

BOTTOM-UP PARSING

A bottom up parse corresponds to construction of a parse tree for an i/p string beginning at the leaves (the bottom) and working up towards the root (the top).

eg: 1 bottom up parser for $id * id$. The grammar is

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

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

$id * id$

$F * id$

$T * id$

$T * F$

E

|

|

|

|

|

id

F

F

id

/

\

|

|

T

$*$

F

|

|

F

id

id

Reduction Reduction is the process of reducing a string 'w' to the start symbol of the grammar.

- At each reduction step, a specific substring matching the body of a production is replaced by the nonterminal at the head of the prod.

- A reduction is the reverse of a step in derivation. In derivations, a nonterminal in sentential form is replaced by the body of one of its prod.

- The goal of bottom-up parsing is to construct a derivation in reverse.

- The key decisions during bottom-up parsing are about when to reduce and about which prod to apply as the parse proceeds.

- The following corresponds to parse in ex 1.

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

This derivation is in fact a rightmost derivation.

Handle Pruning

(The process to construct a bottom-up parse during a left to right scan of the ip constructs a rightmost derivation in reverse.)

- A handle is a substring that matches the body of a prod, and whose reduction, represents one step along the reverse of a rightmost derivation.

parse is called handle pruning

eg: Adding subscripts to the tokens id for clarity. The handles during the parsing of $id_1 * id_2$ according to expr. grammar in eg (1) is given in table below:

Right Senterhial Form	Handle	Reducing Production
$id_1 * id_2$	id_1	$F \rightarrow id$
$F * id_2$	F	$T \rightarrow F$
$T * id_2$	id_2	$F \rightarrow id$
$T * F$	$T * F$	$T \rightarrow T * F$
T	T	$E \rightarrow T$

Consider the senterhial form $T * id$ and the pdn $E \rightarrow T$, Here T is not a handle in the senterhial form $T * id$. If T was indeed replaced by E , we would get the string $E * id$.

Shift Reduce Parsing

It is a form of bottom up parsing in which a stack holds grammar symbols and an i/p buffer holds the rest of the string to be parsed. The handle always appears at the top of the stack. We use $\$$ to mark the bottom of the stack and also the

right end of the i/p. Initially the stack is empty and the string w is on the i/p as follows

STACK	INPUT
\$	w \$

During a left to right scan of the i/p string the parser shift zero or more i/p symbols on to the stack until it is ready to reduce a string β of Φ grammar symbol on top of the stack.

It then reduce β to the head of the appropriate production the parser repeat this cycle until it has detected an error or until the stack contains the start symbol and the i/p is empty.

STACK	INPUT
$\$S$	$\$$

Upon entering this configuration the parser stops and announces successful completion of parsing.

$E \rightarrow E+T / T$

$T \rightarrow T * F / F$

$F \rightarrow (E) / id$

STACK	INPUT	ACTION
\$	id ₁ id ₂ \$	shift
\$ id ₁	* id ₂ \$	reduce $F \rightarrow id$
\$ F	* id ₂ \$	reduce $T \rightarrow F$
\$ T	* id ₂ \$	shift
\$ T *	id ₂ \$	shift
\$ T * id ₂	\$	reduce $F \rightarrow id$
\$ T * F	\$	reduce $T \rightarrow T * F$
\$ T	\$	reduce $E \rightarrow T$
\$ E	\$	accept

Four possible actions a shift reduce parser can make.

1. Shift - shift the next i/p symbol on the top of the stack.
2. Reduce - the right end of the string to be reduced must be at the top of the stack.
3. locate the left end of the string within the stack and decide with what non-terminal to replace the string.

3. accept - announce successful completion of parser.

4. error - discover a syntax error and hold called error recovery routine.

Conflict during shift reduce parsing

There are 2 types of conflict in SRP.

① shift or reduce conflict

The situation in which the parser cannot decide whether to shift or to reduce.

stmt \rightarrow .if expr then stmt / β

- if expr then stmt else stmt /
- other

STACK

INPUT

\$ if expr then stmt

else stmt \$

we cannot tell whether if expression then stmt is the handle. Here there is a shift or reduce conflict

② reduce-reduce conflict

The situation in which the parser cannot decide which of several productions to make

Consider the grammar,

1. $stmt \rightarrow id(parameters_list)$
2. $stmt \rightarrow expr := expr$
3. $parameter_list \rightarrow parameter_list / parameters$
4. $parameter \rightarrow parameter$
5. $parameter \rightarrow id$
6. $expr \rightarrow id(expr_list)$
7. $expr \rightarrow id$
8. $expr_list \rightarrow expr_list, expr$
9. $expr_list \rightarrow expr$

<u>STACK</u>	<u>INPUT</u>
id(id	,id)\$

After shifting, a shift reduce p would be in a configuration $id(id, id)\$$.

There is a conflict that which production either '3' or '7' can be used to reduce. This conflict is called reduce-reduce conflict.

OPERATOR PRECEDENCE PARSING

Operator Grammar:

It is a CFG that has the following properties:

There is no ϵ production on the right side

- there is no adjacent nonterminals on the right side of the pdr.

$$E \rightarrow EAE / (E) / -E / id$$

$$A \rightarrow + / - / * / / / \uparrow$$

This is not an operator grammar.

because $E \rightarrow EAA$ has 3 consecutive nonterminals

This can be converted into an operator grammar

i.e., $E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid E \uparrow E \mid (E) \mid -E \mid id$

Operator precedence parser:

- It is a bottom up parser that interprets an operator grammar.

