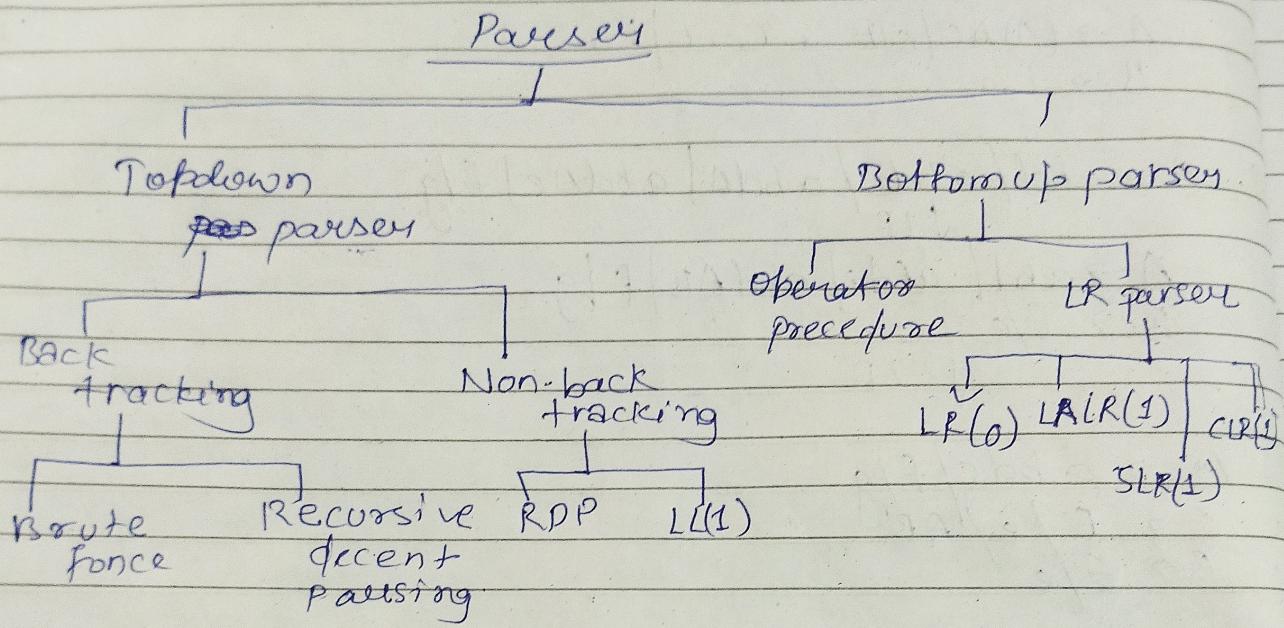


PARSER

Ex:-

⇒ Parser is a piece of Software that can perform Syntax Analysis



□ Top down parser

→ purpose of parser:- check whether this word (w) belongs language generated by a grammar or not.

□ B

→ we start with starting Symbol(s) and gradually move to down until we check the given word is belong to the language generated by grammar or not.

⇒ C

→ For ~~any~~ Top-down parser we want parser ~~free~~ from left recursion, do not have unambiguous grammar

Ambiguous and unambiguous grammar

→ Ambiguous grammar

- * An ambiguous grammar is a grammar where exist more than one derivation for any word, that belongs to the language generated by the grammar (well-formed)
- * Ambiguous means → not specific/confusion

→ Unambiguous grammar

- * where exist only one derivation for any word that belongs to the language generated by the grammar (well-formed)

El1

If grammar ambiguous

An ambiguous grammar is a grammar where exist more than one derivation for any word, that belongs to the language generated by the grammar (well-formed)

Derivation

This ambiguous left

Derivation

Left most

Derivation is equal to Right most derivation (LMDT = RMDT) for ambiguous grammar you have different (LMDT ≠ RMDT)

If grammar is both left recursive and right recursive, then this grammar is ambiguous

(iii) Ambiguous at CFG (context free grammar) is undecidable: means there does not exist an algorithm to determine

PROPERTIES

Properties

Left

Right

Left

Right

Left

Right

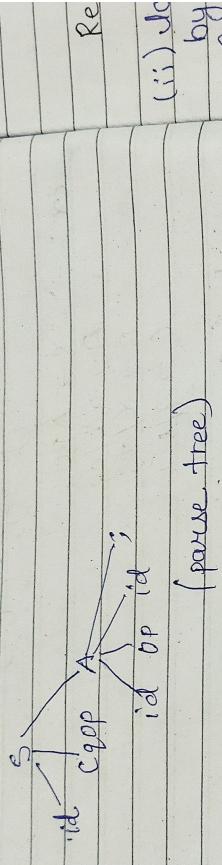
Left

Right

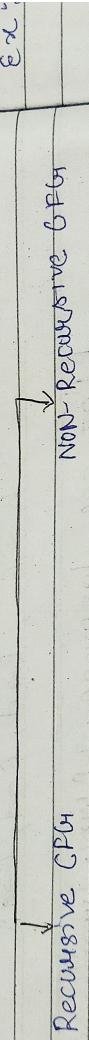
Parse tree :- Parse tree is a pictorial representation at the process of deriving word from this grammar.

Ex :- Take above eq (i) which is derived from the grammar w.

$$S \rightarrow id \ CGOP \ op \ id \ op \ id;$$



Classification of CFG (context free grammar)



Recursive CFG

If there exist at least one production having same variable on L.H.S & R.H.S. Then it is non- recursive CFG.

Ex :- $S \rightarrow Sa/b$
 $; - S \rightarrow aS/b$
 $; - S \rightarrow aSb/E$

$$\begin{array}{ll} S \rightarrow S \text{ or } B & S \rightarrow AB/BA \\ A \rightarrow b/C & A \rightarrow a/B \\ B \rightarrow C/d & B \rightarrow b/c \end{array}$$

Contra
 \Rightarrow For such cases

 1) Ex:-

Lexeme :- A Lexeme, is a sequence of characters in the source program that is matched by pattern by pattern for a token.

Lexical analyzer design :-

specification of token \rightarrow Regular expression
recognizing tokens \rightarrow finite automata.

