# Chapter 1: Course Introduction
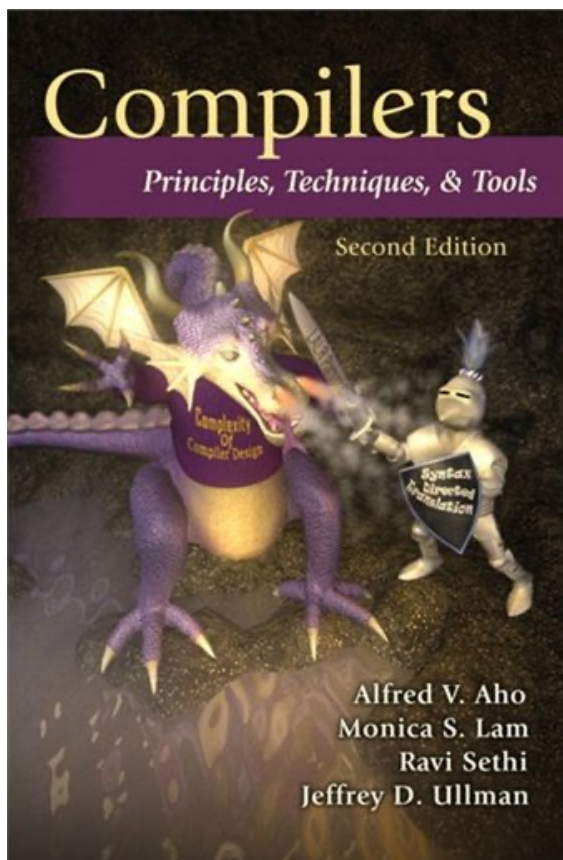
Yepang Liu

liuyp1@sustech.edu.cn

# Outline

- Course Information

- Why Study Compilers?

- The Evolution of Programming Languages

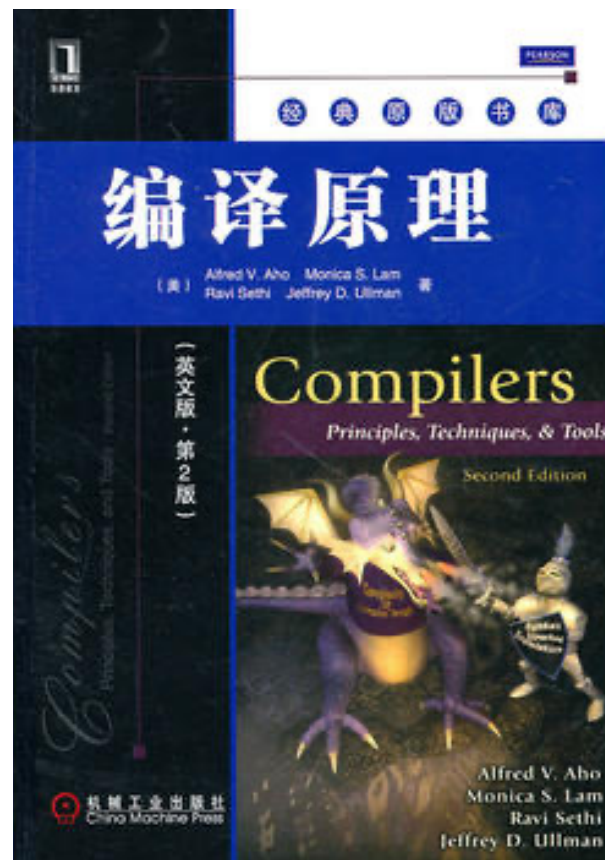- Compiler Structure and Phases

# The Teaching Team

- **Instructor:** Yepang Liu (刘烨庞)
    - Email: liuyp1@sustech.edu.cn
    - Office: Room 609, CoE Building (South)

- **TA:** Yujia Fan (樊宇佳), Yige Chen (陈一戈), Ziming Rui (芮梓铭)
    - Email: 12431253, 12432744, 12432659@mail.sustech.edu.cn
    - Office: Room 650A, CoE Building (South)

- **Communication:**
    - Emails: typically replied within 24 hours
    - Office hour (Yepang Liu): 2:30pm – 4:30pm, every Monday
    - Office hour (Yige Chen): 9:00am – 10:40am, every Friday

# Textbook: The "Dragon Book"



Available at: Yidan Library 2nd Floor ;
TP314 /E7.v1 第3排B面

40 times cheaper than the original edition
~80 ¥ on 京东  (Buy one! It's worth the money ☺.)