- Operator precedence parser is the only one parser that can parse ambiguous grammars.

Disadvantages:

- It is hard to handle tokens like the minus sign, which has two different precedences (depending on whether it is unary or binary).

- Only a small class of grammars can be parsed using operator precedence technique.

Advantage:

- Simple

- Numerous compilers use operator precedence parser for handling expressions.

- In operator precedence parsing, we define 3 disjoint precedence r/n's between certain pair of terminals ($<$, $=$ and $>$).

$a < b$: a yields precedence to b

$a > b$: a takes precedence over b

$a = b$: a has the same precedence as b

• The determination of correct precedence relation b/w terminals are based on the traditional notations of associativity and precedence of operator.

• The intention of precedence relation is to find the handle of a right sentential form.

• The handle can be find out by the following process.

1. scan the symbol from left end until the first \rightarrow symbol is encountered.

2. Then scan backwards over any $=$ until a \leftarrow is encountered.

3. The handle contain everything to the left of the first \leftarrow greater than and to the right of the first \rightarrow less than encountered in step 2.

	id	+	*	\$
id		\rightarrow	\rightarrow	\rightarrow
+	\leftarrow	\rightarrow	\leftarrow	\rightarrow
*	\leftarrow	\rightarrow	\rightarrow	\rightarrow
\$	\leftarrow	\leftarrow	\leftarrow	

Consider the i/p string $id+id*id$ apply the precedence relation
 $E \rightarrow E+E / E * E / id$

~~$w \rightarrow id+id*id$~~

$w \rightarrow id+id*id \rightarrow \$id+id*id\$$

$\$ \leftarrow id \rightarrow + \leftarrow id \rightarrow * \leftarrow id \rightarrow \$$

$\$ \leftarrow + \leftarrow id \rightarrow * \leftarrow id \rightarrow \$$

$\$ \leftarrow + \leftarrow * \leftarrow id \rightarrow \$$

$\$ \leftarrow + \leftarrow * \rightarrow \$$

$\$ \leftarrow + \rightarrow \$$

$\$ \$$

Operator Precedence Relations from Associativity and Precedence

1. If operator \odot_1 has higher precedence than operator \odot_2 , make $\odot_1 \succ \odot_2$ and $\odot_2 \prec \odot_1$

eg: If $*$ is having higher precedence than $+$ make $* \succ +$ and $+ \prec *$

2. If \odot_1 and \odot_2 are operators of equal precedence then make $\odot_1 \succ \odot_2$ and $\odot_2 \prec \odot_1$, If operators are left associative, or make $\odot_1 \prec \odot_2$ and $\odot_2 \succ \odot_1$, if they are right associative

eg: If $+$ and $-$ are left associative, then make $+\succ+$, $+\succ-$, $-\succ+$, and $-\succ-$

If \uparrow is right associative, then make $\uparrow \prec \cdot \uparrow$.

3. make $\alpha \prec \cdot id$, $id \succ \cdot \alpha$, $\alpha \prec \cdot ($, $(\prec \cdot \alpha$, $) \succ \cdot \alpha$, $\alpha \succ \cdot)$,
 $\alpha \succ \cdot \$$ and $\$ \prec \cdot \alpha$ for all operators α .

Also let

(\equiv) $\$ \prec \cdot ($ $\$ \prec \cdot id$ $(\prec \cdot ($
 $id \succ \cdot \$$ $) \succ \cdot \$$ $(\prec \cdot id$ $id \succ \cdot)$ $) \succ \cdot)$

These rules ensure that id will be reduced to ϵ whenever bound and (ϵ) will be reduced to ϵ whenever bound.

- The operator precedence relations for grammar assuming

1. \uparrow is at highest precedence and right associative
2. $*$ and $/$ are at next precedence and left associative
3. $+$ and $-$ are at lowest precedence and left associative is given below:

	+	-	*	/	↑	id	()	\$
+	>	>	<	<	<	<	<	>	>
-	>	>	<	<	<	<	<	>	>
*	>	>	>	<	<	<	<	>	>
/	>	>	>	>	<	<	<	>	>
↑	>	>	>	>	<	<	<	>	>
id	>	>	>	>	>			>	>
(<	<	<	<	<	<	<	=	
)	>	>	>	>	>			>	>
\$	<	<	<	<	<	<	<		

Operator Precedence Parsing Algorithm:

Input: An i/p string w and a table of precedence relations.

Output: If ' w ' is well formed, a skeletal parse tree, with a placeholder nonterminal ϵ labeling all interior nodes, otherwise, an error condn.

Method: Initially the stack contains $\$$ and the i/p buffer contains the string $w\$$.

