

COMPILER DESIGN PROJECT



Visualization of Phases of Compiler in JAVA

Submitted by:

Subham Gupta (13103456)
Aryan Jadon (13103734)
Burhanuddhin (13103735)
Prachi Tiwari (13103742)

Submitted to:

Ms. Aayushee Gupta

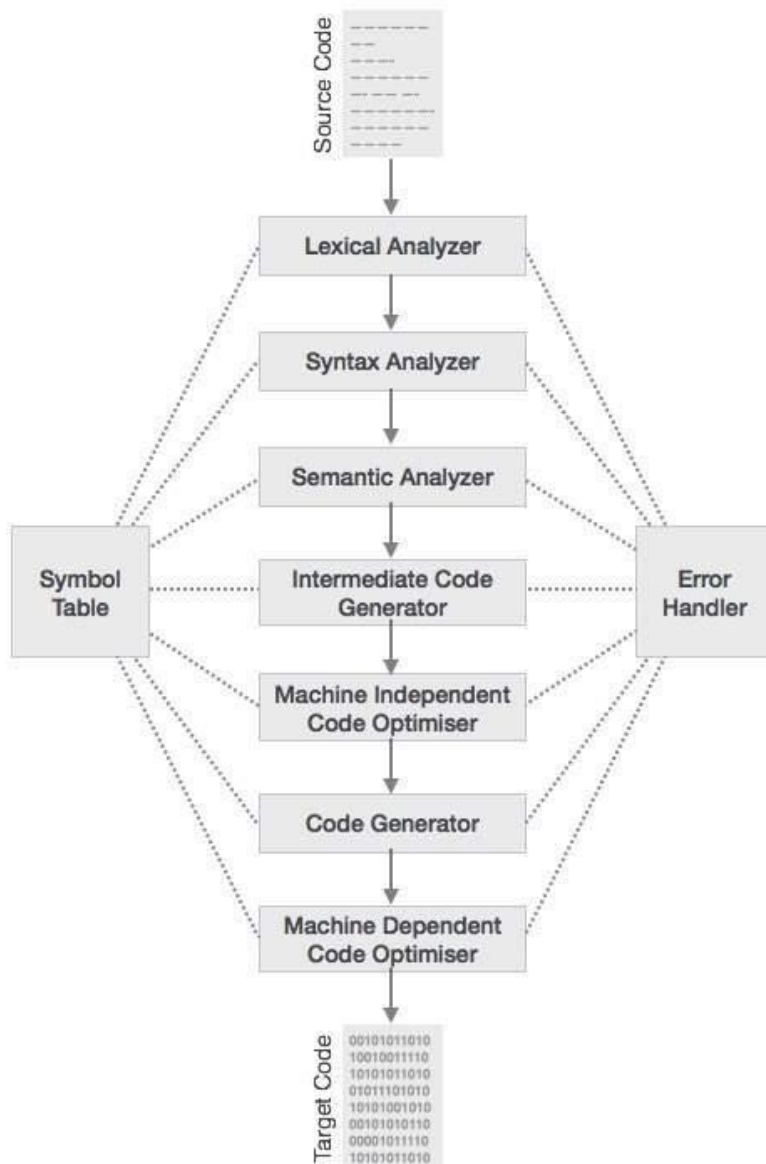
Batch- B10

ABSTRACT

A compiler translates the code written in one language to some other language without changing the meaning of the program. It is also expected that a compiler should make the target code efficient and optimized in terms of time and space.

Compiler design principles provide an in-depth view of translation and optimization process. Compiler design covers basic translation mechanism and error detection & recovery. It includes lexical, syntax, and semantic analysis as front end, and code generation and optimization as back-end.

The compilation process is a sequence of various phases. Each phase takes input from its previous stage, has its own representation of source program, and feeds its output to the next phase of the compiler. Let us understand the phases of a compiler.



Lexical Analysis

The first phase of scanner works as a text scanner. This phase scans the source code as a stream of characters and converts it into meaningful lexemes. Lexical analyzer represents these lexemes in the form of tokens as:

```
<token-name, attribute-value>
```

Syntax Analysis

The next phase is called the syntax analysis or **parsing**. It takes the token produced by lexical analysis as input and generates a parse tree (or syntax tree). In this phase, token arrangements are checked against the source code grammar, i.e. the parser checks if the expression made by the tokens is syntactically correct.

Semantic Analysis

Semantic analysis checks whether the parse tree constructed follows the rules of language. For example, assignment of values is between compatible data types, and adding string to an integer. Also, the semantic analyzer keeps track of identifiers, their types and expressions; whether identifiers are declared before use or not etc. The semantic analyzer produces an annotated syntax tree as an output.