# Textbook: Chinese Version

南京大学赵建华、郑滔、戴新宇译

Available at:

Yidan Library 2[nd] Floor Reading Space；

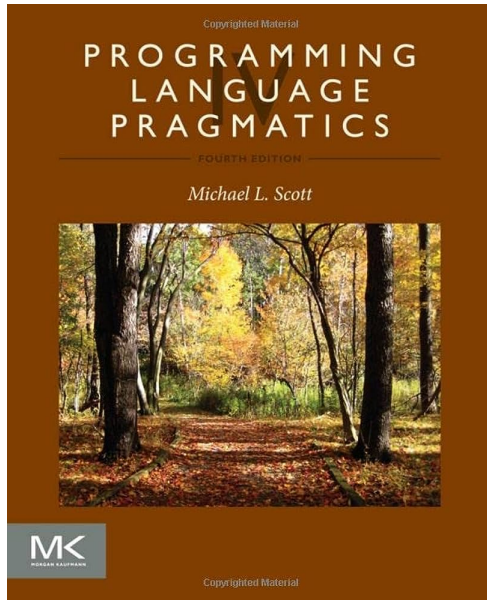TP314 /4.v2 第8排A面

~60 ¥ on 京东

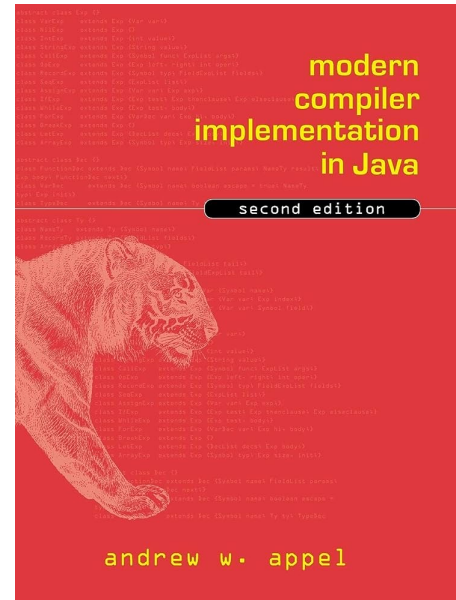# Reference Book for Labs

南京大学许畅、陈嘉、朱晓瑞编著

Available at:

Yidan Library 2nd Floor Reading Space ;

TP314 /18 第8排A面

~100￥on 京东

# Other Reference Books



Available at:

Yidan Library 2nd Floor

TP312 /E 5.v1 第3排B面

Available at:

Yidan Library 1st Floor

TP312JA/E207 第29排A面

Other versions: C, ML

# Course Websites

- Announcements, lecture/lab notes, written assignments, sample answers, etc. are available on Blackboard

  - Registered students will be automatically added to the site

- Other lab/project materials (e.g., project notes, code, test cases) are available on GitHub

  - https://github.com/sqlab-sustech/CS323-2024F

# Marking Scheme (Tentative)



Lecture & Lab Attendance 10%

Project 30%

Research Reports & Presentations 10%

Written Assignments 20%

Final Exam 30%

# Course Content

★ indicates difficulty level, the more the harder

| | |
|---|---|
| Introduction to Compilers (引论) | ★ ☆ ☆ |
| Regular Expressions & Context-Free Grammars (正则表达式与上下文无关文法) | ★ ★ ☆ |
| Lexical Analysis (词法分析) | ★ ★ ★ |
| Syntax Analysis (语法分析) | ★ ★ ★ |
| Syntax-Directed Translation (语法制导的翻译) | ★ ★ ☆ |
| Intermediate-Code Generation (中间代码生成) | ★ ★ ★ |
| Run-Time Environments (运行时刻环境) | ★ ☆ ☆ |
| Code Generation (代码生成) | ★ ★ ☆ |
| Machine-Independent Optimizations (机器无关优化) | ★ ★ ☆ |

# Why Study This Course?

- A fundamental computer science course

- Learn compilation and program analysis techniques

- Learn how to build programming languages

- Course covers both theoretical and practical aspects

  - **Theory:** Lectures and written assignments

  - **Practice:** Lab exercises and projects

- If you want to become a professional programmer

# Outline

- Course Information

- Why Study Compilers?

