

Course: Compiler Construction

Week 4: Compiler Design Phases – Syntax Analysis

Lecturer: Martha Gichuki

Course Description

- The course begins with a coverage of basic Principles in Compiler Design with an introduction to Formal programming language translation and compiler design concepts; compilers and interpreters.
- A coverage of Language theory, Parsing Context-Free Languages (Top - down and bottom - up parsing), translation specifications and machine-independent code optimization will then follow.

Course Description ...

- The main phases of compilation i.e. Lexical, **Syntax** and Semantic analysis will be taught.
- Code generation and Optimization, Symbol table design, Program compilation, Loading and execution will then be covered.
- Compilation techniques, Optimization, Design of a simple complete compiler. will culminate the course.

Learning outcomes Week 4: Compiler

Design Phases – Syntax Analysis

At the end of the lecture, you will be able to:

- (i) Define Parser, Grammar and Ambiguity
- (ii) Describe Context Free Grammar (CFG)
- (iii) Derive strings using parse trees

Syntax Analysis

- ✓ Another term for Syntax Analysis is **Parsing**
- ✓ This is the process of analyzing a sequence of tokens to determine their **grammatical structure** with respect to a **given formal grammar**¹.

Syntax Analysis intro...

- ✓ It is the most important phase of a compiler.
- ✓ A **syntax analyzer (parser)** checks for the correct syntax and builds a data structure (**parse tree**) implicit in the input tokens i.e. it considers the sequence of tokens for possible valid constructs of the programming language².