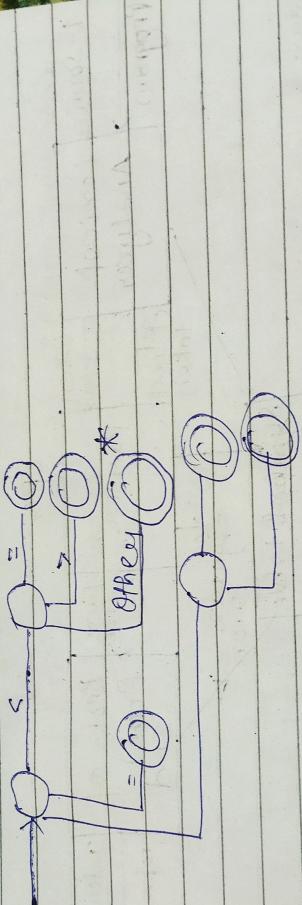
① Finite Automata for keyword and identifiers:-
Others



R.F. :- Identifier (letter + digits)*

in process am. hims the

Transition diagram for identifier :-



is ut

Tokens Ex :-

1. main() \rightarrow 3
2. { \rightarrow 1
3. int x,y; \rightarrow 5
4. f J /* gate, * / 0 at z; \rightarrow 9
5. x = /* exam */ 10; \rightarrow 4
6. y = 20; \rightarrow 1
7. \rightarrow 22

- int (keyword)
- x (identifier)
- = (Assignment operator)
- 10 (literal)
- ; (Delimiter)

Patterns

→ A pattern is a rule or a regular expression
that defines the structure of a token in the source code.

→ The lexical analyzer uses these patterns to match sequences of characters (lexemes) and identify them as tokens.

Input buffering

→ Input buffering is a technique used in lexical analysis to efficiently read and process the source code from an input stream.

→ It minimizes the number of interactions with the input file and optimizes the performance of the lexical analyzer.

Two-buffer Scheme

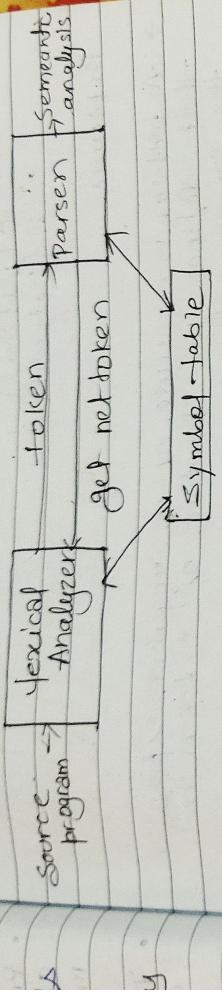
A common approach to input buffering is the two-buffer scheme. Where, the input is divided into two buffers of equal size.

Lexical
choice
is re
or

Lexic

①

The Role of Lexical analysis



(i) Lexical analyzer divides the given program into meaningful word which is known as tokens.

(ii) Tokens are generally identifiers, keywords, operators, constant and special symbols.

(iii) Lexical Analyzer eliminate the comment lines from the given program.

(iv) L.A help in giving error message by providing row no & column no.

Tokens:-

A token is a sequence of characters from the source code that matches a predefined pattern and is related on a single logical unit during lexical analysis. Tokens are the building blocks of a program's syntax.

Ex :- int x = 10;

Code generation phase:-

- It is a program that translates optimized intermediate code into assembly language or target code.

Symbol Table :- It is a data structure that contains information about identifiers and constants present in the program.

Error Handler :- If any phase of compiler detects error then it is stored in error handler.

Front End

Front end comprises of phases which are dependent on the input (source language) and independent on the target language.

Semantic

Semantic Analysis phase:-

1. It is a program that takes parse tree as input and produces annotated parse tree as output.

2. It also detects semantic errors present in the program. (Type checking)

3. Syntax directed Translation is used.

Intermediate code generation phase:-

1. It is a program that translates high level code into intermediate code.

2. Advantages of generating intermediate code is to perform optimization in simple way.

3. Syntax directed Translation is used

Code optimization phase:-

1. It is a program that reduces time and space required by the target machine by removing some unnecessary code.

2. There are two type of optimization performed by compilers known as machine independent optimization and machine dependent optimization

Errors		
Compile-time	Run time	
Lexical phase errors	Syntactic phase error	Semantic errors

- ↳ Infinite Recursive calls;
 - defined by zero
 - infinite loop
 - wild pointers
- ↳ Segmentation

Lexical Analysis Phase :-

1. It is a program that takes high level language as input and produces token as output.
2. It also detects lexical errors present in the program.
3. Regular expression & finite automata used to design.

Syntax Analysis Phase :-

1. It is a program that takes token as input and produces parse tree as output.
2. It detects syntax error present in the program.
3. Context-free grammar & push down automata used.

Seg-S
Semantic

1. It is an input-output process

2. It is a parser.

3. Syntax

3) Writing a Grammar

4) LR Parsers

5) LR Derivation: left-to-right scan, rightmost derivation.

Types

a) SLR (Simple LR): Basic form; uses follow sets for reductions.

b) LALR (Look-Ahead LR): Combines similar states in SLR to reduce table size.

c) Canonical LR: Most powerful but complex.

6) Steps:

a) Construct a DFA for ~~were~~ viable prefixes.

b) Build a parsing table with states and transitions.

10.) Parser Generators (YACC)

7) YACC: Yet another compiler compiler.

8) A tool for automatic parser generation.

9) Input: Grammar specification.

10) Output: C code for a parser.

11) Applications: Building parsers for custom languages.

11. Error Recovery Strategies

12) Goal: continue parsing despite syntax errors.

13) Techniques:

a) Permit mode: skip tokens until a synchronizing token is found.

b) Phrasal-level recovery: Replace / insert / delete tokens locally to fix errors.

c) Error productions: extend grammar with rules to handle common errors.

d) On-the-fly correction: use algorithms to find the minimal changes to the input.

Section on memory.

~~free(header);
return 0;~~

6. Bottom - UP Parsing

Definition: Constructs a parse tree starting from leaves (terminals) and works upward toward the root.

Techniques:

- Shift Reducing Parsing: shift tokens to stack and reduce based on grammar rules.
- Handles both left-reduction and right-reduction.

