

Static array - array whose size is fixed during declaration

Semi-dynamic array - size during runtime.

Dynamic array - array whose size is decided during runtime or data type

int * p = new int [size];

After dynamic memory allocation size cannot change.

Once array is created it cannot change.

- # During execution if type of array can be changed then dynamic array. Type → (data type or size).
- # During a function call when control is passed to a function then the called function uses the same registers as the calling function since number of registers is fixed. During transfer of control saving of data in registers is done. Before transferring control to called function the calling function again puts back the data in registers. The called function stores the contents of registers in main memory, not the calling function. Along with register contents, some local variables are also stored.

When function A calls B, the starting address of B is stored and return address is stored in one of the ^{total} registers of A. There are 16 general purpose registers numbered from 0 - 15.

Registers numbered from 0 - 15.

→ address of instruction → general purpose register

1004 BALR 14, 15 → branching address / address
return address of called unit

[register 14 will store return address of called unit]

BR 14

[branch to content of register 14] last instruction
of called unit.

The responsibility of registers is given to called unit.

#

A → B → C

During return C restores contents of reg B and B
restores contents of A.

Register 13 is used to save address of register save area.

When called unit is returning a single value
that value is stored in register 0.

When more than one value the address of these values is stored in register1.

- # FORTRAN, COBOL does not support recursion → does not have separate code segment and data segment
- # Scope Rules of a language - how a language access local variables.

Block Structured language -

variable - entity

unit — entity .

Each entity will have certain attributes .

Variable's attributes - type, scope, lifetime, lval, rval

decides size

name of a variable is not a property, its an identifier .
lval → address allocated to that name .

If the name appears on the left side of an assignment then its lval is used and if on right side of assignment then its rval is used .

$$a = a + b .$$

$$\text{lval}(a) = \text{rval}(a) + \text{rval}(b) .$$

lval is associated only with identifiers, not with expressions. So,

$$x + y = 3 ;$$

is not a valid statement because LHS is an expression.

len → Unix

Flex → Linux (open source)

Specifications to len:-

\$ cc test.c → for compiling a C program.

a.out

\$ file test.c → for finding type of file.

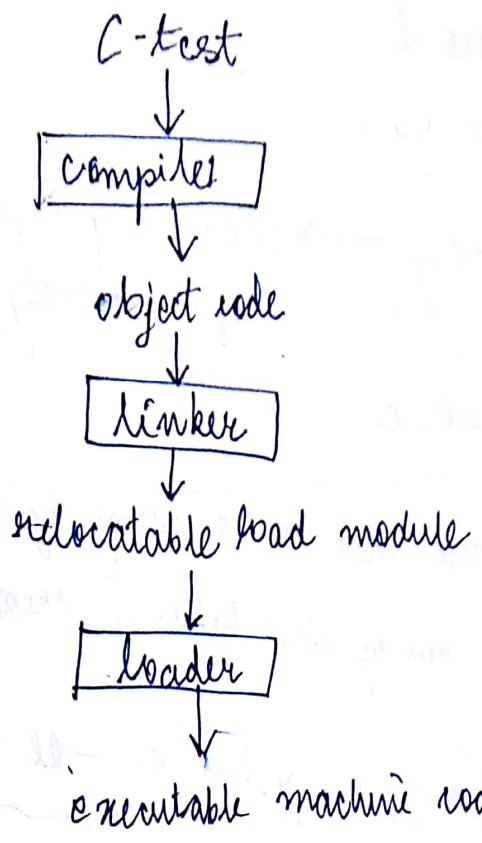
\$ file a.out → pure executable statements

\$ cat a.out → displays junk character.

Difference b/w join and merge is join is vertical and merge is horizontal.

In \$ cc test.c, cc is not just C compiler.

Along with compiling cc also performs linking and loading.



.l is for lex specification

The command for checking lex specification is,

\$ lex test.l

Declaration section → grammar rules.
 Rules section → sections of lexical analyzer.
 codes of rule section

Every rule will contain two parts. The left part is a regular expression (specification) and right part contains the action to be taken.

The lex generates another C-program

\$ lex test.l

lex.y.c

\$ ~~yy~~^{yy-d} test.y → contains specification for parser or
syntax analyzer.
generates
→ y.tab.c

If -d is not used the existing symbol table is modified
otherwise a new symbol table is created.

\$ cc lex.y.c y.tab.c -ll -ly

for including library during
linking.

The library routines are precompiled and stored in
specific locations.

#include <math.h>

.h files contain function prototypes and not the
actual code.

#include "math.h"

↳ we can give location of using
double quotes ↳

The default directory of include statement is
"/~~user~~/include".

- * The name of the directory itself is "include".
* While checking for presence of a.out file, it is checked in default directories along with current directory.

~~lex~~. /a.out → a.out of current directory.

The environment variable PATH has value

$$\text{PATH} = \$\text{PATH}:$$

To this variable we need to add another directory.

A particular user's profile is specified in .profile.

Whenever a profile is created its default .profile is created.

lex/flex - scanner

- lex is a scanner generator

input - set of regular expression and associated actions (written in C)

output - table-driven scanner (lex.yy.c)

- flex - an open source implementation of the original UNIX lex-utility

lex input example -

filename: ex1.l
→ first part empty

% .% → Regular expression

"hello world" printf("GOOD BYE\n"); ;

{} second part

→ prints GOOD BYE whenever string "hello world" is encountered.

→ third part empty does nothing for any other character

Using lex -

1. lex ex1.l → process lex file to generate a scanner (gets saved as lex.yy.c)

2. cc lex.yy.c -ll → compile the scanner and grab main() from the lex library (-ll option)

3. ./a.out hello world → run scanner taking input from std::input

GOOD BYE → standard output

#

- ✓ loweralpha
- ✓ upperalpha
- ✓ num
- ✓ digits
- ✓ words
- ✓ characters
- ✓ whitespaces
- ✓ lines.
- ✓ special character

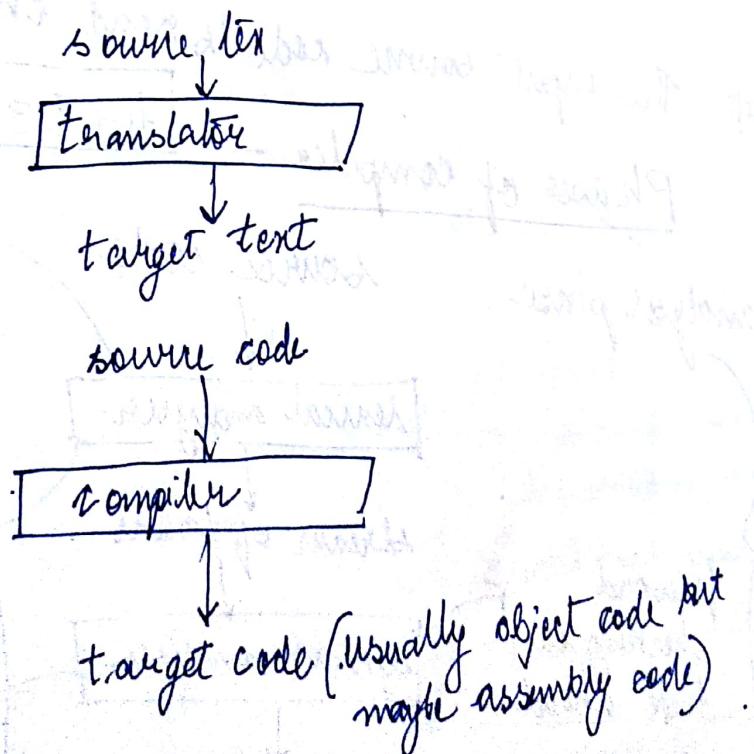
✓ int
✓ float

Anything other than digits, alphabet and whitespace is a special character.

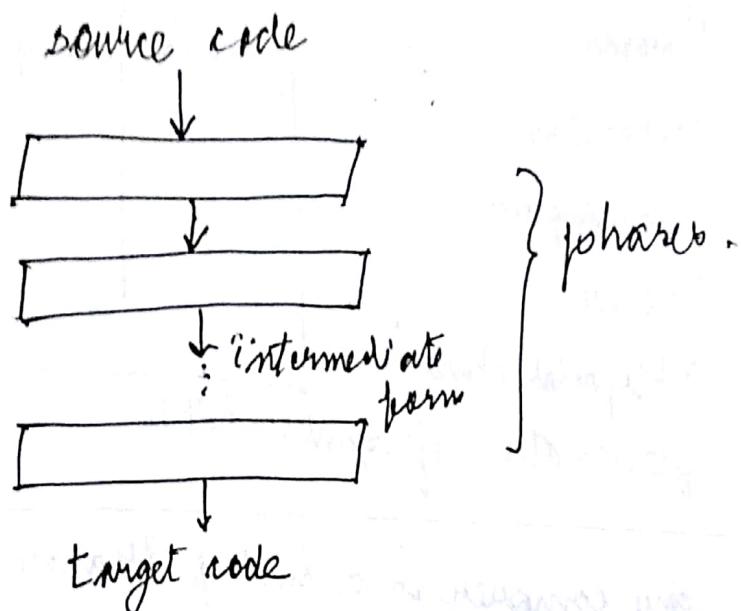
Symbol table -		array or LL
word	frequency	line number
		(a word in one line is not considered only once).

generation of symbol table

Any compiler is a kind of translator

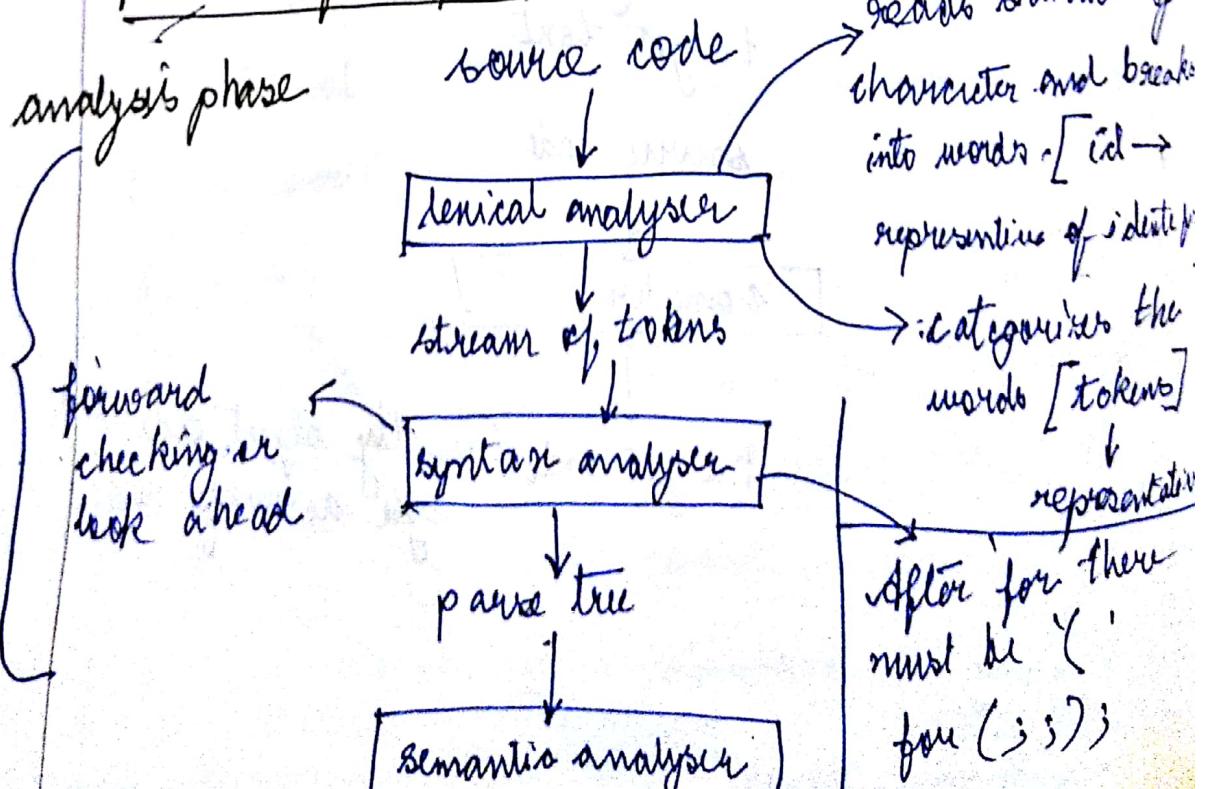


If converting source code to target code is done in steps by converting into some intermediate form, the intermediate form is converted into target form in phases.