- **The Evolution of Programming Languages**

- Compiler Structure and Phases

# Programming Languages

- Notations for describing computations

- All software is written in some programming language

- Nowadays, there are over 700 programming languages (~250 popular ones)[1]

  - Low-level (低级语言): directly understandable by a computer

  - High-level (高级语言): understandable by human beings, need a translator to be understood by a computer

[1] https://en.wikipedia.org/wiki/List_of_programming_languages

# Top Programming Languages

- Winners in the past 20 years according to TIOBE index[1]:

    - Python: 2021, 2020, 2018, 2010, 2007

    - C: 2019, 2017, 2008

    - C#: 2023, 2023

    - Go: 2016, 2009

    - Java: 2015, 2005

    - Objective-C: 2012, 2011

    - C++: 2022; JavaScript: 2014; Transact-SQL: 2013; Ruby: 2006; PHP: 2004

[1] A measure of popularity of programming languages calculated from the number of search engine results for queries containing the name of the language. See the index at https://www.tiobe.com/tiobe-index/.

# When It All Started ...

https://en.wikipedia.org/wiki/ENIAC



The first[*] electronic computer ENIAC appeared in 1946. It was programmed in machine language (sequences of 0's and 1's) by setting switches and cables.

[*] The **Atanasoff–Berry computer** (**ABC**) was the first automatic electronic digital computer. It appeared a few years earlier than ENIAC, but it was neither programmable nor Turing-complete. It was designed only to solve systems of linear equations, not for general purposes.

# Can You Understand This?

```
000010010010111001100110011010010110110001100101000 0100
100100100110110001100101011000110111010001110101011100
100110010100110001001011100110001100100010000010100110 0
1110110001101100011001100100101011110110001101101111011 0
110101110000110100101101100011001010110010000101110001
1101000001010001011100111001101100101011000110111010001
101001011011110110111000001001001001000101110011101000
110010101111000111010000100010000010100000100100101110
011000010110110001101001011001110110111000100000011010
000001010000010010010111001100111011011000110111101100 0
100110000101101100010000001101101011000010110100101101
1100000101000001001001011100111010001111001011100000110
0101000010010010000011011010110000101101001011011 10...
```

# Assembly Language (Early 1950s)

```
save %sp,-128,%sp
mov 1,%o0
st %o0,[%fp-20]
mov 2,%o0
st %o0,[%fp-24]
ld [%fp-20],%o0
ld [%fp-24],%o1
add %o0,%o1,%o0
st %o0,[%fp-28]
mov 0,%i0
nop
```

- 1st step towards human-friendly languages

- Mnemonic names (助记符) for machine instructions

- Macro instructions for frequently used sequences of machine instructions

- Explicit manipulation of memory addresses and content

- Still low-level and machine dependent

# The Move to High-Level Languages

- Disadvantages of assembly language

    - Programming is <span style="color:red">tedious</span> and <span style="color:red">slow</span>

    - Programs are <span style="color:red">not understandable</span> by human beings

    - Programs are <span style="color:red">error-prone</span> and <span style="color:red">hard to debug</span>

- High-level programming languages appeared in the second half of the 1950s

    - <span style="color:red">Fortran:</span> for scientific computation

    - <span style="color:red">Cobol:</span> for business data processing

    - <span style="color:red">Lisp:</span> for symbolic computation

# Fortran: The 1ˢᵗ High-Level Language

- In 1953, John Backus proposed to develop a more practical alternative to assembly language for programming on IBM 704 mainframe computer



John Backus (1924 – 2007)
American Computer Scientist
ACM Turing Award (1997)



IBM 704 mainframe

# Fortran: The 1ˢᵗ High-Level Language

- The 1ˢᵗ Fortran (**For**mula **Tran**slation) compiler was delivered in 1957

- Coding became much faster, 50%+ software was in Fortran in 1958

- **Huge impact**, modern compilers preserve the outline of Fortran I

- Fortran is still used today (No. 10, TIOBE Index August 2024)

```
C---- THIS PROGRAM READS INPUT FROM THE CARD READER,
C---- 3 INTEGERS IN EACH CARD, CALCULATE AND OUTPUT
C---- THE SUM OF THEM.
  100 READ(5,10) I1, I2, I3
   10 FORMAT(3I5)
      IF (I1.EQ.0 .AND. I2.EQ.0 .AND. I3.EQ.0) GOTO 200
      ISUM = I1 + I2 + I3
      WRITE(6,20) I1, I2, I3, ISUM
   20 FORMAT(7HSUM OF , I5, 2H, , I5, 5H AND , I5,
     *    4H IS , I6)
      GOTO 100
  200 STOP
      END
```

Fortran code example



Fortran code on a punch card

http://www.herongyang.com/Computer-History/FORTRAN-Program-Store-on-Punch-Card.html

# Outline

- Course Information

- Why Study Compilers?