7. Handles and viable Prefixes

- ❖ Handles: A substring of the input that matches the RHS of a production and can be reduced to a non-terminal.
- ❖ Viable Prefix: A prefix of the input string that can appear on the stack of a shift-reduce parser.

8. Operator Precedence Parsing

- ❖ Use case: Simplified parsing for expressions with operators like +, -, *, /.

- ❖ Steps
 - a) Define operator precedence and associativity.
 - b) Create a precedence table.
 - c) Parse based on precedence relationships.

Syntax analysis

Module - 3

lexically correct
beginning word from

Syntax analysis translates the sequence of tokens produced by the lexical analyzer into a structural representation (parse tree) using grammatical rules.

1. The Role of a parser

- Parser's job: To check the symbols of a program according to the given grammar (usually a context-free grammar - CFG)
- And create a parse tree.

- Inputs: Tokens from the lexical analyzer. Also
- Outputs: Parse tree or syntax tree.

4 Key tasks:

- Syntax checking: Ensure the code follows grammar rules
- Error detection: Identify and report errors when the symbol doesn't match the grammar.
- Tree structure building: Represent program structure for further phases like semantic analysis.

2. Content - Free grammars (CFG)

- Definition: A set of production rules used to describe the syntax of a language.

- Components:
 - Non-terminals (N): Basic symbols (like keywords, operators, etc).
 - Terminals (T): Variables that represent syntactic constructs.
 - Start Symbol (S): The root of the grammar.
 - Production Rules (P): Rules defining how terminals and non-terminals combine.
- Example: Grammar for arithmetic expressions.

$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow F \mid f \\ F \rightarrow (E) \mid \text{Id.} \end{array}$$

2. Context Free Grammars

5. Recognition of tokens and finite automata to recognize tokens.

- * The lexical analyzer uses finite automata to recognize tokens.
- * Deterministic Finite Automata (DFA): A simple DFA but allows multiple possible states at a time.
- * Non-Deterministic Finite Automata (NFA): Similar to DFA but allows multiple possible states at a time.
- * machine with fixed states that determines if a pattern matches a pattern.

6. Conversion from Regular Expressions to NFA:

- * Regular Expression to NFA:
 - * NFA's are built step by step from regular expressions using like:
 - * concatenation (e.g., ab → a path for a followed by b).
 - * union (e.g., ab → paths for a and b).
 - * Kleene. (e.g., a^* → loops on a).

7. NFA to DFA:

- * convert an NFA into a DFA using a method called subset construction.

8. Lexical analysis generator (Lex)

- * Lex is a tool used to generate lexical analyzers automatically.
- * Lex contains regular expressions and tokens.

9. Design of a lexical analyzer

- * Input: A file containing regular expressions and corresponding action.
- * Output: A C program that performs lexical analysis, identifies, defines patterns for keywords, identifiers, etc.
- * Example: Defining patterns for tokens using regular expression.

Summary:
lexical analysis is about breaking source code into tokens using regular expressions, matching compiler tokens like lexical analyzer, tools like DFA or NFA, DFA is more efficient.

Sanjana Mather

On Roll No. 01, under Abul Kalam Azad University of Technology, West Bengal

2. Key concepts

☞ Tokens: Smallest meaningful elements of a program, like, if, while, +, *, 123.

☞ Patterns: Descriptions of what a token looks like, often defined using regular expressions.

Example: A pattern for an identifier might be $[a-zA-Z][a-zA-Z0-9]^*$.

☞ Lexemes: Actual sequences of characters that match a pattern and form a token. Example: for the pattern

$[a-zA-Z][a-zA-Z0-9]^*$, the lexeme could be x , count, or variable i .

3. Input Buffering

☞ Need: Reading the source code character by character can be slow. Input buffering allows reading in chunks for efficiency.

☞ Mechanisms: It uses two pointers:

1. Start pointer: Marks the beginning of the current lexeme.
2. Forward pointer: Moves ahead to find the end of the lexeme.

☞ Specification of tokens: For example, tokens are defined using regular expressions. For example,

- ☞ A number: $[0-9]^+$
- ☞ An identifier: $[a-zA-Z][a-zA-Z0-9]^*$

Lexemes, tokens, finite automata. From a regular expression to NFA.

Recognizes Finite Automata. From a regular expression to DFA. Design of a

NFA. From a regular expression to DFA. Design of a

From a regular expression to DFA. Design of a

lexical analyzer generator (lex).

lex (91)

Example: In compilers design context

```
int a = 10;  
float b = 20.5;  
float c = a+b;
```

Type Checking:

- a is an integer, and b is a float. The operator a+b is called addition.
- a is automatically promoted to float before addition.

Type Conversion:
• a is automatically promoted to float before addition.

Ex:

```
int a = 10;  
char b = 'C';  
float c = a+b;
```

float c will be.

The type checker will:

- 1) Check b is a char and a is int.
- 2) Convert b into an int (force b into an int before storing the result in c).
- 3) Promote both to float.

Module - 6
Modularization or forward declaration.

Run-time environment: This refers to the setup or memory, which occurs after it is compiled. It involves managing memory, program execution, program flow, etc. It also deals with variables and execution. There are challenges that occur when it comes to programming language itself, like how it handles language issues of scope, flow of control, and memory.

1. Activation records: This is a structure that contains information about the nature of variables, and scope, that shows how the memory handles memory, variables, and activation records. This call each other during the execution of different functions (or procedures) of a program. It helps visualize the order and relationships of calls.

Storage Allocation Strategies

- ↳ Context Stack: A stack used to keep track of active function calls when a function is called, its details are added (pushed) onto the stack, and when it ends, these details are removed (popped).
- ↳ Scope of Declaration: Declares where a variable is accessible in the program.
- ↳ Local Scope: Variable is accessible only inside a specific function or block.
- ↳ Global Scope: Variable is accessible throughout the program.
- ↳ Binding of Names: The process of linking a variable or function name to its static or dynamic binding.
- ↳ compile-time (Static Binding): Fixed before the program compilation.
- ↳ run-time (Dynamic Binding): Decided during program execution.