Intermediate Code Generation

After semantic analysis the compiler generates an intermediate code of the source code for the target machine. It represents a program for some abstract machine. It is in between the high-level language and the machine language. This intermediate code should be generated in such a way that it makes it easier to be translated into the target machine code.

Code Optimization

The next phase does code optimization of the intermediate code. Optimization can be assumed as something that removes unnecessary code lines, and arranges the sequence of statements in order to speed up the program execution without wasting resources (CPU, memory).

Code Generation

In this phase, the code generator takes the optimized representation of the intermediate code and maps it to the target machine language. The code generator translates the intermediate code into a sequence of (generally) re-locatable machine code. Sequence of instructions of machine code performs the task as the intermediate code would do.

It is a data-structure maintained throughout all the phases of a compiler. All the identifier's names along with their types are stored here. The symbol table makes it easier for the compiler to quickly search the identifier record and retrieve it. The symbol table is also used for scope management.

```

graph LR
    START((START)) -- "white space" --> START
    START -- "digit" --> INNUM((INNUM))
    START -- "letter" --> INID((INID))
    START -- "[other]" --> INASSIGN((INASSIGN))
    INCOMMENT((INCOMMENT)) -- "}" --> START
    INCOMMENT -- "{" --> START
    INCOMMENT -- "other" --> INCOMMENT
    INNUM -- "digit" --> INNUM
    INNUM -- "[other]" --> DONE(((DONE)))
    INID -- "letter" --> INID
    INID -- "[other]" --> DONE
    INASSIGN -- "=" --> INASSIGN
    INASSIGN -- "[other]" --> DONE
    INASSIGN -- "=" --> DONE
  
```

```

program  $\rightarrow$  stmt-seq
stmt-seq  $\rightarrow$  stmt stmt-seq`
stmt-seq`  $\rightarrow$  ; stmt-seq`
      stmt-seq`  $\rightarrow$   $\epsilon$ 
stmt  $\rightarrow$  if-stmt
stmt  $\rightarrow$  repeat-stmt
stmt  $\rightarrow$  assign-stmt
stmt  $\rightarrow$  read-stmt
stmt  $\rightarrow$  write-stmt
if-stmt  $\rightarrow$  if exp then stmt-seq
if-stmt` if-stmt`  $\rightarrow$  end

```

$\text{if-stmt} \rightarrow \text{else stmt-seq end}$
 $\text{repeat-stmt} \rightarrow \text{repeat stmt-seq until exp}$
 $\text{assign-stmt} \rightarrow \text{identifier} := \text{exp}$
 $\text{read-stmt} \rightarrow \text{read identifier}$
 $\text{write-stmt} \rightarrow \text{write exp}$
 $\text{exp} \rightarrow \text{simple-exp exp'}$
 $\text{exp'} \rightarrow \text{comparison-op simple-exp exp'} \rightarrow \epsilon$
 $\text{comparison-op} \rightarrow <$
 $\text{comparison-op} \rightarrow =$
 $\text{simple-exp} \rightarrow \text{term simple-exp'}$
 $\text{simple-exp'} \rightarrow \text{addop term}$
 simple-exp'
 $\text{simple-exp'} \rightarrow \epsilon$
 $\text{addop} \rightarrow +$
 $\text{addop} \rightarrow -$
 $\text{term} \rightarrow \text{factor term'}$
 $\text{term'} \rightarrow \text{mulop factor term'}$
 $\text{term'} \rightarrow \epsilon$
 $\text{mulop} \rightarrow *$ $\text{mulop} \rightarrow /$
 $\text{factor} \rightarrow (\text{exp})$
 $\text{factor} \rightarrow \text{number}$
 $\text{factor} \rightarrow \text{identifier}$

First Set of Non-Terminals:

$\text{first}(\text{program}) = \{\text{if, repeat, identifier, read, write}\}$ $\text{first}(\text{stmt-seq}) = \{\text{if, repeat, identifier, read, write}\}$

$\text{first}(\text{stmt-seq}^{\backslash}) = \{;, \epsilon\}$

$\text{first}(\text{stmt}) = \{\text{if, repeat, identifier, read, write}\}$ $\text{first}(\text{if-stmt}) = \{\text{if}\}$

$\text{first}(\text{if-stmt}^{\backslash}) = \{\text{end, else}\}$ $\text{first}(\text{repeat-stmt}) = \{\text{repeat}\}$ $\text{first}(\text{assign-stmt}) = \{\text{identifier}\}$ $\text{first}(\text{read-stmt}) = \{\text{read}\}$ $\text{first}(\text{write-stmt}) = \{\text{write}\}$

$\text{first}(\text{exp}) = \{(\text{, number, identifier}\}$ $\text{first}(\text{exp}^{\backslash}) = \{<, =, \epsilon\}$

