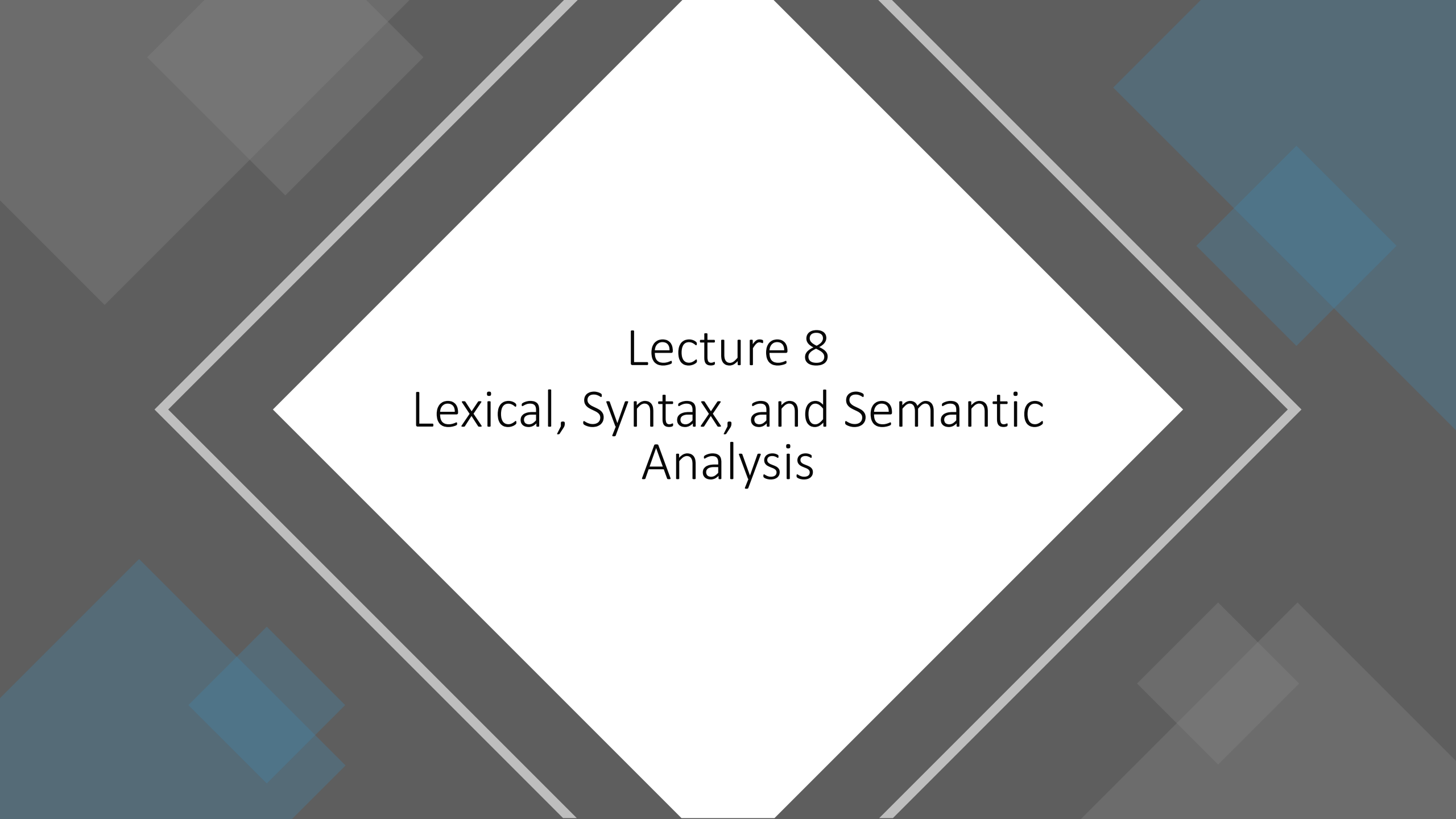CENG204 - Programming Languages Concepts

Asst. Prof. Dr. Emre ŞATIR

# Lecture 8
## Lexical, Syntax, and Semantic Analysis

# Lecture 8 Topics

- Introduction
- Lexical Analysis
- Syntax Analysis
- Semantic Analysis
- Code Optimization
- Code Generation

# Introduction

# Introduction

- **Lexical**, **syntax** and **semantic** analysis steps are also used in "Natural Language Processing (NLP)".

