

Compiler Design

Syntax Analysis

Top-Down Parsing

conf. dr. ing. Ciprian-Bogdan Chirila

chirila@cs.upt.ro

<http://www.cs.upt.ro/~chirila>

Outline

- Recursive-Descent Parsing
- FIRST and FOLLOW
- LL(I) Grammars
- Non-recursive Predictive Parsing
- Error Recovery in Predicting Parsing

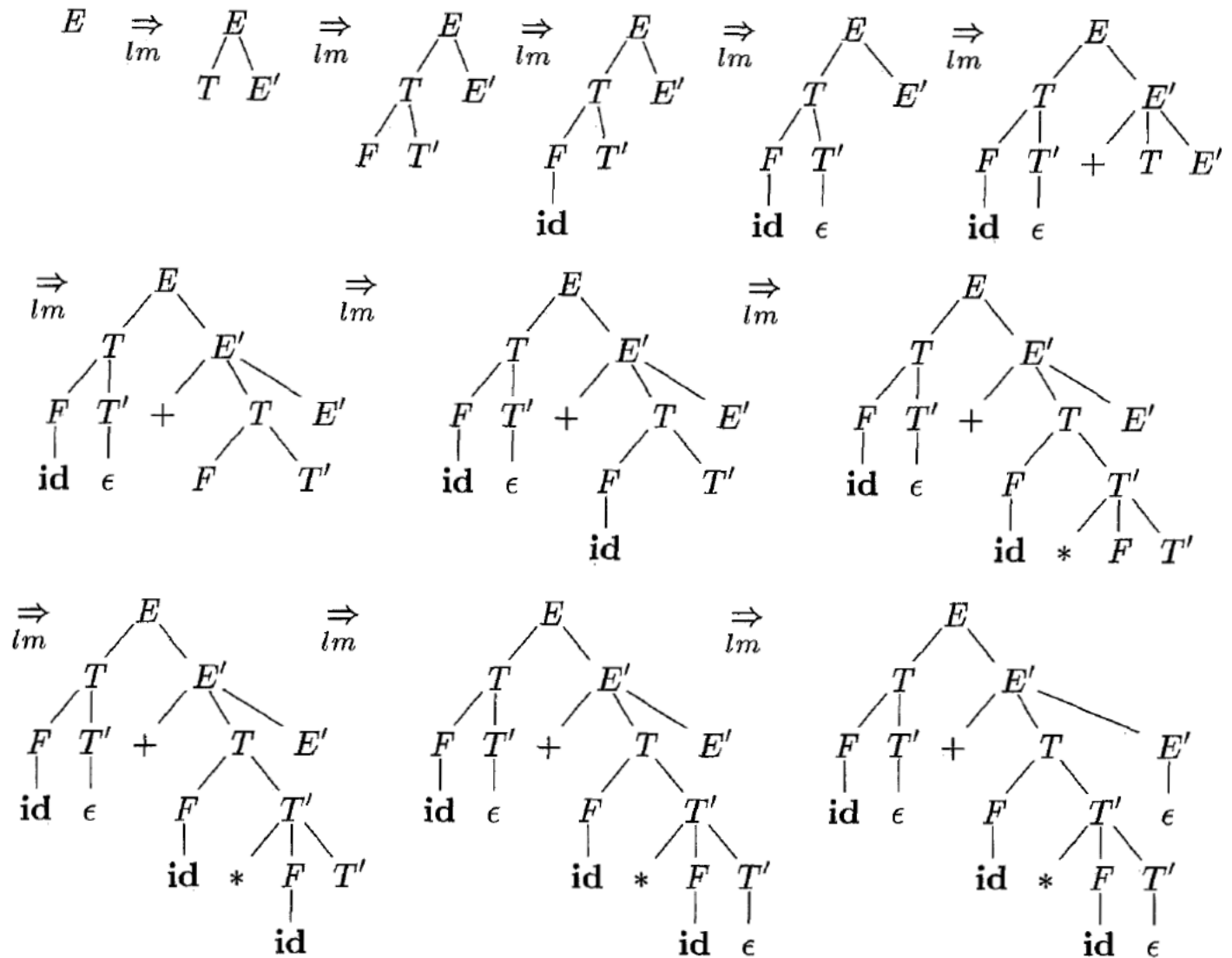
Top Down Parsing

- constructing a parse tree from the input string
 - starting from the root
 - creating the nodes in preorder
- finding the left-most derivation for an input string