- Storage Organization: This is how memory is managed during the execution of a program.
- ↳ Subdivision of Run-time memory:
 - ↳ Memory is divided into sections for different purposes.
 - ↳ Code section: Stores program instructions.
 - ↳ Static section: Stores global variables.
 - ↳ Stack: Stores function call details and local variables.
 - ↳ Heap: Used for dynamic memory allocation.
 - ↳ memory required at run time.
 - ↳ Activation Record: A block of memory created on the stack for each active function call, it stores:
 - ↳ local variables.
 - ↳ local variable addresses (where the program goes after the function ends).
 - ↳ Return address (where the program goes after the function ends).
- the one automatically

c. 3.) Equivalence: If type expressions whether two type expressions are the same.

~~are the same.~~ ~~types are equivalent if they have the same structure, regardless of their names.~~

Example:

type A = float, float y;

type B = float, float y;

A and B are structurally equivalent.

~~4.) Name equivalence~~: Types are equivalent only if explicitly declared as the same.

Example:

type A = int;

type B = int;

A and B are not name-equivalent.

5.) Type conversions: When types don't match but the operation is still valid, the compiler uses type conversion (or type casting) to make them compatible.

Conversion: The compiler automatically converts one type to another.

Implicit type conversion (casting): The compiler automatically converts one type to another.

Example:

int a = 10;
float b = a; // implicitly converts int to float

Explicit type conversion (casting): The programmer explicitly specifies the type to convert to.

Example:

float a = 10.5;
int b = (int)a; // explicitly converts float to int.

~~Conversion can occur if the conversion loses data or changes meaning, e.g., translating a float to an int.~~

Storage Allocation

Module - 5

Type Checking: Is a critical aspect of compiler design, ensuring that the type of variables, expressions, and functions in a program are used consistently and correctly.

1.) Type Systems: A type system defines the rules that assign a type to each variable, expression, and function in a program. It ensures that operations are performed on compatible data types to prevent type-related errors.

- Static Typing: Types are checked at compile-time (e.g., C, Java).
- Dynamic Typing: Types are checked at runtime (e.g., Python, JavaScript).

A type system has the following properties:

- Soundness: Ensures no type errors occur during execution.
- Completeness: Every valid program can be assigned a type.

2.) Specification of a simple type checker

A type checker ensures that the program adheres to the type system rules. Here's how it works step-by-step:

1.) Symbol Table: The compiler generates an abstract symbol table. This table stores information about variables, functions, and other symbols defined in the code.

2.) Type Rules Application: Each node in the AST is checked against predefined rules.

- For example, in an expression $a+b$, the type checker verifies that both a and b are numeric types.

3.) Symbol Table: It uses the symbol table to resolve type information of variables and functions.

4.) Error Reporting: If a type mismatch is found, it reports an error message (in pseudocode):

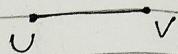
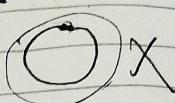
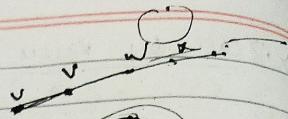
```
if (a+b) &
    check a.type == int & b.type == int;
else &
    report_error("Type Mismatch");
```

for
syntax
batch system

graph has no cycles
no path that starts & ends at same vertex.

DAG (Directed Acyclic Graph)

edges in graph have direction
(they go from one vertex to another)



Properties:-

- 1) Topological sorting
- Linear ordering of vertices
- for every directed edge from vertex U to V
 U comes before V

2) Dependency Resolution

\because no cycle, you can resolve dependencies
in a clear, step-by-step manner.

3) Version control

systems like Git, DAGs used - commit history.

Type Checking

ensure types of values used in a prog. match the expected types.

→ imp. concept in prog. lang.

→ ensure - operations you perform is valid

→ prevents errors caused by mismatched types.

1) Static TC → correctness is checked at compile time

before prog. is run.

→ Type of every variable, fun? & expression is determined & checked by compiler.

C, C++, Java

ex - int n = 5;

n = "Hello"; // error.

less flexibility.

2) Dynamic TC → correctness checked at runtime. (Prog. exec.)

variable can be assigned values of diff. types during execution.

X = 5

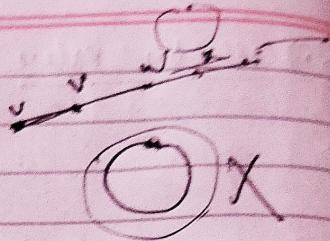
X = "Hello" # No error, X can change type

for
syntax
tree
batch system

graph has no cycles
no path that starts & ends at same vertex.

DAG (Directed Acyclic Graph)

edges in graph have direction
(they go from one vertex to another)



Properties:- 1) Topological sorting

Linear ordering of vertices

for every directed edge from vertex U to V
 U comes before V

2) Dependency Resolution

\because no cycle, you can resolve dependencies
in a clear, step-by-step manner.

3) Version control

systems like Git, DAGs used - commit history

(3)

Compiler

- used in C, C++, Java, etc.
- compile programs at one go.
- CPU utilization is more.
- compilation time is more.
- execution time is less.
- overall better efficiency.

Interpreter

- used in Python, Ruby, etc.
- Interpretation is step by step.
- less.
- Interpretation time is less.
- execution time is more.
- less efficient.

Symbol table

- crucial DS used in compilers
- to store info. about identifiers (variables, functions, objects, etc) in source program.
- features → key info... stored :-
- [Name] - (of identifier)
- [Type] - (data type, return type, etc.)
- [Scope] - (global, local, block level)
- [Memory location] / address
- [Size] - (based on data type)
- [Attributes] - (parameters for func., modifiers, etc.)
- function - facilitates semantic checking (type checking)
- supports code generation.
- detects & handles errors.
- Str. of symbol table → implemented as (hash table / tree-based) structure for efficient lookup, insertion, and deletion.

Symbol table →

<u>Identifier</u>	<u>Type</u>	<u>Scope</u>	<u>Size</u>
x	int	Global	4 bytes
y	float	"	4
func	function	"	-
a	int	local	4
z	int	local	4

⑥ Lexical Analyzer (design?)