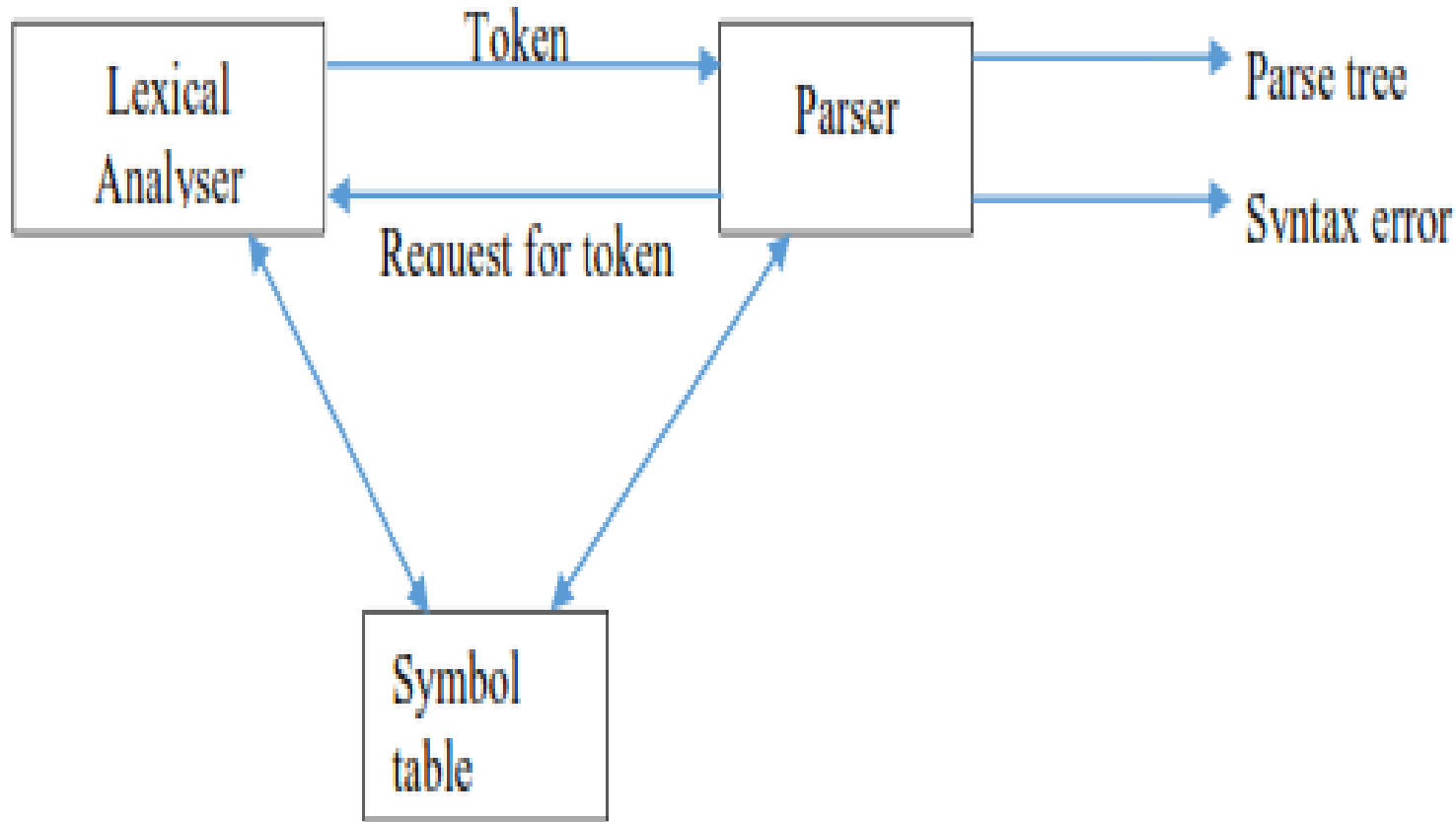
First Role of a Parser

1. To identify language constructs present in a given input program.
- ✓ If the parser determines input to be **valid**, it outputs the presentation of the input in form of a **parse tree**³

Second Role of a Parser

2. If the input is grammatically **incorrect**, the parser declares the detection of **syntax error** in the input.
- ✓ In this case, a parse tree is **not produced**⁴.

Parser Roles Illustration⁵



Grammar

A grammar G is defined as a **four-tuple** with $\langle V_N, V_T, P, S \rangle$ where:-

- I. V_N is the set of **non-terminal symbols** used to write the grammar⁶.

II. V_T is the set of terminals
(set of words of the language, lexicon or dictionary of words).

III. P is the set of production rules that define how a sequence of terminal and non-terminal symbols can be replaced by some other sequence⁷

IV. S: - $S \in V_N$ is a special non-terminal called the **start symbol** of the grammar.

- ✓ The language of the grammar $G = \langle V_N, V_T, P, S \rangle$ denoted by $L(G)$ is defined as all those strings **over** V_T that can be generated by starting with the start symbol S then applying the production rules in P until no more non-terminal symbols are present⁸.

Example: Consider the grammar to generate arithmetic expressions consisting of numbers and operator symbols i.e. +, -, *, /.

Rules of the grammar can be written as follows⁹:-

$$E \rightarrow EAE$$

$$E \rightarrow (E)$$

$$E \rightarrow \textit{number}$$

$$A \rightarrow +$$

$$A \rightarrow -$$

$$A \rightarrow *$$

$$A \rightarrow /$$

- ✓ We can apply these rules to derive the expression **2 * (3 + 5 * 4)** as follows:-

$$\begin{aligned} E &\rightarrow EAE \rightarrow EA(E) \rightarrow EA(EAE) \rightarrow \\ &EA(EAEAE) \rightarrow EA(EAEA 4) \rightarrow \\ &EA(EAE * 4) \rightarrow EA(EA 5 * 4) \rightarrow \\ &EA(E + 5 * 4) \rightarrow EA(3 + 5 * 4) \rightarrow E \\ &* (3 + 5 * 4) \rightarrow 2 * (3 + 5 * 4) \end{aligned}$$

- ✓ In the grammar, ***E*** and ***A*** are **non-terminals** while the rest are **terminals**¹⁰.

Context Free Grammar (CFG)

✓ This is grammar that **defines** Context Free Languages and consist of production rules in which:-

- i. the left hand side contains only **a single non-terminal and no terminals**;
- ii. the right hand side consists of **either terminals, non-terminals or both**.

✓ Note: Most programming language constructs belong to Context Free Languages¹¹

Notational Conventions¹²

The following is the notational convention that would be used when defining grammars:

1. The following will be taken as **terminals**:

- ✓ all operator symbols
- ✓ punctuation symbols including parenthesis ‘()’
- ✓ digits
- ✓ lower case letters of the alphabet such as a, b, c
- ✓ Lexemes such as id, number, while, etc.

2. The following will represent **non-terminals**:

- ✓ Upper case letters of the alphabet such as A, B, C.
- ✓ The letter S will represent the start symbol
- ✓ Lowercase names such as expr, stmt etc¹³.