Grammar Example

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

Derivation Example for $id+id*id$



LL(k) Grammars

- LL(k) – class of grammar for which we can construct predictive parsers looking k symbols ahead
- LL(1)
- FIRST and FOLLOW sets
 - are used to construct predictive parsing tables
 - make explicit the choice of production
 - are useful for bottom-up parsing

Recursive Descendant Parsing Program

- set of procedures
- one procedure for each non-terminal
- the start symbol
 - launches the execution
 - announces success if the body scans it's input string

Recursive Descendant Parsing

```
void A()  
{  
    choose an A-production,  $A \rightarrow X_1 X_2 \dots X_k$ ;  
    for( $i=1$  to  $k$ )  
    {  
        if ( $X_i$  is a non-terminal)  
            call procedure  $X_i()$ ;  
        else if ( $X_i$  equals the current symbol  $a$ )  
            advance the input to the next symbol;  
        else /* an error has occurred */  
    }  
}
```


Recursive Descendent Parsing Pseudocode

- non-deterministic
 - the manner in which the A-production is chosen is not specified
- generally requires backtracking
 - repeated scans over the input
 - rarely needed to parse programming language constructs
 - not very efficient – tabular methods such as dynamic programming are preferred

Allowing Backtracking

```
void A()  
{
```

```
{
```

```
  choose an A-production,  $A \rightarrow X_1 X_2 \dots X_k$ ;
```

```
  for(i=1 to k)
```

```
  {
```

```
    if ( $X_i$  is a non-terminal)
```

```
      call procedure  $X_i()$ ;
```

```
    else if ( $X_i$  equals the current symbol  $a$ )
```

```
      advance the input to the next
```

```
    else /* an error has occurred */
```

```
  }
```

```
}
```

try each
production in
some order

try another
A-production

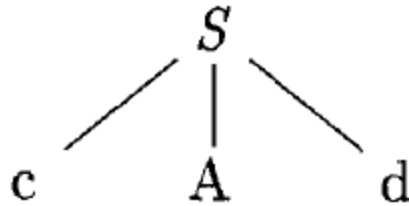
reset the input
pointer

Top-Down Parse Tree

- $S \rightarrow c A d$
- $A \rightarrow a b \mid a$
- $w = cad$

Step 1

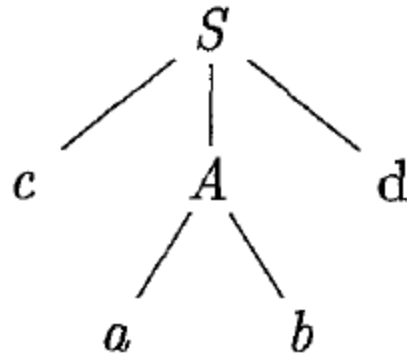
- S has only one production
- we expand S



- first character of input $w=cad$ matches the leftmost leaf in the tree c

Step 2

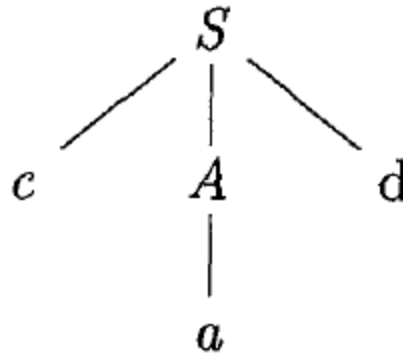
- we expand $A \rightarrow a b$
- we have a match for second input character a



- we go to next symbol d
- b does not match d
- we report failure
- we go back to A to try another alternative
- we reset input pointer to position 2