- 1st phase of the compiler.
- 1st task - read the source code & convert it into a seq. of tokens.
- functions → Tokenization breaks up source code into tokens.
ex: if, int, =
- Removing whitespace and comments → eliminates irrelevant data.
- Symbol table maintenance → add identifiers
- Error detection → identifies invalid tokens and generates appropriate error messages.
- Working of lexical analyzer:
 - Scan the source code character by character.
 - uses patterns (regular expression) to identify tokens & categorizes them based on predefined rules.

Ex:- input:- int x = 10;

Tokens:- int → keyword ; = → assignment operator;
10 → literal ; ; → delimiter.

→ advantage:-

Reduces the complexity of pairing (engineering source) & signifies syntax analysis phase (code into tokens).

Source code → [Lexical Analyzer] → [Tokens] → [Syntax Analyzer]

FF

- (2) Lookahead Operator in RE.
- type of assertion in RE.
- check if something specific comes next without actually grabbing it as part of the match.

→ Ex:- cat(?!s) : Negative lookahead

- Positive lookahead (?!s) : ex: cat(?!s)
- ex: cat(?=s) : find "cat", but only w/ its immediate follower by an "s".
- find "cat" but only if it is immediately followed by an "s".
- Test: cat are rule : cat : It finds cat (cat). Its match: It finds cat (cat). It nor followed by "s".

But "s" is not included in the match.

Identifier

- Name given to variables, functions, etc.
- user-defined? Yes
- Yes, it can be used as long as it's not a keyword.
- ex. count, sum, total
- represent user-defined entities

Error Handling

- detecting & responding to errors, occurs during execution of program.
- Allows program to continue running / fail gracefully, instead of crashing / producing incorrect results.
- helps ensure that prog. behave as expected, even when unexpected events occur.
- provides meaningful feedback to users & developers about issues, making it easier to debug.
- Types of errors →
- o Syntax error : - code is written incorrectly. If the prog. lang. can't understand it.
ex.: print ("Hello world") (missing closing parenthesis)
- o Runtime errors! - Occur while the prog. is running
- ex.: 5 / 0 (dividing by zero, accessing an undefined variable).
- 3) Logic errors : when prog. runs without crashing, but output is incorrect due to flaws/bugs
ex - using wrong formula to calculate area of circle.

Regular Expressions

- Reserved word with pre-defined meaning in the lang.
- NO
- NO
- powerful tool for pattern matching
- previous to search, specific patterns for use in tasks like
- different characters
- Metacharacters → An asterisk → matches any character
- Quantifiers → A plus → matches one or more characters

Lexical Analysis

- in context of LL(1) parser
- refers to a state
- o Valid State
- Once the rule is triggered will produce useful output
- Error

LL(1) Parser

- lexical analysis
- tool used +

YACC

- (Not YACC)

Thompson's Algorithm

- used +

Algo. w/ NFA

- used +

Steps

- helps !

Common Error Handling Techniques:

1. Try & catch (exception handling).
2. Throw codes
3. Assertions
4. Logging

Register Allocation and Assignment

Efficient use of CPU registers → crucial for performance.
→ determine which variable should use registers during program execution.

→ Assign specific registers to variables.

Techniques :-

(1) Graph colouring → Variable nodes in graph. Variable

if variable interfere:— edge connects 2 nodes

assign registers (color) to nodes such that

(2) Spilling →

If there are more variables than registers,

store some variables in memory.

Challenges :- Optimization → balance compilation speed

- 2) handles diverse architecture (RISC vs CISC)
- 3) manage limited register availability.

Processor

Code Optimizations

Bottom-up Parsing

→ starts from the root of the parse tree & goes down to the leaves.

attempt to construct tree down to the leaves.

Code Generation

- (1) final phase
- (2) Intermediate \rightarrow target m/c code.
- (3) focus — producing efficient & correct code.

Issues \rightarrow

- (1) convention \rightarrow at b/c pressure demanded by source program.
Order of vars. $\&$ variables \rightarrow source program.
must remain unchanged. \rightarrow time \downarrow space \downarrow memory usage \downarrow minimize resource consumption.
- (2) efficiency \rightarrow execution time \downarrow space \downarrow minimize
 \uparrow time
- (3) error handling \rightarrow manage runtime errors
(e.g. divide by zero)
- (4) Handling complex str. \rightarrow array \downarrow pointer \downarrow function calls \downarrow generate code effectively
- (5) Register Allocation \rightarrow efficiently manage limited no. of registers in target m/c.

Register Allocation or
determining which
variables do
not require
register allocation.

Techniques :-

- (1) Graph colouring

- (2) Spilling -
3 > 2 > 1

Challenges :-

Code Generation

\rightarrow start / TELL how
attempt to config

- Abstract Code Generator
- Input — Intermediate representation (3-address code)
 - Output — Target m/c code / assembly instruction
 - Steps :-
 - (1) Instruction selection \rightarrow each IR statement \rightarrow m/c instruction.
 - (2) Register Allocation \rightarrow Assign variables to register while minimizing memory.
 - (3) Optimization \rightarrow Apply peephole optimization to improve efficiency.
 - (4) Final code \rightarrow Output — the final instructions.

- IR code :-
- ```

t1 := a + b
 LOAD R1, a
 ADD R1, b
 STORE R1, t1
t2 := t1 * c
 MUL R2, t1, c
 STORE R2, d

```

## ① Three Address code Generation (TAC)

- intermediate representation used in compilers
- to simplify the process of code generation.
- breaks down complex expressions into simpler statements.
- where each statement involves at most 3 operands.
- each instruction has at most 1 operator and 3 operands.

2 Source  
operands      1 destination

Ex:- High-level code:-  $a = b + c * d;$

Step by Step TAC:-  
 $t_1 = c + d$  //multiple c and d & store result  
 $t_2 = b + t_1$  //add b and  $t_1$ , store result  
 $a = t_2$  //assign  $t_2$  to a.

→ components of TAC:-

i) Operators:- Arithmetic (+, -, \*, /)

Relational (<, >, <=, >=, ==, !=)

Logical (And, Or).

Assignment (=)

ii) Operands:- variables (a, b)

constants (10, 3.14)

Temporary Var. ( $t_1, t_2$ ).

→ Benefits → simplifies code generation

→ improves Optimization

→ mc - independent

[YACC]

Yet Another Compiler Compiler  
prog that generates parser for a given grammar.

