

COMPILER DESIGN NESO ACADEMY

Introduction

0, 1 → Binary

Paper Tape, Punched card → hole indicates 0 and absence means 1.

Punched card

P - 1010000

u - 1110101

n - 1101110

c - 1100011

h - 1101000

e - 1100101

d - 1100100

C - 1000011

a - 1100001

g - 1110010

d - 1110011

ASCII

Language Translator

Language Translators:

(i) Assembler

MOV R1, 02H

MOV R2, 03H

ADD R1, R2

STORE X, R1

[Assembly language]

Assembler →

0110100101010

0010101010010

0100111100101

0101010101010

[Machine code]

(ii)

Interpreter:

one line at a time

↓

eg Php, Python, Ruby, JS

(iii)

Compiler:

all at a time

↓

⊙

(speed)

eg C, C++, Erlang

C Language : Middle Level Language

- ✓ Direct Memory Access through pointers,
- ✓ Bit Manipulation through Bitwise operator
- ✓ Writing Assembly code within "C" code

But initially: C - High Level Language

```
#include <stdio.h> // Header file for printf  
int main() // Main function
```

```
{  
    int x, a = 2, b = 3, c = 5;  
    x = a + b * c;  
    printf("The value of x is %d", x);  
    return 0;
```

```
}  
Source code / HLL code  
[High Level Language code]
```

Source code / HLL code $\xrightarrow[\text{Translator}]{\text{Language}}$ Machine code

Internal Architecture of Language Translator

- ① Preprocessor
- ② compiler
- ③ Assembler
- ④ Linker / Loader

① Preprocessor

- * removes preprocessing directives
- * removes all comments

O/P: stdio.h

```
int main()
```

```
{
```

```
}
```

↓
pure HLL

[comments are removed]

②

Compiler

Generates the equivalent assembly code for the given HLL code.

Pure HLL $\xrightarrow{\text{compiler}}$ Assembly Language

③

Assembler

Generates relocatable machine code for the given assembly language code.

Assembly Language $\xrightarrow{\text{Assembler}}$ Relocatable Machine Code

Explanation of "relocatable":

A program turns into a process only during execution. only during execution, it is allocated to the main memory. Before execution, we cannot know where exactly in the RAM the process will be allocated. so for the time-being, the sequence of machine codes is given a relocatable address like $i, i+1, i+2, \dots$

④

Linker / Loader

Generate the Absolute machine code that is actually executable.

$E+0 : 001010101001 \quad 0x000004B8 : 001010101001$

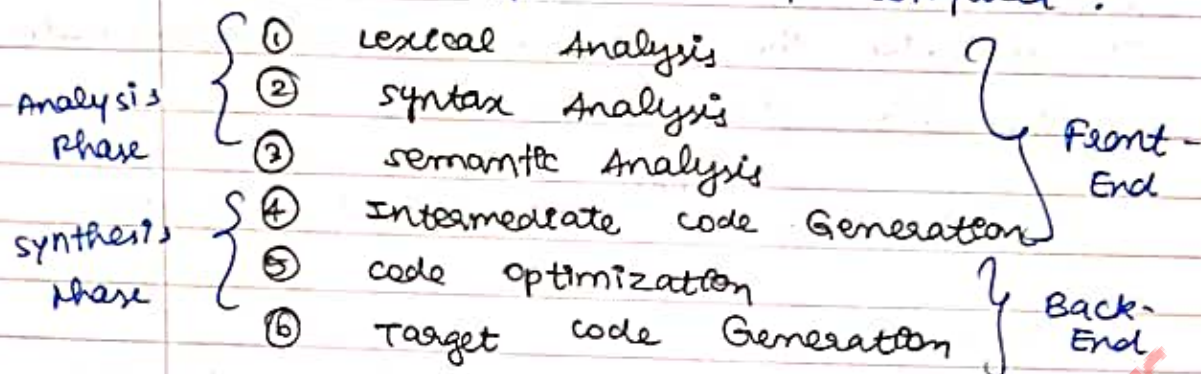
$E+1 : 0101101001100 \rightarrow 0x000004B9 : 0101101001100$

Relocatable
Machine code

Absolute
Machine code

process loaded into RAM
and ready for execution

Internal Architecture of compiler :



Backend : varies with OS.

eg if a C compiler is already designed for windows, and we want to build a new Mac compiler \rightarrow only backend needs to be modified!

other components :

- ① Symbol Table Manager
- ② Error Handler

- ① Information gathered in analysis are stored in the symbol table and used by the synthesis phase.
- ② The Error Handler deals with error detection and recovery.

