

Name Manas Vinod Darbelli

Roll 27

Div A

Tybsc(cs)

Role of parser

Context free grammars

Top down parsing

Bottom up parsing

Parser generators

The role of parser

Lexical

Analyzer Parser Source

program

token

getNext

Token

Symbol

table

Parse tree Rest of

Front

End

Intermediate  
representation

Uses of grammars

Error handling

Common programming errors

Lexical errors

Syntactic errors

Semantic errors

Lexical errors

Error handler goals

Report the presence of errors clearly and accurately

Recover from each error quickly enough to detect  
subsequent errors

Add minimal overhead to the processing of correct  
programs

Error-recover strategies

Panic mode recovery

Discard input symbol one at a time until one of  
designated set of synchronization tokens is found

Phrase level recovery

Replacing a prefix of remaining input by some string  
that allows the parser to continue

Error productions

Augment the grammar with productions that generate  
the erroneous constructs

Global correction

Choosing minimal sequence of changes to obtain a  
globally least-cost correction

Context free grammars

Terminals

Nonterminals

Start symbol

productions

$\text{expression} \rightarrow \text{expression} + \text{term}$

$\text{expression} \rightarrow \text{expression} - \text{term}$

$\text{expression} \rightarrow \text{term}$

$\text{term} \rightarrow \text{term}^* \text{factor}$

$\text{term} \rightarrow \text{term} \text{ factor}$

$\text{term} \rightarrow \text{factor}$

$\text{factor} \rightarrow (\text{expression})$

$\text{factor} \rightarrow \text{id}$

Derivations

Productions are treated as rewriting rules to generate

a

string

Rightmost and leftmost derivations

$E \rightarrow E + E \mid E^* E \mid -E \mid (E) \mid id$

Derivations for  $-(id+id)$

$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(id+E) \Rightarrow -(id+id)$

Parse trees

$-(id+id) \quad E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(id+E) \Rightarrow -(id+id)$

Ambiguity

For some strings there exist more than one parse tree

Or more than one leftmost derivation

Or more than one rightmost derivation

Example:  $id + id^* id$

Elimination of ambiguity

Elimination of ambiguity (cont.)

Idea:

A statement appearing between a `then` and an `else`  
must be matched

Elimination of left recursion

A grammar is left recursive if it has a non-terminal  $A$   
such that there is a derivation  $A \Rightarrow A\alpha$  Top down parsing methods can't  
handle left recursive grammars

A simple rule for direct leftrecursion elimination: For a rule like:  $A \rightarrow$

A

Left recursion elimination (cont.)

There are cases like following

$S \rightarrow Aa \mid b$

$A \rightarrow Ac \mid Sd \mid \epsilon$  □ Left recursion elimination algorithm:

Arrange the nonterminals in some order  $A_1, A_2, \dots, A_n$ .

For (each  $i$  from 1 to  $n$ ) {

For (each  $j$  from 1 to  $i-1$ ) {

Replace each production of the form

are all current  $A_j$  productions

} Eliminate left recursion among the  $A_i$ -productions

}

Left factoring

Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive or top-down parsing.

Consider following grammar:

$\text{Stmt} \rightarrow \text{if expr then stmt else stmt}$

$\mid \text{if expr then stmt}$

On seeing input if it is not clear for the parser which production to use We can easily perform left factoring:

Left factoring

(cont.)

Algorithm

For each non-terminal  $A$ , find the longest prefix  $\alpha$  common to two or more of its alternatives.

Example:  $S \rightarrow IEtS \mid iEtSeS \mid a$   
 $E \rightarrow b$

Introduction

A Top-down parser tries to create a parse tree from the root towards the leafs scanning input from left to right  
It can be also viewed as finding a leftmost derivation for an input string

Recursive descent parsing

Consists of a set of procedures, one for each nonterminal

Execution begins with the procedure for start symbol

A typical procedure for a non-terminal

void  $A()$  {

choose an  $A$ -production,  $A \rightarrow X_1X_2..X_k$

for ( $i=1$  to  $k$ ) {

if ( $X_i$  is a nonterminal

call procedure  $X_i()$ ;

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

advance the input to the next

symbol:

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

### Recursive descent parsing (cont)

General recursive descent may require backtracking

The previous code needs to be modified to allow  
backtracking

In general form it can't choose an A-production easily.