1. set ip to point to the first symbol of $w\$$

repeat forever

if $\$$ is on top of stack and ip points to $\$$ then accept and return.

else begin

let ' a ' be the topmost terminal symbol on the

stack and let b be the symbol pointed to by ip
 if $a < b$ or $a = b$ then begin
 push b onto the stack
 advance ip to the next ip symbol

end

else if $a > b$ then // reduce

 repeat

 pop the stack

 until the top stack element terminal is
 related by $<$ to the terminal most recently
 popped

else

 error();

end

eg:

stack	ip	ip buffer	Action
\$	*	$a+(b * c)$ \$	$\$ < a$ push
\$a		$+(b * c)$ \$	$a > +$ pop.
\$		$+(b * c)$ \$	$\$ < a$ noop
\$		$+(b * c)$ \$	$\$ < +$ push
\$+		$(b * c)$ \$	$+ < c$ push
\$+c		$b * c)$ \$	$c < b$ push
\$+cb		$* c)$ \$	$b > *$ pop
\$+c		$* c)$ \$	$(< b$ noop
\$+c		$* c)$ \$	$(< * $ push
\$+(*		$c)$ \$	$* < c$ push

$\$ + (* ($	$) \$$	$(\>) \text{ pop}$
$\$ + (*$	$) \$$	$* < (\text{ noop}$
$\$ + (*$	$) \$$	$* \geq) \text{ pop}$
$\$ + ($	$) \$$	$(< * \text{ noop}$
$\$ + ($	$) \$$	$(\doteq) \text{ push}$
$\$ + ()$	$\$$	$) \geq \$ \text{ pop}$
$\$ + ($	$\$$	$(\doteq) \text{ pop}$
$\$ +$	$\$$	$+ \geq (\text{ noop}$
$\\$ +$	$\\$	$\\$ < + \text{ no}$
$\$ +$	$\$$	$\$ < + \text{ noop}$
$\$$	$\$$	Accept.

Compilers using operator precedence parsers need not store the table of precedence relations. In most cases the table can be encoded by two precedence functions, f and g , which map terminal symbols to integers.

1. $f(a) < g(b)$ whenever $a < b$
2. $f(a) = g(b)$ whenever $a \doteq b$
3. $f(a) > g(b)$ whenever $a \geq b$

Thus the precedence r/n b/w a and b can be determined by a numerical comparison between $f(a)$ and $g(b)$.

LR parser is a bottom-up parser which is used to parse a large class of context-free grammars. In LR(k) parser; 'L' is for left to right scanning of the input, the 'R' for constructing a rightmost derivation in reverse, and 'k' for the number of input symbols of lookahead that are used in making parsing decisions.

LR parser is attractive because:

- > LR parsers can be constructed to recognize virtually all programming language constructs for which context free grammars can be written.
- > LR parsing method is the most general nonbacktracking shift-reduce parsing method known.
- > The class of grammars that can be parsed using LR method is a proper superset of the class of grammars that can be parsed with predictive parsers.
- > An LR parser can detect a syntactic error as soon as it is possible to do so on a left to right scan of the input.

The principal drawback of the method is that it is too much work to construct an LR parser by hand for a programming language.

The LR Parsing Algorithm.

The schematic form of an LR parser is shown below. It consists of an input, an output, a stack, a driver program, and a parsing table that has two parts, 'action' and 'goto'. The parsing program reads characters from input buffer one at a time. The stack store a string of the form $S_0 X_1 S_1 X_2 S_2 \dots X_m S_m$, where S_m is on top. Each X_i is a grammar symbol and S_i is state of DFA.

function action and a goto function goto. The parser determines s_m , the state currently on top of the stack, and a_i , the current input symbol, it then consults $\text{action}[s_m, a_i]$, which can have one of four values:

1. Shift s , where s is a state.
2. reduce by a grammar production $A \rightarrow \beta$,
3. accept,
4. error.

The goto takes a state and grammar symbol as argument and produces a state.

The parser table can be constructed using 3 methods:

1. Simple LR (SLR)
2. Canonical LR
3. Lookahead-LR (LALR).

Algorithm:

Input: An input string w and an LR parsing table with functions action and goto for a grammar G .

Output: If w is in $L(G)$, a bottom-up parse for w ; otherwise, an error indication.

Method. Initially, the parser has S_0 on its stack, where S_0 is the initial state, and $w\$$ in the input buffer. The parser then executes the program below until an accept or error action is encountered.

set ip to point to the first symbol of $w\$$;

repeat forever begin

let s be the state on top of the stack and

a the symbol pointed to by ip ;

if $\text{action}[s, a] = \text{shift } s'$ then

push a then s' on top of the stack,
advance ip to the next input symbol,

end

else if action $[s, a] = \text{reduce } A \rightarrow \beta$ then

begin

pop $2 * |\beta|$ symbol off the stack; $\dots \rightarrow$ eg: $2 * |E+T|$
 $= 2 * 3 = 6$

let s' be the state now on top of the stack;

push A then goto $[s', A]$ on top of the stack;

output the production $A \rightarrow \beta$.

end

else if action $[s, a] = \text{accept}$ then

return

else error()

end

Example

Consider the grammar for arithmetic expressions with binary operators $+$ and $*$:

(1) $E \rightarrow E+T$

(2) $E \rightarrow T$

(3) $T \rightarrow T * F$

(4) $T \rightarrow F$

(5) $F \rightarrow (E)$

(6) $F \rightarrow id$

The figure below shows the parsing action and goto functions of LR parsing table for the above grammar. The code for actions are:

1. s_j means shift and stack state j ,
2. r_j means reduce by production numbered j ,
3. acc means accept
4. blank means error.

state	action						goto		
	id	+	*	()	\$	E	T	F
0	S ₅			S ₄			1	2	3
1		S ₆				acc			
2		r ₂	S ₇		r ₂	r ₂			
3		r ₄	r ₄		r ₄	r ₄			
4	S ₅			S ₄			8	2	3
5		r ₆	r ₆		r ₆	r ₆			
6	S ₅			S ₄				9	3
7	S ₅			S ₄					10
8		S ₆			S ₁₁				
9		r ₁	S ₇		r ₁	r ₁			
10		r ₃	r ₃		r ₃	r ₃			
11		r ₅	r ₅		r ₅	r ₅			

On input $id * id + id$, the sequence of stack and input content shown below.

	stack	input	Action.
1)	0	$id * id + id \$$	shift.
2)	0 id 5	$* id + id \$$	reduce by $F \rightarrow id$
3)	0 F 3	$* id + id \$$	reduce by $T \rightarrow F$
4)	0 T 2	$* id + id \$$	shift
5)	0 T 2 * 7	$id + id \$$	shift
6)	0 T 2 * 7 id 5	$+ id \$$	reduce by $F \rightarrow id$
7)	0 T 2 * 7 F 10	$+ id \$$	reduce by $T \rightarrow T * F$
8)	0 T 2	$+ id \$$	reduce by $E \rightarrow T$
9)	0 E 1	$+ id \$$	shift
10)	0 E 1 + 6	$id \$$	shift
11)	0 E 1 + 6 id 5	$\$$	reduce by $F \rightarrow id$
12)	0 E 1 + 6 F 3	$\$$	reduce by $T \rightarrow F$
13)	0 E 1 + 6 T 9	$\$$	$E \rightarrow E + T$
14)	0 E 1	$\$$	accept.

There are three methods to construct an LR parsing table from a grammar.

1. Simple LR or SLR :- is the weakest of the three in terms of the number of grammars for which it succeeds, but is the easiest to implement.
2. Canonical LR :- is the most powerful and most expensive.
3. Lookahead LR (LALR) :- is intermediate in power and cost between the other two.

SLR Parsing Tables

The parsing table constructed by this method is an SLR table, and the LR parser using an SLR parsing table is called an SLR parser. A grammar for which an SLR parser can be constructed is said to be an SLR grammar.

The central idea in the SLR method is first to construct a DFA from the grammar to recognize viable prefixes [The set of prefixes of right sentential form that can appear on the stack of a shift-reduce parser are called viable prefixes].

To construct the DFA we make use of augmented grammar. If G_1 is a grammar with start symbol S , then G_1' the augmented grammar for G_1 , is G_1 with a new start symbol S' and production $S' \rightarrow S$. The purpose of this new starting production is to indicate to the parser when it should stop parsing and announce acceptance of the input.

Consider the following grammar:

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

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

↳ The augmented grammar is

$$E' \rightarrow E$$

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

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

LR(0) items

An LR(0) item of a grammar G is a production of G with a dot at some position of the right side. Thus production $A \rightarrow XYZ$ yields the four items

$$A \rightarrow \cdot XYZ, A \rightarrow X \cdot YZ, A \rightarrow XY \cdot Z, A \rightarrow XYZ \cdot$$

A production $A \rightarrow E$ generate only one item $A \rightarrow \cdot$.

If the length of the right side of the production is 'n', then there are $(n+1)$ items can be generated.

The items that can be generated using the above productions are $T \rightarrow T * F$ are

$$T \rightarrow \cdot T * F, T \rightarrow T \cdot * F, T \rightarrow T * \cdot F, T \rightarrow T * F \cdot$$

Canonical collection

A set of items that corresponds to the states of DFA that recognizes viable prefixes is called canonical collection.

Construction of a DFA involves finding canonical collections. Canonical collection provide basis for SLR parsers. To construct canonical collection make use of the two functions.

1. Closure

2. goto.

Closure operation.

If I is a set of items for a grammar G , then $\text{closure}(I)$ is the set of items constructed from I by the two rules:

1. Initially, every item in I is added to $\text{closure}(I)$
2. If $A \rightarrow \alpha \cdot B \beta$ is in $\text{closure}(I)$ and $B \rightarrow \gamma$ is a production, then add the item $B \rightarrow \cdot \gamma$ to I , if it is not already there. We apply this rule until no more new items can be added to $\text{closure}(I)$.

In our example $\text{closure}(E' \rightarrow \cdot E) = \left\{ \begin{array}{l} E' \rightarrow \cdot E, \\ E \rightarrow \cdot E + T, \\ E \rightarrow \cdot T, \\ T \rightarrow \cdot T * F, \\ T \rightarrow \cdot F, \\ F \rightarrow \cdot (E), \\ F \rightarrow \cdot id \end{array} \right\}$

The goto operation.