Step 3

- the second alternative for A is $A \rightarrow a$



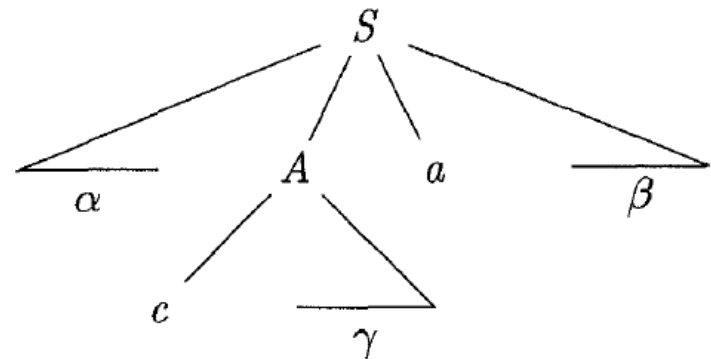
- leaf **a** matches second symbol
- leaf **d** matched the third symbol
- we halt with successful parsing message

FIRST and FOLLOW Functions

- two functions useful in creating parsers for both
 - top-down
 - bottom-up
- helps which production to apply based on next input symbol
- in panic mode error recovery tokens produced by FOLLOW are used for synchronization

The FIRST Function

- $\text{FIRST}(\alpha)$
 - set of terminals that begin strings derived from α
 - α is any string of grammar symbols
 - if $\alpha \xRightarrow{*} \varepsilon$ then ε is in $\text{FIRST}(\alpha)$
- $A \xRightarrow{*} c\gamma$
 - c is in $\text{FIRST}(A)$



How FIRST function works ?

- $A \rightarrow \alpha | \beta$
- $FIRST(\alpha)$ and $FIRST(\beta)$ are disjoint sets
- input symbol a can be in one of the two sets
- if a is in $FIRST(\beta)$ we can choose the production $A \rightarrow \beta$

The FOLLOW Function

- FOLLOW(A)
 - the set of terminals a that can appear immediately to the right of A in some sentential form
 - the set of terminals a such that $S \xRightarrow{*} \alpha A a \beta$ for some α and β

How to compute FIRST ?

- if X is terminal then $\text{FIRST}(X) = \{X\}$
- if X is non-terminal $X \rightarrow Y_1 Y_2 \dots Y_k$ is a production for some $k \geq 1$
 - place a in $\text{FIRST}(X)$ if for some i
 - a is in $\text{FIRST}(Y_i)$ and
 - ϵ is in $\text{FIRST}(Y_1) \dots \text{FIRST}(Y_{i-1})$
 - if ϵ is in all $\text{FIRST}(Y_j)$ $j=1, \dots, k$
 - then add ϵ to $\text{FIRST}(X)$
- if $X \rightarrow \epsilon$ is a production
 - then add ϵ to $\text{FIRST}(X)$

How to compute FIRST ?

- input string $X_1X_2\dots X_n$
- add to $\text{FIRST}(X_1X_2\dots X_n)$
 - all non- ϵ symbols of $\text{FIRST}(X_1)$
 - all non- ϵ symbols of $\text{FIRST}(X_2)$ if ϵ is in $\text{FIRST}(X_1)$
 - all non- ϵ symbols of $\text{FIRST}(X_3)$ if ϵ is in $\text{FIRST}(X_1)$ and in $\text{FIRST}(X_2)$
 - ...
 - ϵ , if ϵ is in all $\text{FIRST}(X_i)$ $i=1,\dots,n$

How to compute FOLLOW ?

- place \$ in FOLLOW(S)
 - S is the start symbol
 - \$ is the right end-marker
- if there is a production $A \rightarrow \alpha B \beta$
 - everything in FIRST(β) *except* ϵ is in FOLLOW(B)
- if there is a production $A \rightarrow \alpha B$ or $A \rightarrow \alpha B \beta$ where first(β) contains ϵ
 - everything in FOLLOW(A) is in FOLLOW(B)

Example

- $\text{FIRST}(F) = \{ (, id \}$
- $\text{FIRST}(T) = \text{FIRST}(F) = \{ (, id \}$
- $\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{ (, id \}$

- $\text{FIRST}(E') = \{ +, \epsilon \}$
- $\text{FIRST}(T') = \{ *, \epsilon \}$