So we need to try all alternatives

If one failed the input pointer needs to be reset and  
another alternative should be tried

Recursive descent parsers can't be used for left-recursive grammar

Example

$S \rightarrow cAd$

$A \rightarrow ab / a$  Input: cad

ScA d ScA d a b ScA da

### First and Follow

First() is set of terminals that begins strings derived from  
and First are disjoint sets then we can select  
appropriate A-production by looking at the next input

Follow, for any nonterminal A, is set of terminals a that  
can appear immediately after A in some sentential

form

If we have  $S \Rightarrow a$  for some  $a$  and then  $a$  is in

$\text{Follow}(A)$

If  $A$  can be the rightmost symbol in some sentential form,  
then  $\$$  is in  $\text{Follow}(A)$

Computing First

To compute  $\text{First}(X)$  for all grammar symbols  $X$ , apply  
following rules until no more terminals or  $\epsilon$  can be  
added to any First set:

1. If  $X$  is a terminal then  $\text{First}(X) = \{X\}$ .
2. If  $X$  is a nonterminal and  $X \Rightarrow Y_1 Y_2 \dots Y_k$  is a production  
for some  $k \geq 1$ , then place  $a_i$  in  $\text{First}(X)$  if for some  $i$   $a_i$  is  
in  $\text{First}(Y_i)$  and  $\epsilon$  is in all of  $\text{First}(Y_1), \dots, \text{First}(Y_{i-1})$  that  
is  $Y_1 \dots Y_{i-1} \Rightarrow \epsilon$ . if  $\epsilon$  is in  $\text{First}(Y_j)$  for  $j=1, \dots, k$  then add  
 $\epsilon$  to  $\text{First}(X)$ .
3. If  $X \Rightarrow \epsilon$  is a production then add  $\epsilon$  to  $\text{First}(X)$

Example

Computing follow

To compute  $\text{First}(A)$  for all nonterminals  $A$ , apply  
following rules until nothing can be added to any  
follow

set:

1. Place \$ in Follow(S) where S is the start symbol

2. If there is a production A->

then everything in

First except  $\epsilon$  is in Follow(B).

3. If there is a production A->B or a production

where First contains then everything

in Follow(A) is in Follow(B)

Example!

LL(1) Grammars

Predictive parsers are those recursive descent parsers needing no backtracking

Grammars for which we can create predictive parsers are called

LL(1)

The first L means scanning input from left to right

The second L means leftmost derivation

And 1 stands for using one input symbol for lookahead

A grammar G is LL(1) if and only if whenever A-> $\alpha/\beta$  are two distinct productions of G, the following conditions hold:

For no terminal  $a$  do  $\alpha$  and  $\beta$  both derive strings beginning with a

At most one of  $\alpha$  or  $\beta$  can derive empty string

\* If  $\alpha \Rightarrow \epsilon$  then  $\beta$  does not derive any string beginning with a terminal in Follow(A).

Construction of

predictive

parsing table

For each production in grammar do the following:

1. For each terminal  $a$  in  $\text{First}(A)$  add  $A \rightarrow a$  in  $M[A,a]$
2. If  $\epsilon$  is in  $\text{First}(A)$ , then for each terminal  $b$  in  $\text{Follow}(A)$  add  $A \rightarrow \epsilon$  to  $M[A,b]$ . If  $\epsilon$  is in  $\text{First}(A)$  and  $\$$  is in  $\text{Follow}(A)$ , add  $A \rightarrow \epsilon$  to  $M[A,\$]$  as well

If after performing the above, there is no production in  $M[A,a]$  then set  $M[A,a]$  to error

Example

$E \rightarrow TE'$

$E' \rightarrow +TE' /$

$T \rightarrow FT'$

$T' \rightarrow *FT' /$

$F \rightarrow (E) / \text{id}$

$FTE$

$E'$

$T'$

$\text{First}$

Follow

{(id)}

{(id)}

{(id)}

{+, ε}

{\*, ε}

{+, \*, )\$}

{+, )\$}

{+, )\$}

{), \$}

{), \$}

E

E'