$\text{goto}(I, X)$ where I is set of items and X is grammar symbol, is defined to be the closure of the set of all items $A \rightarrow [\alpha X \beta]$ such that $[A \rightarrow \alpha \cdot X \beta]$ is in I

In our example consider $I = \{ E' \rightarrow E \cdot, E \rightarrow E \cdot + T \}$, then

$\text{goto}(I, +) = \left\{ \begin{array}{l} E \rightarrow E + \cdot T, \\ T \rightarrow \cdot T * F, \\ T \rightarrow \cdot F, \\ F \rightarrow \cdot (E), \\ F \rightarrow \cdot id \end{array} \right\}$

Now we can construct the canonical collection of sets of LR(0) items for an augmented grammar G' ,

The algorithm for canonical collection is:

algorithm for canonical collection is :

procedure items (G)

begin

$C := \text{closure}(S' \rightarrow \cdot S)$

repeat

for each set of items I in C and each grammar symbol X , such that $\text{goto}(I, X)$ is not empty, and not in C do

add $\text{goto}(I, X)$ to C

until no more set of items can be added to C .

end.

Eg: $E \rightarrow E+T/T$
 $T \rightarrow T \times F/F$
 $F \rightarrow (E)/id$

Its augmented grammar is:

$E' \rightarrow \cdot E$
 $E \rightarrow \cdot E+T/T$
 $T \rightarrow \cdot T \times F/F$
 $F \rightarrow \cdot (E)/id$

- reductio no.
- 1) $E \rightarrow E+T$
 - 2) $E \rightarrow T$
 - 3) $T \rightarrow T \times F$
 - 4) $T \rightarrow F$
 - 5) $F \rightarrow (E)$
 - 6) $F \rightarrow (id)$

Now, closure of $E' \rightarrow \cdot E = \{$

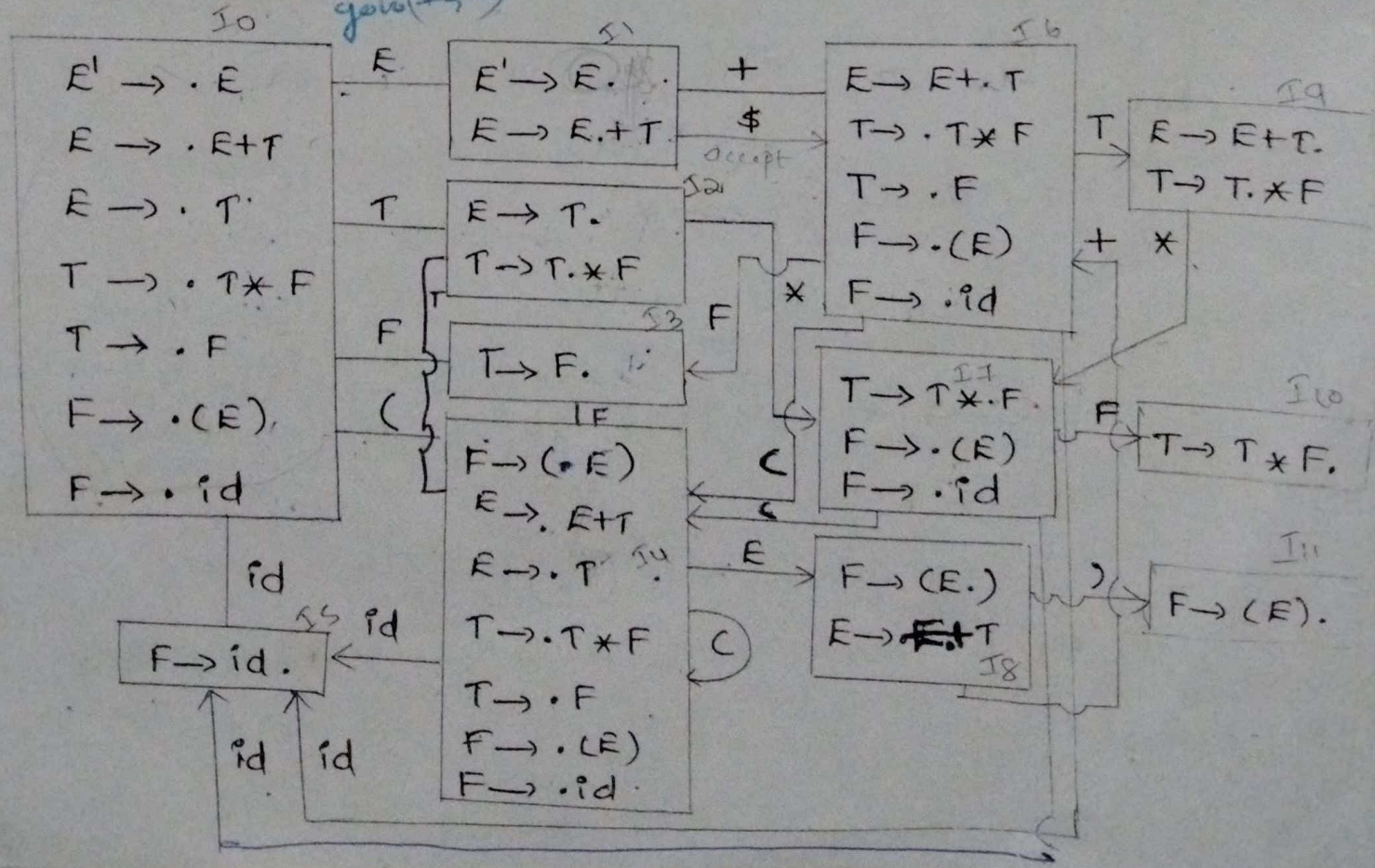
- $E' \rightarrow \cdot E$
- $E \rightarrow \cdot E+T$
- $E \rightarrow \cdot T$
- $T \rightarrow \cdot T \times F$
- $T \rightarrow \cdot F$
- $F \rightarrow \cdot (E)$
- $F \rightarrow \cdot id$

} I_0

So,

DFA

$goto(I_0, E)$



for constructing an SLR parsing table

Input. An augmented grammar G' .