→ I/P to YACC → context-free grammar (CFG)

→ automatically assigns no. for tokens, but can be overridden.

→ generates LALR(1) parser.

→ parser generated through this has something to do with syntax definitions of prog. lang.

→ detect syntax errors.

→ able to continue parse after errors have occurred.

Left to Right most derivation with 1 token.  
look ahead (what type of parser YACC produces)

## ② Quadruples, Triples, Indirect

Aspect

Represent.

Quadruples

Each inst. has 4

(operator, Arg1, Arg2)

Fields

1. Operator

2. Argument 1

3. Argument 2

4. Result (stores name)

Ex. code

$a = b + c$

operator

+

Arg 1

b

Result field

Explicit

use

Easy to implement  
straightforward handling.

(i) call by value → actual value of argument is passed to the function.

→ any changes made to parameter inside the function don't affect the original variable.

ex: void func(int n){  
 ? ~~int~~ x = 10;  
 } // not affect the original variable.  
  
int main(){  
 int a = 5;  
 func(a); // a remains 5 after fun. call.  
 ?  
}

(ii) call by reference → address (reference) of argument is passed to the fun.

→ changes → affect original

b/c both fun. & caller refer to same memory location.

void func(int &x){  
 ?  
}

(iii) call by name → more commonly seen in lang. like

ALGOL

→ rarely used in modern prog. lang.  
→ arg. not evaluated before fun. is called.  
→ no textual substitution.

procedure func(a){  
 print(a);  
 ?  
}

begin  
 func(2 + 3); // 2 + 3 is substituted directly.  
end

$$\begin{aligned}
 & a) \{0, 1, 2\} \quad b) b * ab * ab \\
 & = (0, 1, 2) \quad c) (aabb)^*(aab|bbb)(ab|b)^* \\
 & \text{ex: } (aabb)^* = \underbrace{(aabb)^*}_{(aabb)^*} \cdot \underbrace{(aabb)^*}_{(aabb)^*} \cdot \underbrace{(aabb)^*}_{(aabb)^*} \cdot \underbrace{(aabb)^*}_{(aabb)^*} \cdot \underbrace{(aabb)^*}_{(aabb)^*} \\
 & \text{Ques: Constant folding} \\
 & \text{Calculating constant expressions during compilation instead of at runtime.}
 \end{aligned}$$

$\Rightarrow$  Before optimization: -  $x = 2 + 3$

After optimization (fold constant): -  $y = x * 4$

Compiler precomputes  $2+3$ , replacing it with  $5$ .

b) Common Sub-expression elimination (CSE)

Identifying & removing duplicate calculations by reusing the result of a previous calculation.

Before optimization: -  $t1 = a+b$   
 $t2 = a+b$   
 $t3 = t1 * c$

After optimization (eliminate duplication): -  $t4 = a+b$   
 $t3 = t4 * c$

c) Loop Unrolling

This expands the body of a loop to produce the no. of iterations which can make the prog. run faster.

Before optimization:

```

for (i=0; i<4; i++) {
 x[0] = 0+2;
 x[1] = i+2;
 x[2] = 2+2;
 x[3] = 3+2;
}

```

After optimization:

|                                                                       |                                                         |
|-----------------------------------------------------------------------|---------------------------------------------------------|
| loop w/ automatic, main local var., for                               | loop w/ automatic, main local var., for                 |
| replaced with direct calculation, size-limited, no for each iteration | direct calculation, size-limited, no for each iteration |
| speed-decreased due to deallocation of stack                          | deallocation of stack                                   |

d) Loop Jamming

When combines multiple loops into one when they share the same iteration range, reducing overhead.

Ex: - Before

```

for (i=0; i<n; i++) {
 for (j=0; j<m; j++) {
 a[i][j] = b[i][j] + c[i][j];
 d[i][j] = e[i][j] + f[i][j];
 }
}

```

After:

$\{$

$\}$

4

(2) Quadruples, Triples, Indirect Triples

Aspect

Quadruples

Triples

Indirect Triples

Represents:  
"Each inst" has 4 fields

3 fields

uses a table of pointers  
to triples.

(operator, Arg1, Arg2, Result)  
(operator, Arg1, Arg2)

operator

points to triple index

Argment 1  
Argument 2

Argument 1  
Argument 2

Argment 1  
Argment 2

Argument 1  
Argument 2

Result (stores temp. Variable  
name)

As in Triple.

Ex. code  
 $a = b + c$

$a = b + c$

operator  
+

pointer to triple of +

Arg 1  
b

b

Arg 2  
+

+

Result field  
Explicit

No result field  
(implicitly)

uses pointer table  
to reference triples.

use  
easy to implement

More compact repres.

Best for

straightforward result  
handling

pointer to manage  
code reordering  
during optimization.

advanced  
optimization.

### State transitions.

State  $I_1 \rightarrow F \rightarrow id$ .

State  $I_2 \rightarrow T \rightarrow F$ .  
     $T \rightarrow T * F$ .  
     $T \rightarrow T + F$ .

State  $I_3 \rightarrow T * F$ .  
 $F \rightarrow id$

State  $I_4 \rightarrow T \rightarrow T * F$ .

State  $I_5 \rightarrow E \rightarrow T$ .  
 $E \rightarrow E + T$

$T \rightarrow T * F$

$T \rightarrow T + F$

$F \rightarrow id$

State  $I_6 \rightarrow T \rightarrow T$ .

State  $I_7 \rightarrow E \rightarrow E + T$ .

### (5) $\Psi = (a+b)*c$

i. Understand the "operator's" precedence.

( ) → higher precedence.  
+ → higher precedence.  
\* → lowest precedence.

ex: void fo

?  
int r

?  
i) calls by Re

### Follow

#### (6) First

Set of terminals that can appear immediately after a non-terminal at the beginning of any string derived from non-terminal.

determines which terminal to choose  
when parsing.