- For example, consider the sentence "John ate an apple." **The lexical analysis** <u>provides the words</u> ("John", "ate", "an", "apple").

- **The syntax rules** <u>define the structure of the sentence</u>, with the word "ate" serving as the verb.

- **Semantic analysis** helps to <u>determine the meaning of the sentence</u> by looking at the context of the words. In this case, the meaning is that John is eating an apple (past tense).
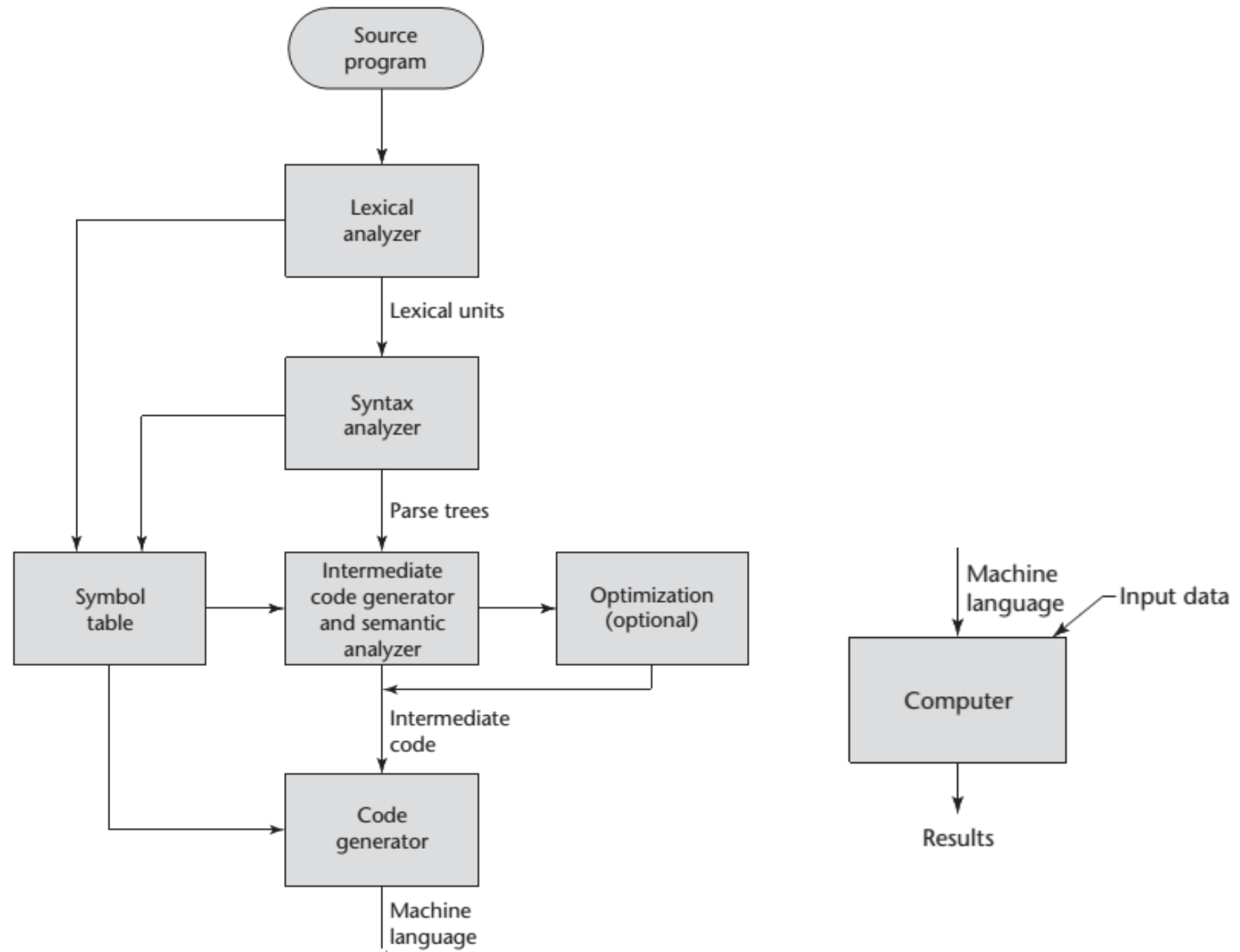
# Introduction