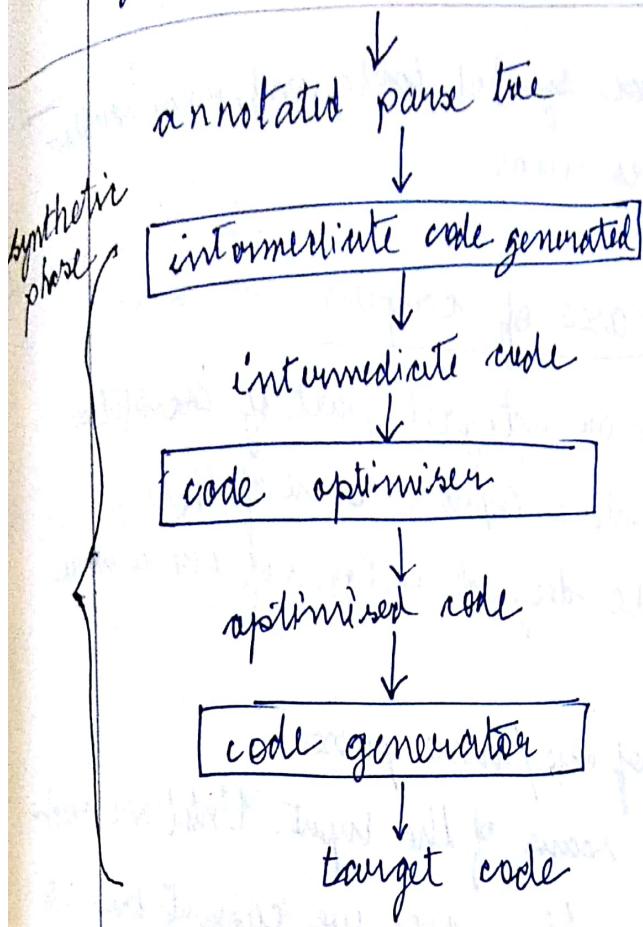


The input source code is read character by character.

Phases of compiler -



Syntax analysis generates parse tree → so also known as parser.



Syntax analyzer checks grammar rules and checks whether the string is an yield of a parse tree.

(membership problem)

(order of derivation)

(if the string is not a member then check for error)

Semantic analysis -

Type checking is done

$$i = n * 10 \\ \downarrow \\ \text{float int}$$

semantic → means 'logical meaning' of a statement.

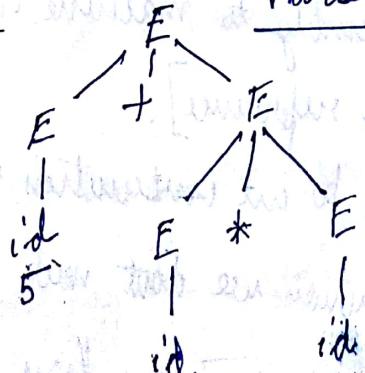
i) operational semantics
(imperative lang.)

ii) denotational semantics

iii) axiomatic semantics

$$\begin{array}{l} E \rightarrow E + E / E * E / (E) \\ E \rightarrow id / num \end{array}$$

Annotated Parse tree



Parse tree is travelled in bottom up in last order.

We get a dependency graph.

Topographical ordering gives scenario of evaluation.

We have tools to describe lexical and syntax analyzer, not for semantic analyzer

These 6 phases use symbol table and ~~avtar~~ ^{Scan Scan} handles finds and notifies error.

Phases of compiler

Phase - a tool or an integral part of translation.

Pass - scan the entire input. Some of these translations can be designed in terms of one or more passes.

A phase consists of one/more passes.

We need multiple scans of the input. Until we ~~read~~ read the input more than once we cannot translate properly.

some instructions require two passes for proper translation from assembly to machine language.

[Eg - macros, forward reference].

If we are branching to an instruction which is behind current instruction we don't need two passes.

In case of macros we require more than two passes.

The preprocessing is done even before lexical analysis.

In compiler design we look at phases, not passes.

The lexical analysis is the only phase where input is read character by character. After this no other phase reads input. It uses output of previous phase.

In code optimisation phase we reduce size of code by modifying intermediate code. Thus we remove unnecessary instructions and reduce overall execution time.

lexical analyzer - (scan input)

Find longest word that matches a particular rule for tokenization.

Eg. id: alpha(alpha/digit)*. [rule for identifier].

That is lexical analyzer matches longest applicable rule classifying the input string and returns a representation for classified set (token).

a
ab } all these strings are identifiers
abc
ab2 } (lexical analyzer returns a representative
of the "identifier set".

$\rightarrow \quad a = b + c$ [5 tokens]
id aop id bop id [identifiers and operations]

\rightarrow The type of variables is to be defined before using it.
int a, b, c;

When this definition is read the information about symbols used in the program is stored in symbol table. All attributes associated with an identifier are stored in symbol table [real, eval, lifetime, scope, type].

The symbol table consists of pointers, variables. A pointer of the symbol table points to an array of characters containing the name of the variable.

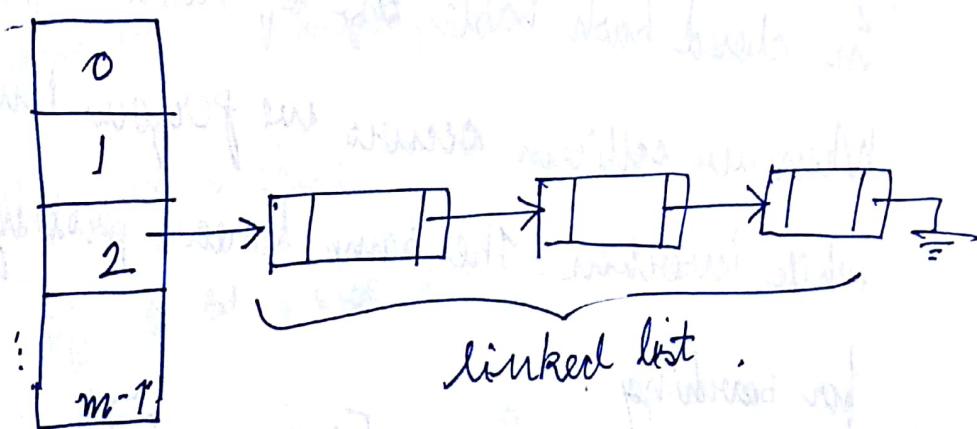
Operations on symbol table -

- (i) accessing the variable
- (ii) checking whether variable is already present
- (iii) inserting a variable in symbol table
- (iv) deleting a variable from symbol table when its scope finishes

In symbol table for each node we store more than one fields. [AVL tree \rightarrow height balanced tree
B tree \rightarrow balanced tree (leaf nodes are at same level)].

Hash table can be used to store symbol table .
[Two methods \rightarrow open hash table
closed hash table].

In open hash table we convert a string into an integer, x which is an index of table.



If there are n variables and table size is m then each linked list is assumed to be of size $\frac{n}{m}$ on an average. A particular rule is followed

for hashing - "abc" \rightarrow 6 ? both strings hashed to same location
"bac" \rightarrow 6
 $x \rightarrow h \rightarrow i$
 $k \leq m$

For finding k ,

$$3 * \underline{32^0} + 2 * \underline{32^1} + 1 * \underline{32^2} =$$

The base is taken as a power of 2, so that we don't need to multiply, we can just perform left shift.

But in open hash table the size of linked list is fixed. So we use closed hash table.

In closed hash table size of table is fixed.

Whenever collision occurs we perform linear probing while insertion. The same linear probing is used for searching.

$$\text{linear probing} \rightarrow [h(a) + i]$$

The problem of linear probing is formation of clusters. In order to avoid that we go for quadratic probing.

$$[h(a) + i^2] \rightarrow \begin{array}{l} \text{we may not be able to} \\ \text{reach all locations in} \\ \text{table.} \end{array}$$

If load factor < 0.5 , then free spaces are occupied using quadratic probing.

$$\text{load factor} = \frac{\text{occupied space}}{\text{total space}}$$

* consider an arithmetic expression,

$x = a * b + 10 ;$ incomplete in C without semi-colon

The ; in Pascal is a separator of statements and a terminator in C-language.

Pascal - if e then s1 else s2

C - if e s1 ; else s2 ;

~~alpha~~ alpha (alpha | digit)* → rule applied to recognise identifiers.

The text must match one

of the rules:

$x \rightarrow \text{lexeme}$ $x = a * b + 10$

$id \rightarrow \text{token}$ $id = id * id + num$

[stored in symbol table].

lex analyzer returns token and actual text through interface

Symbol table

	Token	Type	Scope		
x	id	float			
a		float			
b		float			
10	num	int			

The information related to each lexical entry is stored into symbol table. The symbol table is used by subsequent phases.

Stream → delimited

string → insequeut; no delimiter.

The next phase checks a derivation sequence for $id = id + id + num$.

We can construct a derivation tree for this string. Derivation tree is sometime called parse tree. In syntax analysis, again we are checking membership problem. In lexical analysis also we are checking membership of a string in a language of some grammar rules.

LEXICAL ANALYSIS

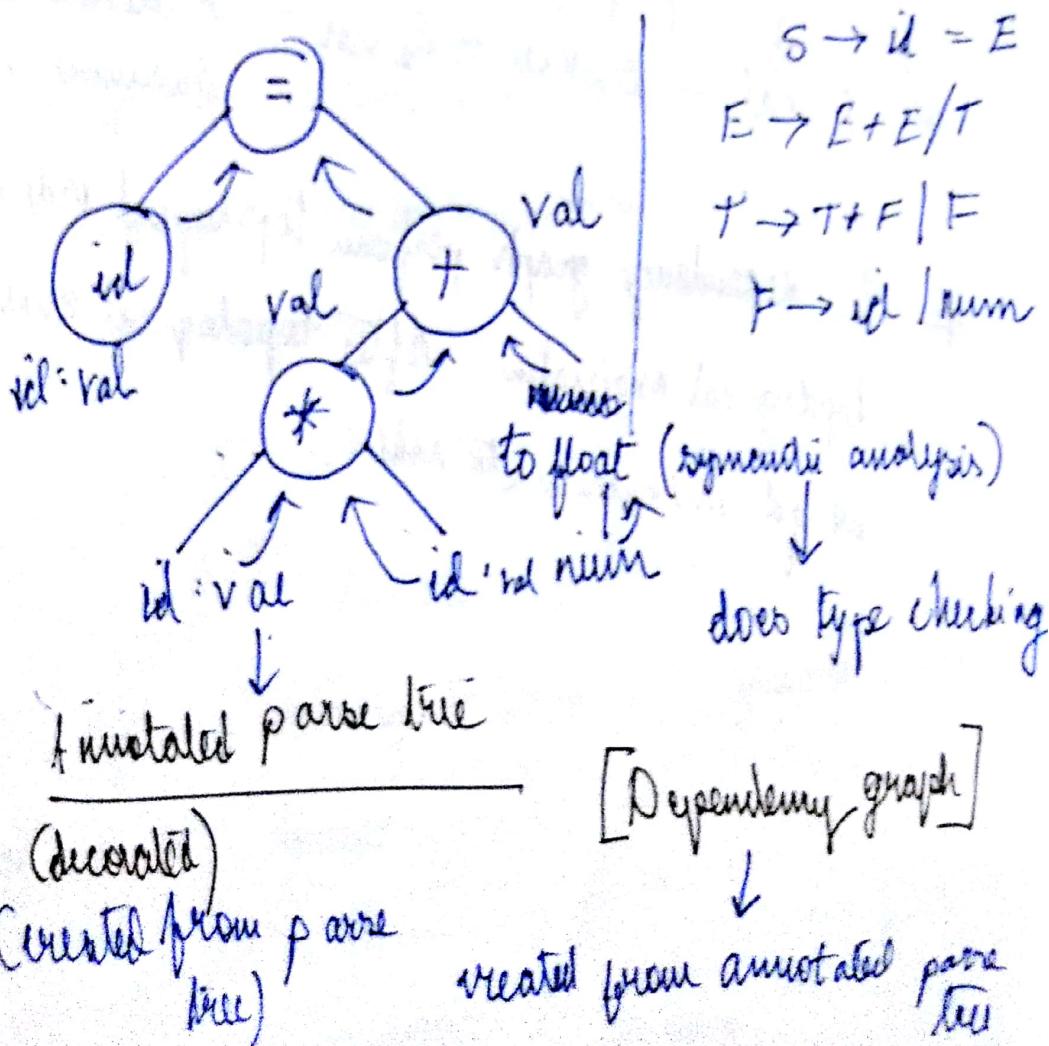
Lexical analyser uses recursive descent and syntax analyser uses recursive constraint.