Symbol:  $\text{First}(X)$

$\text{First}(\epsilon) \rightarrow \epsilon$

Ex:-  $A \rightarrow \alpha\beta$   
 $\text{First}(B) \rightarrow \text{First}(\beta\epsilon) \cup \text{Follow}$

$\{\text{Follow}(X)\}$

$\{\text{Follow}(start)\} \rightarrow \$$

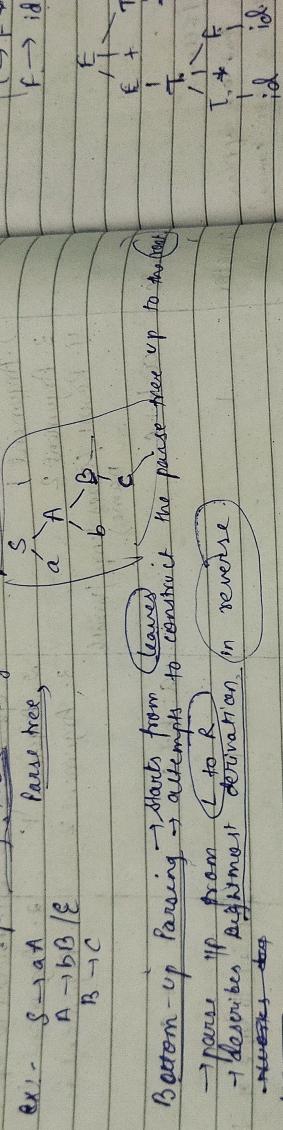
i)  $A \rightarrow \alpha\beta$   
 $\text{Follow}(B) \rightarrow \text{First}(\beta\epsilon) \cup \text{Follow}$

Top Down → Predictive Parsing  
 Non-deterministic  
 Bottom Up → LR(0) Parsing

$S \rightarrow A$   
 $A \rightarrow B$

- (1) Top Down Parsing → starts from root of the parse tree to the leaves  
 → parse if it attempts to construct the tree down  
 → derives the leftmost derivation of the string  
 ↓  
 (2) Recursive descent parsing → involves recursion procedure for each  
 → attempts to parse non-terminal in the grammar  
 → part of the grammar has to be or ambiguous

- 2) Predictive Parsing → A: non recursive version of top down  
 parsing that uses lookahead symbol to decide which prod. to use  
 Reg. gramm. to be LLL(1)  
 • uses a priority table to guide parsing process.



Types

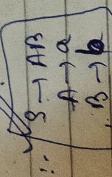
- 1) Shift-Reduce Parsing → uses a stack to store symbols and a pointer to consider  
 → alternates between shift (push LR to stack) and reduce (replace stack symbols with non-terminal based on prod. rule)  
 ↗ (top to right)

2) LR Parsing →

- LR(0)'s → basic version with no lookahead.  
 SLR(1) → Single LR → use lookahead to improve decision making

- CLR(1) → (Lookahead LR) → a more memory-efficient variant of SLR.

CLR(1) → Canonical LR → most powerful but complex



1. Shift a onto the stack.
2. Reduce a to A.
3. Shift b onto the stack.
4. Reduce b to B.
5. Reduce A to S.

Now,  $E \rightarrow T$   
 $E \rightarrow E + T$   
 $E \rightarrow T + T$   
 $F \rightarrow T F$   
 $F \rightarrow id$

$T \rightarrow i \theta$ .

$F \rightarrow i \theta$ .

$E \rightarrow E T$

$F \rightarrow T F$

$T \rightarrow F F$

$F \rightarrow id$

$E \rightarrow E + T$

$F \rightarrow T + T$

$T \rightarrow i \theta$

$F \rightarrow id$

$E \rightarrow E T$

$F \rightarrow T F$

$T \rightarrow F F$

$F \rightarrow id$

$E \rightarrow E + T$

$F \rightarrow T + T$

$T \rightarrow i \theta$

$F \rightarrow id$

$E \rightarrow E T$

$F \rightarrow T F$

$T \rightarrow F F$

$F \rightarrow id$

$E \rightarrow E + T$

$F \rightarrow T + T$

$T \rightarrow i \theta$

$F \rightarrow id$

|                         |       |                         |       |
|-------------------------|-------|-------------------------|-------|
| $S \rightarrow S, C$    | $T_0$ | $S \rightarrow S, C$    | $T_1$ |
| $C \rightarrow CC$      | $T_0$ | $C \rightarrow CC$      | $T_1$ |
| $C \rightarrow d$       | $T_0$ | $C \rightarrow d$       | $T_1$ |
| $C \rightarrow .d, \$$  | $T_0$ | $C \rightarrow .d, \$$  | $T_1$ |
| $S \rightarrow .CC, \$$ | $T_2$ | $S \rightarrow .CC, \$$ | $T_2$ |
| $C \rightarrow .CC, \$$ | $T_2$ | $C \rightarrow .CC, \$$ | $T_2$ |
| $C \rightarrow .d, \$$  | $T_2$ | $C \rightarrow .d, \$$  | $T_2$ |
| $C \rightarrow d, \$$   | $T_2$ | $C \rightarrow d, \$$   | $T_2$ |

Thus,  $T_0 := [S \rightarrow S, C]$

Compute  $GOTO(T_0, S) \rightarrow$  create  $T_1$

(More dot over  $S$  in  $S'$ )

Thus,  $T_1 := [S' \rightarrow S, C]$

Similarly,  $S \rightarrow C, C, \$$

$T_2 := [S \rightarrow C, C, \$]$

$C \rightarrow .CC, \$$

$C \rightarrow .d, \$$

$GOTO(T_0, C) \rightarrow$  create  $T_2$

$[C \rightarrow .CC, \$]$

$C \rightarrow .CC, \$$

$C \rightarrow .d, \$$

$GOTO(T_0, d) \rightarrow$  create  $T_4$

$T_4 := [d, \$]$

Final Table of LCL() pass : Items

|       |                                                                                                    |
|-------|----------------------------------------------------------------------------------------------------|
| $T_0$ | $S' \rightarrow S, \$$ , $S \rightarrow .CC, \$$ , $C \rightarrow .CC, \$$ , $C \rightarrow d, \$$ |
| $T_1$ | $S' \rightarrow S, \$$                                                                             |
| $T_2$ | $S \rightarrow C, C, \$$ , $C \rightarrow .CC, \$$ , $C \rightarrow d, \$$                         |
| $T_3$ | $C \rightarrow .CC, \$$ , $C \rightarrow .CC, \$$ , $C \rightarrow d, \$$                          |
| $T_4$ | $C \rightarrow d, \$$                                                                              |