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

Example

- $\text{FOLLOW}(E) = \{), \$\}$
 - E is the start symbol so it must include \$
 - the body (E) tells that the) symbol must be included
- $\text{FOLLOW}(E') = \{), \$\}$
 - $E \rightarrow TE'$ so what follows after E will follow after E'
 - $\text{FOLLOW}(E') = \text{FOLLOW}(E)$

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

Example

- FOLLOW(T)={+,),}\$}
 - E->TE' so FOLLOW(T) includes FIRST(E')={+} (except ϵ)
 - E->TE' and E' includes ϵ , so FOLLOW(E)={),}\$} is included in FOLLOW(T)
- FOLLOW(T')={+,),}\$}
 - T->FT' so FOLLOW(T) is included in FOLLOW(T')

E->TE'

E'->+TE'| ϵ

T->FT'

T'->*FT'| ϵ

F->(E) | id

Example

- FOLLOW(F) = {+, *,), \$}
 - T' → *FT' so FOLLOW(F) includes FIRST(T') = {*} (except ε)
 - T → FT' and T' → ε so FOLLOW(F) includes FOLLOW(T) = {+,), \$}

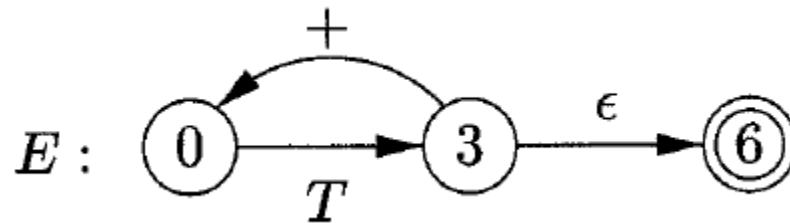
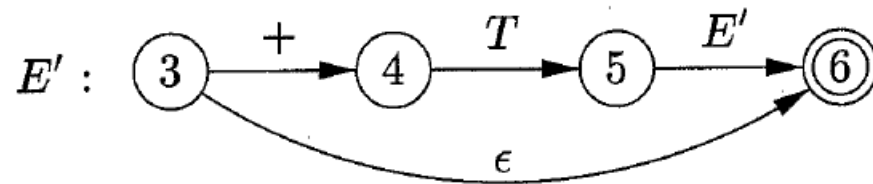
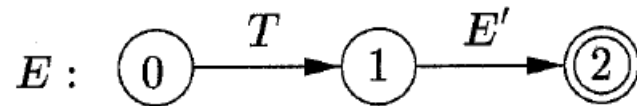
E → TE'
E' → +TE' | ε
T → FT'
T' → *FT' | ε
F → (E) | id

LL(I) Grammars

- predictive parsers
 - recursive descendant with no backtracking
- can be constructed for LL(I) grammar class
 - first L stands for scanning the input **from left to right**
 - second L for producing **leftmost derivation**
 - the I is for using one input symbol of lookahead at each step to make parsing actions decisions

Transition Diagrams for Predictive Parsers

- useful for visualizing predictive parsers
 - $E \rightarrow TE'$
 - $E' \rightarrow +TE' \mid \epsilon$



Building a Transition Diagram

- eliminate left recursion
- left factor the grammar
- for each non-terminal
 - create an initial and a final state
 - for each production $A \rightarrow X_1 X_2 \dots X_k$
 - create a path from initial state to final state with edges labeled X_1, X_2, \dots, X_k
 - if $A \rightarrow \epsilon$ the path is an edge labeled ϵ
- label of edges can be tokens or non-terminals
- ϵ -transitions are the default choice

LL(I) Grammar Definition

- rich enough to cover most programming constructs
- a grammar G is LL(I) iff $A \rightarrow \alpha | \beta$
 - 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 $\beta \xRightarrow{*} \epsilon$
 - α does not derive any string beginning with a terminal in $\text{FOLLOW}(A)$

LL(I) Grammar Definition

- $\text{FIRST}(\alpha)$ and $\text{FIRST}(\beta)$ are disjoint sets
- If ϵ is in $\text{FIRST}(\beta)$ then $\text{FIRST}(\alpha)$ and $\text{FOLLOW}(A)$ are disjoint sets
- vice versa if ϵ is in $\text{FIRST}(\alpha)$

