

# **Course: Compiler Construction**

## **Week 9: Intermediate Code Generation**

**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 9: Intermediate Code Generation

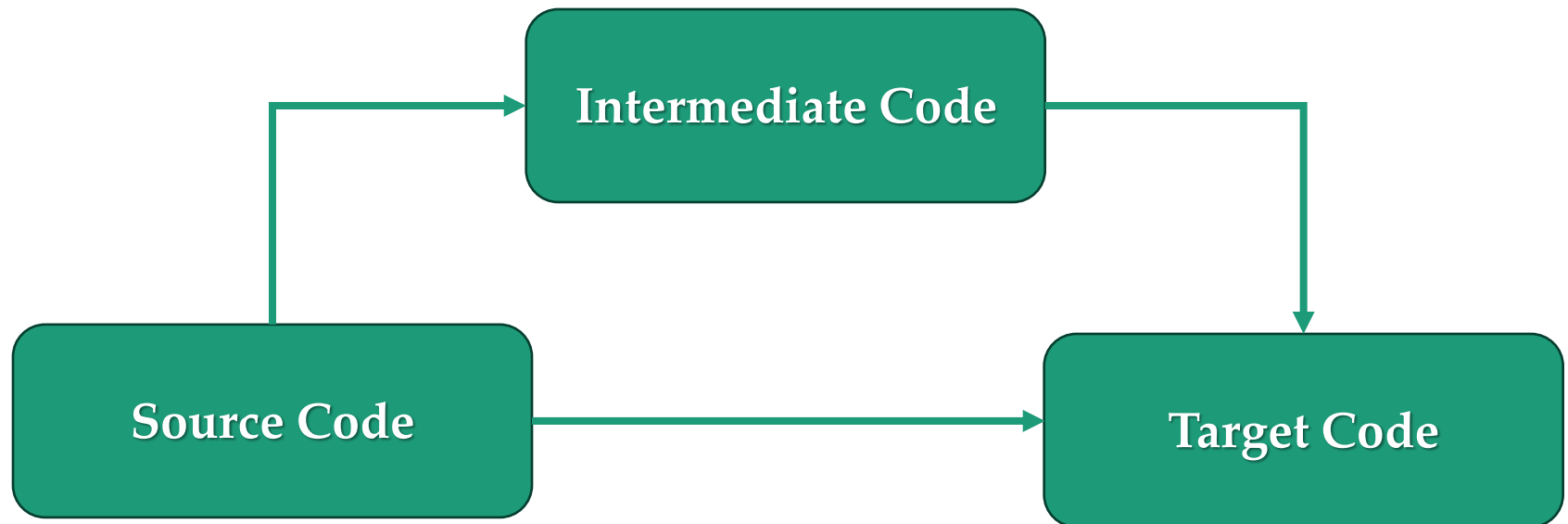
At the end of the lecture, you will be able to:

- i. Differentiate Source code from intermediate output/code and Target Code
- ii. Differentiate High Level and Low Level Intermediate Representation
- iii. Describe the Intermediate Code Generation methods

# INTERMEDIATE CODE GENERATION

- ✓ Compilers are designed to produce target code after working on some source code.
- ✓ When the input program is represented in some hypothetical language or data structure known as **intermediate code** an intermediate output is generated<sup>1</sup>.

- ✓ This representation between the source language and the target machine language program is depicted below<sup>2</sup>