Syllabus :

Pre-requisite :

- ① Introduction
- ② Syntax Analyser
- ③ Top Down Parsers
- ④ Bottom Up Parsers
- ⑤ Syntax Directed Translation Scheme
- ⑥ Intermediate code Generation
- ⑦ Runtime Environment & code optimization

Different phases of compiler

outcome:

- ① overview of various phases of compiler
- ② tools to implement different phases

$x = a + b * c;$

visualize how this statement is converted into assembly language through all ⑥ phases of the compiler.

①

lexical analyzer

$x = a + b * c;$

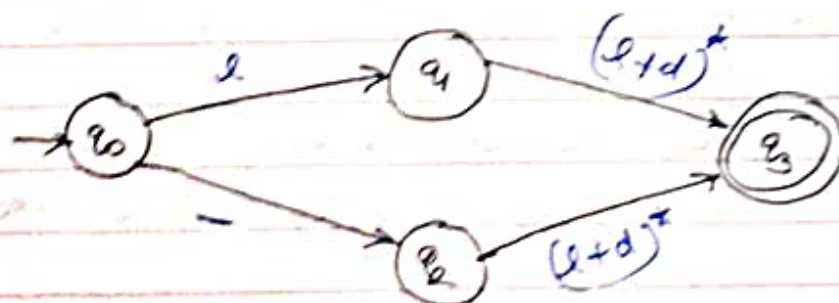
Lexical Analyzer

Tokens
↓
Meaning of the
tokens

Lexemes	tokens
x	identifier
=	operator
a	identifier
+	operator
b	identifier
*	operator
c	identifier

Lexemes: similar to words with only one small difference words have individual meanings. But group of lexemes convey meaning in their entirety

token recognition is done through regex.
Regex for identifier: $[l(l+d)^+ | (l+d)^+]$
 $l \rightarrow$ letter $d \rightarrow$ digit $- \rightarrow$ underscore



②

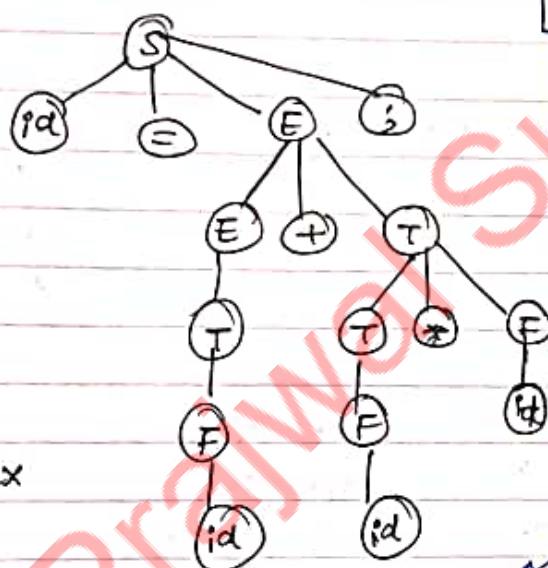
syntax Analyzer

depends on context-free grammar

$x = a + b * c;$

These are the rules, ^{production} that the parser will use to form the parse tree

$S \rightarrow id = E ;$
 $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow id$



$id, =, ;, +, *$
terminal

S, E, T, F
non-terminals

Traversal:
Top to Bottom,
Left to Right

Yield of parse tree:

$id = id + id * id ;$

(If not, there is some syntax error in the statement.)

since the yield of the parse tree and the expression are the same, the syntax analyzer will not produce any errors.