3. A set of **productions** with the same left hand side such as: $A \rightarrow \alpha_1, A \rightarrow \alpha_2 \dots$
 $A \rightarrow \alpha_n$ will be written as
 $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$.
4. If no start symbol is represented, then the non-terminal appearing on the left hand side of the production rule will be considered as the **start symbol**¹⁴

Derivation

- ✓ This is the process of generating a sequence of intermediary strings to expand the start symbol of the grammar to the desired string of terminals.
- ✓ The **graphical depiction of a derivation** is a parse tree¹⁵.

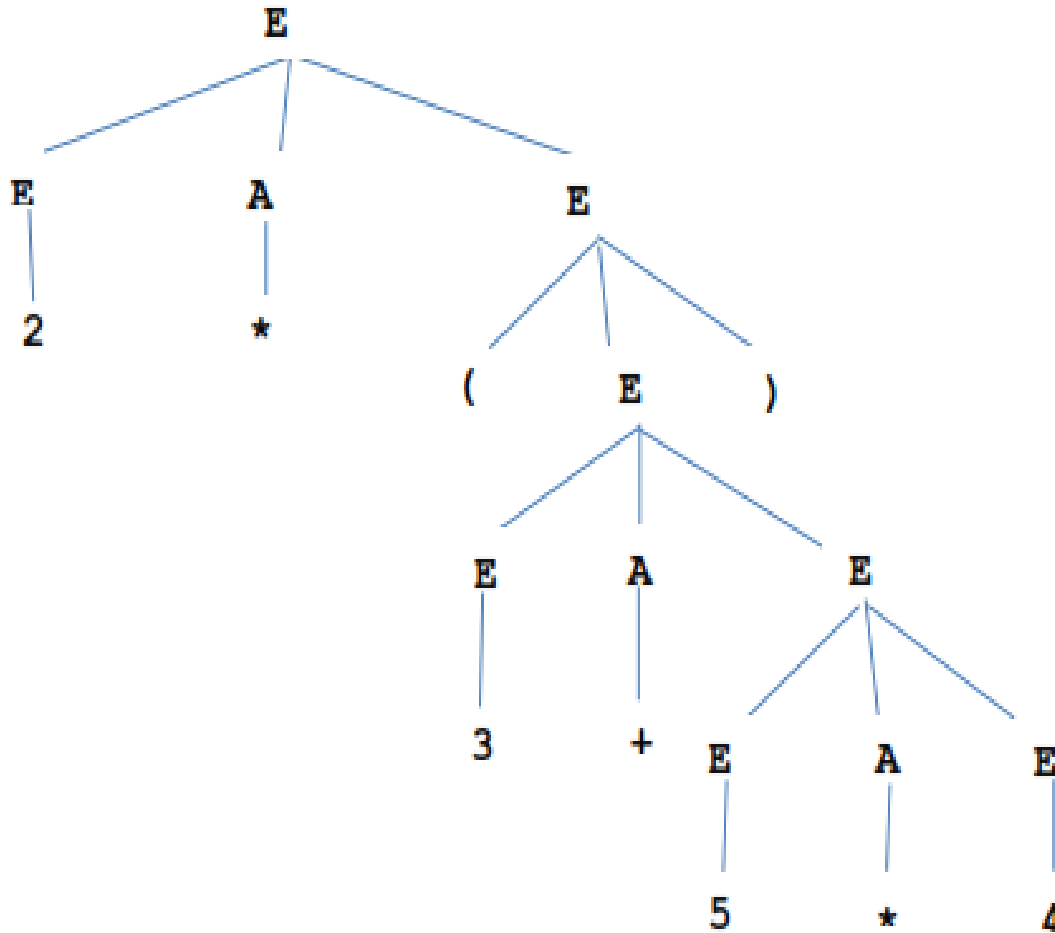
Derivation ...

- ✓ Each step of the derivation modifies an intermediary string to a new one by replacing a substring of it that matches the **left hand side** of the **production rule** by a string on **the right hand side of the rule**¹⁶.