- Why? compiler "sees" at a small part of the code
- ⇒ Peephole Optimization → focuses on improving small chunks rather than the entire program.
- simple & quick method used to improve small part of the code during compilation.
- compiler looks at a small "window" (or peephole) for of a few instructions at a time to improve them without changing what the code does.
- Goal → make the code run faster
- use less memory
  - reduce redundant unnecessary instruction

Top-down and Bottom-up parsing

Remove left recursion from

$$G: S \rightarrow Aa | b, A \rightarrow Ac | Sd$$

Formula:  $A \rightarrow \alpha \beta \gamma$

$$\begin{array}{l} A \rightarrow \beta \alpha \\ A^1 \rightarrow \alpha | \epsilon \end{array}$$

$$\text{Now, } S \rightarrow A^{\alpha} | b^{\beta} \quad \left\{ \begin{array}{l} \text{Given} \\ \text{Left Recur} \end{array} \right.$$

Removing Left Recur

$$\begin{array}{l} S \rightarrow bA^1 \\ A^1 \rightarrow CA^1 | \epsilon \end{array}$$

~~Recursion~~

$$\begin{array}{l} (x+y)^2 = x^2 + 2xy + y^2 \\ (a+bx)^2 = \\ (2+3)^2 = \end{array}$$

$$(x+y)^2 = x^2 + 2xy + y^2$$

$$(a+bx)^2 =$$

$$(2+3)^2 =$$

### 2) Syntax-Directed Translation (SDT)

→ SDT is like giving instructions for a recipe where each step tells how to combine ingredients ~~to make~~ to make something (SDP).

Similarly, SDT is a technique used in ~~SDT~~ to associate semantic rules with grammar of a progr. lang.

→ uses syntax directed definitions (SDD).

→ specifies rules like type checking, code generation, intermediate rep, construction.

### 3) Attribute

Info moves sideways from parent to child, "use this results from child to parent."

ex- A child tells parent "Here's my answer."

i) easier - only one direction (up)

ii) Bottom-up parsers

iii) Synthesized attributes

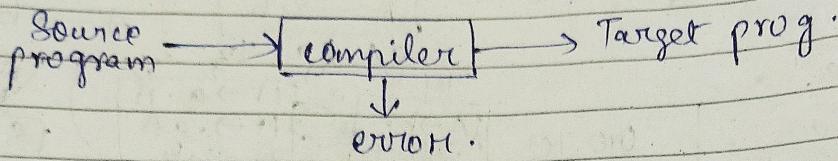
iv) Suitable for LR parsers

- i) A parent tells child, "use this info."
- ii) siblings share info.
- iii) A bit more complex - info can flow in multiple directions
- iv) Top-down parsers
- v) Generalized & inherited attributes.
- vi) LL parsers

## COMPILERS

What is a Compiler?

- A prog. that reads a prog. written in one lang. & translates it into an equivalent prog. in another lang.
- IIP lang. → Source lang. OIP lang. → Target lang.
- A compiler also reports error present in the source code program as a part of its translation process.



Types of compilers → Single Pass compiler.

Multi-Pass compiler.

The 2 parts of compilation process

Analysis

It takes the IIP source prog. & breaks it into parts. It then creates an intermed. represent. of source program.

Synthesis

It takes the intermediate represent. as the IIP. And creates the desired target program.

Cousins of a compiler

Source Prog.

Preprocessor  
directives

Pre PROCESSOR

#define  
#include

Modified Source Program.

COMPILER

MOV R1, R2  
JMP

Target Assembly Program

ASSEMBLER

Relocatable M/C code

In addn. to compiler, many other prog. are needed to create absolute m/c code.

PREPROCESSOR → macro expansion  
→ file inclusion

ASSEMBLER → takes assembly code that has been generated by compiler & converts it into relocatable m/c code.

LINKER / LOADER → takes the

library routines, relocatable object files, are linked main memory for execution

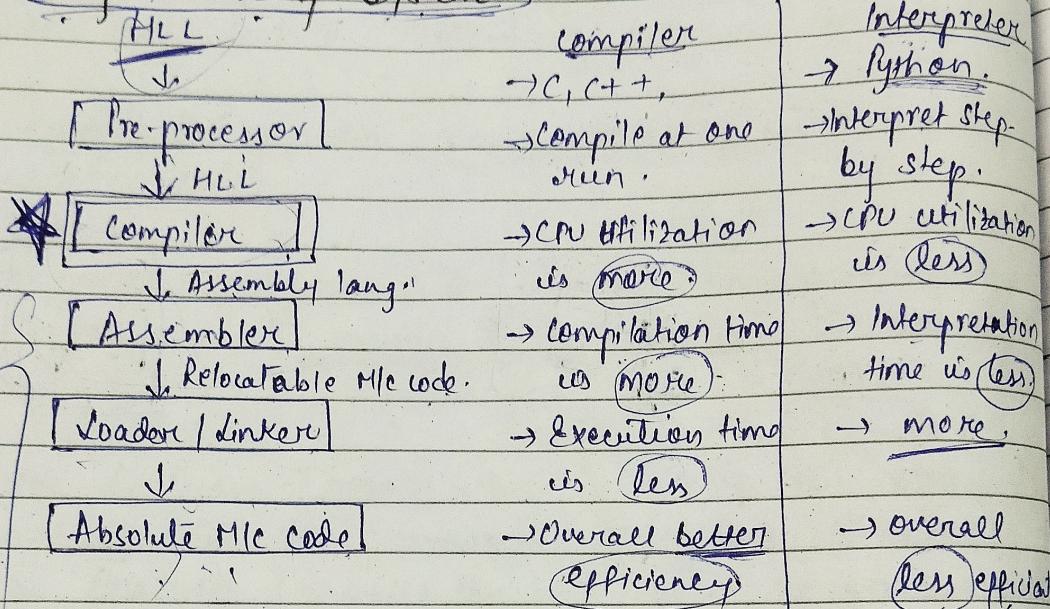
## Module 2

### Lexical Analysis

→ fundamental phase of the compilation process where the source code

28/12/22

### (1) Overview of Lang. Processing System



### (2) Phases of compiler.