- The Evolution of Programming Languages

- **Compiler Structure and Phases**

# The Structure of a Compiler

Source Code

| Lexical Analyzer |
| --- |
| Syntax Analyzer |
| Semantic Analyzer |
| Intermediate Code Generator |

Frontend
(the *analysis* part)

| Machine-Independent Code Optimizer |
| --- |
| Code Generator |
| Machine-Dependent Code Optimizer |

Backend
(the *synthesis* part)

Machine Code

# The Frontend (前端) of a Compiler

| | |
|---|---|
| Source Code → | **Lexical Analyzer** |
| | **Syntax Analyzer** |
| | **Semantic Analyzer** |
| | **Intermediate Code Generator** → IR |

- Breaks up the source program into constituent pieces and imposes a grammatical structure on them

- Uses the grammatical structure to create an intermediate representation (IR) of the source program

- Collect the information about the source program and stores it in a data structure called *symbol table* (will be passed to backend with IR)

# The Backend (后端) of a Compiler

```
                    ┌─────────────────────────┐
                    │  Machine-Independent    │
  ┌──────────┐      │     Code Optimizer      │         ┌──────────┐
  │          │      ├─────────────────────────┤         │          │
  │    IR    │ ───▶ │     Code Generator      │  ───▶   │ Machine  │
  │          │      ├─────────────────────────┤         │   Code   │
  └──────────┘      │   Machine-Dependent     │         └──────────┘
                    │     Code Optimizer      │
                    └─────────────────────────┘
```

- Constructs the target program (typically, in machine language) from the IR and the information in the symbol table

- Performs code optimizations during the process[*]

[*] Lexing and parsing are most complex and expensive in the early days, while in today, optimization dominates all other phases and lexing and parsing are very cheap.

# Lexical Analysis (Scanning, 词法分析)

Character stream ⟶ | Lexical Analyzer | ⟶ Token stream

- The lexical analyzer (lexer/tokenizer/scanner) breaks down the source code into a sequence of "lexemes" (词素) or "words"

- For each lexeme, produce a "token" (词法单元) in the form:

`<token-name, attribute-value>`

An abstract symbol that is used during syntax analysis

Points to an entry in the symbol table. Info in the table entry is for semantic analysis and code generation.

# Lexemes vs. Tokens

- A **lexeme** is a string of characters that is a lowest-level syntactic unit in the programming language

  - "words" and punctuation of the programming language (**instance**)

- A **token** is a syntactic category representing a class of lexemes

  - **In English:** Noun, Verb, Adjective…

  - **In programming language**: Identifier, Keyword, Whitespace… (**pattern**)

https://courses.cs.vt.edu/~cs1104/Compilers/Compilers.070.html

# Lexical Analysis (Example)

`position = initial + rate * 60`

Lexical Analyzer

`<id, 1>`

`<=>`

`<id, 2>`

`<+>`

`<id, 3>`

`<*>`

`<60>`

SYMBOL TABLE

| | | |
|---|---|---|
| 1 | position | ... |
| 2 | initial | ... |
| 3 | rate | ... |
| | | |

Note: <=>, <+>, <*>, <60> are not in the defined form. This is for notational convenience. <=> could have been <assign, -> and <60> could have been <number, 4>.

# Lexical Analysis (Analogy)

position = initial + rate * 60          SUSTech is a great university

| Lexical Analyzer |                        | Human brain |

<**id**, 1>                              <**noun**, "SUSTech">

<=>                                     <**verb**, "is">

<**id**, 2>                              <**article**, "a">

<+>                                     <**adjective**, "great">

<**id**, 3>                              <**noun**, "university">

<*>

<60>                                    Example adapted from Aiken's notes (Stanford CS143)

# Syntax Analysis (Parsing, 语法分析)

Token stream ⟶ | Syntax Analyzer | ⟶ Syntax tree

- The syntax analyzer (parser) uses the token names produced by the lexer to create an intermediate representation that depicts the grammar structure of the token stream, typically a *syntax tree*

- Each interior node represents an operation and the children of the node represent the arguments of the operation

# Syntax Analysis (Example)

`<id, 1> <=> <id, 2> <+> <id, 3> <*> <60>`

Syntax Analyzer

```
            =
          /   \
    <id, 1>    (+)  ← Interior node (operation)
              /   \
        <id, 2>    *
          ↓       /  \
children node  <id, 3>  60
(argument)
```