$$E \rightarrow E + E \quad (\text{recursive})$$

$$E \rightarrow AB \quad (\text{non-recursive})$$

Lexical and syntax analysis is done simultaneously but for ease of understanding we take them as separate phases. Lexical analyser is an integral part of syntax analyser.

Derivation Tree ($id = id \neq id + num$)



The grammar rules must be in Backus-Naur Form, or in EBNF.

↓
some special characters represent many variables
other characters

\$, um *, → represent "all" → meta-characters

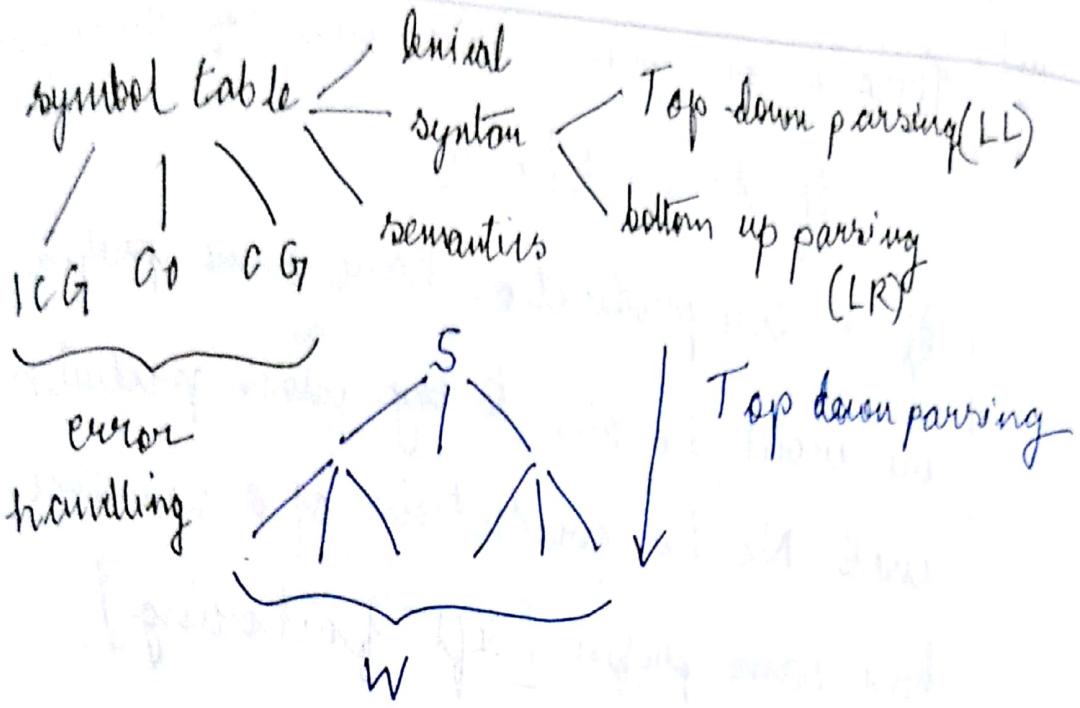
\$, um - i, *
[enquiries before deleting].

" is a meta-character in yacc.

' is a meta-character in yacc.

$E = E + E$
 $E.\text{val} = E_1.\text{val} + E_2.\text{val} \rightarrow \text{attributed grammar rule}$

For dependency graph perform topological ordering or
topological evaluation. After topological evaluation
we get intermediate rule.



Most compilers can't perform top down parsing.

Expression grammars can be parsed using top down parsing.

Most general parsing technique is bottom-up.

LLParser is a topdown parser using left-derivation sequence.

$A \rightarrow A\alpha$ (immediate left recursion)

$A^* \rightarrow Ad$ (left recursion)

① To apply LL parsing technique the grammar should not be left recursive.

- (a) There is one more issue with top down parsing.
- If $A \rightarrow \alpha_1 \beta_1 / \alpha_1 \beta_2$
- i.e. two productions have same prefix α_1 , we won't be able to say which production is used. No two productions of a variable should have same prefix [left factoring].

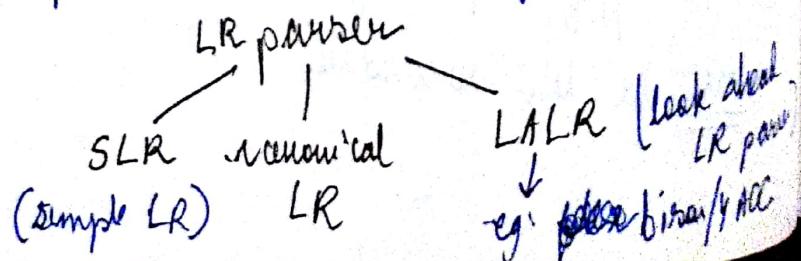
To perform LL parsing - remove left recursion and perform left factoring.

LR parsing → right most derivation sequence in reverse order.

i/p read from left
to right.

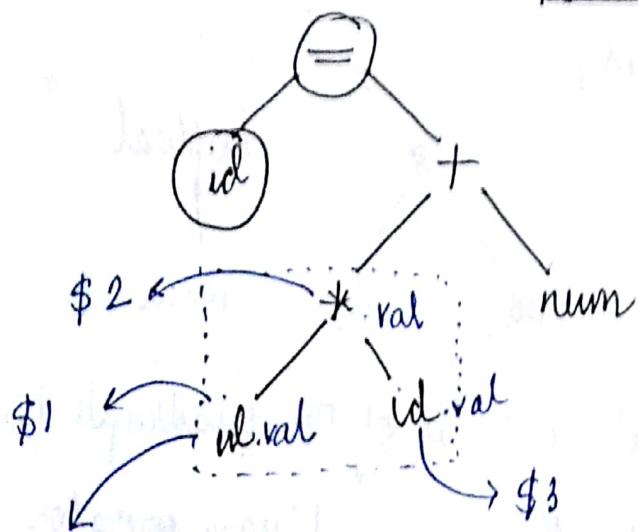
[first we get right most sentential form and finally we get S].

A lookup table is created using LR parsing!



Illustration

$$id = id \neq id + num.$$



suppose production rule used is $E \rightarrow E * E$.

[then directly generates intermediate code from

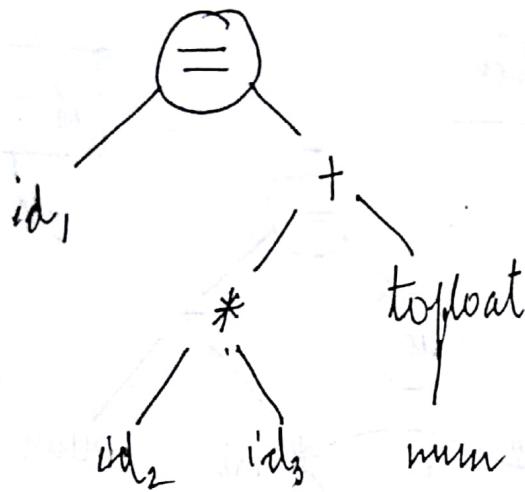
symbol analysis]

$$\begin{aligned} \$\$ &= \$1 * \$3 && \text{value associated with 1st token} \\ &&& \text{on right side} \\ &&& \text{value of attribute on left side} \longrightarrow \text{intermediate} \\ &&& \text{code} \\ E.\text{val} &= E_1.\text{val} * E_2.\text{val} \end{aligned}$$

(syntax directed translation)

[goto instructions are essential for intermediate code generation].

#



Intermediate code is of the quadruple form.

$ln - a = b$ → binary operator
 $a = b \oplus c$ → binary operator

'?:' is a ternary operator which is right associative.

Ternary operators are converted into binary operators during ICG. → temporary identifier names.

$$t_1 = id_2 \cdot val * id_3 \cdot val$$

$$t_2 = t1float (num, levval)$$

$$t_3 = t_1 + t_2$$

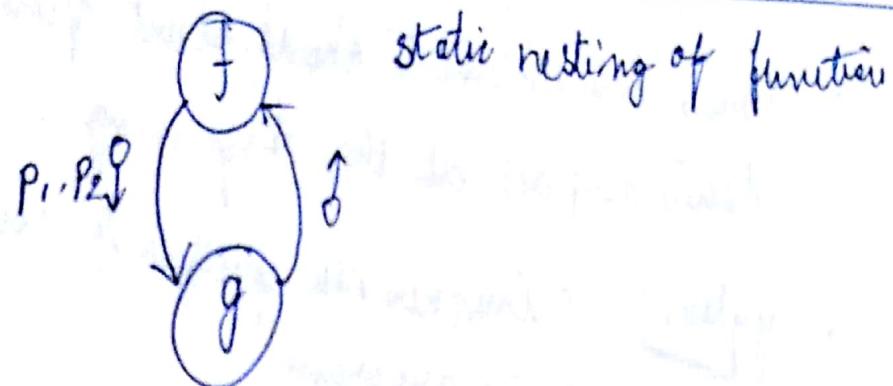
typical,

$$id_1 \oplus t_3 = t_3$$

Most of the temporary identifier names generated by system begin with '_' (underscore).

We have to avoid memory references and number of instructions during code optimisation.

Eg- $id_1 = t_1 + t_2$.



Lex specification - (structure).

%{ Declaration section %}

%%

Rules (Regular Ex. Action) → functions can be called here which are written in auxiliary function section

%%

Auxiliary functions

%{ C-definitions (copied verbatim to lex.yy.c by lexer)

% regular definitions

% symbolic naming for a class of character

Example of regular definition -

digit [0-9]

digits digit+

- Every lex statement should start from column 1. Do not leave a space at the beginning.
- `yylex()` → invokes `lex` compiler to convert the rules in C-program.
- `lex` and you do not do lexical analysis.
- `yylex` → generates a program in C (`lex.yyc.c`) which is the lexer and that C program creates tokens.
- `you` → generates ~~lexer~~ parser program.
- `lex` matches longest input string.
Suppose input string is "abc" and rule is `[a-z]` then only "a" is matched.
 $[a-z]^+$ → also called a pattern.
"abc" → lexeme (the actual text which is matched against a rule)
`yyltext()` returns the lexeme

$[a-z]^+$ "ab"
 $[a-z]^*$

every rule returns a token

$[a-z]^+$ $\xrightarrow{\text{action}}$ returns tokens

return (10) \curvearrowright manifest constants

\downarrow
value is returned to yylex(). The system again and again calls yylex().