- Lexical, syntax and semantic analysis are <u>essential components</u> of natural language processing.

- Lexical analysis is the process of <u>breaking down a large text into smaller parts</u>, such as words, phrases or symbols, while syntax analysis is the process of <u>understanding how these parts fit together</u> to form meaningful sentences.

- Semantic analysis helps to <u>determine the meaning</u> of a sentence or phrase.

- By <u>combining these three components</u>, computers can understand natural language.
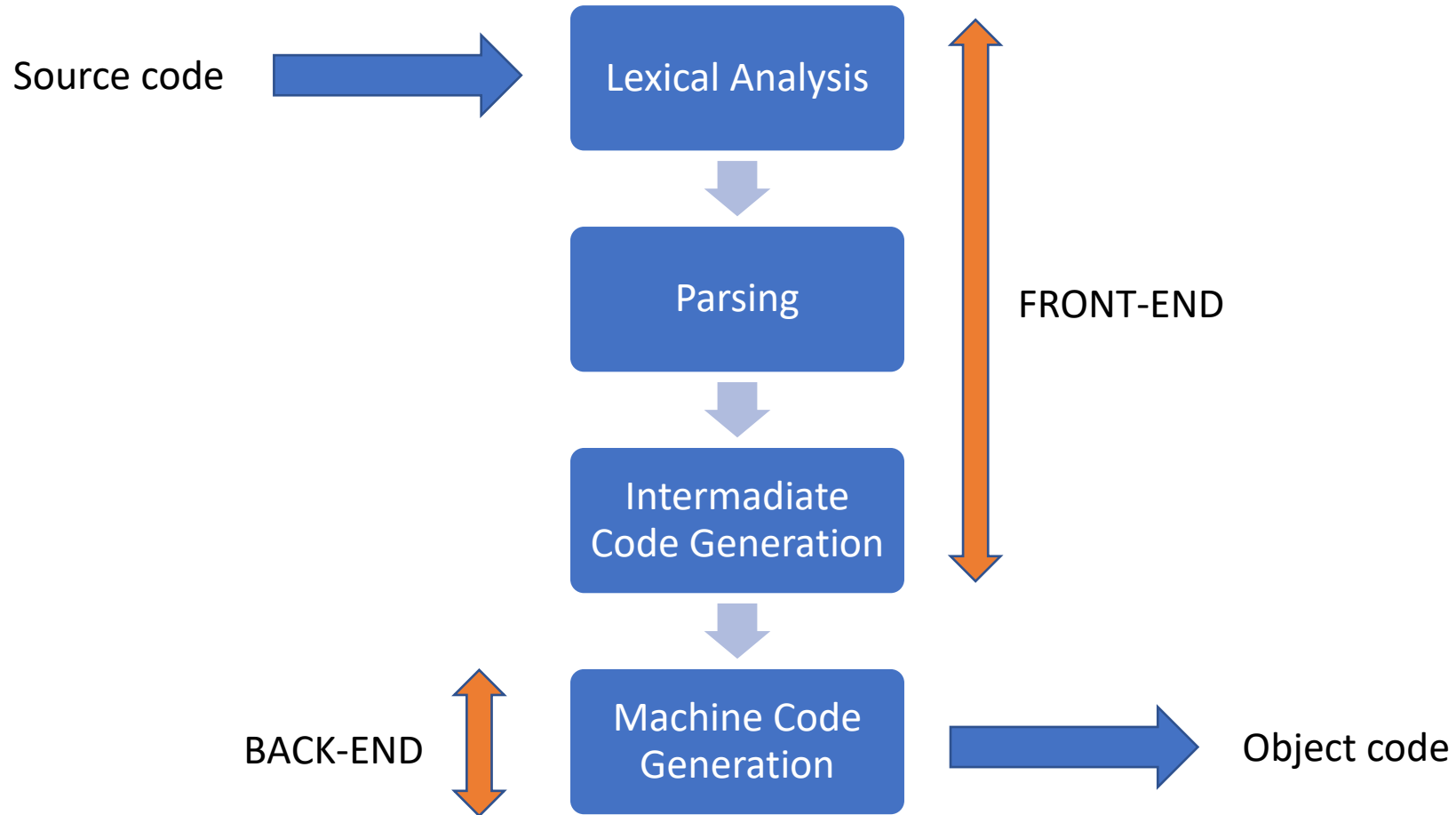
# Introduction