Output. The SLR parsing table functions action and goto for G' .

Method.

1. Construct $C = \{I_0, I_1, I_2, \dots, I_n\}$, the collection of sets of LR(0) items for G' .
2. State i is constructed from I_i . The parsing actions for state i is determined as follows:
 - a) If $[A \rightarrow \alpha \cdot a \beta]$ is in I_i and $\text{goto}(I_i, a) = I_j$, then set $\text{action}[i, a]$ to 'shift j '. Here a must be a terminal.
 - b) If $[A \rightarrow \alpha \cdot]$ is in I_i , then set $\text{action}[i, a]$ to "reduce $A \rightarrow \alpha$ ". for all a in $\text{FOLLOW}(A)$; here A may not be S' .
 - c) If $[S' \rightarrow \cdot S]$ is in I_i , then set $\text{action}[i, \$]$ to 'accept'.if any conflicting actions are generated by the above rules, we say the grammar is not SLR(1).
3. The goto transitions for state i are constructed for all nonterminals A using the rule: If $\text{goto}(I_i, A) = I_j$ then $\text{goto}[i, A] = j$.
4. All entries not defined by rules (2) and (3) are made 'error'.
5. The initial state of the parser is the one constructed from the set of items containing $[S' \rightarrow \cdot S]$.

SLR(1) parse table for our example grammar is as shown below.

State	Action						Jobs		
	id	+	*	()	\$	E	T	F
0	S5			S4			1	2	3
1		S6				acc			
2		02	S7		02	02			
3		04	04		04	04			
4	S5			S4			8	2	3
5		06	06		06	06			
6	S5			S4				9	3
7	S5			S4					10
8		S6			S11				
9		01	S7		01	01			
10		03	03		03	03			
11		05	05		05	05			

CANONICAL LR (CLR OR LR)

In SLR, the FOLLOW set may contain more elements than actual lookahead (which can be drawback)

- In CLR, the lookahead symbol is calculated based on the pdn. Every pdn has a lookahead associated with it.

- CLR requires more memory consumption as it generates more no. of states as the lookahead symbols are associated with the pdn. And for different states lookahead symbols will vary.

- The extra info of lookahead is incorporated into the state by redefining items to include a terminal symbol as a second component.

The general form of an item becomes

$$A \rightarrow \alpha \cdot \beta, a$$

where $A \rightarrow \alpha \beta$ is a pdn, and a is a terminal or the right endmarker $\$$.

Such an object is called as LR(1) item. The 1 refers to the length of the 2nd component, called the lookahead of the item.

- The lookahead has no effect in an item of the form $A \rightarrow \alpha \cdot \beta, a$ where β is not ϵ .

But an item of the form $A \rightarrow \alpha \cdot, a$ calls for reduction by $A \rightarrow \alpha$ only if the next i/p symbol is a .

Constructing LR(1) set of items:

- The method for building the collection of set of valid LR(1) items is essentially the same as that of canonical collection of sets of LR(0) items. The two procedures CLOSURE and GOTO has to be modified.

Set of items CLOSURE (I)

{

repeat

for (each item $[A \rightarrow \alpha \cdot B \beta, a]$ in I)

for (each pdn $B \rightarrow \gamma$ in G')

for (each terminal b in $\text{FIRST}(\beta a)$)

add $[B \rightarrow \cdot \gamma, b]$ to set I;

until no more items are added to I;

return I;

}

Set of items GOTO (I, x)

{

initialise J to the empty set;

for (each item $[A \rightarrow \alpha \cdot x \beta, a]$ in I)

add item $[A \rightarrow \alpha x \cdot \beta, a]$ to set J;

return CLOSURE (J);

}

void items(a')

{

initialize C to {CLOSURE({[S' → • S, \$])}};

repeat

for (each set of items I in C)

for (each grammar symbol x)

if (GOTO(I, x) is not empty and not in C)

add GOTO(I, x) to C;

until no new sets of items are added to C;

}

eg: Suppose given $s \rightarrow cc$

$c \rightarrow cc$

$c \rightarrow d$

step 1: Create Augmented Grammar

$s' \rightarrow \cdot s$

$s \rightarrow \cdot cc$

$c \rightarrow \cdot cc$

$c \rightarrow \cdot d$

In canonical LR, along with the pdn, lookahead symbol is passed.

step 2: Compute closure:

We begin by computing the closure of $s' \rightarrow \cdot s, \$$ here $\$$ is the lookahead symbol as it is the accepting state.