Only the first applicable rule is considered.

'.' dot is a metacharacter \rightarrow matches every character other than newline characters.

space $[\t \n]^*$

alpha $[a-zA-Z]$

ws $[\text{space}]^+$

$[]$ \rightarrow set of characters
 $[^ab]$ \rightarrow all characters except "a" or "b".
 $[A]$ \rightarrow checks for "A".
blank

lower $[a-z]$

upper $[A-Z]$

alpha $\left[\cancel{\text{lower | upper}} \right] \{ \text{lower} \} \{ \text{upper} \}$

∴ $\{ \text{alpha} \}^+ \{ \text{action}, ; \text{action}, ; \cdot \}$

$\{ \text{alpha} \}^+ \{ \text{action}, ; \text{action}, ; \cdot \}$

{alpha} + {word ++; ch + = yytext ;
lower + = int-lower(yytext);
}

%
int int-lower(char * s) {
=====
}

#	<u>Assignment</u> - ① Postfix evaluation ② Postfix to infix	AWK → Unix tool
---	---	-----------------

compiler for language lisp - (lots of irrelevant
and silly programming) → functional programming
fortran → formula translation

* There are only two data types in lisp →
① atom ② list

lisp is not an imperative language -

In an imperative language we have Assignment
operation and they follow Content Free programs

- In lisp we give input to a function and the output to this function is given to another function as input. Lisp was parenthesis notation for denoting expressions. Parenthesis is unambiguous with respect to precedence of operators.
- Functional languages do not have a compiler. Lisp has a translator interpreter
 - ↓
translates line by line.
- The lisp was developed using 5 basic functions and these functions are used to make other functions. Lisp is inbuilt in Linux -
- C-language is a low level language in a way that it has less level of abstraction. Using C program we can invoke Linux commands.

`system(command);`

↓
control goes to OS. OS executes command and control goes back to program. C language is close to kernel -

Writing a compiler for a language -

find CFG for the language.

Nobody writes compiler from scratch. The tools existing compilers are used to build new compilers. Can a constraint for language Y be written using constraint of language X? If yes then tools of X are used.

Bootstrapping - the technique used is called

bootstrapping. It is used to reuse compiler X to get compiler Y, reuse compilers.

Source language → language in which implementation language → language in which compiler is written

Target language

$S_I T$

language → $L_M N$ → machine code

language of compiler

The compiler for language L written in language

M used for machine N .

$$L_s N + S_M M \rightarrow L_M N.$$

concept of cross compiler intermediary
cross compiler

$$L_L N + L_M M \rightarrow L_M N. \quad \} \text{Porting a}$$

$$L_L N + L_N N \rightarrow L_N N. \quad } \text{compiler.}$$

The program $L_L N$ is again given as input to intermediary cross compiler $L_M N$ and then we are getting $L_N N$.

$M \rightarrow$ dirty compiler

$N \rightarrow$ improved or optimised compiler

The process bootstrapping is "not reinventing the wheel".

Sem. Assignment - [operators in 2nd line have higher precedence than 1st].

Data types: int, unsigned int, boolean

Assignment: $=, +=, -=, *=, /=$

Binary operators: $, +, -, *, /, ^$ (exponentiation)
left association right association

Bitwise: $, \&, \sim, \wedge$

logical: ||, &&

Relational op: ==, !=, <, <=, >, >=

Assignment stat: $S \rightarrow id \ AOP \ expression$

Iterative

conditional: if, if else, switch case

Repetitive: while

id : simple ids (no special characters).

Expressions in infix notation :-

Input: $x = a + b * c ;$

Output: $t = b * c ; \quad \left. \begin{array}{l} \text{intermediate} \\ x = a + t ; \end{array} \right\} \text{code}$

If you have,

$E \rightarrow E + T$

write semantic action associated with it in
face. Semantic action: $E = E_{\text{val}} + T_{\text{exp}}$

Traverse parse tree in bottom up manner
using face.

generally semantic action will be

$$\$\$ = \$1 + \$3.$$

printf ("\\$\\$" = " \\$1 " "+ " "\\$3");

So next when you compile it as C program - you can generate $\$\$ = \$1 + \$3$.

∴ You have to break an expression in simple binary expression → 1st part.

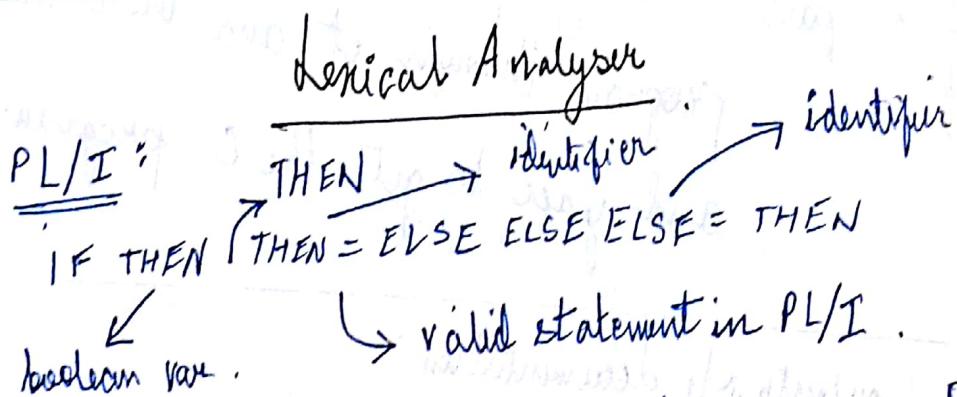
2nd part → using specification given write a program through it and then use lex

and you'll get the C program -

Contents of documentation :-

- Problem description
- Decomposition
- Approach
- Algorithm proposed
- Description of functional units
- Relationship among functions - I/O.

- Static nesting tree
- Implementation details of main / functional units
- use of each for thorough comments for all non-trivial statements ~~trivial~~.
- sample i/o for testing
- test results
- major limitations / assumptions



Keywords in the programming language are not reserved. [THEN → keyword as well as identifier]

- To make keywords reserved -

store them in symbol table and mark it as keyword.

PL/I language was a failure.

The lexical analyser works as a subordinate of you due to simplicity of design.
If lexical analyser works as an independent entity then error detection is easy.

For statement,

$f_i(a == n)$.

f_i is not a keyword, so it's interpreted as an identifier. On seeing opening and closing parenthesis its interpreted as a function.

Error Recovery - Try to find as many errors as possible.

① $\text{int } i, n \rightarrow \text{error}$
② $\boxed{\text{float } a; } \rightarrow \text{ignored}$
 $a = i \rightarrow \text{error}$

One of the methods of error recovery is panic mode error recovery.

For statement ① compiler expecting , or ; so it ignores all subsequent characters until it gets , or ;

This statement ① is ignored in parser code.

≠ Designing a lexical analyzer - (3 days)

- ① Use a tech for → lexical analyzer part
- ② Develop lexical analyzer using a system level program or any low level language like C which allows reading or writing from a buffer character by character.

get() → buffered input

unget() → push back into buffer.

getch() → not buffered input. One not

want be pushed back into buffer.

- ③ Use assembly language.

Input buffer pairs - use two buffers. Fill one

buffer and process it. While processing

fill ^{other} input buffer. There are two pointers for

identifying keyword. One stays at beginning

and other one moves forward. Example when
lexer decides "for" it may be a possible
keyword.

optional $[^]$ set of chs other than specified
 $a? \rightarrow \backslash/a$

$\backslash \rightarrow$ escape character

$[]$ \rightarrow set of characters

$\{ \}$ - range

Number representation in scientific notation.

$atof(1.23 E + 2)$ $atof('5678')$
 \downarrow
 $1.23 * 10^2$ $n = n * 10 + (\text{ch} - '0')$

digits $\rightarrow [0-9]^*$

opt-fraction \rightarrow digits | E

opt-exp $\rightarrow (E (+/-) digits) | E$

num \rightarrow digits opt-fraction opt-exp

111 calc. calc.

Q- Lex program to simulate hex calc. [Assignment]
0-9 a-f +, -, *, /

$| -5A | \rightarrow 5A$ } absolute
 $\text{def } 5A \rightarrow 5A$ } own choice

$\% \rightarrow$ both integers (final remainder)
(probably two integers)

logical: OR, AND, ^, XOR $45 \mid 76$

bitwise OR

Both input and output \rightarrow hexadecimal.

Q- English Text encryption [case sensitive]

12
— special words — different approach
— other words — normal approach
— \boxed{spt} — left rotate by 2 positions (cyclic)

— [print in reverse order]

— other words — right rotate by 3 positions (print rev order)

- Number → divide into two halves and swap.
[$1945 \rightarrow 4519$]
- other characters - retain
- Words like can't should also be accepted.

B- SubC program — Stream of tokens.

- B Word frequency [read text file]
[Print word frequency and line numbers]
[form a symbol table]

E Infix expression evaluation

Output — evaluation order.

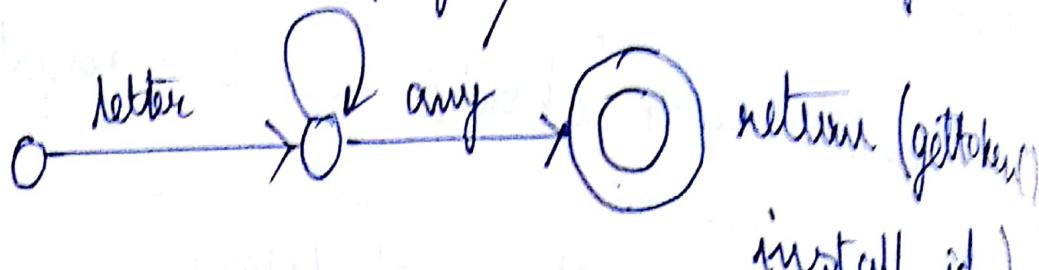
[~~the~~ precedence, associativity to be specified].

$$\begin{array}{c} \downarrow \\ 3 + 4 + 5 \\ = 7 + 5 = 12 \end{array}$$

Draw a transition diagram for recognizing
a number / identifier .

identifies, letter / digit

when you get any character other than letter or digit

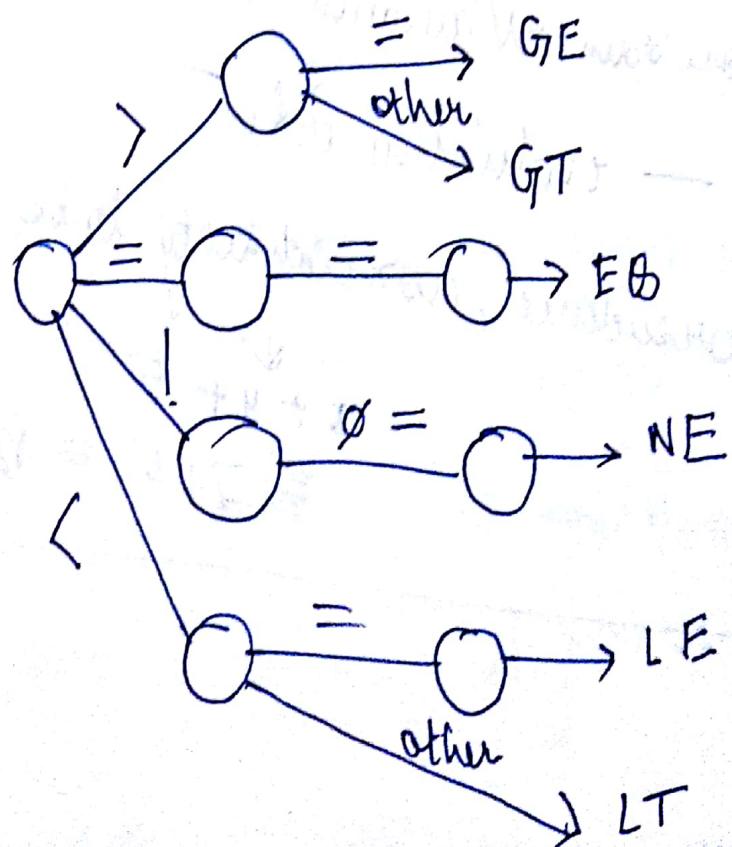


```
return (gettoken()  
install_id);
```

install in symbol table
along with attributes.