- In order for a source code <u>written in a high-level language</u> to be executed, it must be <u>translated into machine language</u>, the only language that the <u>computer directly recognizes</u>.

- The source code translated into machine language is called **object code**.

- <u>Three different approaches</u> for implementing programming languages were introduced before:
  - compilation
  - pure interpretation
  - hybrid implementation

- When a compiler converts source code to object code, it first examines whether this source code was built <u>according to the rules</u>. This job includes "**lexical**" and "**syntax**" analysis (All three of the implementation approaches just discussed use both lexical and syntax anayzers).

The
Compilation
Process

Source
program

↓

Lexical
analyzer

↓ Lexical units

Syntax
analyzer

↓ Parse trees

Symbol
table

Intermediate
code generator
and semantic
analyzer

→

Optimization
(optional)

↓ Intermediate
code

Code
generator

↓ Machine
language

Machine
language → Computer ← Input data

↓

Results

# The Simplified Compilation Process

# The Compilation Process

**FRONT-END:** In this section, the <u>syntax</u> and <u>sematics</u> of the program are analyzed according to the rules of the language, and an <u>intermediate code</u> is created. It is <u>independent</u> of the machine.
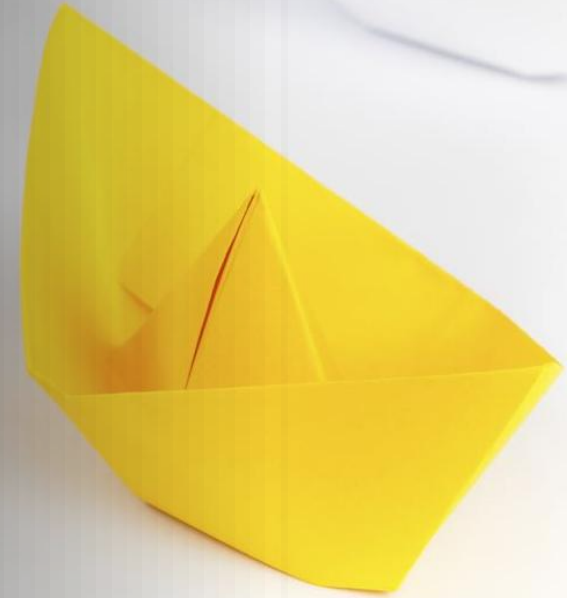
**BACK-END:** In this part, the <u>machine code is created</u> by applying the optimization (optional) from the intermediate code representation of the program. These parts are known as the **synthesis** step. It is <u>machine dependent</u>.

# Introduction

- Nearly all compilers separate the task of <u>analyzing syntax</u> into two distinct parts, <u>lexical analysis</u> and <u>syntax analysis</u> for simplicity and efficiency.