Consider the items for CLOSURE (I)

here I contains $s' \rightarrow \cdot s, \$$, which is similar to $A \rightarrow \alpha \cdot B \beta, a$

$$\text{Here, } A = s'$$

$$\alpha = \epsilon$$

$$B = s$$

$$\beta = \epsilon \text{ and}$$

$$a = \$$$

Function closure tells to add $B \rightarrow \cdot \gamma, b$, for each pdn $B \rightarrow \gamma$ and terminal b in $\text{FIRST}(\beta a)$. Thus, must add $s \rightarrow \cdot c, c$ and here b is $\$ \because \beta = \epsilon$ and $a = \$$, b may be only $\$$

Illy, compute closure for all items.

$$I_0 : s' \rightarrow \cdot s, \$$$

$$s \rightarrow \cdot c, c, \$$$

$$c \rightarrow \cdot c, c, c/d$$

$$c \rightarrow \cdot d, c/d$$

$$\left. \begin{array}{l} c \rightarrow \cdot c, c, c/d \\ c \rightarrow \cdot d, c/d \end{array} \right\} \Rightarrow \begin{array}{l} c \rightarrow \cdot c, c, c \\ c \rightarrow \cdot c, c, d \\ c \rightarrow \cdot d, c \\ c \rightarrow \cdot d, d \end{array}$$

step 3:

$$\text{GOTO}(I_0, s) \Rightarrow s' \rightarrow s \cdot, \$ = I_1,$$

$$\text{for GOTO}(I_0, c) \Rightarrow s \rightarrow c \cdot, c, \$$$

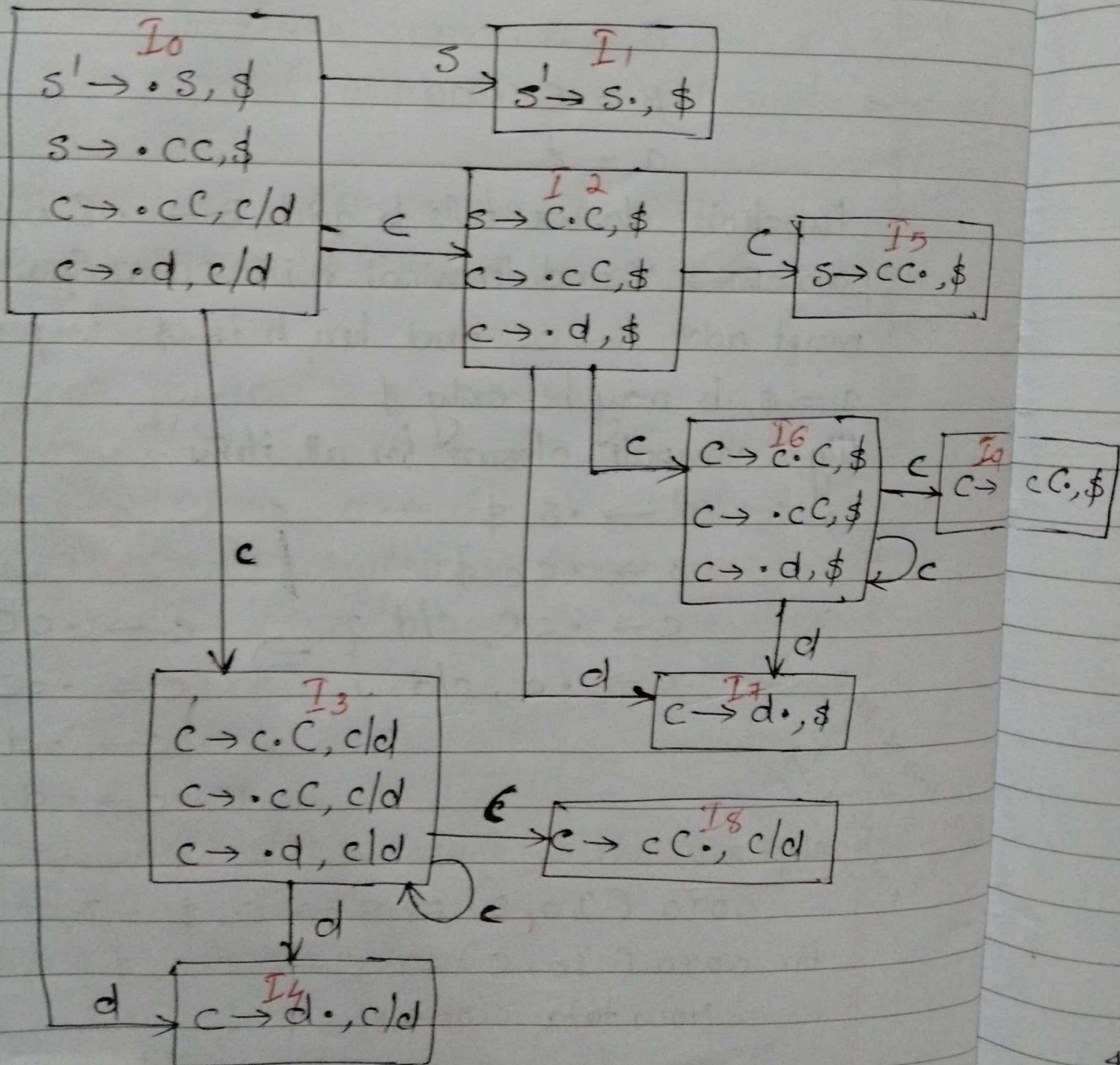
Now take CLOSURE

$$c \rightarrow \cdot c, c, \$$$

$$c \rightarrow \cdot d, \$$$

ie, I_2 : $S \rightarrow c \cdot C, \$$
 $c \rightarrow \cdot cC, \$$
 $c \rightarrow \cdot d, \$$

||ly carry out for the remaining states, and generate DFA or LR(1) automaton.