Interior node (operation)

children node (argument)

# Syntax Analysis in English

<article, "SUSTech">  <verb, "is">  <article, "a">

<adjective, "great">  <noun, "university">

```
Human brain
```

sentence

subject          object

noun    verb    article    adjective    noun

SUSTech    is    a    great    University

# Semantic Analysis (语义分析)

Syntax tree ⟶ | Semantic Analyzer | ⟶ Syntax tree

- The semantic analyzer uses the syntax tree and the information in the symbol table to check the source program for <span style="color:red">semantic consistency</span> with the language definition

- Also gathers <span style="color:red">type information</span> for type checking, type conversion, and intermediate code generation

# What is Semantics?

- The syntax of a programming language describes the proper form of its programs


- The semantics of a programming language describes the meaning of its programs, i.e., what each program does when it executes

# Semantic Analysis in English

Jack said Jerry left his assignment at home.

*What does "his" refer to? Jack's or Jerry's?*

Jack said Jack left his assignment at home.

*How many Jacks? Which one left the assignment?*

Examples are from Aiken's notes (Stanford CS143)

# Semantic Analysis in Programming

- Understanding the meaning of a program is very hard ☹

- Compilers perform only very limited analysis to catch semantic inconsistencies.

```
1. {
2.    int Jack = 3;
3.    {
4.       int Jack = 4;
5.       print Jack;
6.    }
7. }
```

*Which value will be printed?*

**Programming languages define strict rules to avoid ambiguities.** Compiler will bind Jack at line 5 to its inner definition at line 4.

# Type Checking (类型检查)
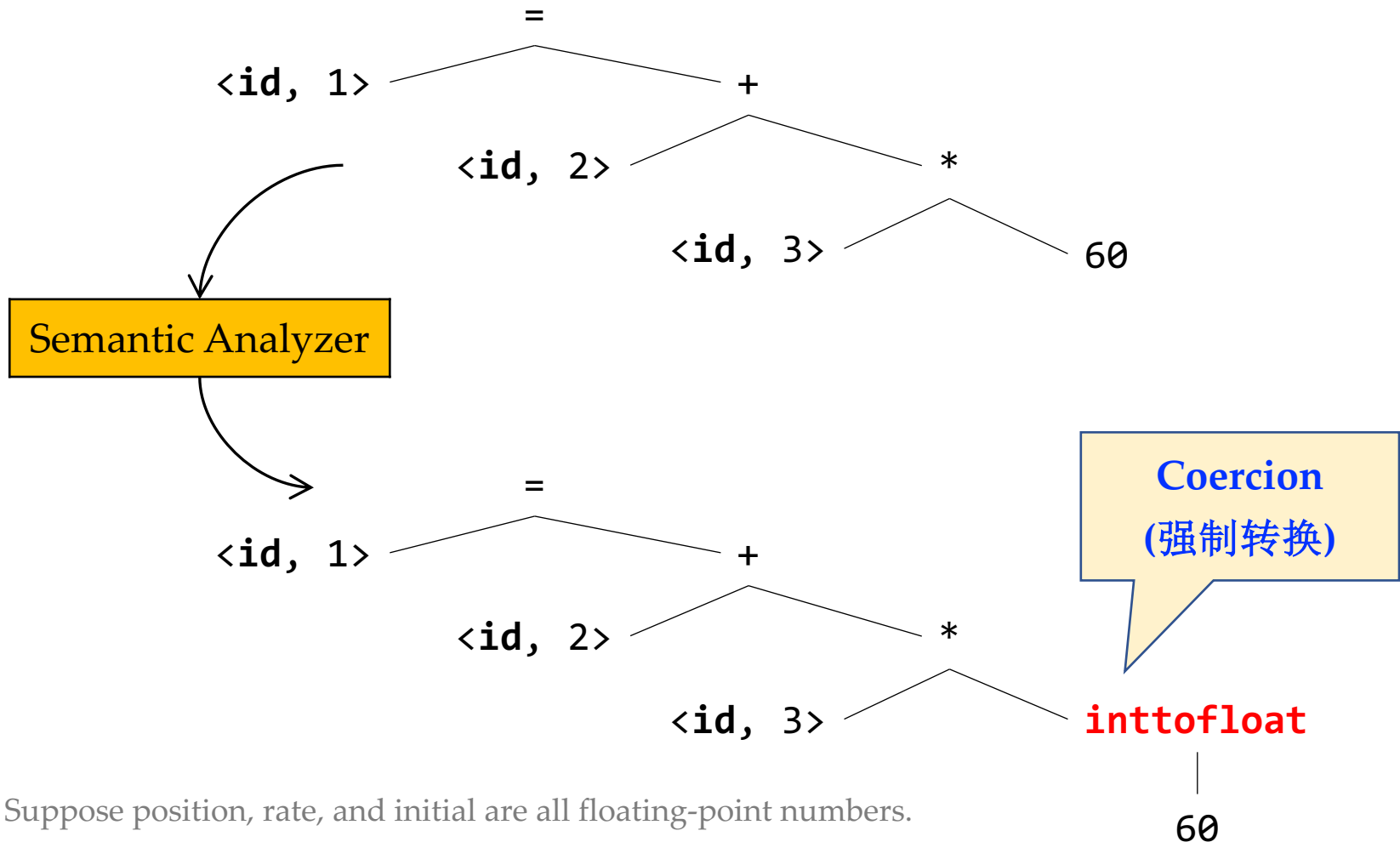
- An important part of semantic analysis is type checking