- The lexical analyzer deals with <u>small-scale</u> language constructs, such as <u>names</u> and <u>numeric literals</u> (mathematically, a finite automaton based on a regular grammar).

- The syntax analyzer deals with the <u>large-scale</u> constructs, such as <u>expressions</u>, <u>statements</u>, and <u>program units</u> (Syntax analyzers are also called as "**parsers**") (mathematically, a push-down automaton based on a context-free grammar, or BNF).

# Lexical Analysis

# Lexical Analysis

- A lexical analyzer is essentially a "***pattern matcher***" for character strings. It identifies and distinguishes the lowest level units of the source program (lexemes).

- A lexical analyzer serves as the <u>front end</u> of a <u>syntax analyzer</u>.

- Technically, lexical analysis is a <u>part of syntax analysis</u>.

- A lexical analyzer performs syntax analysis at the <u>lowest level</u> of program structure.

- An input program appears to a compiler as a <u>single string of characters</u>.

- The lexical analyzer collects characters <u>into logical groupings</u> and <u>assigns internal codes to the groupings</u> according to their structure.

# Lexical Analysis

- In chapters before, these logical groupings are named **lexemes**, and the internal codes for categories of these groupings are named **tokens**.

- Lexemes are recognized by matching the input character string against <u>character string patterns</u>.

- Although tokens are usually <u>represented as integer values</u>, for the sake of <u>readability</u> of lexical and syntax analyzers, they are often referenced through <u>named constants</u>.

# Lexical Analysis

- Consider the following example of an assignment statement:
  ```
  result = oldsum - value / 100;
  ```

- Following are the lexemes and tokens of this statement:

| Lexemes | Tokens |
|---------|-----------|
| result  | IDENT |
| =       | ASSIGN_OP |
| oldsum  | IDENT |
| -       | SUB_OP |
| value   | IDENT |
| /       | DIV_OP |
| 100     | INT_LIT |
| ;       | SEMICOLON |

# Lexical Analysis

- Lexical analyzer identifies the language keywords (`for,` `if,` `while,` `...`), variable names (`i, j, counter, ...`), literals (`3, 5, 4.54,` …), punctuation (`'(',` `')',';'`) and operators (`+, *,` …) and classify them as tokens.

- Lexical analyzer removes non-programming language symbols such as spaces and comments.

# Lexical Analysis

- One of the approaches to building a lexical analyzer is to design a "state transition diagram" that <u>describes the token patterns of the language</u> and write a program that implements the diagram.

- State diagrams of the form used for lexical analyzers are representations of a class of mathematical machines called "finite automata".

- Finite automata can be designed to recognize members of a class of languages called "regular languages".

- "Regular grammars (Regular expressions)" are generative devices for regular languages.

- The tokens of a programming language are a regular language, and a lexical analyzer is a finite automaton.

# Syntax Analysis

# The Parsing

- The part of the process of <u>analyzing syntax</u> that is referred to as "**syntax analysis**" is often called **parsing**. We will use these two interchangeably.

- Syntax analysis is performed to determine whether the words that make up the source program are <u>in an order</u> in accordance with the <u>grammatical rules</u> of the programming language.

- Parsers for programming languages construct **parse trees** for given programs. The parse tree is used as <u>the basis for translation</u>.

# The Parsing

- Nearly all syntax analysis is <u>based on a formal description of the syntax</u> of the source language (BNF).

- Goals of the parser, given an input program:
  - Find all <u>syntax errors</u>; for each, produce an appropriate <u>diagnostic message</u> and recover quickly.
  - <u>Produce the parse tree</u> for the program.

# Semantic Analysis

# Semantic Analysis

- **Semantic analysis** is the creation of a code in an <u>abstract programming language</u> (intermediate code) using the parse tree created during syntax analysis.
  - Lexical Analysis → RGs / REs – DFA / NFA
  - Syntax Analysis → CFG / BNF - PDA
  - Semantics Analysis → Semantic features of the language
- Intermediate languages are similar to an <u>assembly language</u>.
- This abstract language forms an <u>intermediate step</u> between the compiler's source and object languages, designed to be <u>compatible</u> with the source language's <u>data types and operations</u>.