③

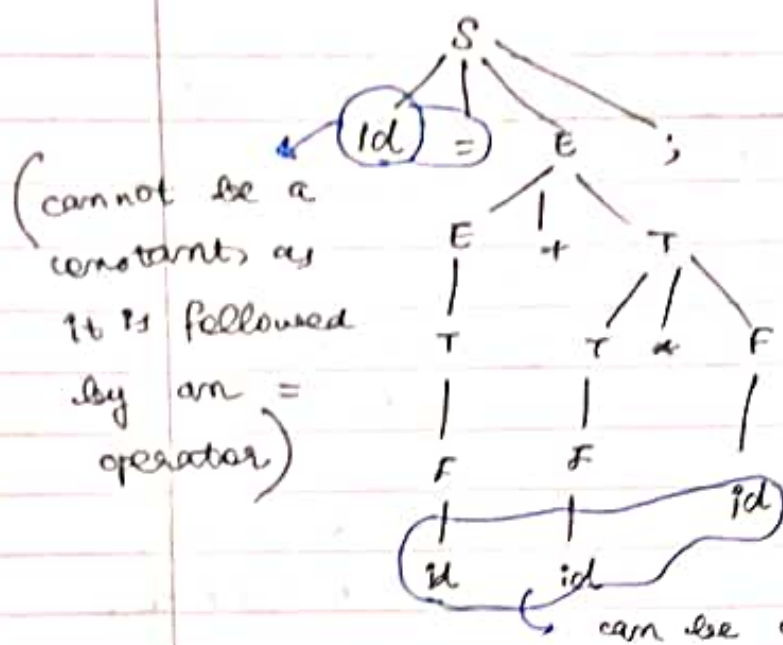
semantic Analyzer

Parse Tree \xrightarrow{SA}

semantically verified parse tree

semantic Analyzer is responsible for:

- * type checking
- * error bound checking
- * correctness of scope resolution



- Semantic analyzer detects
- ✓ type mismatch errors
 - ✓ undeclared variable
 - ✓ misuse of reserved words
 - ✓ multiple declaration of variable within a single scope
 - ✓ accessing an out of scope variable
 - ✓ mismatch between actual & formal parameters

Semantic Analyzer checks for "meaningfulness" of the parse tree and verifies that.

④

Intermediate code Generator:

Semantically verified parse tree → Intermediate code

$id = id + id * id;$

↕

$x = a + b * c;$

① $t_0 = b * c;$

② $t_1 = a + t_0;$

③ $x = t_1;$

Till this phase

↓

FRONT-END

↓

After this backend : To generate target code specific to a platform, we need to only modify the next 2 phases.

Three Address code
[At most 3 addresses of variables are referenced in one line]

TAC

⑤

code optimizer

Intermediate code $\xrightarrow[\text{optimizer}]{\text{code}}$ optimized code

code optimization can be

① Machine dependent

(OR)

② Machine independent

① $t_0 = b * c;$

② $\{ t_1 = a + t_0; \}$

③ $\{ x = t_1; \}$ merge

① $t_0 = b * c;$

② $x = a + t_0;$

⑥

Target code Generator

optimized code $\xrightarrow[\text{Generator}]{\text{Target code}}$ Assembly code segment

3 Address code

① $t_0 = b * c;$

② $x = a + t_0;$ \rightarrow assembly level code

$\left\{ \begin{array}{l} \text{mov } \text{eax}, \text{DWORD PTR}[1bp-8] \\ \text{imul } \text{eax}, \text{DWORD PTR}[2bp-12] \\ \text{mov } \text{edx}, \text{eax} \\ \text{push } \text{eax}, \text{DWORD PTR}[2bp-4] \\ \text{add } \text{eax}, \text{edx} \\ \text{mov } \text{DWORD PTR}[2bp+16], \text{eax} \end{array} \right.$

$\left[\begin{array}{l} \text{mov} \\ \text{moving} \end{array} \right.$

eax: extended version of ax register, a combination of ah and al registers

eax stores 32 bits $\left\{ \begin{array}{l} \text{ah: stores 16 higher order bits} \\ \text{al: stores remaining 16 lower order bits} \end{array} \right.$
eax: accumulator register

DWORD PTR data word pointer, scaling factor -8

a : ebp - 4

b : ebp - 8

c : ebp - 12

x : ebp - 16

imul \rightarrow mnemonic for
signed multiplication

edx \rightarrow dh + dl
(another accumulator register)

Tools for Practical Implementation

LEX \rightarrow Lexical Analysis phase
standard Lexical Analysis Generator in
UNIX Based systems. It reads an input
stream specifying the lexical analyzer
and writes the source code which
implements the lexical analyzer for the
same C programming language.

YACC \rightarrow Yet Another compiler compiler
 \rightarrow Syntax Analysis Phase

It is an LA-LR parser generator. It is
commonly used with LEX to implement
syntax analysis.

LANCE compiler \rightarrow The entire front-end for a C
language for an embedded processor.

Summary

- ① overview of various phases of compiler
- ① Tools to implement different phases.