In a parse tree

- ✓ The **start symbol** of the derivation becomes the **root**.
- ✓ All **leaf nodes** are **terminals**
- ✓ All **interior nodes** are **non-terminals**
- ✓ **In-order traversal** gives original **input string**
- ✓ A parse tree depicts **associativity and precedence of operators**¹⁷

Example: The following is the derivation of the string $2*(3+5*4)$ ¹⁸



Associativity & Precedence

A. Associativity

- ✓ If an operand has operators on both sides, left associativity dictates that the operand be taken by the left operator
- ✓ If the operation is right-associative, the right operator will take the operand¹⁹

B. Precedence

- ✓ To decrease the chances of ambiguity in a language or its grammar, **operands in brackets have the highest precedence**, followed by division, multiplication, addition and **subtraction comes last**²⁰

Types of Derivation

There are two types of derivation:-

i) **Left-Most Derivation**

- ✓ Here, the left-most non-terminal is replaced at each step.
- ✓ The intermediate strings are called **left-sentential forms**. They consist of grammar symbols (both terminal & non-terminals) ²¹

ii) Right-Most Derivation

- ✓ In this derivation the right-most non-terminal is replaced at each step.
- ✓ The intermediary strings are called **right-sentential forms**.
- ✓ Right-most derivation is often referred to as **canonical representation**²².

Ambiguity

- ✓ A grammar is said to be **ambiguous** if there exists **more than one parse tree for the same sentence**²³

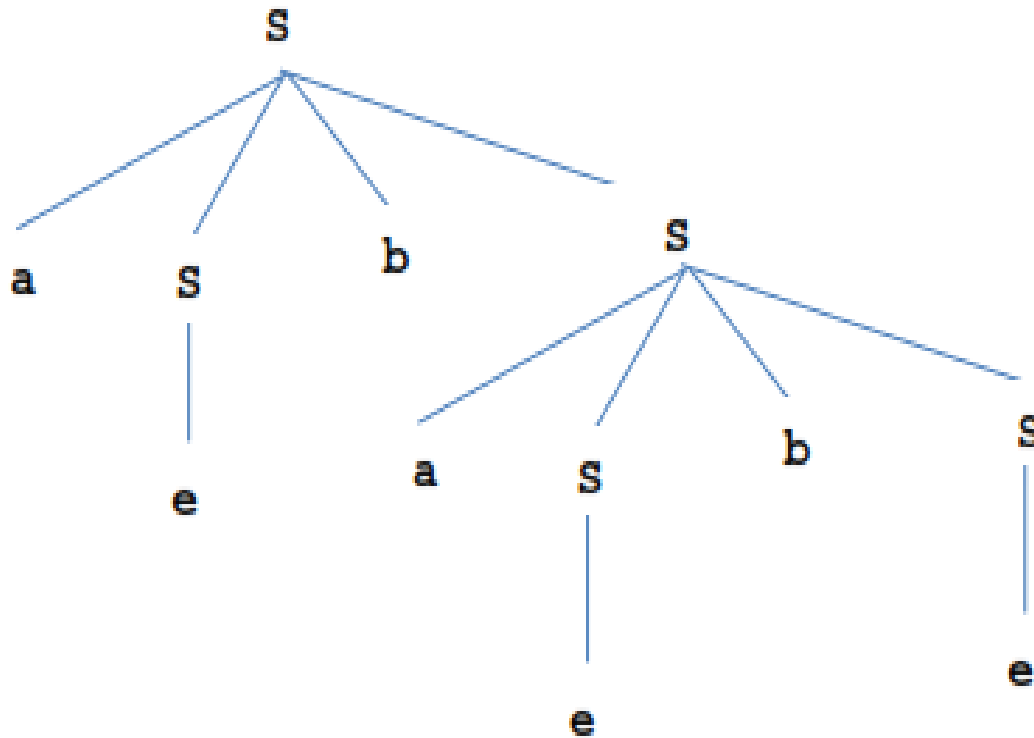
Example:- Consider the following grammar: $S \rightarrow aSbS \mid bSaS \mid e$

- ✓ To show that the grammar is ambiguous, we present two different parse trees for the string “*abab*”.