T

T'

F

Non -

terminal

Input Symbol

id + \* ( )\$

$E \rightarrow TE'E \rightarrow TE'$

$E' \rightarrow +TE'E' \rightarrow \epsilon E' \rightarrow \epsilon$

$T \rightarrow FT'T \rightarrow FT'$

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

$F \rightarrow id F \rightarrow$

(E)

Another example

$S \rightarrow iEtSS' / a$

$S' \rightarrow eS / \epsilon$

$E \rightarrow b$

$S$

$S'$

$E$

Non -

terminal

Input Symbol

a b e i t \$

$S \rightarrow aS \rightarrow iEtSS'$

$S' \rightarrow \epsilon$

$S' \rightarrow eS$

$S' \rightarrow \epsilon$

$E \rightarrow b$

Non-recursive predicting parsing

a + b

\$

Predictive

parsing

program

output

Parsing

Table

M

stack XYZ\$

Predictive parsing algorithm

Set ip point to the first symbol of w;

Set X to the top stack symbol;

While ( $X \neq \$$ ) { /\* stack is not empty \*/

if (X is a) pop the stack and advance ip;

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

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

else if ( $M[X, a] = X \rightarrow Y_1 Y_2 \dots Y_k$ ) {

output the production  $X \rightarrow Y_1 Y_2 \dots Y_k$ ;

pop the stack;

push  $Y_k, \dots, Y_2, Y_1$  on to the stack with  $Y_1$  on top;

}

set X to the top stack

symbol:

}

Example

□  $id + id^* id \$$

Matched Stack Input Action

$E \$ id + id^* id \$$

Error recovery in predictive parsing

□ Panic mode

□ Place all symbols in  $Follow(A)$  into synchronization set for nonterminal A: skip tokens until an element of  $Follow(A)$  is seen and pop A from stack.

□ Add to the synchronization set of lower level construct the symbols that begin higher level constructs

□ Add symbols in  $First(A)$  to the synchronization set of nonterminal A

□ If a nonterminal can generate the empty string then the production deriving can be used as a default

□ If a terminal on top of the stack cannot be matched, pop the terminal, issue a message saying that the terminal was inserted

Example E

$E'$

T

$T'$

F

Non

-  
terminal

Input Symbol

$id^* ( \ $$

$E \rightarrow TE' E \rightarrow TE'$

$E' \rightarrow +TE' E' \rightarrow \epsilon E' \rightarrow$

$T \rightarrow FT' T \rightarrow FT'$

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

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

synch synch

synch synch synch

synch synch synch synch

Stack Input Action

$E \$ id^* + id \$ \text{ Error, Skip } )$

$E \$ id^* + id \$ \text{ id is in First}(E)$

$TE' \$ id^* + id \$$

$FT'E' \$ id^* + id \$$

$idTE' \$ id^* + id \$$

$TE' \$ ^* + id \$$

$*FT'E' \$ ^* + id \$$

$FT'E' \$ + id \$ \text{ Error, } M[F, +] = \text{synch}$

$TE' \$ + id \$ \text{ F has been popped}$

Introduction

- Constructs parse tree for an input string beginning at the leaves (the bottom) and working towards the

root

(the top)

□ Example:  $id^* id$

$E \rightarrow E + T / T$

$T \rightarrow T^* F / F$

$F \rightarrow (E) / id id$

$id^* id F^* id T^* id$

$id$

$F$

$T^* F$

$id$

$F id T^* F$

$id$

$F id$

$F$

$T^* F$

$id$

$F id$

$FE$

Shift-reduce parser

The general idea is to shift some symbols of input to the stack until a reduction can be applied

At each reduction step, a specific substring matching the body of a production is replaced by

the

nonterminal at the head of the production

The key decisions during bottom-up parsing are about  
when to reduce and about what production to apply

A reduction is a reverse of a step in a derivation

The goal of a bottom-up parser is to construct a  
derivation in reverse:

$E \Rightarrow T \Rightarrow T^* F \Rightarrow T^* id \Rightarrow F^* id \Rightarrow id^* id$