Symbol Table

outcome :

- ① Usage of symbol table by various phases
- ② Entries of symbol table
- ③ operations on symbol table

Symbol Table :

- ① data structure which is created and maintained by compilers in order to store information about the occurrences of various entities such as :
 - variables and function names
 - objects
 - classes
 - interfaces

Symbol Table - usage by phases

1)

Lexical Analysis

↳ creates entry for identifiers

Lexical Analyzer is a scanner. It scans the entire source code line by line. During the scanning, whenever it encounters any identifier, it creates an entry for that in the symbol table.

2)

Syntax Analysis

↳ Add information regarding attributes like type, scope, dimension, line of reference, line of usage, etc

- 3) Semantic Analysis
→ using available information, checks semantics and updates the symbol table.
- 4) Intermediate code Generation
→ Available information helps in adding temporary variables' information.
- 5) code optimization
→ Available information is used in machine dependent optimization.
- 6) Target code Generation
→ Generates the target code using address information of identifiers.

Symbol Table - Entries

- ① Name - stores name of identifiers
- ② Type - stores datatype of identifiers
- ③ size - specifies size of the identifier
- ④ Dimension - for 1D and multi-dimensional arrays, it stores dimension; for primitive datatypes, value = 1
- ⑤ Line of Declaration - line number of the source code where the identifier has been declared is stored
- ⑥ Line of Usage - line number of the source code where the identifier has been used is stored

[If used in many places, a linked list representation is used]

- ⑦ Address - stores address info of the identifier

(platform dependent)

```
int count ;
char x[] = "NESO ACADEMY";
```

Name	Type	Size	Dimension	LOD	LOU	Address
count	int	2	0	-	-	-
x	char	12	1	-	-	-

As none of the attributes are of a fixed size, it is not possible to know how much space is required for the symbol table before it is actually created.

If we choose a size which turns out to be smaller, we will be able to save the space wastage, however we won't be able to store all the entities.

on the other hand, a bigger size will lead to wastage of space.

considering these circumstances, the best soln would be to dynamically allocate the size of the symbol table during compile time

Symbol Table - operations

- ① Non-Block structured Language: eg Fortran
 - contains single instance of the variable declaration.
 - operations: Insert() Lookup()
 - Discontinued particularly because of the use of unstructured control flow using goto statements.

- _ / _ / _
- ② Block structure Languages
- variable declaration may happen multiple times
 - operations
- Insert() Lookup() set() Reset()

summary

- ① usage of symbol table by various phases
- ② Entries of symbol table
- ③ operations on symbol table

symbol table - solved PYQs

outcome :

- ① GATE 2021 Q on ST
- ② ISRO 2016 Q on ST

GATE
2021

In the context of compilers, which of the following is/are NOT an intermediate representation of the source program?

- (A) Three Address code
- (B) Abstract Syntax Tree (AST)
- ✓ (C) symbol table
- (D) control flow Graph (CFG)

Intermediate code can be of 2 forms :

linear form
↓
3 AC

AND the Tree form
↓
AST

CFG : It is a representation using Graph notation of all the paths that might be travelled through a program during its execution.

ISRO
2016

Access Time of the symbol table will be logarithmic if it is implemented by:

- (A) Linear List (C) Hash Table
✓ (B) Search Tree (D) None of these

Symbol Table - various Implementations

Implementation	Insertion time	Lookup time	Disadvantage
(A) Linear List			
(i) ordered lists			i) For ordered lists, every insertion is preceded by lookup operation
a) Arrays	$O(n)$	$O(\log n)$	
b) Linked Lists	$O(n)$	$O(n)$	
(ii) unordered lists	$O(1)$	$O(n)$	ii) Access time & table size
(B) Search tree	$O(\log_m n)$	$O(\log_m n)$	Always needs to be balanced
(C) Hash table	$O(1)$	$O(1)$	Too many collisions increases time/complexity to $O(n)$ (linked list)

Summary

- ✓ GATE 2021 ✓ ISRO 2016

Introduction to Lexical Analyzer

outcome: ✓ working principle of Lexical Analyzer

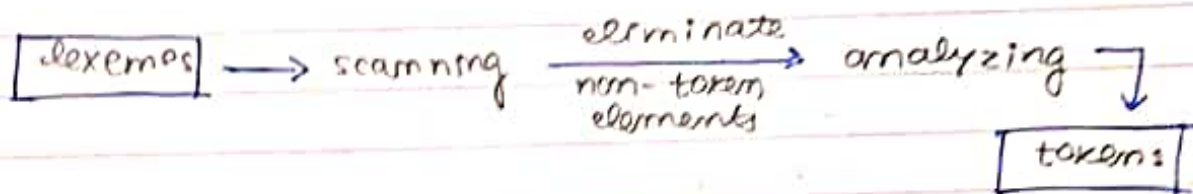
Lexical Analyzer:

(uses DFA for pattern matching) ✓ scans the pure HLL code line by line
✓ takes lexemes as I/p and produces tokens

Tokens

- 1) identifier 2) operator 3) constants 4) keyword
5) literals 6) punctuators 7) special character

2 functions take place in the lexical analyzer : ① scanning ② analyzing

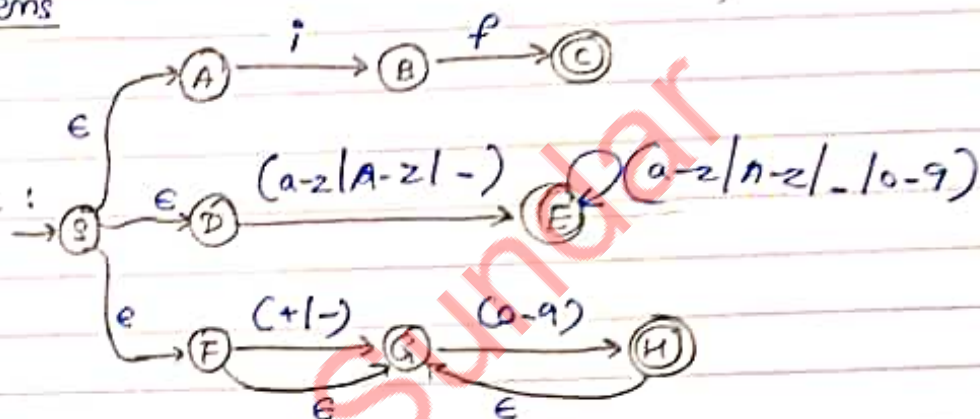


C- tokens

if :

Identifiers :

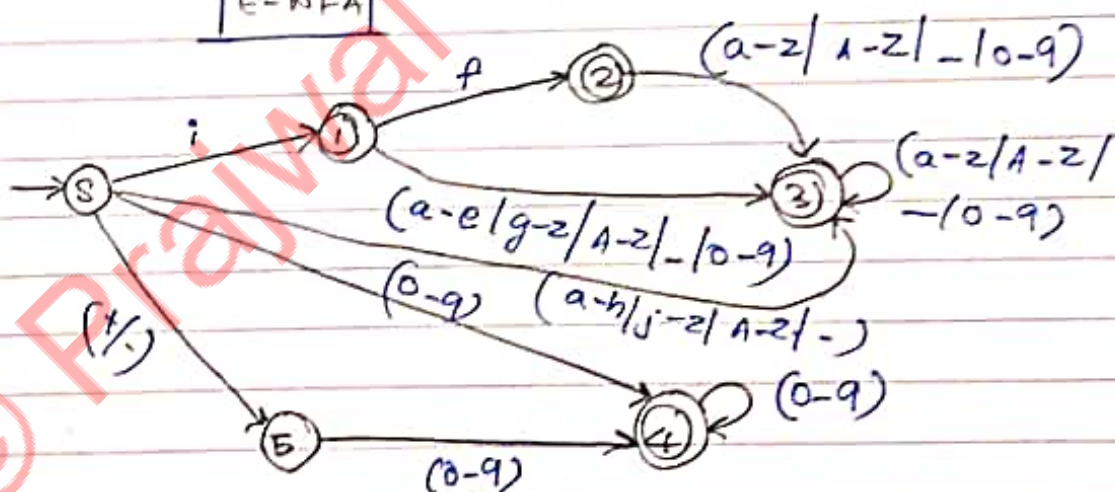
Integer



(DFA is implemented in the Analyzing phase)

E-NFA

equivalent DFA



- ② identifies if
- ① and ③ identify identifiers
- ④ identifies integers

NFA \rightarrow DFA \rightarrow mDFA (minimized DFA)

note features of Lexical analyzer

- ✓ removes comments and whitespace from the pure HLL code.