- Compilers check that each operator has matching operands (of correct types)

Example: Many language definitions require an array index to be an integer.

```
double x = 3.2;
int[] nums = new int[5];
nums[x] = 6;
```

Compilers should report an error!

# Semantic Analysis (Example)

```
            =
  <id, 1>        +
         <id, 2>        *
                 <id, 3>        60
```

Semantic Analyzer

```
            =
  <id, 1>        +
         <id, 2>        *
                 <id, 3>        inttofloat
                                    |
                                   60
```

**Coercion**
**(强制转换)**

Suppose position, rate, and initial are all floating-point numbers.

# Intermediate Code Generation
# (中间代码生成)

Syntax tree ⟶ | Intermediate Code Generator | ⟶ IR
(in three-address code)
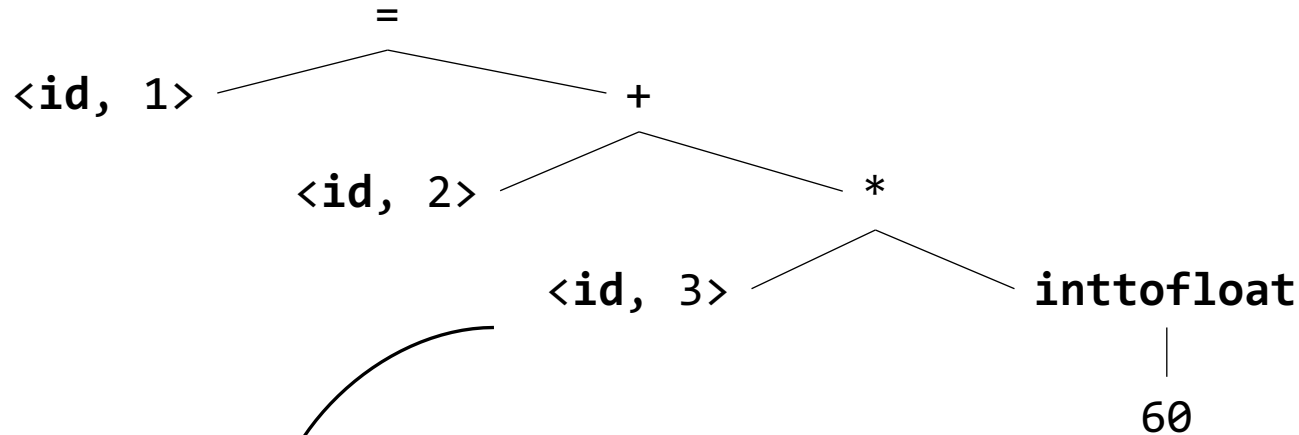
- After semantic analysis, compilers generate an intermediate representation, typically *three-address code (三地址码)*

  - Assembly-like instructions with three operands per instruction

  - Each operand acts like a register

  - Each assignment instruction has at most one operator on the RHS

  - Easy to translate into machine instructions of the target machine

# Three-Address Code Example

```
                          =
        <id, 1>                      +
                    <id, 2>                   *
                                <id, 3>            inttofloat
                                                        |
                                                        60
```

Intermediate Code Generator

```
t1 = inttofloat(60)
t2 = id3 * t1
t3 = id2 + t2
id1 = t3
```