Example

- control flow constructs having distinguishable keywords generally satisfies the LL(1) constraints
- $\text{stmt} \rightarrow \mathbf{if} (\text{expr}) \text{ stmt } \mathbf{else} \text{ stmt}$
| $\mathbf{while}(\text{expr}) \text{ stmt}$
| $\{\text{stmt_list}\}$
- keywords like: if, while, { tells which alternative to take in order to succeed in finding a statement

The Construction of a Predictive Parsing Table

- to collect information from FIRST and FOLLOW
- to store them into a predictive parsing table $M[A,a]$ – two dimensional array
 - A – non-terminal
 - a – terminal or the \$ end marker
- main idea
 - $A \rightarrow \alpha$ is chosen if the next input symbol a is in $FIRST(\alpha)$
 - if $\alpha \Rightarrow \epsilon$ or $\alpha^* \Rightarrow \epsilon$ production $A \rightarrow \alpha$ is chosen when the current input symbol or \$ is in $FOLLOW(A)$

The Construction Algorithm

- Input
 - Grammar G
- Output
 - Parsing table M
- Method
 - for each production $A \rightarrow \alpha$
 - for each terminal a in $FIRST(A)$ add $A \rightarrow \alpha$ to $M[A,a]$
 - if ϵ is in $FIRST(\alpha)$ then for each terminal b in $FOLLOW(A)$ add $A \rightarrow \alpha$ to $M[A,b]$
 - if ϵ is in $FIRST(\alpha)$ and $\$$ is in $FOLLOW(A)$ then add $A \rightarrow \alpha$ to $M[A,\$]$
 - after filling the table if there is no production in $M[A,a]$ then set $M[A,a]$ to error, represented by an empty entry in the table

Example

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

	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	$F \rightarrow id$			$F \rightarrow (E)$		

Example

- $E \rightarrow TE'$
 - $\text{FIRST}(TE') = \text{FIRST}(T) = \{ (, \text{id} \}$
 - added to $M[E, (]$ and $M[E, \text{id}]$
- $E' \rightarrow +TE'$
 - $\text{FIRST}(+TE') = \{ + \}$
 - added to $M[E', +]$
- $E' \rightarrow \epsilon$
 - $\text{FOLLOW}(E') = \{), \$ \}$
 - added to $M[E',)]$ and $M[E', \$]$