// single line comment
 /* Multi-line
 comment */

int NE /* a comment */ so;
 int NE SO;

lexical analyzer removes the comment & replaces it with a blank whitespace. During semantic analysis, an error will be generated → so will be reported as an undeclared variable. However, for the lexical analyzer, NE and SO will be treated as 2 different tokens.

whitespaces :

' ' → space 't' → horizontal tab
 '\n' → newline '\v' → vertical tab
 '\f' → form feed '\r' → carriage return

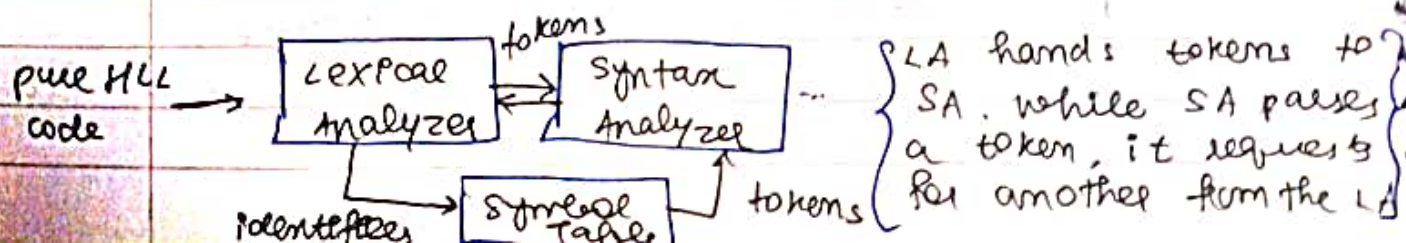
['\t' = 6 x '\n', '\f' → page breaking]
 '\r' → space created by the ENTER key
 ASCII character

All these are recognized and eliminated during the scanning phase of the lexical analyzer.

Summary
 working principle of LA

Feature

① helps in macro expansion in the pure HLL code.



LA hands tokens to SA. While SA parses a token, it requests for another from the LA.

Lexical Analyzer - Tokenization

outcome

- ① count the no. of tokens in a given code segment

```
int main()
{
    int x, a = 2, b = 3, c = 5;
    x = a + b * c;
    printf("The value of x is %d", x);
    return 0;
}
```

Tokens:

①	keyword :	int	→	①
②	Identifier :	main	→	①
③	Punctuation :	() { }	→	③
④	keyword :	int	→	①
⑤	Identifier :	x	→	①
⑥	Punctuation :	,	→	①
⑦	Identifier :	a	→	①
⑧	operator :	=	→	①
⑨	constant :	2	→	①

counts

→ 39

summary :

- | | | |
|---|---------------|-----------------------------|
| ① | keywords : | int ,
return |
| ② | Identifier : | main,
x, a, b, c, printf |
| ③ | Punctuation : | () { }
, ; |
| ④ | operator : | = + * |
| ⑤ | constant : | 2, 3, 5, 0 |
| ⑥ | Literal : | "The value of
x is %d" |

Note : Every occurrence of the token is counted by the lexical analyzer, no matter how many times it occurs in the pure HLL code.

summary : ① count no. of tokens

Lexical Analyzer - solved Problems - set 1

outcome

① 3 solved questions on Lexical Analyzer

Q1
GATE 2011

The lexical analysis for a modern computer language such as Java needs the power of one of the following machine models in a necessary and sufficient sense?

- (A) Finite State Automata
- (B) Deterministic Push Down Automata
- (C) Non-Deterministic Push Down Automata
- (D) Turing Machine

D-PDA and ND-PDA are the acceptors of D-CFL and ~~ND~~-CFL [context free language]

↓

CFLs are generated from CFGs [context free grammars]. During the different phases of compilers, we have observed that in the syntax analysis phase, the syntax analyzer makes use of the CFGs in the construction of the parse trees.

N-PDAs have more expressing power than PDAs

↓

Also, since compilers are implemented on physical systems, ATM can be designed that can map the HLL string into assembly language target code. So Turing Machine has the power to implement the entire compiler itself.