Consider the following
grammar: $S \rightarrow aSbS \mid bSaS \mid e$

1. First Parse Tree for the string “ $abab^{24}$ ”

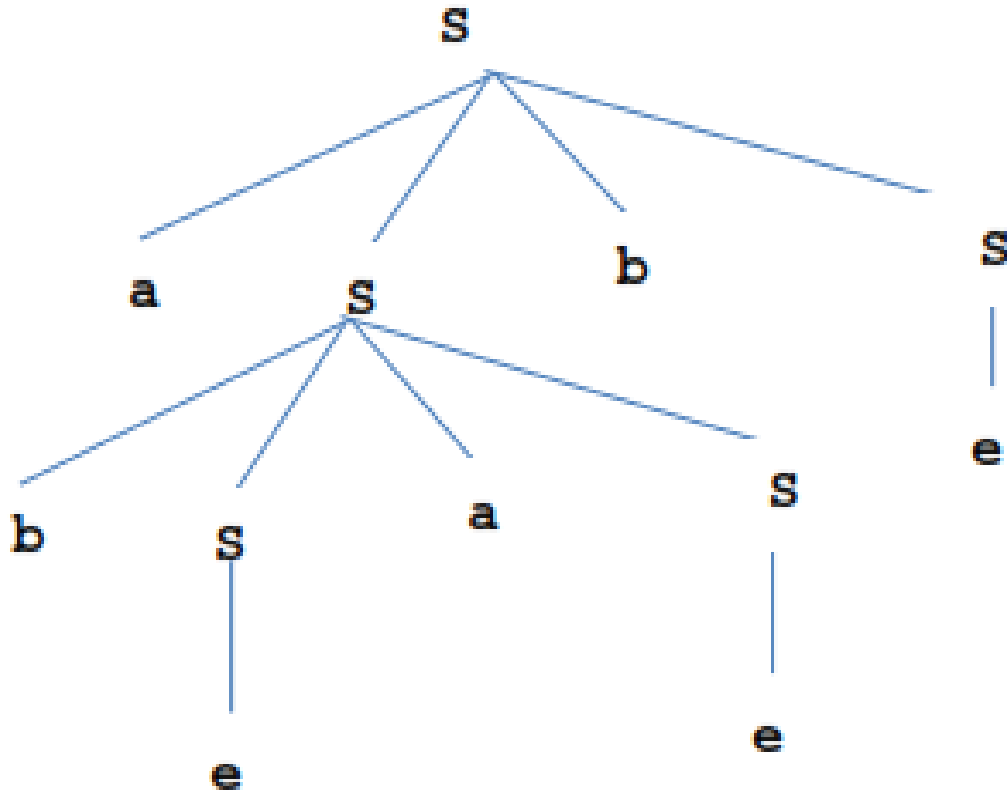
First parse tree



Consider the following grammar:

$S \rightarrow aSbS \mid bSaS \mid e$

2. Second Parse Tree for the string "abab²⁵"



Dangling Else Ambiguity

- ✓ Most programming languages have both if... then and if... then ... else versions of the statement.
- ✓ The grammar rules are as follows:-