Example 2

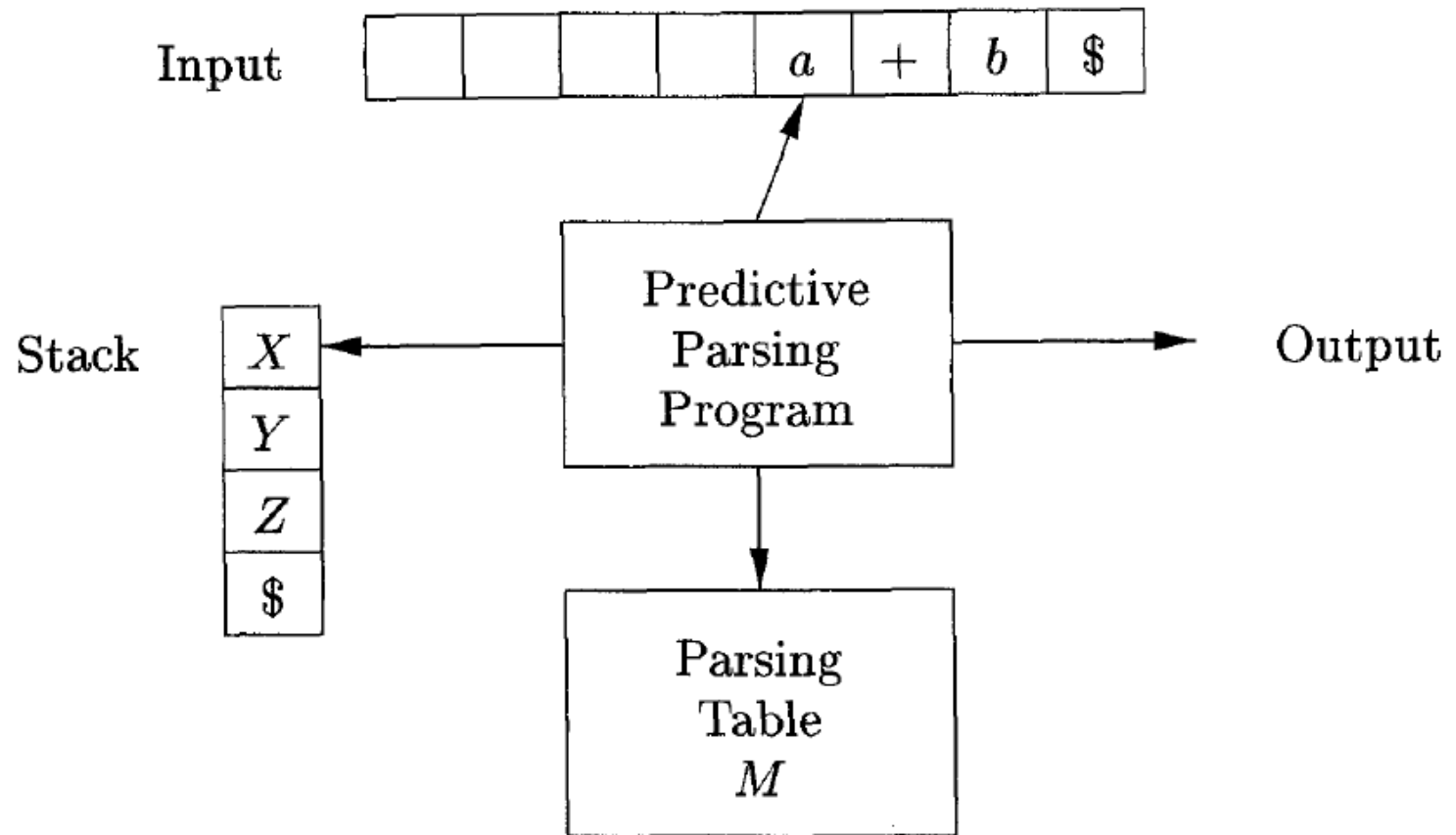
- $S \rightarrow iEtSS' \mid a$
- $S' \rightarrow eS \mid \epsilon$
- $E \rightarrow b$

	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

Non-recursive Predictive Parsing

- to maintain a stack explicitly
- rather than implicitly by recursive calls
- the parser simulates the leftmost derivation
- if w is the input matched so far
 - then the stack holds a sequence of grammar symbols α such that $S \xRightarrow{*}_{lm} w\alpha$

Model of a Table Driven Predictive Parser



Model of a Table Driven Predictive Parser

- input buffer
 - string to be parsed
 - end marker \$
- stack containing grammar symbols
 - it's bottom is marked by \$
- parsing table
- output stream

Model of a Table Driven Predictive Parser

- **X** is the symbol on top of the stack
- **a** is the current input symbol
- if **X** is non-terminal
 - the parser chooses a production by consulting the entry $M[X,a]$
 - semantic actions can be added to build a node in the parse tree
- if **X** is a terminal
 - a match is checked between **X** and input symbol **a**

Model of a Table Driven Predictive Parser

- parser configurations
 - stack content
 - remaining input

Table Driven Predictive Parsing

- Input
 - a string w
 - parsing table M for a grammar G
- Output
 - if w is in $L(G)$ then
 - a leftmost derivation of w
 - otherwise error indication
- Method
 - initially the parser has
 - $w\$$ in the input buffer
 - start symbol S of G on the stack top, above $\$$

Predictive Parsing Algorithm

```
set ip to point the first symbol a of w
set X to the top stack symbol
while(X!=$)
{
  if (X is a) then pop the stack and advance ip
  elseif (X is a terminal) error();
  elseif (M[X,a] is an error entry) error();
  elseif (M[X,a]=X->Y1,Y2,...,Yk)
  {
    output the production X->Y1,Y2,...,Yk
    pop the stack
    push Yk,Yk-1,...,Y1 onto the stack with Y1 on top
  }
  set X to the top stack symbol
}
```

