



- Quit on first error. Considered unacceptable. But, with fast parse-edit-compile-cycles, why not?
- Repair the error by inserting/deleting tokens, to turn the erroneous program into a legal one.
- Repair the error by guessing what the programmer meant.

_____ Simple error handling: _____

- Skip tokens until a *synchronizing* token is found.
- The tokens in FOLLOW(A) are in the SYNCH set.
- If the language has an hierarchical structure
 prog \rightarrow proc \rightarrow stat \rightarrow expr, then add FIRST sets of higher
 constructs to the SYNCH sets of lower ones.

```

PROCEDURE F();
  CONST SYNCH = FOLLOW(F)  $\cup$  FIRST(stat);
  IF cur_tok =    THEN
    cur_tok := next_token(); E();
  IF cur_tok =    THEN
    cur_tok := next_token()
  ELSE WHILE cur_tok  $\notin$  SYNCH DO
    cur_tok := next_token() ENDDO
  ENDIF
  ELSIF cur_tok = id THEN cur_tok := next_token();
  ELSE WHILE cur_tok  $\notin$  SYNCH DO
    cur_tok := next_token() ENDDO
  ENDIF;

```

Example I

$Z \rightarrow d$	$Y \rightarrow c$
$Z \rightarrow XYZ$	$X \rightarrow Y$
$Y \rightarrow \epsilon$	$X \rightarrow a$

round	FIRST(Z)	FIRST(Y)	FIRST(X)
1	d	c	a
2	d	c, ϵ	a, c
3	d, a	c, ϵ	a, c, ϵ
4	d, a, c	c, ϵ	a, c, ϵ
5	d, a, c	c, ϵ	a, c, ϵ

Example I...

$Z \rightarrow d$	$Y \rightarrow c$
$Z \rightarrow XYZ$	$X \rightarrow Y$
$Y \rightarrow \epsilon$	$X \rightarrow a$

round	FOLLOW(Z)	FOLLOW(Y)	FOLLOW(X)
0	$\$$		
1	$\$$	d, a, c	c
2	$\$$	d, a, c	$c, d, a, \$$

PRODUCTION $Z \rightarrow d$:
 $\text{FIRST}(d) = \{d\}$
 $\Rightarrow M[Z, d] = Z \rightarrow d$

PRODUCTION $Z \rightarrow X, Y, Z$:
 $\text{FIRST}(XYZ) = \{a, c, d\}$
 $\Rightarrow M[Z, a] = Z \rightarrow XYZ$
 $\Rightarrow M[Z, c] = Z \rightarrow XYZ$
 $\Rightarrow M[Z, d] = Z \rightarrow XYZ$

PRODUCTION $Y \rightarrow c$:
 $\text{FIRST}(c) = \{c\}$
 $\Rightarrow M[Y, c] = Y \rightarrow c$

PRODUCTION $Y \rightarrow \epsilon$:
 $\text{FIRST}(\epsilon) = \{\epsilon\}$
 $\text{FOLLOW}(Y) = \{d, a, c\}$
 $\Rightarrow M[Y, d] = Y \rightarrow \epsilon$
 $\Rightarrow M[Y, a] = Y \rightarrow \epsilon$
 $\Rightarrow M[Y, c] = Y \rightarrow \epsilon$

PRODUCTION $X \rightarrow a$:
 $\text{FIRST}(a) = \{a\}$
 $\Rightarrow M[X, a] = X \rightarrow a$

PRODUCTION $X \rightarrow Y$:
 $\text{FIRST}(Y) = \{c, \epsilon\}$
 $\text{FOLLOW}(Y) = \{d, a, c\}$
 $\Rightarrow M[X, c] = X \rightarrow Y$
 $\Rightarrow M[X, d] = X \rightarrow Y$
 $\Rightarrow M[X, a] = X \rightarrow Y$
 $\Rightarrow M[X, \$] = X \rightarrow Y$

- Note that this grammar is not LL(1): it has multiple entries in some table cells.
- The grammar is, in fact, ambiguous. The string "d" can be derived in more than one way.

$S \rightarrow ABC$	$B \rightarrow b$
$A \rightarrow aA C$	$C \rightarrow c$

#	FIRST(S)	FIRST(A)	FIRST(B)	FIRST(C)
0		a	b	c
1	a	a, c	b	c
2	a, c	a, c	b	c
#	FOLLOW(S)	FOLLOW(A)	FOLLOW(B)	FOLLOW(C)
0	\$	b	c	
1	\$	b	c	\$, b

PRODUCTION $S \rightarrow ABC$:
 $\text{FIRST}(ABC) = \{a, c\}$
 $\Rightarrow M[S, a] = S \rightarrow ABC$
 $\Rightarrow M[S, b] = S \rightarrow ABC$

PRODUCTION $A \rightarrow aA$:
 $\text{FIRST}(aA) = \{a\}$
 $\Rightarrow M[A, a] = A \rightarrow aA$

PRODUCTION $A \rightarrow C$:
 $\text{FIRST}(C) = \{c\}$
 $\Rightarrow M[A, c] = A \rightarrow C$

PRODUCTION $B \rightarrow b$:
 $\text{FIRST}(b) = \{b\}$
 $\Rightarrow M[B, b] = B \rightarrow b$

PRODUCTION $C \rightarrow c$:
 $\text{FIRST}(c) = \{c\}$
 $\Rightarrow M[C, c] = C \rightarrow c$

	a	b	c	\$
S	$S \rightarrow ABC$	$S \rightarrow ABC$		
A	$A \rightarrow aA$		$A \rightarrow C$	
B		$B \rightarrow b$		
C			$C \rightarrow c$	

Readings and References

- Read Louden, pp. 143–196.
- Or, the Dragon Book:
 - Top-Down Parsing 181–190
 - Error Recovery 192–195
 - Recursive Descent Parsing 40–55, 75–76