$\text{first}(\text{comparison-op}) = \{<, =\}$

$\text{first}(\text{simple-exp}) = \{(\text{, number, identifier}\}$ $\text{first}(\text{simple-exp}^{\backslash}) = \{+, -, \epsilon\}$

$\text{first}(\text{addop}) = \{+, -\}$

$\text{first}(\text{term}) = \{(\text{, number, identifier}\}$ $\text{first}(\text{term}^{\backslash}) = \{*, /, \epsilon\}$

$\text{first}(\text{mulop}) = \{*, /\}$

$\text{first}(\text{factor}) = \{(\text{, number, identifier}\}$

Follow Set of Non-Terminals:

$\text{follow}(\text{program}) = \{\$ \}$ $\text{follow}(\text{stmt-seq}) = \{\$, \text{end, else, until}\}$ $\text{follow}(\text{stmt-seq}^{\backslash}) = \{\$, \text{end, else, until}\}$ $\text{follow}(\text{stmt}) = \{;, \$, \text{end, else, until}\}$ $\text{follow}(\text{if-stmt}) = \{;, \$, \text{end, else, until}\}$ $\text{follow}(\text{if-stmt}^{\backslash}) = \{;, \$, \text{end, else, until}\}$

$\text{follow}(\text{repeat-stmt}) = \{;, \$, \text{end, else, until}\}$ $\text{follow}(\text{assign-stmt}) = \{;, \$, \text{end, else, until}\}$ $\text{follow}(\text{read-stmt}) = \{;, \$, \text{end, else, until}\}$ $\text{follow}(\text{write-stmt}) = \{;, \$, \text{end, else, until}\}$ $\text{follow}(\text{exp}) = \{\text{then, }, \$, \text{end, else, until,)}\}$ $\text{follow}(\text{exp}^{\backslash}) = \{\text{then, }, \$, \text{end, else, until,)}\}$ $\text{follow}(\text{comparison-op}) = \{(\text{, number, identifier}\}$

$\text{follow}(\text{simple-exp}) = \{<, =, \text{then, }, \$, \text{end, else, until,)}\}$ $\text{follow}(\text{simple-exp}^{\backslash}) = \{<, =, \text{then, }, \$, \text{end, else, until,)}\}$ $\text{follow}(\text{addop}) = \{(\text{, number, identifier}\}$

$\text{follow}(\text{term}) = \{+, -, <, =, \text{then, }, \$, \text{end, else, until,)}\}$ $\text{follow}(\text{term}^{\backslash}) = \{+, -, <, =, \text{then, }, \$, \text{end, else, until,)}\}$ $\text{follow}(\text{mulop}) = \{(\text{, number, identifier}\}$

$\text{follow}(\text{factor}) = \{*, /, +, -, <, =, \text{then, }, \$, \text{end, else, until,)}\}$

Implementation:

I. Tokens are divided as follows:

0. Number

1. ID

2. Left Parenthesis

3. Right Parenthesis

4. Plus

5. Minus

6. Divide

7. Multiplication

8. Equal

9. Assignment

10. Error

11. SEMICOLON

12. *SMALLERTHAN*
13. *GREATER THAN*
20. *IF*
21. *THEN*
22. *ELSE*
23. *END*
24. *REPEAT*
25. *UNTIL*
26. *READ*
27. *WRITE*

Then it prints "Scanning is complete" if the scanning is done successfully.

II. The Parser receives the tokens from the scanner using the function *getToken()*; reading from the file "*test.txt*".

According to the grammatical rules in EBNF:

```

program -> stmt-sequence
stmt-sequence -> stmt { ; stmt }
stmt -> if-stmt | repeat-stmt | assign-stmt | read-stmt | write-stmt
if-stmt -> if exp then stmt-sequence [ else stmt-sequence ] end
repeat-stmt -> repeat stmt-sequence until exp
assign-stmt -> identifier := exp
read-stmt -> read identifier
write-stmt -> write exp
exp -> simple-exp [ comparison-op simple-exp ]
comparison-op -> < | =
simple-exp -> term { addop term }
addop -> + | -
term -> factor { mulop factor }
mulop -> * | /
factor -> (exp) | number | identifier

```

The Parser checks whether the input text matches the grammatical rules or not.

If yes, it prints "Parsing complete".

If not, it prints "Parsing error" with hints about the error place of the received token.

It also prints the matched token values.

III. Created an HTML Page that takes input as a source code and on pressing 'Start Button' generates Visualization Of Token generated. This is done using Javascript and d3(Document Data Viewer Library). The user is also given a feature to be able to control the speed of the scanning Phase of the Compiler.