CANONICAL LR(1) PARSING TABLE CONSTRUCTION

LR(1) ACTION and GOTO bns are constructed from LR(1) items. The bns are represented by a table.

ALGORITHM: Constr of LR(1) Parsing Table

INPUT: Augmented Grammar, G'

OUTPUT: Canonical LR parsing table bns ACTION and GOTO for G'

Method:

1. Construct $C' = \{I_0, I_1, \dots, I_n\}$, the collection of sets of LR(1) items for G' .

2. state i of the parser is constructed from I_i . The parsing fraction for state i is determined as follows.

a) If $[A \rightarrow \alpha \cdot a \beta, b]$ is in I_i and $GOTO(I_i, a) = I_j$ then set ACTION $[i, a]$ to "shift j ". Here, a must be a terminal.

b) If $[S' \rightarrow s \cdot, \$]$ is in I_i , then set ACTION $[i, \$]$ to "accept".

If any conflicting actions from the above rules, we say the grammar is not LR(1).

3. The GOTO transitions for state i are constructed for all nonterminals A using the rule: If $GOTO(I_i, A) = I_j$, then $GOTO[i, A] = j$.

4. All entries not defined by rules (2) and (3) are made "error".

5. The initial state of the parser is the one constructed from the state set of items containing $[s' \rightarrow \cdot s, \$]$

So the table formed for the grammar

1. $s \rightarrow cc$

2. $c \rightarrow cc$

3. $c \rightarrow d$ is given below

STATE	ACTION			GOTO	
	c	d	\$	s	c
0	S3	S4		1	2
1			acc		
2	S6	S7			5
3	S3	S4			8
4	r3	r3			
5			r1		
6	S6	S7			9
7			r3		
8	r2	r2			
9			r2		

LALR (LOOKAHEAD LR) PARSER.

Date _____

Page _____

Memory consumption is very high for CLR as the no. of states generated is more. (because of lookaheads included in each pair)

- This can be reduced, if some states are merged which are having same items with same or different lookaheads.

LR parsers using LALR parsing table is called LALR parser. The table obtained by it are considerably smaller than CLR tables.

- The SLR and LALR tables of a grammar always have the same no. of states.
- Consider the example discussed in the previous section.

When we look at the states I_4 and I_7 , we'll find that the 1st component of the items (core) are same but the lookaheads are different (c/d for I_4 and $\$$ for I_7).

~~Since~~

Similarly I_3 and I_6 are forms another pair with core

$\{ C \rightarrow c.C, C \rightarrow .cC, C \rightarrow d. \}$

One more pair I_8 and I_9 is common with

$\{ C \rightarrow cC. \}$

LALR Parsing table

ALGORITHM:

I/p: An augmented grammar G'

O/p: LALR parsing table b/w ACTION & GOTO for G'

Method:

1. construct $C = \{I_0, I_1, \dots, I_n\}$ the collection set of LR(1) items.
2. For each core present among the set of LR(1) items, bind all sets having that core, and replace these sets by their union.
3. let $C' = \{J_0, J_1, \dots, J_m\}$ be the resulting sets of LR(1) items. The parsing action and state i are constructed from J_i in the same manner as canonical LR parsing table items. If there is a parsing action conflict, the algorithm fails to produce a parser and the grammar is said not to be LALR(1).
4. The GOTO table is constructed as follows:
If J is the union of one or more sets of LR(1) items, i.e., $J = I_1 \cup I_2 \cup \dots \cup I_k$, then the core of $\text{GOTO}(I_1, x)$, $\text{GOTO}(I_2, x) \dots \text{GOTO}(I_k, x)$ are same. since $I_1, I_2 \dots I_k$ have the same core.
let K be union of all sets of items having the

same core as $\text{GOTO}(I_i, x)$. Then $\text{GOTO}(J, x) = k$.
 - The collection of sets of items constructed in step 3 is called LALR(1) collections.

eg: Consider the grammar:

$$s' \rightarrow s$$

$$s \rightarrow cc$$

$$c \rightarrow cC/d$$

and canonical parsing table.

Here I_3 and I_6 have the same core. So we can merge these two states

$I_{36} : c \rightarrow c \cdot C, c/d/\$$
 $c \rightarrow \cdot cc, c/d/\$$
 $c \rightarrow \cdot d, c/d/\$$

Similarly I_4 and I_7 are replaced by their union

$I_{47} : c \rightarrow d \cdot, c/d/\$$

and I_8 and I_9 are replaced by their union

$I_{89} : c \rightarrow cc \cdot, c/d/\$$

The LALR parsing table is shown below:

State	ACTION.			GOTO	
	c	d	\$	s	C
0	S36	S47		1	2
1			acc		
2	S36	S47			5
36	S36	S47			89
47	r3	r3	r3		
5			r1		
89	r2	r2	r2		