ISRO 20/7

The output of a lexical analyzer is :

- (A) A parse tree (B) Intermediate code
(C) Machine code ✓ (D) stream of Tokens

GATE 2000

The number of tokens in the following statement is

```
printf ("i = %d, &i = %x", i, &i);
```

- (A) 3
↓
no. of arguments
- (B) 26
↓
counting each token in string
- ✓ (C) 10
↓

correct ans

" "
is one token
- (D) 21
↓
counting permit as one string + each token in string
new

Summary :

- ① 3 solved questions on Lexical Analyzer

Lexical Analyzer (solved Problems) - set 2

outcome

- ① 2 solved questions on Lexical Analyzer

GATE 2017
NIELIT
scientist-B
2017

In a compiler, keywords of a language are recognized during

- (A) parsing of the program
(B) the code generation
✓ (C) the lexical analysis of the program
(D) dataflow analysis

DF analysis is performed during code optimization on the program flow graph. So, it will not recognize tokens, rather it analyses the flow control of the tokens.

Q2 :
GATE 2018

A lexical analyzer uses the following patterns to recognize 3 tokens T_1, T_2, T_3 over the alphabet $\{a, b, c\}$.

$T_1: a?(b|c)^*a$

$T_2: b?(a|c)^*b$

$T_3: c?(b|a)^*c$

[note that $x?$ means 0 or 1 occurrence of the symbol x . Note also that the analyzer outputs the token that matches the longest possible prefix.]

GATE 2018

which one of the following is the sequence of tokens output by the analyzer, if the string "bbaacabc" is processed?
(A) $T_1 T_2 T_3$ (B) $T_1 T_1 T_3$ (C) $T_2 T_1 T_3$ (D) $T_3 T_3$

b b a a c a b c
 T_1 T_2 T_3

✓

b b a a c a b c
 T_1 T_1 T_3

✓

b b a a c a b c
 T_2 T_1 T_3

✓

b b a a c a b c
 T_3 T_3

✓

↓
(longest possible prefix)

summary

✓ 2 solved Q on LA

Errors and Error-Recovery in Lexical Analysis

outcome :

- ✓ Role of Error handler, especially for Lexical Analysis.
- ✓ Types of Error
- ✓ Different Types of Lexical Errors
- ✓ Error Recovery in Lexical Analysis

Error Handler:

- ① Error detection
- ② Error reporting [generating error reports to the user]
- ③ Error recovery [implementation of some recovery strategy for handling errors]

Types of Errors:

- ① Run-Time Errors: These take place during the execution of the program code. Eg: insufficient ^{available} memory space, unexpected errors, etc. all the examples of runtime errors.
- ② Compile-Time Errors: These occur during compile time, i.e. before execution of the program. There are 3 types of compile-time errors:
 - ① lexical errors
 - ② syntax errors
 - ③ semantic errors

Lexical Errors:

- ① Identifiers that are way too long
eg C: ANSI ^{allows} 6 significant characters from the identifier's names. external identifiers are the identifiers that are used in macros.

Anyways ANSI allows:

- 31 significant characters → internal identifiers
- ② which are basically declared in a function block.

on the other hand, microsoft c compilers allow internal ^{identifiers} ~~variables~~ to have ^{upto} 247 significant characters.

C++

For C++, both Intel and Microsoft compilers allow upto 2048 significant characters for the identifiers.

Python

In python, identifiers are allowed to have upto 79 significant characters.

- ① Exceeding length of numeric constants.

int i = 4567891;
size 2 Bytes : - 32768 to 32767 only

- ② Numeric constants which are ill-formed.

int i = 4567 \$91;

- ③ Illegal characters that are absent from the source code.

char x[] = "NESO ACADEMY"; \$

Lexical - Error - Recovery

- ① Panic - Mode Recovery

→ It is the most basic way of recovery. In this, once an error has been encountered, successive characters are skipped until a valid delimiter is found.

int 4 NESO) → delimiter
error skipped

while (condition)
{

} → delimiter

For blocks, } is the delimiter.

Ignore in b/w → "panic mode" recovery

//_

② Transpose of 2 adjacent characters

union test

```
{  
    int x;  
    float y;  
}  
TI;
```

union wrong positions

union Ⓢ valid keyword



Error handler transposes o↔i
and forms the correct keyword
-taken union.

③ Insert a missing character

it NESO; → int NESO;

④ deleting an unknown or extra character

int NESO; → int NESO;

⑤ Replacing one character with another

itt NESO; → int NESO;

therefore, these are the various error recovery strategies implemented by the error handler during the lexical analysis phase.

summary:

- ✓ Role of Error Handler, especially for LA.
- ✓ Types of Errors.
- ✓ Different Types of Lexical Errors.
- ✓ Error Recovery in LA.

— X —