1

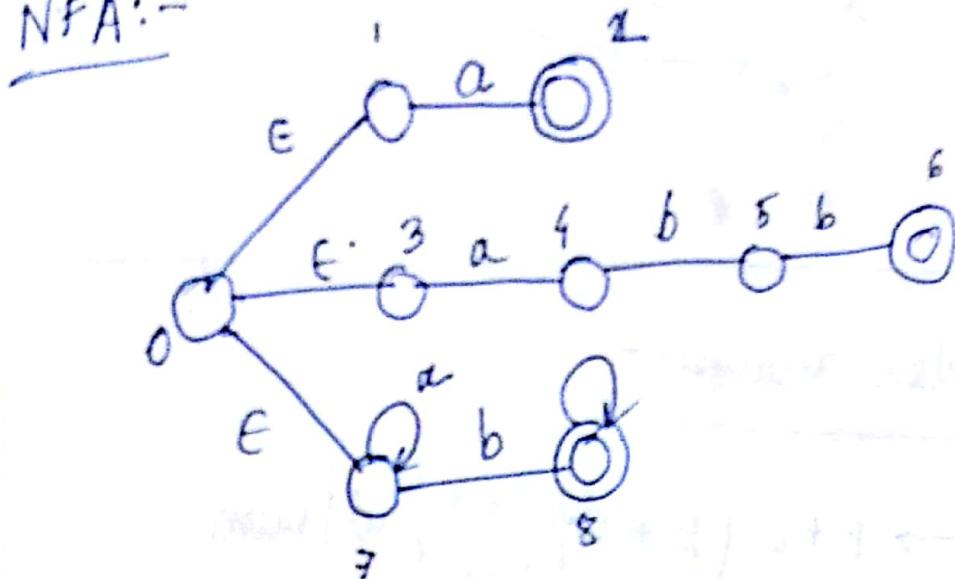
Transition diagram to recognise relational spectra



Consider three patterns,

- $a \{ \text{action}_1 \}$
- $abb \{ \text{action}_2 \}$
- $a^+b^+ \{ \text{action}_3 \}$

NFA:-



Converting into tabular form,

$[0137]$ (start state)

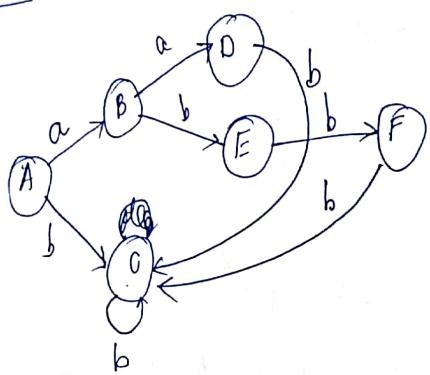
↓ a

↓ b

$[58]$ not an
acceptable
word

	a	b	token
$[0137]_A$	$[247]_B$	$[8]_C$ none	
$[247]_B$	$[7]_D$	$[58]_E$ (1) 'a'	
$[8]_C$	\emptyset	$[8]_F$ (2) a^+b^+	
$[7]_D$	$[7]_E$	$[8]_G$ (3) a^+b^+	reject
$[58]_E$	\emptyset	$[8]_H$ (4) abb	
$[68]_F$	\emptyset		

DFA :-



Syntax analyzer -

$$E \rightarrow E+E \mid E*E \mid (E) \mid id \mid num.$$

We give precedence order of operators to remove ambiguity. The least precedent rule is written first.

$$\begin{aligned} E &\rightarrow E+T \mid T \mid E-T \\ T &\rightarrow T.*F \mid F \mid /F \\ F &\rightarrow (E) \mid id \mid num \end{aligned}$$

should be GEF

In yacc specification

- token + -
- token * /

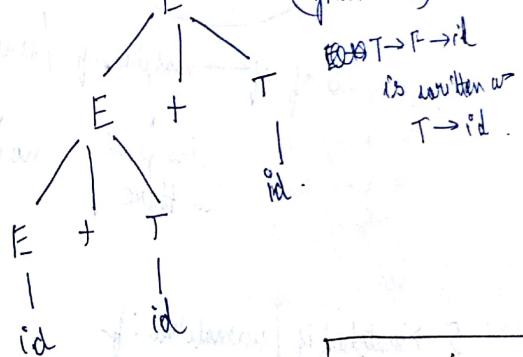
• / left + - * / → left associative
• right @ → right associative.

Derivation of $a+b+c$:-

$$\begin{aligned} E &\rightarrow E+T \rightarrow E+T+T \\ &\rightarrow T+T+T \rightarrow id+id+id. \end{aligned}$$

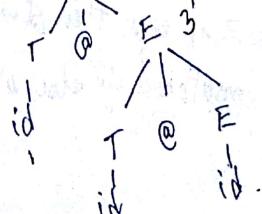
$a+b+c$

(Syntax tree)



$$\begin{aligned} E &\rightarrow E op T \quad (op is left) \\ E &\rightarrow T op E \quad (op is right) \end{aligned}$$

$$2 @ 3 @ 4$$



Only fortran language uses exponentiation as $a^{**} b = a^b$.

In C we write as function $\text{pow}(a,b)$.

$S \rightarrow iEtsS \mid (\text{if then else statement})$.

iETS.

if E_1 , then if E_2 then S_1 , else S_2 (ambiguous)

Two kinds of if → matched if | unmatched if
↓ ↓
else part no else part
is there

$S \rightarrow \text{matched if } \mid \text{unmatched if}$

matched if → if expr then statement else stmt.

unmatched if → if expr then stmt1 | if expr then
matched if else unmatched if

If error is not detected in lexical analysis phase then it will be caught in syntax analysis.

If caught in lex then rest is less.

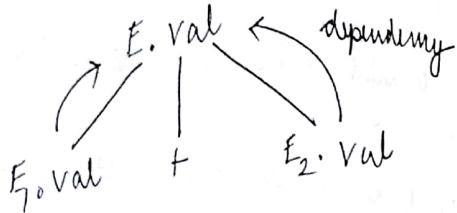
Top down parser → not left recursive but left factored.

Top Down Parsing

Checking membership
Pairing
Identify derivation sequence
Construct parse tree
Attributed parse tree
Dependency graph
Evaluation order

- In expression grammar there's a type associated with every non-terminal.
- Attributes associated with each non-terminal.
Not same for every non-terminal. Identifying dependency of parent-child node using attributes.

These attributes are synthesized attributes or inherited attributes.



The value associated with an attribute of a child is inherited from parent.

We need to traverse dependency graph to get evaluation order (intermediate code generation).

Traversal - Topological ordering

- For every grammar production we associate an action

$$E \rightarrow E + E \quad \{ E.\text{val} = E_1.\text{val} + E_2.\text{val} \}$$

[syntax directed translation].

Left Derivation Sequence (for top-down parsing)

$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow (E) \mid \text{id} \end{array} \left. \begin{array}{c} \\ \\ \end{array} \right\} \text{Left recursive.}$$

For a production like,

$$A \rightarrow A \alpha$$

leftmost symbol of RHS is same as LHS: then immediate left recursion.

$$A \xrightarrow{*} A \alpha \text{ (left recursive).}$$

I Elimination of (immediate) left recursion

$$A \rightarrow A \alpha \mid \beta.$$

- ① separate all productions of A having immediate left recursion.

$$\begin{aligned} A &\rightarrow A \alpha_1 \mid A \alpha_2 \dots \mid A \alpha_m \\ &\rightarrow \beta_1 \mid \beta_2 \dots \mid \beta_n \end{aligned}$$

We have to modify our grammar rules and show that strings in both grammar rules are same.

$$\bullet \frac{A \rightarrow A \alpha_1 \beta \rightarrow A \alpha_2 \alpha_1 \dots \rightarrow \beta_j \alpha_2 \alpha_1}{A \rightarrow \beta_j B} \quad \frac{A \rightarrow \beta_j B \rightarrow \beta_j \alpha_2 \alpha_1 \dots}{B \rightarrow \alpha_j B \mid E} \quad \frac{A \rightarrow \beta_j B \rightarrow \beta_j \alpha_2 \alpha_1 \dots}{B \rightarrow \beta_j \alpha_2 \alpha_1 B \rightarrow \beta_j \alpha_2 \alpha_1}$$

Example - removing immediate left recursion
expression grammar.

$$\begin{aligned} E &\rightarrow T F' \\ F' &\rightarrow +TE \mid \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid E \\ F &\rightarrow (E) \mid id \end{aligned}$$

II. Left Factoring - (merging productions having
same prefix)

$$A \rightarrow \alpha P_1 \mid \alpha P_2 \mid \dots \mid \alpha P_k$$

some prefix

(helps in checking unmatched if)

$$A \rightarrow d A'$$

$$A \rightarrow \beta_1 \beta_2 \dots \beta_k$$

- Left recursion removed and left factored.

Removal of left recursion and left factoring given
a CFG. [you may use it].

Top Down Parsing

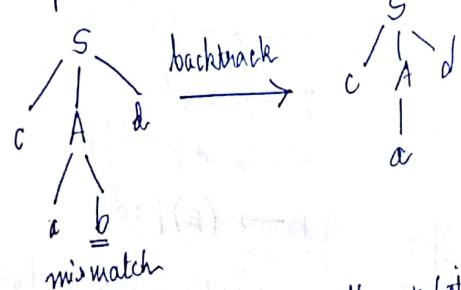
[Recursive descent parsing]

— requires backtracking [inefficient algorithm]

$$S \rightarrow c A d$$

$$A \rightarrow ab/a$$

consider input : cad



Whenever you get mismatch, check for other production
of A. Keep two pointers increase first pointer as
long as you are getting match.

