# *Compiler Module 1 Important Questions*

# 1. Describe input buffering scheme in the lexical analyzer.

**Input Bufferring**

- Input buffering in a lexical analyzer involves reading a block of characters from the source program into a buffer to improve efficiency.
- It employs two buffers of the same size N (Size of a disk block) to efficiently handle input. While one buffer is being processed, the other buffer is filled with characters from the input file
- Input buffering allows for lookahead, which is essential for identifying certain lexemes that require examining characters beyond the current position, such as identifiers and compound operators.
  - Examples of identifiers: `variableName`, `function_name`, `ClassName123`, etc.
  - Examples of compound operators: `+=`, `-=`, `<=`, `>=`, `&&`, `||`, etc.
- The lexical analyzer maintains two pointers:
  - `lexemeBegin`, marking the start of the current lexeme
  - `forward`, which scans ahead to identify lexemes.
- Once a lexeme is identified, the `forward` pointer is updated, and the lexeme is recorded as a token. The `lexemeBegin` pointer is then positioned to the character immediately following the lexeme.
  - Let's say we have the following line of code in a programming language:
    - `int variableName = 42;`
  - Here's how the input buffering scheme in a lexical analyzer would identify the lexeme "variableName":
    - **Initial State**:
      - The `forward` pointer and the `lexemeBegin` pointer are both pointing to the beginning of the line.
    - **Scanning**:
      - The lexical analyzer starts scanning characters from the beginning of the line.
      - It encounters the characters "i", "n", "t", which form the keyword "int". This is recognized as a token representing the data type.
        - The `forward` pointer moves to the character immediately after "t".
      - `lexemeBegin` pointer is typically updated after a complete lexeme has been identified and recorded as a token. It marks the beginning of the next lexeme to be analyzed in the input stream.
    - **Identifying the Identifier**:

- - - The lexical analyzer encounters the characters "v", "a", "r", "i", "a", "b", "l", "e", "N", "a", "m", "e", which form the identifier "variableName".
  - - As it scans these characters, it recognizes that they follow the rules for identifiers in the programming language.
  - - The `forward` pointer moves to the character immediately after "e".
  - - Once the lexeme "variableName" is identified, it is recorded as a token representing an identifier.
  - - The `lexemeBegin` pointer is positioned to the character immediately following "e".
  - The lexical analyzer continues scanning characters until the end of the line, identifying and recording tokens for other lexemes like the assignment operator "=", the integer literal "42", and the semicolon ";"
  - 
- When the `forward` pointer reaches the end of a buffer, the lexical analyzer reloads the other buffer from the input file and moves the `forward` pointer to the beginning of the newly loaded buffer.