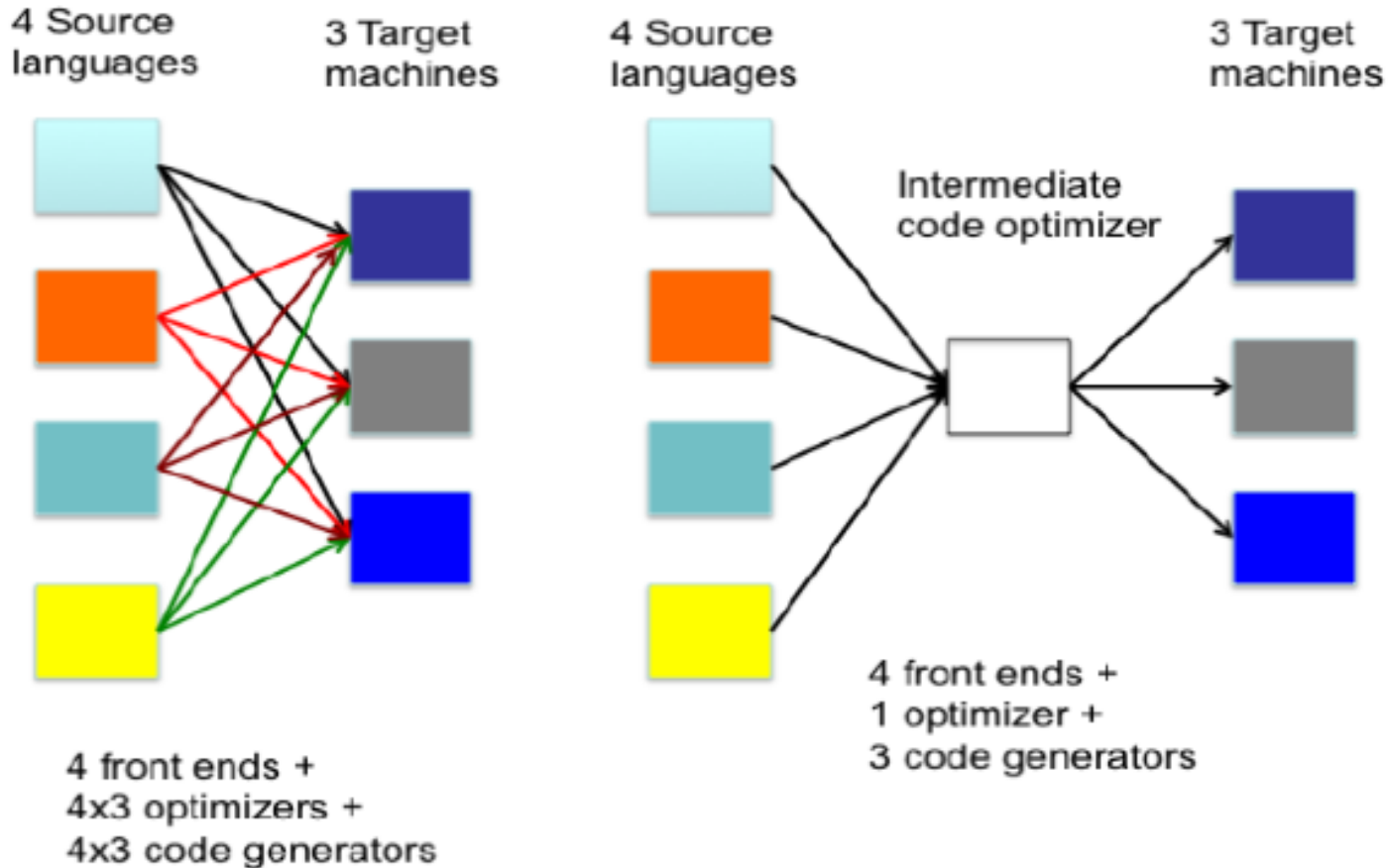


- ✓ Generally, *while it is possible to generate machine code directly from source code*, it entails **two problems**: -
  - i. With **m languages** and **n target machines**, we need to write m front ends,  **$m \times n$  optimizers**, and  **$m \times n$  code generators**<sup>3</sup>

- ✓ Second problem of *generating machine code directly from source code*
- ii. The code optimizer which is one of the largest and very-difficult-to-write components of a compiler, cannot be reused<sup>4</sup>

- ✓ By converting source code to an intermediate code, a machine-independent code optimizer may be written
- ✓ This means just *m front ends*, **n code generators** and **1 optimizer**<sup>5</sup>

# Example of a machine-independent code optimizer<sup>6</sup>



# Intermediate Representation (IR)

- ✓ Intermediate codes can be represented in a variety of ways; each of which have their own benefits.
- ✓ In this lecture we look at two ways of representing Intermediate Code<sup>7</sup>

## I. High Level IR

## II. Low Level IR

# I. High Level IR

- ✓ High-level intermediate code representation is *very close to the source language itself*.
- ✓ They can easily be generated from the source code and we can easily apply *code modifications to enhance performance*.
- ✓ But for target machine optimization, *it is less preferred*<sup>8</sup>.

## II. Low Level IR

- ✓ Low level Intermediate Representation is *close to the target machine*, which makes it suitable for register and memory allocation, instruction set selection, etc.
- ✓ It is good for machine-dependency<sup>9</sup>.

- ✓ Intermediate code can either be:
  - a) **Language-specific** (e.g., Byte Code for Java) or
  - b) **Language-independent** (three-address code)<sup>10</sup>.