- Predictive parser - (a top down parser which does not require backtracking but might use recursion).
- Non-recursive predictive parser :- (LL(1) parser)
 - If predictive parsing table does not have multiple entries then LL(1) parsing mechanism a tree (leftmost derivation).

$$E \rightarrow TE'$$

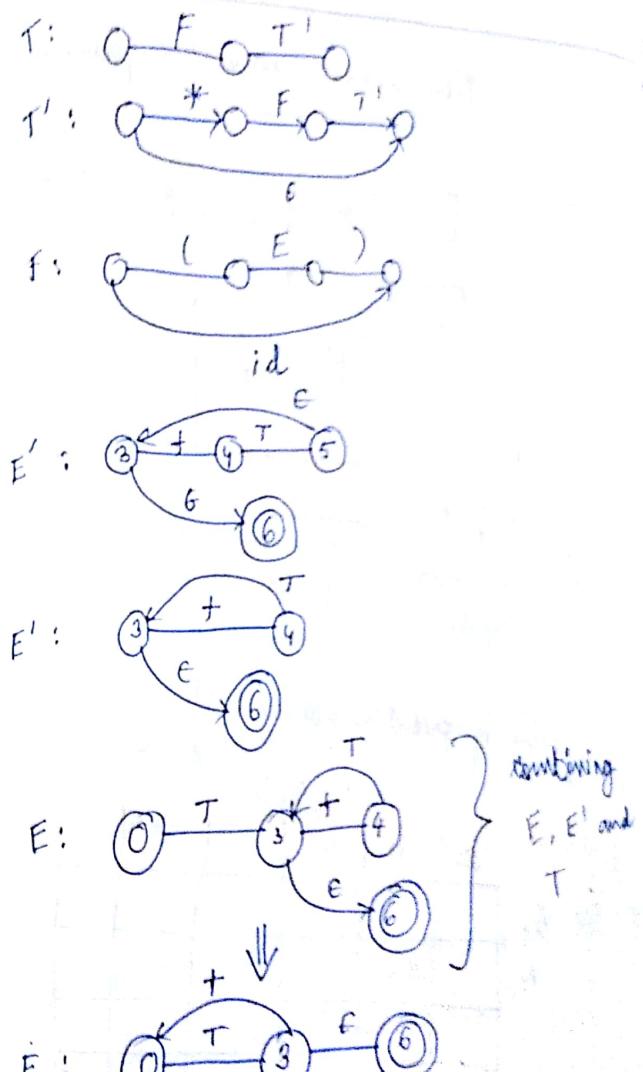
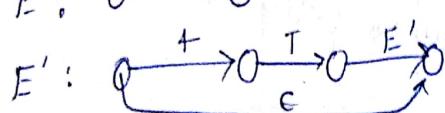
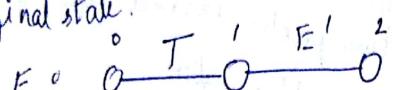
$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

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

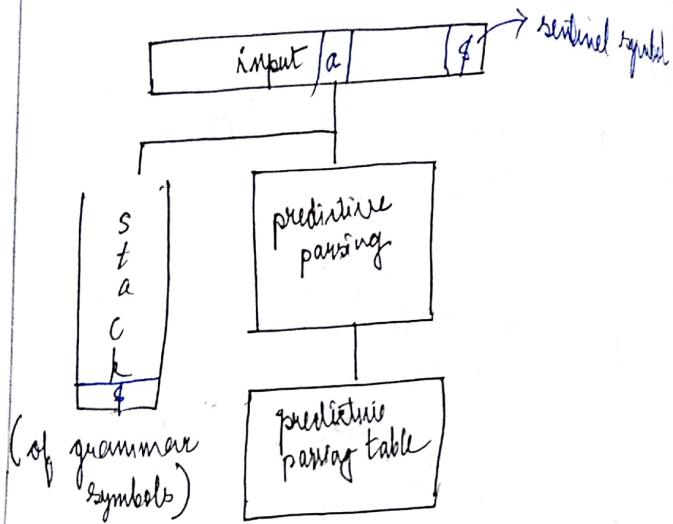
$$F \rightarrow (E) \mid id$$

For every variable we have a starting state and a final state.



A predictive parser can comfortably parse any binary expression.

Non recursive predictive parser



Entries of predictive parsing table

$$M[A, a] = A \rightarrow \alpha$$

a	b	c	\$
A			
$A \rightarrow \alpha$			

check current input symbol and current top of stack.

If x is a terminal, and matches $s.top()$,

$s = a$, $.pop(a)$, advance i/p.

If x is a non-terminal, A

look for entry in cell $[A, a]$

If $M[A, a] : A \rightarrow \alpha$

pop A, push α such that leftmost symbol is at top.

A	$A \rightarrow \alpha$		$A \rightarrow \alpha_2$	
---	------------------------	--	--------------------------	--

If entry corresponding $M[A, a]$ is empty then error condition. That particular symbol is not expected at that point.

Algorithm ends when entire input is exhausted i.e we are seeing a sentinel symbol '\$'.

\$ is also pushed as bottom of stack.

If input symbol = \$ and $s.top = '$'$ then

stop -

- not left recursive - operators left associative

Parsing Table M :-

E	$E \rightarrow TE'$		$E \rightarrow TE'$		
E'	$E' \rightarrow HE'$			$E' \rightarrow E$	$E' \rightarrow E$
T	$T \rightarrow FT'$		$T \rightarrow FT'$		
T'	$T' \rightarrow E$	$T' \rightarrow FT'$		$T' \rightarrow E$	$T' \rightarrow E$
F	$F \rightarrow E$		$F \rightarrow (E)$		

Example:- \downarrow id + id * id \$

$$\$E \Rightarrow \$E'T \Rightarrow \$E'FTF \Rightarrow \$E'Tit$$

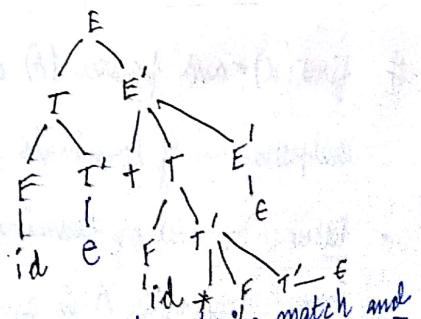
↓
match

$\$E' \leftarrow \$E'E \leftarrow \$E'T'$
 \downarrow
 $\$E'T^+ \xrightarrow{\text{match (pop)}} \$E'T \rightarrow$

Parsing action :-

production to be applied

<u>Stack</u>	<u>Input</u>	<u>Action</u>
\$E	t + id + id \$	E → TE'
\$E'T	id + id + id \$	T → FT'
\$F'T'F	id + id + id \$	F → id
\$E'T'(id)	(id + id + id \$	match and pop.
\$E'T'	+ id + id \$	T' → E
\$E'	+ id + id \$	E' → +TE'
\$E' T(+)	(id + id \$	match and pp.
\$E'@T	id + id \$	T → FT'
\$E' T'F	id * id \$	F → id
\$E'T'(id)	(id + id \$	match and pop.
\$E'T)	+ id \$	T' → *FT'
\$E'T'F(*)	(id \$	match and pop.
\$E'T'F	id \$	f → id
	Parse tree	



[No change in parse tree during match and
pop]

- In CYK algorithm we cannot find the derivation sequence but in top down parser we get parse tree. (eg - exponentiation)
- On including another operator there will not be any change in method.
- C-language does not have exponentiation due to left right associativity. Also exponentiation is repeated multiplication by shift operation. User has to write his/her own exponentiation function.

Construction of parsing table:

- # $\text{first}(x)$ and $\text{follow}(A)$ are two functions outputting set of terminals.
- $\text{follow}(A) \rightarrow$ set of terminals which can appear immediately after A in any sentential form.

In follow(A) we will not include ϵ . (VUT)

$\text{first}(x)$: consider any string derived from x .
 $x \xrightarrow{*} a\beta$. so $a \in \text{first}(x)$.

$\text{first}(x)$ is the set of strings which appears the leftmost symbol in any string derived from x .

[x can be null as well]

Consider $A \rightarrow Y_1 Y_2 \dots Y_n$

so $\text{first}(A) = \text{first}(Y_1)$ if $\text{first}(Y_i)$ does not contain ϵ .

If $\text{first}(Y_1)$ contains ϵ ,

$A \rightarrow Y_2 Y_3 \dots Y_n$

so $\text{first}(A) = \text{first}(Y_2)$.

$\text{first}(x)$ can be ϵ as well unlike $\text{follow}(A)$ which does not contain ϵ .

$\text{first}(x)$:

If x is a terminal, it is in $\text{first}(x)$.

If $X \rightarrow C$ is in P , then add C to $\text{first}(x)$.

If $X \rightarrow Y_1 Y_2 \dots Y_n$ and X is non-terminal,
place a in $\text{first}(x)$, if a is in $\text{first}(Y_i)$
and a is in $\text{first}(Y_1 Y_2 \dots Y_{i-1})$ each.

If C is in $\text{first}(Y_i)$ for all $i = 1, 2, \dots, k$ then
add C to $\text{first}(x)$.

$\$$ is in $\text{follow}(S)$.

If $A \rightarrow \alpha B \beta$ is in P , the symbols in $\text{first}(\beta)$
except ϵ are in $\text{follow}(B)$.

If $A \rightarrow \alpha B$ or $A \rightarrow \alpha B \beta$ with $\beta \not\rightarrow \epsilon$ then
every symbol in $\text{follow}(A)$ is in $\text{follow}(\beta)$.

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow + \quad TE'/\epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow * \quad FT'/\epsilon \\ F &\rightarrow (\epsilon) \mid id \end{aligned}$$

$$\text{first}(E) = \{(, id\} \quad \text{follow}(E) = \{, \$\}$$

$$\text{first}(T) = \{(, id\} \quad \text{follow}(T) = \{+,), \$\}$$

$$\text{first}(F) = \{(, id\} \quad \text{follow}(F) = \{+,), \$\}$$

$$\text{first}(E') = \{+, E'\} \quad \text{follow}(E') = \{, \$\}$$

$$\text{first}(T') = \{*, E'\} \quad \text{follow}(T') = \{+, +, \$\}$$

For a production, $A \rightarrow \alpha B \beta$, $\text{first}(\beta)$ is in $\text{follow}(A)$

$A \rightarrow \alpha B$ or $A \rightarrow \alpha B \beta$, ϵ is in $\text{first}(\beta)$

then everything in $\text{follow}(A)$ is in $\text{follow}(\beta)$.

For $E \rightarrow TE'$, $E' \not\rightarrow \epsilon$, so everything in $\text{follow}(E')$
will be in $\text{follow}(T')$.

construction of parsing-table :-

For each production $A \rightarrow \alpha$ in G_1 ,

- for each terminal a in $\text{first}(\alpha)$
add $A \rightarrow \alpha$ to $m[A, a]$
- if ϵ is in $\text{first}(\alpha)$
add $A \rightarrow \epsilon$ to $M[A, b]$ where b is in $\text{follow}(A)$
for each b .
- if ϵ is in $\text{first}(\alpha)$ add $A \rightarrow \alpha$ to
 $M[A, \$]$ if $\$$ is in $\text{follow}(A)$.

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

All the blank entries correspond to certain error conditions.

construct predictive parsing table and parse a given input. [lab assignment] \rightarrow C++.

(given in book for +, *) $[+ - * / ^]$

No unary minus.

Consider integral division.

Use exponentiation as well.

In input, each operator should be used twice.
[remove left {recursion and left {rule}].

$S \rightarrow ; E \{ S \mid ; E \} S \quad E \rightarrow b$.

~~$\text{first}(S) = \{ ; \}$~~

$\text{follow}(S) = \{ \$, e \}$

~~$\text{first}(E) = \{ \}$~~

$\text{follow}(E) = \{ +, ^ \}$

~~$\text{first}(D) = \{ \}$~~

	i	b	e	b	\$
S	$S \rightarrow ; E \{ S \mid ; E \} S$				
S	$S \rightarrow ; E \{ S \mid ; E \} S$				
E					

$S \rightarrow iE + tSeS \mid iEt \not\in a \quad E \rightarrow b.$

$S \rightarrow iE + Sx \mid a$

$X \rightarrow eS \mid E$

$E \rightarrow b.$

$\text{first}(S) = \{i, a\}$

$\text{first}(X) = \{e, E\}$

$\text{first}(E) = \{b\}$

$\text{follow}(S) = \{e, \$\}$

$\text{follow}(X) = \{e, \$\}$

$\text{follow}(E) = \{t\}$

	a	b	i	t	e	\$
S	$S \rightarrow a$		$S \rightarrow iE + Sx$			
X			$X \rightarrow eS$		$X \rightarrow E$	
E		$E \rightarrow b.$				

We get multiple entries due to ambiguity.

Or without multiple parsers in PPT is called LL(1).