## 2. Define a token. Identify the tokens in the expression a := b + 10.

Tokens

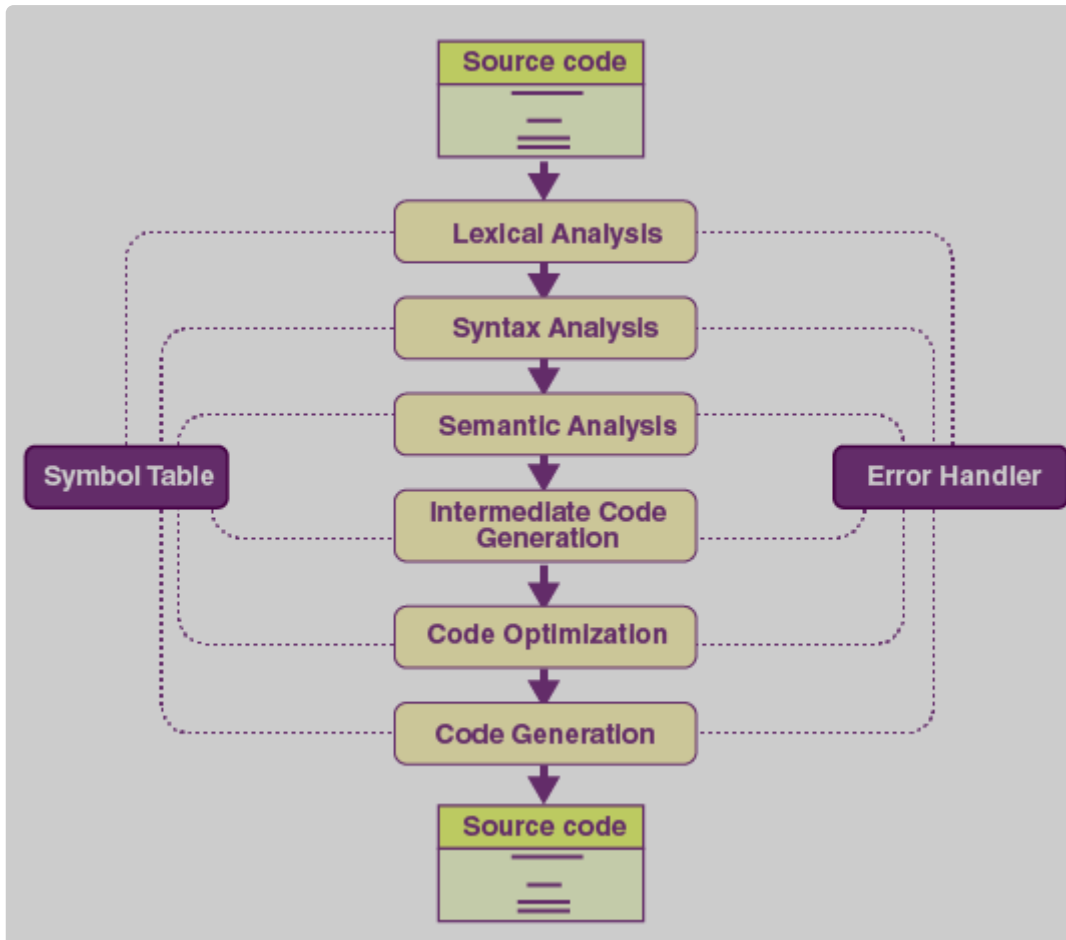- Tokens are basically a sequence of characters that are treated as a single unit as they cannot be further broken down
- Tokens are the building blocks of a program and are used to convey information to the compiler or interpreter
- Examples of tokens include
  - identifiers (user-defined names)
  - keywords (int, float, goto, continue, break)
  - operators (+,-,/, * )
  - literals
  - punctuation symbols,

- the expression `a := b + 10` consists of five tokens: two identifiers (`a` and `b`), an assignment operator (`:=`), an addition operator (`+`), and a numeric literal (`10`).

# 3. Construct a regular expression to denote a language L over = {0,1} accepting all strings of 0's and 1's that do not contain the substring 011.

# 4. Explain the different phases in the design of a compiler.



- Position := initial + rate * 60

**1. Lexical Analyzer**

- Dividing the statement to tokens
- We get the identifiers as position, initial, rate
- We create a symbol table

| **1. Position** | |
|---|---|
| 2. Initial | |

| 1. Position | |
| --- | --- |
| 3. rate | |

- We convert it to the identifier form
- id1 = id2 + id3 * 60

## 2. Syntax Analyzer

- Syntax analyzer converts to parse tree

```
          =
         / \
       id1  +
           / \
         id2  *
             / \
          id3   60
```

## 3. Semantic Analyzer

- Type checking and does modification

```
        ┌───┐
        │ = │
        └───┘
       ↙      ↘
  ┌─────┐   ┌───┐
  │ id1 │   │ + │
  └─────┘   └───┘
           ↙      ↘
      ┌─────┐   ┌───┐
      │ id2 │   │ * │
      └─────┘   └───┘
               ↙      ↘
          ┌─────┐   ┌─────────────┐
          │ id3 │   │ int_to_real │
          └─────┘   └─────────────┘
                          ↓
                      ┌──────┐
                      │  60  │
                      └──────┘
```

### 4. Intermediate code generation

- temp1 = int_to_real(60)
- temp2 = id3 * temp1
- temp3 = id2 + temp2
- id1 = temp3

### 5. Code Optimizer

- temp1 = id3 * 60.0
- id1 = id2 + temp1

### 6. Code generator

```
MOVF id3, R2
MULF #60.0, R2
MOVF id2,R1
ADDF R2,R1
MOVF R1,id1
```

```
; Move the value of id3 into register R2
MOVF id3, R2

; Multiply the value in R2 by 60.0
MULF #60.0, R2

; Move the value of id2 into register R1
MOVF id2, R1

; Add the value in R2 to the value in R1
ADDF R2, R1

; Move the result in R1 to id1
MOVF R1, id1
```

## 5. Explain any four compiler writing tools.

Compiler writers use software development tools and more specialized tools for implementing various phases of a compiler.