# Intermediate Representation Techniques

## a) High level representation

- ✓ The following are the various ways of representing the input program in high level intermediate representation<sup>11</sup>:

*i) Abstract Syntax Trees (AST)*

*ii) Directed Acyclic Graphs (DAG)*

*iii) P-code*

## i) Abstract Syntax Trees (Syntax Trees)

- ✓ An *abstract syntax tree* is a *compacted form of a parse tree* that represents the hierarchical structure of the program.
- ✓ The *nodes* represent *operators*
- ✓ The *children of a node* represent the *operands* on which the operators operate<sup>12</sup>.

# Example:

✓ Consider the following piece of code in the source language<sup>13</sup>:

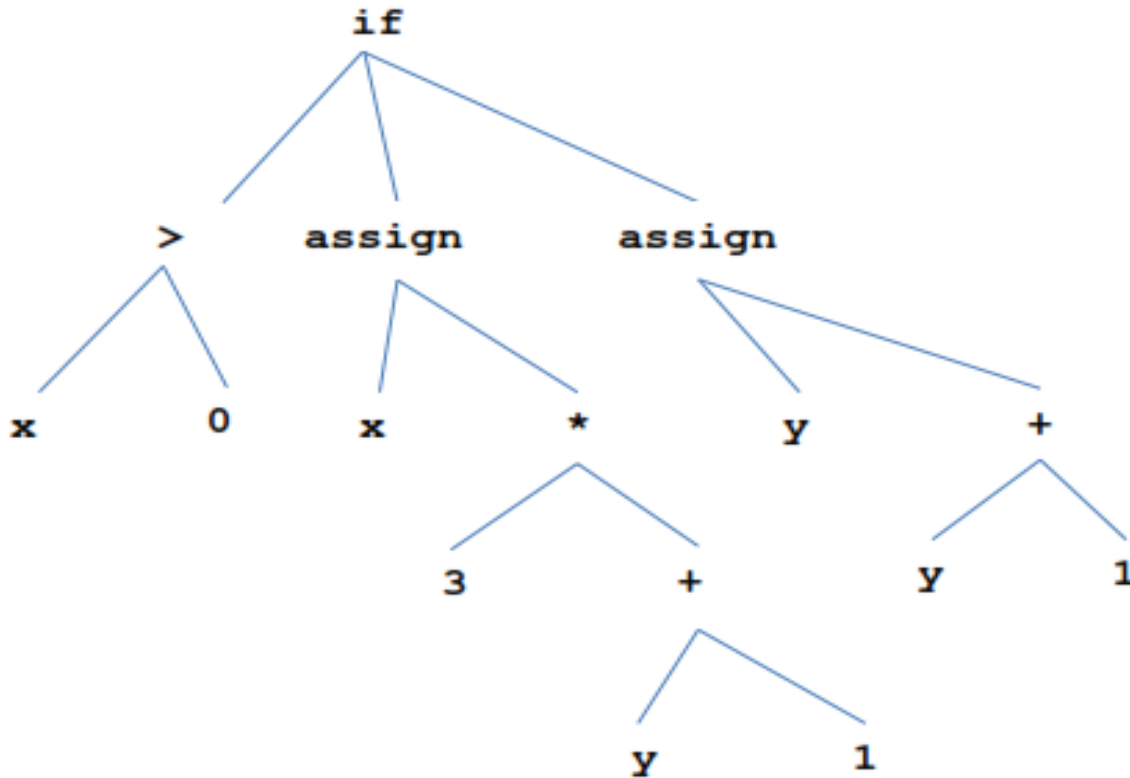
if  $x > 0$

then  $x = 3 * (y + 1)$

else  $y = y + 1$

✓ The **Abstract Syntax Trees (AST)** is as shown below<sup>14</sup>:

*if  $x > 0$  then  $x = 3 * (y + 1)$  else  $y = y + 1$*



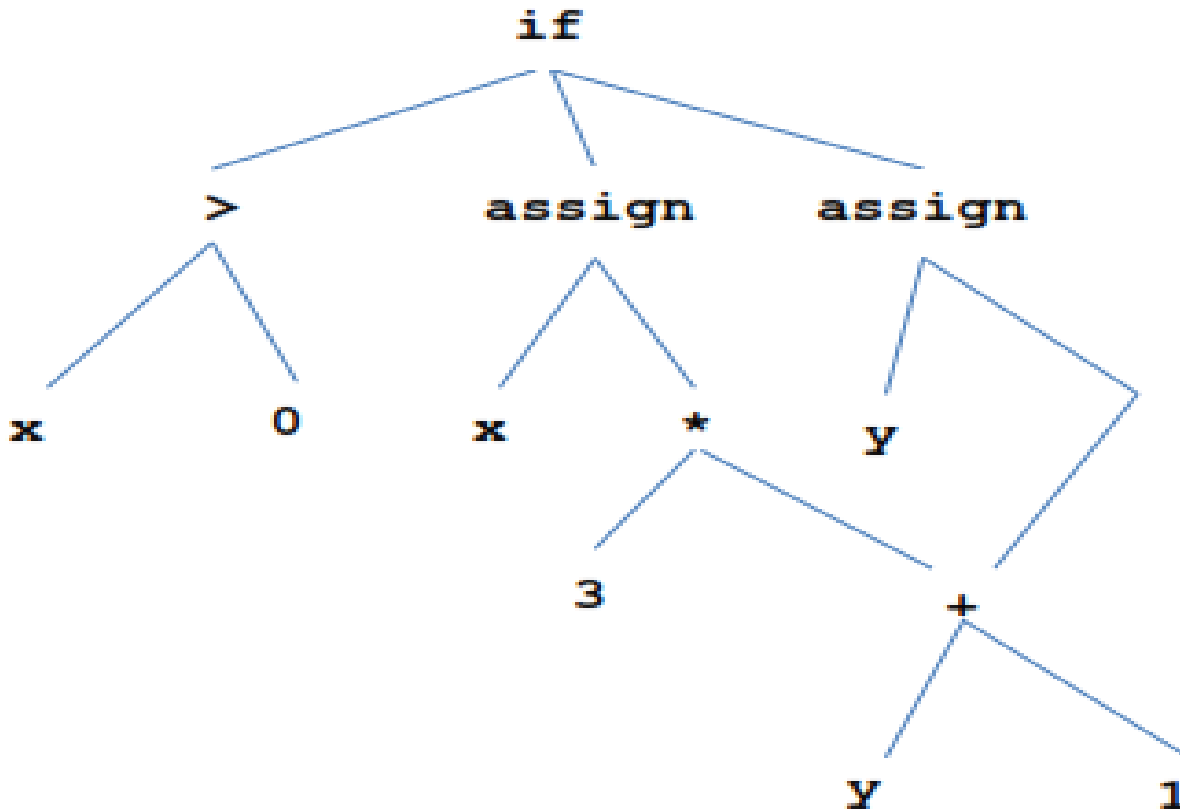
## ii) Directed Acyclic Graphs (DAG)

✓ These are similar to AST except that common sub-expressions are represented by a **single node**<sup>15</sup>.

## Example:

✓ Consider the following piece of code whose DAG is described below<sup>16</sup>

*if  $x > 0$  then  $x = 3 * (y + 1)$  else  $y = y + 1$*



### iii) P-code

- ✓ This representation is used for stack based virtual machines.
- ✓ In this machine, the operands for the operators are always found on top of the stack<sup>17</sup>.

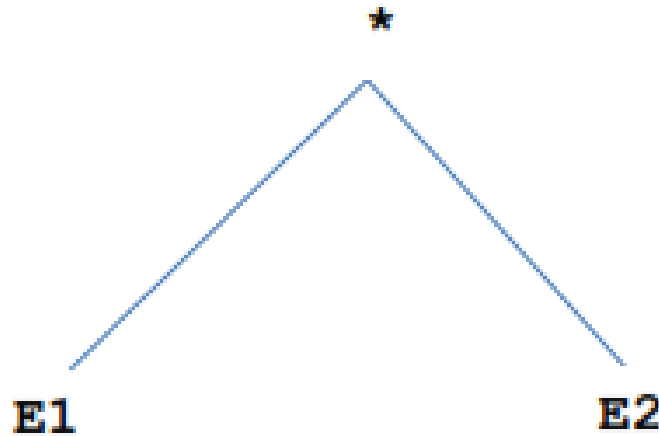
# Example:

- ✓ For the following AST, the code is as follows<sup>18</sup>:

*Code to evaluate E1*

*Code to evaluate E2*

*Mult*



## **b) Low Level Representation**

### **i) Three Address Code**

- ✓ This is a form of low level intermediate representation that permits only one operator to the Right Hand Side (RHS) <sup>19</sup>

# Low Level Representation

## Example:

Consider the following source language statement<sup>20</sup>:-

$$x = y * z + w * a$$

is translated as follows:-

$$t1 = y * z$$

$$t2 = w * a$$

$$x = t1 + t2$$

# Type Checking

- ✓ Type checking is verifying that each operation executed in a program respects the type system of the language, i.e., that all operands in any expression are of appropriate types and number.
- ✓ It is used in the following ways:
  - i). Allow programmers to limit what types may be used in certain circumstances.
  - ii). Assign type to values.
  - iii). Determine whether these values are used in an appropriate manner<sup>21</sup>.