# Machine-Independent Code Optimization
# (机器无关的代码优化)

IR ⟶ | Machine-Independent Code Optimizer | ⟶ IR

- Akin to article editing/revising in English

- Improve the intermediate code for better target code

  - Run faster

  - Use less memory

  - Shorter code

  - Consume less power …

# Code Optimization (Example)

```
t1 = inttofloat(60)

t2 = id3 * t1

t3 = id2 + t2

id1 = t3
```

1. 60 is a constant integer value. Its conversion to floating-point can be done once and for all at compile time

2. t2 and t3 are only used for value transmitting

Optimization
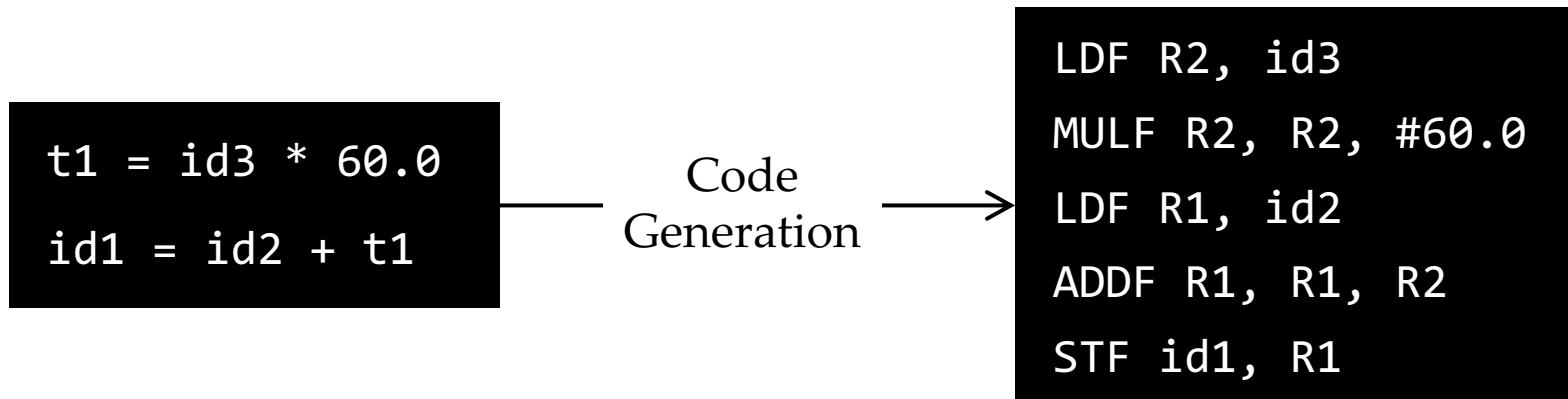
```
t1 = id3 * 60.0

id1 = id2 + t1
```

# Code Generation (代码生成)

IR $\longrightarrow$ [ Code Generator ] $\longrightarrow$ Target-machine code

- Map IR to target language, analogous to human translation

- It is crucial to <span style="color:red">allocate register and memory</span> to hold values

```
t1 = id3 * 60.0
id1 = id2 + t1
```

Code Generation $\longrightarrow$

```
LDF R2, id3
MULF R2, R2, #60.0
LDF R1, id2
ADDF R1, R1, R2
STF id1, R1
```
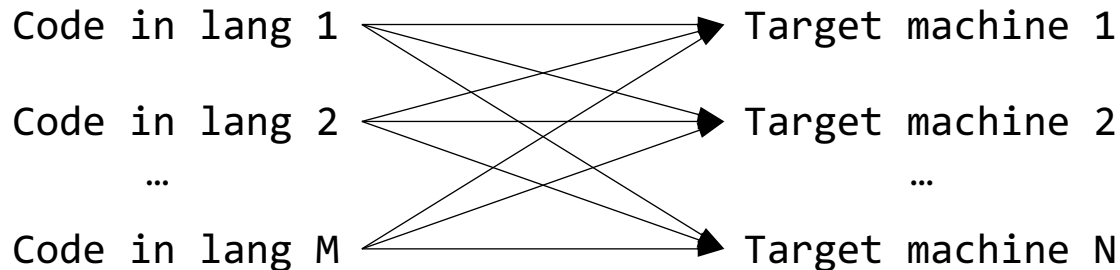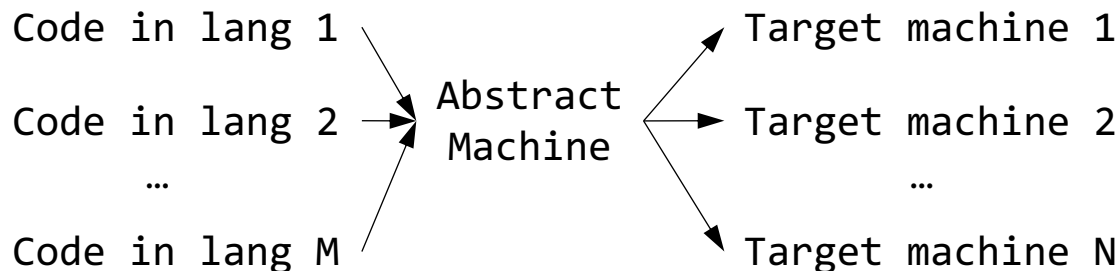
# Symbol Table Management