The tools include

- **Scanner Generators**:
  - Generate lexical analyzers (scanners) from a description of the tokens of a programming language using regular expressions.
  - Lexical analyzers produced by scanner generators recognize patterns in the input source code and generate tokens for further processing by the compiler.
  - Example LEX
- **Parser Generators**:
  - Parser generators, such as Yacc (Yet Another Compiler Compiler) automatically generate syntax analyzers (parsers) from a grammatical description of a programming language.
  - Syntax analyzers produced by parser generators parse the input source code according to the specified grammar rules, generating parse trees or abstract syntax

trees (ASTs) that represent the syntactic structure of the code.

- **Syntax-Directed Translation Engines**:
  - A Syntax-Directed Translation Engine is a tool used during the compilation process of a programming language. Its main job is to help convert code from one form to another.
  - set of instructions or rules that guide the translation process. These rules are closely tied to the grammar of the programming language being compiled.
  - Produce collections of routines for walking a parse tree and generating intermediate code.

- **Automatic code generators:**
  - This generator takes an intermediate code as input, and converts each operation of the intermediate code into the equivalent machine language.

- **Data flow analysis:**
  - Data-flow analysis is a technique used in compilers and programming language processing to gather information about how values (such as variables or expressions) flow or propagate through a program.
  - It helps analyze how data is used and manipulated within different parts of the program.
  - They track how values are defined, assigned, and used across different parts of the program. By doing so, they create a map or model of how data moves through the program's execution paths.
  - Data-flow analysis helps identify opportunities for optimizing code by eliminating redundant computations, identifying variables that can be safely removed or optimized, and improving memory usage.

- **Compiler construction toolkits:**
  - Compiler construction toolkits, such as GCC (GNU Compiler Collection), provide an integrated set of routines, libraries, and tools for constructing various phases of a compiler.
  - Provides an integrated set of routines for constructing various phases of a compiler.

# 6. Trace the output after each phase of the compiler for the assignment statement: a = b + c * 10, if variables given are of float type.

- Refer Question Number 4, it uses example Position := initial + rate * 60

# 7. Regular Expressions

## What are regular expressions?

- Imagine you are looking for specific patterns inside text.
- Regular expressions (regex) is a special tool that helps you define exactly what you're looking for.
  - Think of it like a recipe for finding text. You can use letters, numbers, and symbols to build your instructions.
  - For example,the pattern "letter (letter|digit) * " means it starts with a letter and then has zero or more letters or digits.
- Some common rules
  - The empty string (ε) is a pattern that matches nothing.
  - A single character like "a" matches itself.
  - You can combine patterns with "|" (OR) to find one or another pattern. For instance, "a|b" matches either "a" or "b"
  - Putting patterns together without any symbol means matching them one after another. So "ab" finds the letter "a" followed by "b".
  - The asterisk after a pattern means you can match it zero or more times. a * matches any number of "a"s,
  - Round brackets can be used to group parts of the pattern.
- Examples
  - Write regular expression for language accepting all strings containing any number of a's and b's
    - The answer is
      - (a+b)*
  - Write regular expression for the language accepting all the string which are starting with 1 and ending with 0 over {0,1}
    - The answer is
      - 1(0+1) * 0
      - Here We start with 1, so one is put first
      - After 1 we can have any combination of 0s and 1s so we do (0+1) *
      - At the end we need 0 so putting 0 at the end

- Write regular rexpression for the language L over {0,1} such that all the string do not contain the substring 01
  - Our language is
    - L = [ε, 0, 1, 00, 11, 10, 100, ...]
  - The expression is
    - (1 *0* )
- Design Regular expression to accept all possible combination of a's and b's

# 8. Write a note on bootstrapping.

Bootstrapping in the context of compiler design refers to the process of using a simpler language or tool to create a more complicated one, which can then be used to handle even more complex tasks

- Suppose we want to write a compiler for a new programming language called X. We need to consider three languages involved in this process:
  - **Source Language (X)**: The language for which we are writing the compiler.
  - **Object Language (or Target Language) (Z)**: The language that the compiler will produce code for.
  - **Implementation Language (Y)**: The language in which the compiler itself is written.

---

# 9. What is the relevance of input buffering in lexical analysis?

Input buffering in a lexical analyzer involves reading a block of characters from the source program into a buffer to improve efficiency.

- It employs two buffers of the same size N (Size of a disk block) to efficiently handle input. While one buffer is being processed, the other buffer is filled with characters from the input file
- Input buffering allows for lookahead, which is essential for identifying certain lexemes that require examining characters beyond the current position, such as identifiers and compound operators.
  - Examples of identifiers: `variableName`, `function_name`, `ClassName123`, etc.
  - Examples of compound operators: `+=`, `-=`, `<=`, `>=`, `&&`, `||`, etc.