Handle pruning

A Handle is a substring that matches the body of a  
production and whose reduction represents one step  
along the reverse of a rightmost derivation

Right sentential form Handle Reducing production

$id^* id id F \Rightarrow id$

$F^* id F$

$id$

$T \Rightarrow F$

$T^* id F \Rightarrow id$

$T^* F T^* F E \Rightarrow T^* F$

Shift reduce parsing

A stack is used to hold grammar symbols

Handle always appear on top of the stack

Initial configuration:

Stack

Input

\$ w \$

Acceptance configuration

Stack Input

\$ S \$

Shift reduce parsing (cont.)

Basic operations:

Shift

Reduce

Accept

Error

Example: id\*id

Stack Input Action

\$

\$ id

id\*id\$ shift

\*id\$ reduce by F->id

\$ F \*id\$ reduce by T->F

\$ T \*id\$ shift

\$ T \*id\$ shift

\$ T \*id \$ reduce by F->id

\$ T \*F \$ reduce by T->T\*F

\$ T \$ reduce by E->T

\$ E \$

accept

Handle will appear on top of  
the stack

Conflicts during shift reduce  
parsing

Two kind of conflicts

Shift/reduce conflict

Reduce/reduce conflict

Example:

Stack Input

... if expr then stmt else ...\$

Reduce/reduce conflict

stmt  $\rightarrow$  id(parameter\_list)

stmt  $\rightarrow$  expr:=expr

parameter\_list  $\rightarrow$  parameter\_list, parameter

parameter\_list  $\rightarrow$  parameter

parameter  $\rightarrow$  id

expr  $\rightarrow$  id(expr\_list)

expr  $\rightarrow$  id

expr\_list  $\rightarrow$  expr\_list, expr

expr\_list  $\rightarrow$  expr Stack Input

... id(id id) ...\$

LR

## Parsing

The most prevalent type of bottom-up parsers

$LR(k)$ , mostly interested on parsers with  $k \leq 1$

Why  $LR$  parsers?

Table driven

Can be constructed to recognize all programming language constructs

Most general non-backtracking shift-reduce parsing method

Can detect a syntactic error as soon as it is possible to do so

Class of grammars for which we can construct  $LR$  parsers are superset of those which we can construct  $LL$  parsers

States of an  $LR$  parser

States represent set of items

An  $LR(0)$  item of  $G$  is a production of  $G$  with the dot at some position of the body:

For  $A \rightarrow XYZ$  we have following items

$A \rightarrow .XYZ$

$A \rightarrow X.YZ$

$A \rightarrow XY.Z$

$A \rightarrow XY.Z.$

In a state having  $A \rightarrow .XYZ$  we hope to see a string derivable from  $XYZ$  next on the input.

◻ What about  $A \rightarrow X.YZ$ ?

Constructing canonical

$LR(0)$

item sets

Augmented grammar:

$G$  with addition of a production:  $S'' \rightarrow S$

Closure of item sets:

If  $I$  is a set of items,  $\text{closure}(I)$  is a set of items constructed from  $I$  by the following rules:

Add every item in  $I$  to  $\text{closure}(I)$

If  $A \rightarrow a.B\beta$  is in  $\text{closure}(I)$  and  $B \rightarrow V$  is a production then add the item  $B \rightarrow .V$  to  $\text{closure}(I)$ .

Example:  $E' \rightarrow E$

$E \rightarrow E + T / T$

$T \rightarrow T^* F / F$

$F \rightarrow (E) / \text{id}$

$I_0 = \text{closure}(\{[E' \rightarrow .E]\})$

$E' \rightarrow .E$

$E \rightarrow .E + T$

$E \rightarrow .T$

$T \rightarrow .T^* F$

$T \rightarrow .F$

$F \rightarrow .(E)$

$F \rightarrow .\text{id}$

Constructing canonical  $LR(0)$

item sets

(cont.)

Goto  $(I, X)$  where  $I$  is an item set and  $X$  is a grammar symbol is closure of set of all items  $[A \rightarrow \alpha X \beta]$  where  $[A \rightarrow \alpha X \beta]$  is in

