# Intermediate Code

- If the compiler directly translates source code into the machine code without generating intermediate code then a <u>full native compiler</u> is required for <u>each new machine</u>.

- The intermediate code keeps the <u>analysis portion same</u> for all the compilers that's why it <u>doesn't need a full compiler for every unique machine</u>.

- Intermediate code generator receives input from its predecessor phase (semantic analyzer phase). It takes input in the form of an "annotated syntax tree".

- Using the intermediate code, the <u>second phase</u> of the compiler (the synthesis phase) produces the object code (<u>changed according to the target machine</u>).

# Intermediate Code Generation Example in C Programming Language

```
int main(){
        int dizi1[5], dizi2[5], carpim, i;
        int *p;
        int *r;
        carpim = 0;
        p = dizi1;
        r = dizi2;
        for (i=0; i<10; i++){
                carpim += *p++ * *r++;
        }
        return 0;
}
```

```
carpim=0;
p=&dizi1;
r=&dizi2;
i=0;
S1: if(i>=10) goto S2;
        T3=*p;
        T4=p+1;
        p=T4;
        T5=*r;
        T6=r+1;
        r=T6;
        T7=T3*T5;
        T8=carpim+T7;
        carpim=T8;
        T9=i+1;
        i=T9;
        goto S1;
S2:
```

# Code
# Optimization

# Code Optimization

- **Optimization** is the name given to the <u>collection of techniques</u> that compilers may use to <u>decrease the size</u> and/or <u>increase the execution speed</u> of the code they produce.

- Example:
```
x = 5;
y = 8;
if (x > y)
    …                    → dead code
else
    …
```

- Disposal of dead code will <u>save space</u>.

# Code Optimization

- Example:

```
do
{
    counter = 10;
    sum = sum + counter;
}
while (sum < 100);
```

→

```
counter = 10;
do
{
    sum = sum + counter;
}
while (sum < 100);
```

- Since it is enough to assign the value to the "counter" variable once, taking it out of the loop will <u>increase the running speed</u>.

# Code Optimization

- Optimization is an <u>optional</u> step (it <u>can be turned off</u> in the compiler settings).

- If little or no optimization is done, compilation can be done much <u>faster</u> than if a significant effort is made to produce optimized code.

- The statement order of the code resulting from the optimization may differ from the original source code, which can make <u>debugging</u> (locating errors) <u>difficult</u>.

- Remember that most optimization is done on <u>the intermediate code</u>.

# Code
# Generation

# Code Generation

- At this stage, the conversion from the machine independent "intermediate code" to the "machine code" of the target hardware is performed.

- This is known as **the synthesis step**.

- The best machine instructions that will perform intermediate code operations must be selected.

# Code Generation

- Intermediate code generation has <u>some advantages</u> over directly generating machine code.

- In case no intermediate code is used, if there are x different target machines for a source language, x number of optimizations and code generators will be needed.

- The optimization part written for one machine <u>cannot</u> be used for another machine and will have to be rewritten, which is one of <u>the most difficult parts to write in compiler design</u> (We have to mentioned that many kinds of optimization are <u>difficult to do on machine language</u>).

# Code Generation

- In general, the benefit of the "FRONT-END" and "BACK-END" studies being independent of each other is that:
  - if a new machine (platform) emerges, it enables the production of solutions using the existing "FRONT-END" (and optimization part),
  - and if a new programming language is produced (that uses the same intermediate code), using the existing "BACK-END".

# Code Generation

- The machine code produced in the code generation phase <u>contains reference sections to library files</u>.

- In this case, it turns out that a complete machine code has <u>not</u> been produced and there are <u>missing parts</u>.

- At this stage, the **linker** will come into play and add machine codes to the places where references are left in the library files, and the final executable version of the file will be produced.