- Performed by the frontend, symbol table is passed along with the intermediate code to the backend

- Record the variable names and various attributes
  - storage allocated, type, scope

- Record the procedure names and various attributes
  - the number and type of arguments
  - the way of passing arguments (by value or by reference)
  - the return type

# Intermediate Language (IL)

- Intermediate code is in IL (e.g., three-address code)

- A good IL eases compiler implementation

```
Code in lang 1 ─────────────────→ Target machine 1

Code in lang 2 ─────────────────→ Target machine 2
      …                                  …
Code in lang M ─────────────────→ Target machine N
```

**M * N** compilers
without a good IL

```
Code in lang 1 ╲
                ╲        ╱→ Target machine 1
Code in lang 2 ──→ Abstract ──→ Target machine 2
      …          Machine         …
Code in lang M ╱        ╲→ Target machine N
```
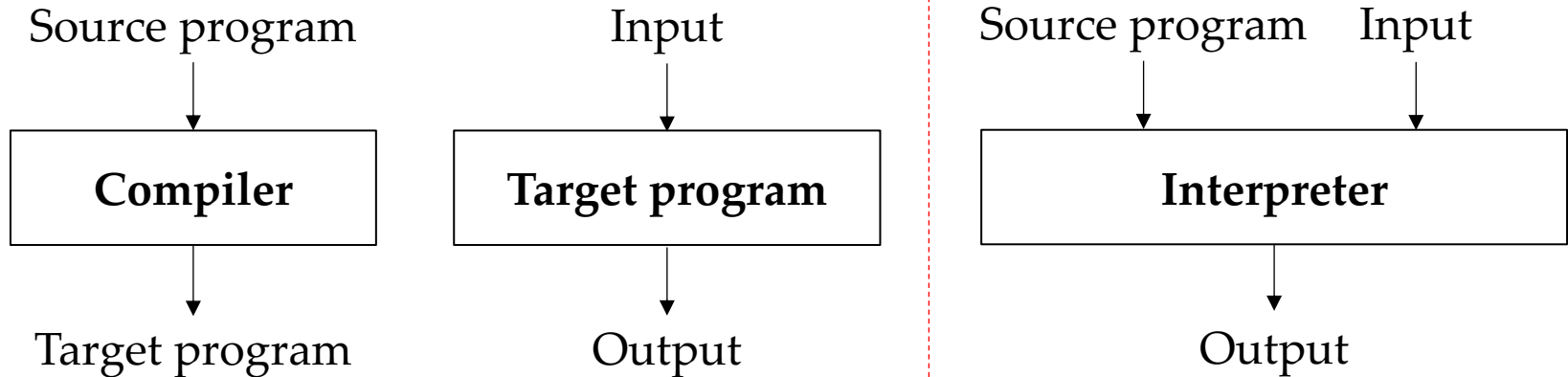
**M + N** compilers
with a good IL

# Compilers vs. Interpreters

> A complier translates source programs written in high-level languages into machine codes that can run directly on the target computer.
>
> An interpreter directly executes each statement in the source code, without requiring the program to have been compiled into machine codes.

**1**

Source program

↓

**Compiler**

↓

Target program

Input

↓

**Target program**

↓

Output

Source program    Input

↓              ↓

**Interpreter**

↓

Output

# Compilers vs. Interpreters

**2**

Interpreters often take less time to analyze the source code: they simply parse each statement and execute it (e.g., Python code).

In comparison, compilers typically analyze the relationships among statements (e.g., control and data flows) to enable optimizations.

**3**

Interpreters continue executing a program until the first error is met, in which case they stop.

For compiled languages, programs are executable only after they are successfully compiled.

# Reading Tasks

- Chapter 1 of the Dragon book

    1.1 Language Processors

    1.2 The Structure of a Compiler

    1.3 The Evolution of Programming Languages