Example

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

	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	$F \rightarrow id$			$F \rightarrow (E)$		

Moves Made by a Predictive Parser on $id+id*id$

MATCHED	STACK	INPUT	ACTION
	$E\$$	$id + id * id\$$	
	$TE' \$$	$id + id * id\$$	output $E \rightarrow TE'$
	$FT' E' \$$	$id + id * id\$$	output $T \rightarrow FT'$
	$id T' E' \$$	$id + id * id\$$	output $F \rightarrow id$
id	$T' E' \$$	$+ id * id\$$	match id
id	$E' \$$	$+ id * id\$$	output $T' \rightarrow \epsilon$
id	$+ TE' \$$	$+ id * id\$$	output $E' \rightarrow + TE'$
$id +$	$TE' \$$	$id * id\$$	match $+$
$id +$	$FT' E' \$$	$id * id\$$	output $T \rightarrow FT'$
$id +$	$id T' E' \$$	$id * id\$$	output $F \rightarrow id$
$id + id$	$T' E' \$$	$* id\$$	match id
$id + id$	$* FT' E' \$$	$* id\$$	output $T' \rightarrow * FT'$
$id + id *$	$FT' E' \$$	$id\$$	match $*$
$id + id *$	$id T' E' \$$	$id\$$	output $F \rightarrow id$
$id + id * id$	$T' E' \$$	$\$$	match id
$id + id * id$	$E' \$$	$\$$	output $T' \rightarrow \epsilon$
$id + id * id$	$\$$	$\$$	output $E' \rightarrow \epsilon$

Error Recovery in Predictive Parsing

- error recovery refers to the stack of the table driven predictive parser
- it makes explicit the terminals and non-terminals the parser hopes to match
- the techniques can be used with recursive-descendant parsing
- an error is detected when:
 - stack top terminal does not match the next input symbol
 - $M[A,a]$ is error (empty)
 - A is the non-terminal on the top of the stack
 - a is the next input symbol

Panic Mode

- skipping input symbols until a set of synchronizing symbols appear
- effectiveness depend on the chosen set
- the sets should be chosen so the parser recovers quickly from errors that are likely to occur in practice

Some Heuristics

- all symbols in $\text{FOLLOW}(A)$ will be added to the synchronizing set for A non-terminal
- skip tokens until an element of $\text{FOLLOW}(A)$ is seen
- pop A from the stack
- the parsing is likely to continue

Some Heuristics

- only FOLLOW(A) set is not enough
- because semicolons terminate statements in C
- keywords that begin statements may not appear in the FOLLOW set for expression non-terminal
- a missing semicolon after an assignment may result in the keyword beginning next statement to be skipped
- expressions appear within statements
- we need to add
 - **to** the synchronizing symbols of lower level constructs
 - the synchronizing symbols of higher level constructs
- we can add symbols that begin statements to the synchronizing sets for the non-terminals generating expressions

Some Heuristics

- all symbols from $\text{FIRST}(A)$ will be added to the synchronizing set for A non-terminal
 - it is possible to resume parsing according to A
 - if a symbol from $\text{FIRST}(A)$ appears
- if a non-terminal can generate the empty string
 - then the production deriving in ϵ can be used by default
 - we may postpone some error detection
 - cannot cause an error to be missed
 - reduces the number of non-terminals to be considered during error recovery

Some Heuristics

- if a terminal on the top can not be matched
 - pop the terminal
 - issue a message
 - continue parsing
 - the synchronization set of a token consists in all other tokens

Phrase Level Recovery

- filling in the blank cells pointers to error routines
 - change, insert, delete symbols
 - pop from the stack
- stack alteration is questionable
 - modifying the stack might not enable derivation at all
 - risk of infinite loop
 - to check the stack size after modifying it
 - it should decrease

Bibliography

- Alfred V. Aho, Monica S. Lam, Ravi Sethi, Jeffrey D. Ullman – Compilers, Principles, Techniques and Tools, Second Edition, 2007