stmt \rightarrow if condition then stmt else stmt

| if condition then stmt²⁶

Consider the following code segment

if a > b then

if c > d then x = y

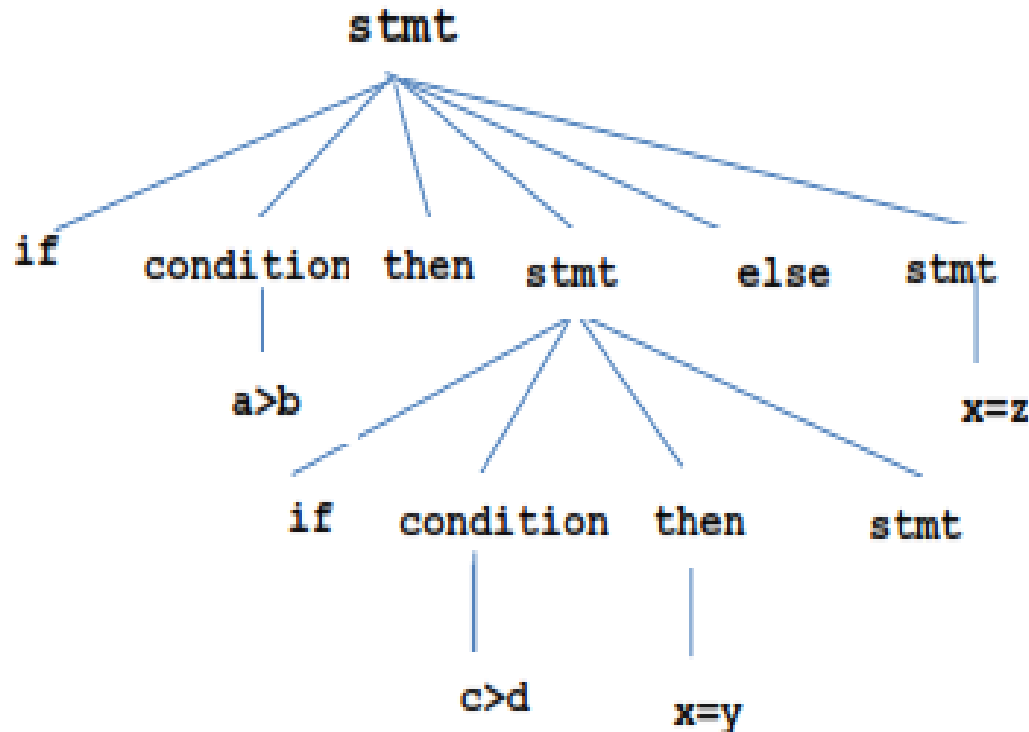
else x = z

✓ Two parse trees can be generated by the grammar as shown below²⁷:

i) the else is taken with the outer if statement²⁸. *if a > b then*

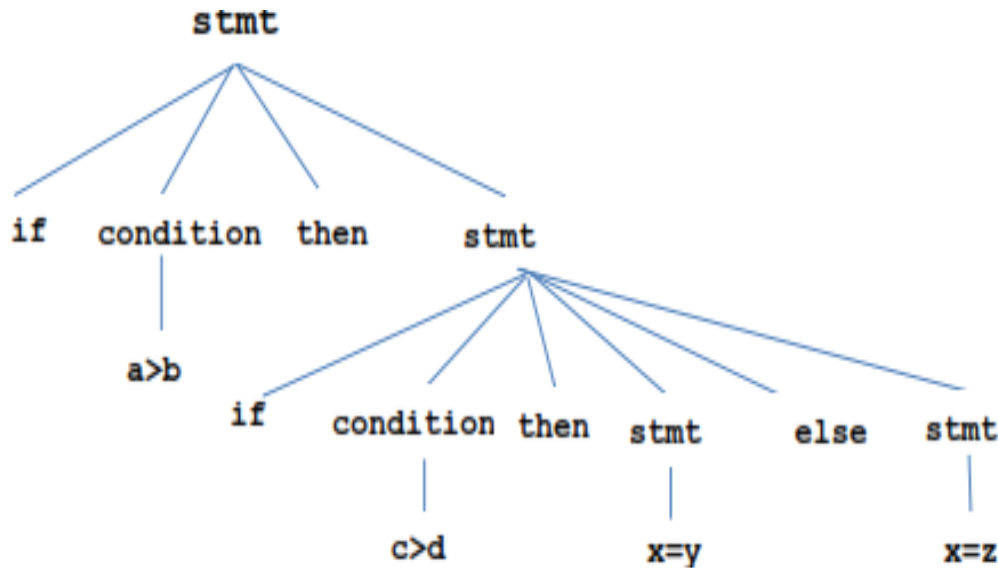
if c > d then x = y

else x = z



ii) the else is taken with the inner if statement²⁹. *if a > b then*

if c > d then x = y
else x = z



✓ **Note:** most programming languages accept this second parse tree as the correct syntax.

Eliminating Ambiguity³⁰

Ambiguity may be eliminated by:-

- i. Rewriting the grammar
- ii. Modifying the grammar e.g. Many programming languages require that “if” should have a matching “end if ” as shown below:-

stmt → *if condition then stmt else stmt*
endif

| *if condition then stmt endif*

Content Covered in Week 4: Compiler Design Phases – Syntax Analysis

At the end of the lecture, we were able to:

- i. Define Parsers, Grammar and Ambiguity
- ii. Describe Context Free Grammar (CFG)
- iii. Derive strings using parse trees

Course Text Books

1. Compilers: Principles, Techniques and Tools, Aho, Lam, Sethi, and Ullman; *Addison- Wesley Pub Co*, ISBN: 0201100886 (2007))
2. Compiler Construction: Principles and practices; Kenneth C. Loudon; *Cengage Learning*; 1st edition, ISBN-10 : 0534939724 (1997)
3. Basics of Compiler Design: Torben Mogensen; *DIKU University of Copenhagen Universitetsparken 1 DK-2100 Copenhagen DENMARK* (2007)
4. Engineering a Compiler: 2nd edition; Keith D. Cooper and Linda Torczon; *Morgan Kaufmann Publishers* ISBN: 978-0-12-088478-0 (2003)
5. Compiler Design: Santanu Chattopadhyay; *PHI Learning Publishers*, ISBN 812032725X, 9788120327252 (2005)