- The lexical analyzer maintains two pointers:
    - `lexemeBegin`, marking the start of the current lexeme
    - `forward`, which scans ahead to identify lexemes.
- Once a lexeme is identified, the `forward` pointer is updated, and the lexeme is recorded as a token. The `lexemeBegin` pointer is then positioned to the character immediately following the lexeme.
    - Let's say we have the following line of code in a programming language:
        - `int variableName = 42;`
    - Here's how the input buffering scheme in a lexical analyzer would identify the lexeme "variableName":
        - **Initial State**:
            - The `forward` pointer and the `lexemeBegin` pointer are both pointing to the beginning of the line.
        - **Scanning**:
            - The lexical analyzer starts scanning characters from the beginning of the line.
            - It encounters the characters "i", "n", "t", which form the keyword "int". This is recognized as a token representing the data type.
            -
                - The `forward` pointer moves to the character immediately after "t".
            - `lexemeBegin` pointer is typically updated after a complete lexeme has been identified and recorded as a token. It marks the beginning of the next lexeme to be analyzed in the input stream.
        - **Identifying the Identifier**:
            -
                - The lexical analyzer encounters the characters "v", "a", "r", "i", "a", "b", "l", "e", "N", "a", "m", "e", which form the identifier "variableName".
            -
                - As it scans these characters, it recognizes that they follow the rules for identifiers in the programming language.
            -
                - The `forward` pointer moves to the character immediately after "e".
            -
                - Once the lexeme "variableName" is identified, it is recorded as a token representing an identifier.

- - The `lexemeBegin` pointer is positioned to the character immediately following "e".
  - The lexical analyzer continues scanning characters until the end of the line, identifying and recording tokens for other lexemes like the assignment operator "=", the integer literal "42", and the semicolon ";"
    - 
- When the `forward` pointer reaches the end of a buffer, the lexical analyzer reloads the other buffer from the input file and moves the `forward` pointer to the beginning of the newly loaded buffer.

---

# 10. With an example source language statement, explain tokens, lexemes, and patterns.

**Tokens**

- Tokens are basically a sequence of characters that are treated as a single unit as they cannot be further broken down
- Tokens are the building blocks of a program and are used to convey information to the compiler or interpreter
- Examples of tokens include
  - identifiers (user-defined names)
  - keywords (int, float, goto, continue, break)
  - operators (+,-,/, * )
  - literals
  - punctuation symbols,
- In the statement `int num = 10;` , the tokens are:
  - Keyword token: `int`
  - Identifier token: `num`
  - Assignment operator token: `=`
  - Integer literal token: `10`
  - Semicolon token: `;`

**Lexemes**

- A lexeme is a sequence of characters in the source program that matches the pattern for a token
- Lexemes are the actual occurrences of the tokens in the source code.

Example:

- In the statement `int num = 10;`, the lexemes are:
  - `int` : Keyword lexeme
  - `num` : Identifier lexeme
  - `=` : Assignment operator lexeme
  - `10` : Integer literal lexeme
  - `;` : Semicolon lexeme

**Pattern**

- A pattern is a set of rules that the scanner or the lexical analyzer follows to create a token
- For example, in case of keywords, the pattern is just the sequence of characters that form the keyword.\
- Regular expression pattern for integer literals in a programming language: `[0-9]+`
  - - This pattern specifies that an integer literal lexeme consists of one or more digits (0-9) in sequence.

---

# 11. Scanning of source code in compilers can be sped up using input buffering. Explain.

- Input buffering is a crucial technique in compiler design that significantly speeds up the process of scanning source code.
- The Problem:
  - Traditionally, the lexical analyzer (scanner) in a compiler reads the source code one character at a time
  - This can be slow and inefficient, especially for large programs.
  - Every time the scanner needs to check for a pattern (like a keyword or identifier), it has to access the disk or file system to get the next character.
  - This constant back-and-forth between memory and storage creates overhead and slows down the process.

- The Solution: Input Buffering
  - Input buffering solves this issue by reading the source code in larger chunks (blocks) into a designated memory area called a buffer.
  - This buffer acts as a temporary storage for the source code characters.
  - The scanner then processes the characters within the buffer instead of constantly accessing the source file.
- Benefits of using Input Bufferring
  - Reduced I/O operations:
    - By reading the source code in larger chunks,the scanner minimizes the number of times it needs to access the disk or file system
  - Improved processing efficiency:
    - The scanner can efficiently scan through the characters in the buffer without the delays involved in reading individual characters.
  - Faster scanning
    - As a result,input buffering leads to a noticeable improvement in the overall speed of the scanning phase in a compiler.