On removing $X \rightarrow t$ from PPT, the grammar becomes LL(1).

Error recovery in Top/Down LL parsing -

Panic mode error correction. (Element at top of stack is not expected).

identify → through tables.

recover

proceed for finding further error

} parser.

→ in case of panic mode recovery ignore those symbols which were not expected.

Example: If T is in top and input symbol is not then ignore.

For every non-terminal symbol, find a set of terminals which can be present (synchronizing set).

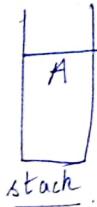
$M[A, a] = \text{error synch.}$

If $M[A, a] = E \rightarrow TE'$ (valid)

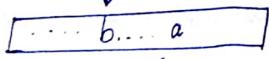
The symbol in $\text{follow}(A)$ will be definitely in synchronizing set of A.

Illustration:

$\text{first}(A) = \{a, \dots\}$



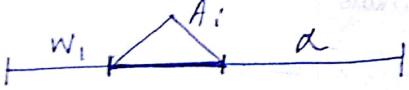
input pointer



ignored by parser as long
as a symbol belonging to $\text{first}(A)$
is not obtained.

Another way of recovery is to invert symbols in
input but that leads to many other errors and
confuses programmer. So many commercial
compilers use panic mode recovery.

Bottom Up Parsing



terminal string.

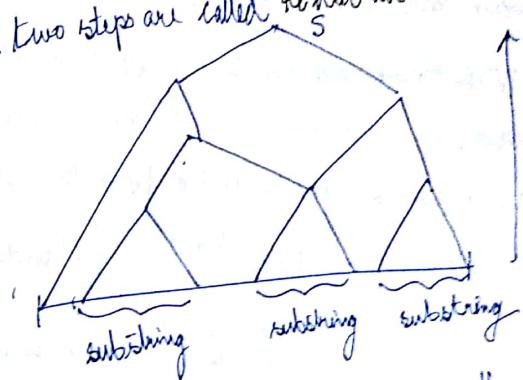
① Identify a substring of this terminal string which
is right side of a production of a grammar rule.

$a, a_2, \dots, a_i, \boxed{a_i}, \dots, a_n$
reduced to $A:$

$a, a_2, \dots, a_i, A_i, \dots, a_n$

② Replace that substring with a variable in LHS of
this production.

These two steps are called reduction.



- In any substring form we can find more than one substring which matches RHS of a production.
- The derivation sequence should be rightmost derivation sequence in reverse order.

Example:

$S \rightarrow aABe$
 $A \rightarrow Abc \mid b$
 $B \rightarrow d$
 $S \xrightarrow{a} aABe \xrightarrow{B} aAde \xrightarrow{A} aAbcd \xrightarrow{b} abbcde$

- In Bottom-up parsing we perform rightmost derivation sequence.
 - A handle is a substring in a given sentential form which matches RHS of a production rule.

When this homerule can be reduced to a RHS of a production rule such that the new sentential form is previous right sentential form then that

handles is called proper handle. [Handle pruning]

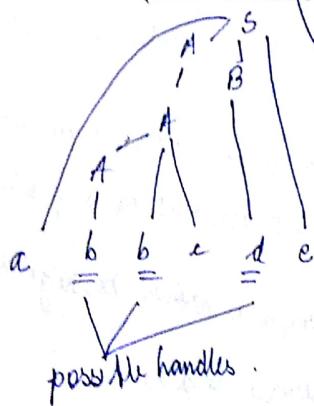


abb ad

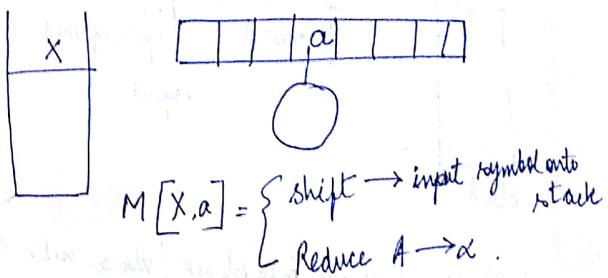
a ABCde

a ABe

S

Shift-Reduce Parsing - (bottom up parsing)

→ has a stack of grammar symbols.
(terminals & non-terminals).

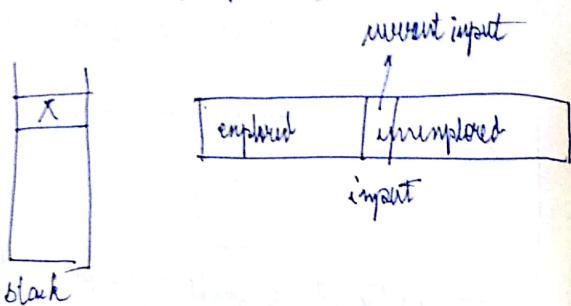


how many symbols of α has to be popped off?

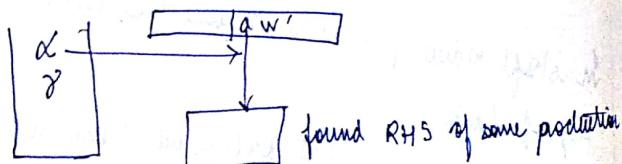
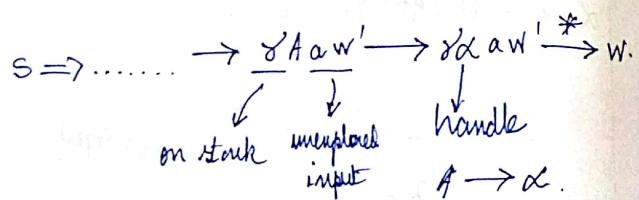
$$\frac{\beta}{\beta + w'} \frac{\alpha}{w'} \frac{w'}{(\text{unexplored input})}$$

- In shift reduce parsing we assume that initial α is on top of stack.
 - A viable prefix is that prefix which does not cross the stack.
 - If $\beta A w'$ is previous rightmost substring form then w' is a viable prefix.

Operator Precedence Parsing - (makes sure handle is on top of stack).



shift current input onto stack (along with few user-defined symbols).



If a reduce operation is performed on a right sentential form we must get previous right sentential form.

Viable prefix - prefix of a right sentential form.
Not every prefix of a right sentential form is viable.
A viable prefix is one that appears entirely on stack.

A prefix of a right sentential form that does not move past the right end of rightmost handle.

Identifying sequence of viable prefixes -

(or identifying right end of handle).

Only after finding left end we can pop off α .
There is one language for identifying left end of handle. (Limited language)

→ Operator Precedence grammar -
(parses expressions containing operators)

$E \rightarrow EAE \quad id \quad \} \text{ not OPG}$
 $A \rightarrow +|-|*|/| \uparrow \}$

Any grammar whose right-hand side does not have two or more consecutive non-terminals.

$$E \rightarrow E+E/E-E/E \cdot E/E/E/F \uparrow E/\text{id}/(E)$$

Operator precedence relationships

$$\theta_1 > \theta_2$$

$$\theta_1 < \theta_2$$

$$\theta_1 = \theta_2$$

If $\theta_1 > \theta_2$ and $\theta_2 > \theta_3$ it may happen that θ_1 and θ_3 are not related at all. Transitivity, reflexivity does not apply to operator precedence.

If $\theta_1 > \theta_2$ it may happen that ~~both~~ $\theta_1 < \theta_2$.

Example: $a [b + c. \quad a f b - c.$
 ↓ ↓
 more precedence higher precedence.

$\theta_1 > \theta_2$: θ_1 takes precedence over θ_2 .

$\theta_1 < \theta_2$: θ_1 yields precedence over θ_2

$\theta_1 = \theta_2$: θ_1 and θ_2 have same precedence.

If θ_1 has higher precedence over θ_2 make
 $\theta_1 > \theta_2$ and $\theta_2 < \theta_1$.

If θ_1 and θ_2 have equal precedence,

if they are left associative
 $\theta_1 > \theta_2$ and $\theta_2 > \theta_1$

if they are right associative
 $\theta_1 < \theta_2$ and $\theta_2 < \theta_1$.

[between bi-ary operators]

From operator precedence relationships we make operator precedence functions f and g using operator precedence functions f and g using which we identify precedence relationships.

[Pg-219, 220]
 List

⑦.

	id	+	*	\$
id	→	→	→	
+	<	→	<	→
*	<	→	→	→
\$	<	<	<	

Operator precedence relationship.

θ < id	(. =)
id → θ	(< (
θ < ((< id
(< θ	\$ < C
) > θ	id → \$
θ > θ)	id >)
θ > \$	\$ < id
\$ < θ) > \$
) >)

⑧.

Create operator precedence table from operator precedence relationships.

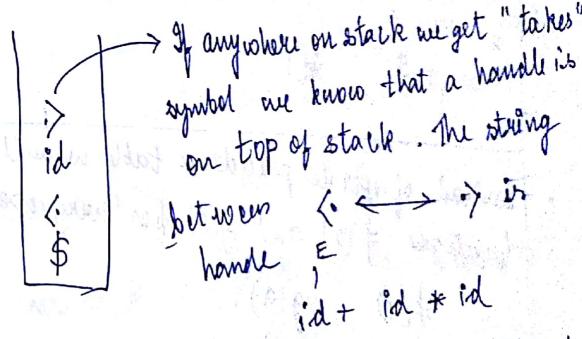
Relationship
between every pair of non-terminal symbols give an operator symbol.

Suppose string is "id + id * id \$" to indicate beginning of input.

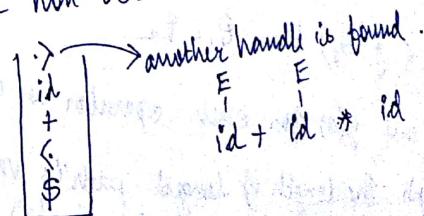
\$ < id → + < id → * < id → \$ "

Note: Nowhere in the parsing we ~~without~~ have used non-terminal symbol E till now.

Push all symbols onto stack including \$.



Between every pair of consecutive operators there is a non-terminal. Non-terminals are just placeholder



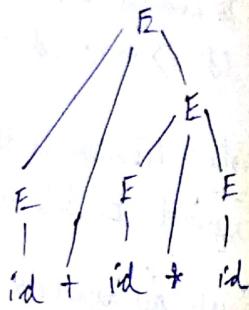
>
 id
 <
 *
 <
 +
 <
 ;
 \$
 handle

E E E
 → id + id + id

id + id + id → E + E + E

The stack contents are,

>
 *
 <
 +
 <
 ;
 handle
 id



- Instead of operator precedence table we will have functions $f(a)$ and $g(a)$ for every operator a .

$f(a)$ $g(a)$
 ↓ ↓

used when operator used when operator a is on left side right side.

If $f(\theta_1) \leq g(\theta_2)$ then $\theta_1 \leq \theta_2$

- Using $f(a)$ and $g(a)$ for each operator we draw a directed graph. The length of longest path is value of $f(a)$ and $g(a)$.

If there are n operators, we have $2n$ nodes. Some of these nodes will be grouped into one.

if $a \cdot = b$, group f_a and g_b .

if $c \cdot = b$, group f_c and g_b .

In same group we have here f_a, f_c and g_b .

Even if f_a and f_c are in same group then they might have no relationship between them.

$a \cdot = b$
 $\cdot c \cdot = b$
 $c \cdot = d$

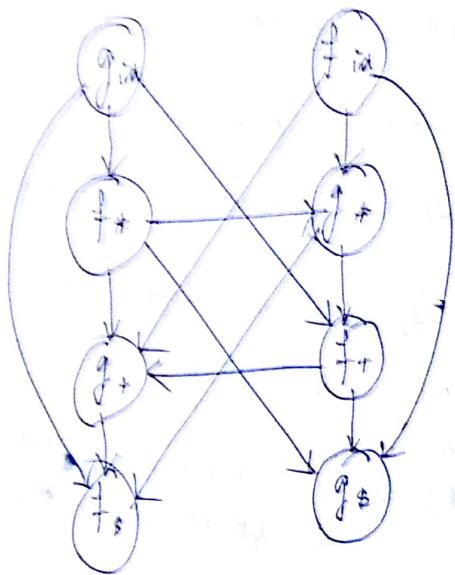
} f_a and g_d are in same group
 but have no relationship.