# Type Checking ...

- ✓ Apart from *verifying the code to be correct*, like checking if a function call has the correct number and types of parameters, type checking also helps in *deciding which code is generated* as in the case of arithmetic expressions.
- ✓ There are many variants of type checking, it may be used to check the type of objects and report a type error in case of violation or the incorrect type may be corrected (coercing) <sup>22</sup>.

# Static Checking and Dynamic Checking

## i) Static Checking

- ✓ In modern languages, type checking is done before execution (compile-time). It is known as **static checking** since properties are verified before the program is run.
- ✓ Two important checks are: -
  - i. **Scope checking**- ensures that all variables and functions used within a given scope are *declared correctly*.
  - ii. **Type checking**: ensures that an operator or function is applied to the correct number of arguments of the correct types<sup>23</sup>

# Advantages of Static Checking<sup>24</sup>:

- a) Can catch many common errors.
- b) It is desirable when speed is important  
i.e. it can result in faster code that does not perform any type checking during execution.

## ii) **Dynamic checking (run-time checking)**

This is performed during program execution.

### **Advantages of Dynamic Checking:**

- i) It permits programmers to be less concerned with **types**
- ii) It may be required in some cases like array bound checks which can only be performed during execution.
- iii) May generate more robust code by ensuring thorough checking of values for program identifiers during execution<sup>25</sup>

# Type systems: Strong Typing and Weakly typed Languages

## A. Strongly typed language

- ✓ Language implementation ensures that whenever an operation is performed, the arguments to the operation are of a type that the operation is defined,
- ✓ **Example:-** do not concatenate a string and a floating-point number. This is independent of whether this is ensured **statically** (*prior to execution*) or **dynamically** (*during execution*)<sup>26</sup>.

# Type systems ...

## B. Weakly typed language

- ✓ There is no guarantee that operations are performed on arguments that make sense for the operation.
- ✓ The archetypical weakly typed language is machine code where operations are just performed with no checks.
- ✓ Weakly typed languages are mostly used for system programming, where data needs to be manipulated, moved, copied, encrypted or compressed without regard to what the data represents<sup>27</sup>.

	Static	Dynamic
Strong	Java	Lisp
Weak	C	Perl(1-5)

# Type Expressions

- ✓ These are used to represent the type **of language constructs**.
- ✓ Here, we discuss **SIX** different kinds of type expressions:-
  - i) **Basic type** which includes: **integer, real, character, Boolean**, and other atomic types with no internal structures
  - ii) **Arrays** which are specified as **array(I,T)** where **T** is a **type expression**; **I** is an **integer or a range of integers** e.g. **array(I,T)** is an array of type T with T elements. **Int a [100]** is **array type a with 100 elements**<sup>28</sup>

## Type Expressions...

iii) If  $T1$  and  $T2$  are two typed expressions, then the **Cartesian product  $T1 \times T2$**  is also a **type expression**.

Products denote anonymous records (i.e. field names are absent) and functions argument lists. E.g. functions argument list passed to *function funct* with first argument as *integer*, second as *real* has the associated type as **integer  $\times$  real**.

iv) **Named records** are products, but with **named elements** e.g. a record structure with two named fields: length, an integer and word which is **type array (10, character)** <sup>29</sup>

## Type Expressions...

v) If **T** is a type expression, then **pointer(T)** is a type expression representing objects which are pointers to objects of type **T** e.g. **int \*p**; may be represented as **pointer(integer)**

vi) Functions map a collection of types to another.

✓ They are represented by the type expression  $D \rightarrow R$ , where **D** is the **domain** and **R** is the **range**.

✓ **D** and **R** are type expressions

✓ The type expression **integer × integer → character** represents a function that takes two integers as arguments and returns a character value<sup>30</sup>

# Content Covered in Week 9: Intermediate Code Generation

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

- i. Differentiate Source code from intermediate output/code and Target Code
- ii. Differentiate High Level and Low Level Intermediate Representation
- iii. Describe the Intermediate Code Generation methods

## 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)