Number of nodes becomes $< 2n$.

If there is a relationship $a \leq b$ (a yields to b), add an edge from g_b to f_a .

If $a \geq b$, add edge from f_a to g_b .

graph :-



- Define $f_a = \text{length of longest path from group containing } f_a \text{ to any leaf node}$

	id	+	*	\$
$f(a)$	4	2	4	0
$g(b)$	5	1	3	0

g_a : ... path from group containing g_a to my leaf node.

LR Parsing :- (bottom up parsing)

left to right

LR parsers are also table driven

shift-reduce is also done here

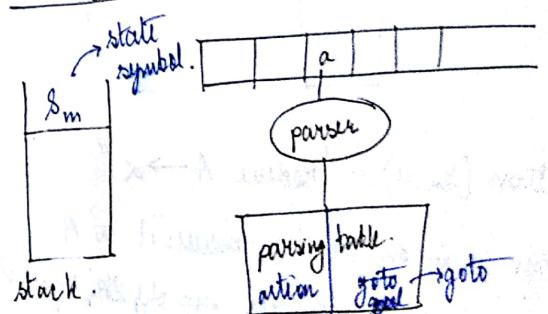
there are 3 different approaches for constructing LR parsing table

i) SLR

ii) Canonical LR

iii) LALR (your or bison is similar to LALR)
use LALR table.

Parsing action -

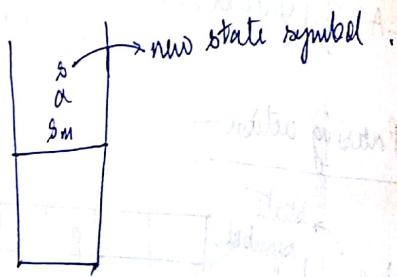


The stack contains grammar symbol and some state symbols. The t.o.b always contain state symbol. State symbol is representation of the configuration or the contents of stack.

The parsing table contains two parts,
action and goal.

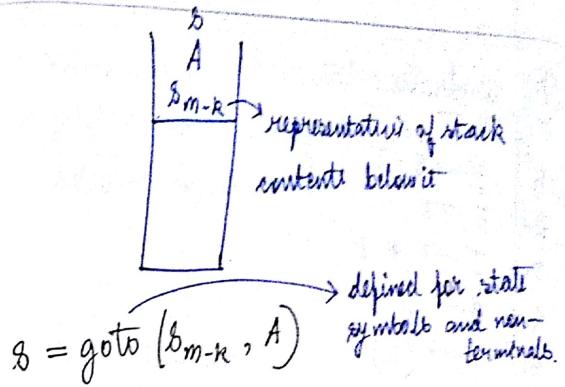
- $\text{action}(s_m, a) = \text{"shift } s\text{"}$ → state symbol
 ↓ input
 shift current symbol
 on top of stack
 [a is pushed on stack]

When a is pushed, the stack configuration changes.
 So, s is new state configuration which is pushed on stack.



- $\text{action}(s_m, a) = \text{"Reduce } A \rightarrow \alpha"$

When α is on top, then reduce it to A.
 Suppose α has k symbols, pop off 2k symbols from stack and push A, [k symbols of α and k state symbols] and push state symbol s.



- $\text{action}(s_m, a) = \text{"accept"}$

If none of the action can be done i.e. we have blank then,

- $\text{action}(s_m, a) = \text{"error"}$

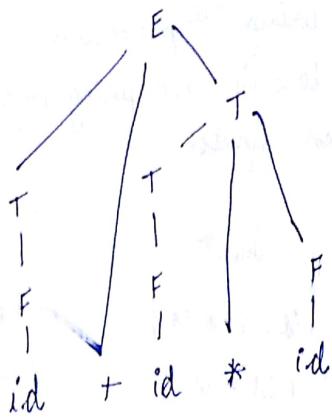
④ Illustration :-

1. $E \rightarrow E + T$
2. $E \rightarrow T$
3. $T \rightarrow T * F$
4. $T \rightarrow F$
5. $F \rightarrow (E)$
6. $F \rightarrow id$

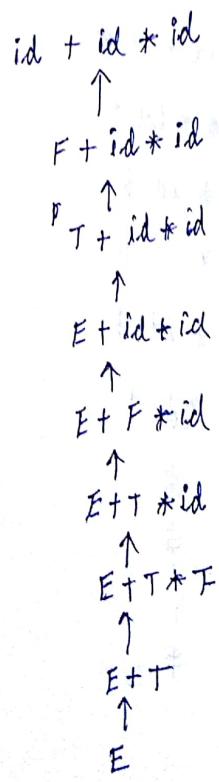
option							get
state	id	+	*	()	\$	E T
0	s_5			s_4			1 2
1		s_6				acc	
2	s_2	s_7		s_2	s_2		
3	s_4	s_4					
4	s_5		s_4			s	2
5	s_1	s_2	s_2	s_2	s_2		
6	s_5		s_4				
7	s_5		s_4				
8	s_6			s_{11}			
9	s_1	s_7		s_3	s_1		
10	s_3	s_3		s_3	s_3		
11	s_5	s_5		s_5	s_5		

$s_i \rightarrow$ reduce using i^{th} production
 $s_i \rightarrow$ shift to state i and pushing corresponding column symbol.

Stack	Input	Action
0	$id + id * id \$$	"shift 5"
$0 id 5$	$+ id * id \$$	"reduce $F \rightarrow id$ "
$0 F 3$	$+ id * id \$$	"reduce $T \rightarrow F$ "
$0 T 2$	$+ id * id \$$	"reduce $E \rightarrow T$ "
$0 E 1$	$+ id * id \$$	"shift 6"
$0 E 1 + 6$	$id * id \$$	"shift 5"
$0 E 1 + 6 id 5$	$* id \$$	"reduce 6"
$0 E 1 + 6 F 3$	$* id \$$	"reduce $T \rightarrow F$ "
$0 E 1 + 6 T 9$	$+ id \$$	"shift 7"
$0 E 1 + 6 T 9 + 7$	$id \$$	"shift 6"
$0 E 1 + 6 T 9 + 7 id 5$	$\$$	"reduce $F \rightarrow id$ "
$0 E 1 + 6 T 9 + 7 F 10$	$\$$	"reduce $T \rightarrow F$ "
$0 E 1 + 6 T 9$	$\$$	"reduce $E \rightarrow T$ "
$0 E 1$	$\$$	accepted



Rightmost derivation :-



Construction of SLR Parsing Table :-

LR(0) Item

Item is a production by placing a '·' before, before or after any symbol.

Itemset $\left\{ \begin{array}{l} A \rightarrow \cdot X Y Z \\ A \rightarrow X \cdot Y Z \\ A \rightarrow X Y \cdot Z \\ A \rightarrow X Y Z \end{array} \right\}$ a production with n symbols gives $(n+1)$ items

* If stack is $s_0 s_1 s_2 s_3 \dots s_m$

① and $M[s_m, a] = \text{push } s$.

then stack becomes,

$s_0 s_1 s_2 s_3 \dots s_m s a$.

② and for $M[s_m, a] = \text{pop}(A \rightarrow \alpha)$ and $\alpha = \underbrace{x_{k+1} \dots x_m}_{(m-k) \text{ symbols}}$.

stack becomes,

$s_0 s_1 s_2 s_3 \dots s_n A$

where $\text{goto}[s_k, A] = s$.

④

LR(0) Item :-

It's a production in G₁ with . placed in any position on the right side.

$A \rightarrow XYZ$ yields

$A \rightarrow X.YZ$

$A \rightarrow X.Y.Z$

$A \rightarrow XY.Z$

$A \rightarrow XY.Z.$

} LR(0) items depend on number of β symbols in production.

⑤

Sets of LR(0) items :- Each set of LR(0) item

represents a state in SLR table.

The common property of items in a set comes from position \neq in production

$A \rightarrow X.YZ$ means we have seen a string derivable from X and we are expecting a string derivable from YZ to reduce XYZ to A.

Therefore $A \rightarrow \alpha.\beta$ means

string derivable from α is seen and expecting to see a string derivable from β so that $\alpha\beta$ can be reduced to A.

Closure of a set of items :-

Suppose itemset I = { $A \rightarrow \alpha.\beta$ },

Closure(I) = {every element in I is in closure(I).
Suppose $A \rightarrow \alpha.B\beta$ is in I and $B \rightarrow \gamma$ is in P then $B \rightarrow \gamma$ is in closure(I).
 \Rightarrow we are expecting string deriving from $B\beta$.

- if a substring is reduced to start symbol then we can't say that parsing is complete.
That's why we add a production (augmenting) an extra production $E' \rightarrow E$ to grammar. Such a grammar is augmented grammar. When E is reduced to E' we say parsing is complete.

Closure($E' \rightarrow E$) = { $E' \rightarrow .E$,
 $E \rightarrow .E + T$
 (I_0) ↓ item
 $E \rightarrow .T$
 $T \rightarrow .T * F$
 $T \rightarrow .F$
 $F \rightarrow .id$
 $F \rightarrow .(E)$ }

goto operation -

$$g(I, x) \quad \begin{cases} I \rightarrow \text{item} \\ x \rightarrow \text{any grammar symbol.} \end{cases}$$

If I_0 contains $A \rightarrow \alpha \cdot X \beta$ then $g(I, x)$

contains closure ($A \rightarrow \alpha \cdot X \cdot \beta$).

We can define for any symbol after.

$$g(I_0, E) = \text{closure } \{ E \xrightarrow{\cdot} E \cdot, E \rightarrow E \cdot T \}$$

we do not have a non-terminal after. no more items can be added to g this set

$$g(I_0, C) = \{ F \rightarrow (\cdot, E), F \rightarrow \cdot (E) \\ E \rightarrow E \cdot T, F \rightarrow \cdot id \\ E \rightarrow \cdot T \\ T \rightarrow T \cdot F \\ T \rightarrow \cdot F \}$$

Canonical collection of LR(0) item set -

We start from I_0 (closure of augmented grammar)

We define all possible goto items I_0 .

$$I_1 = g(I_0, E) = \{ E \xrightarrow{\cdot} E \cdot \\ E \rightarrow E \cdot T \}$$

$$I_2 = g(I_0, T) = \{ E \rightarrow T \cdot \\ T \rightarrow T \cdot F \}$$

$$I_3 = g(I_0, F) = \{ T \rightarrow F \cdot \} \quad [No \text{ goto can be} \\ \text{defined on } I_3]$$

$$I_4 = g(I_0, C) = \{ F \rightarrow (\cdot, E) \\ E \rightarrow E \cdot T \\ E \rightarrow \cdot T \\ T \rightarrow T \cdot F \\ T \rightarrow \cdot F \\ F \rightarrow \cdot (E) \\ F \rightarrow \cdot id \}$$

$$I_5 = g(I_0, id) = F \rightarrow id.$$

$$I_6 = g(I_1, +) = \{ E \rightarrow E \cdot T \\ T \rightarrow \cdot T \cdot F \\ T \rightarrow \cdot F \\ F \rightarrow \cdot (E) \\ F \rightarrow \cdot id \}$$

$$I_7 = g(I_2, *) = \{ T \rightarrow T \cdot F \\ F \rightarrow \cdot (E) \\ F \rightarrow \cdot id \} .$$