



南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

Chapter 1: Introduction

Yepang Liu

liuyp1@sustech.edu.cn

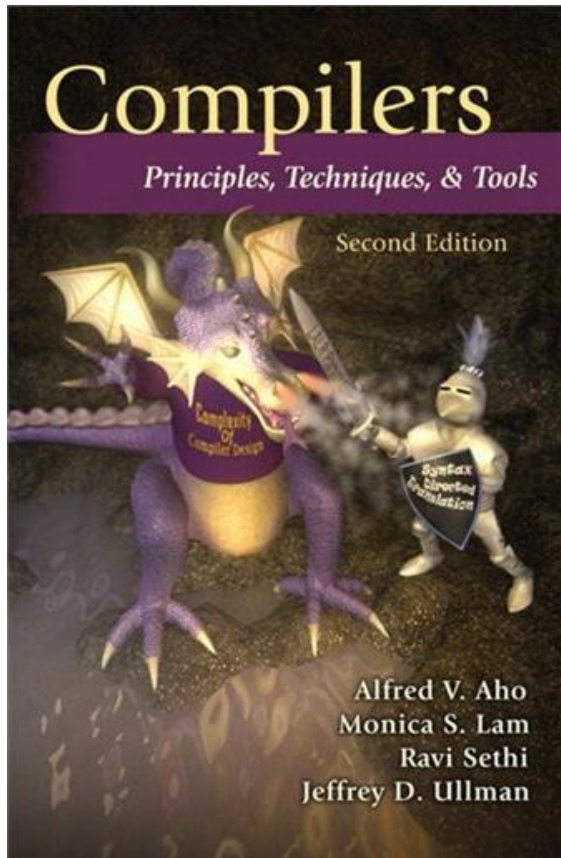
Outline

- Course Information
- Why Study Compilers?
- The Evolution of Programming Languages
- Compiler Structure and Phases

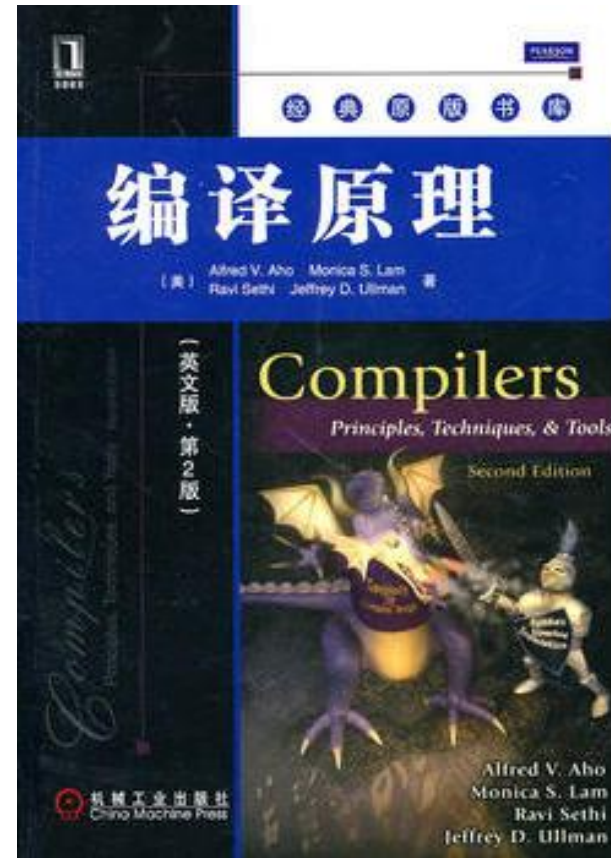
The Teaching Team

- **Instructor:** Yepang Liu (刘烨庞)
 - Email: liuyp1@sustech.edu.cn
 - Office: Room 609, CoE Building (South)
- **TA:** Pengkun Jiang (姜朋坤), Shangwen Wang (王尚文)
 - Email: 12132334@mail.sustech.edu.cn, 1453505201@qq.com
 - Office: Room 650A, CoE Building (South)
- **Communication:**
 - Emails: typically replied within 24 hours
 - Office hour: 4pm – 5pm, every Monday

Textbook: The “Dragon Book”



Available at: Lynn Library 2nd Floor ;
TP314 /E 7 第20排B面



40 times cheaper than the original edition
~70 ¥ on 京东 (Buy one! It's worth the money ☺.)

Textbook: Chinese Version



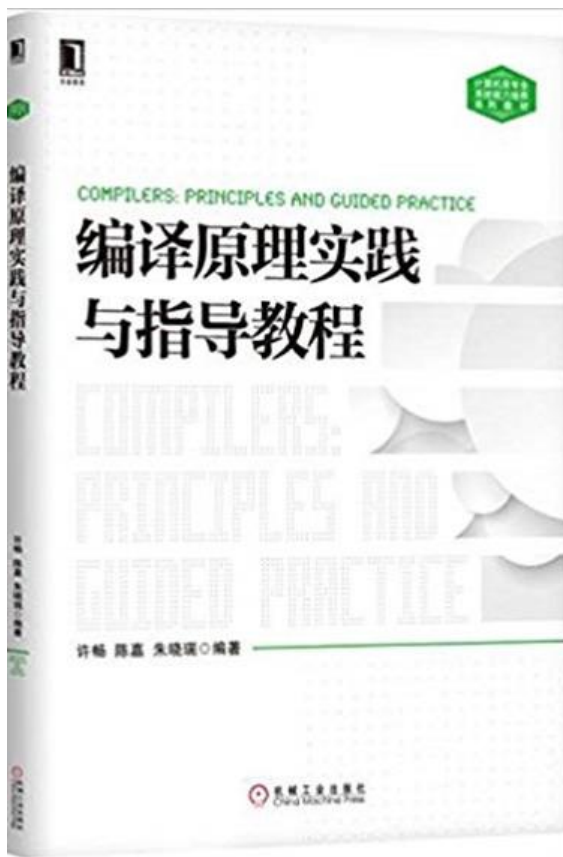
南京大学赵建华、郑滔、戴新宇译

Available at:

Yidan Library 2nd 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 ;

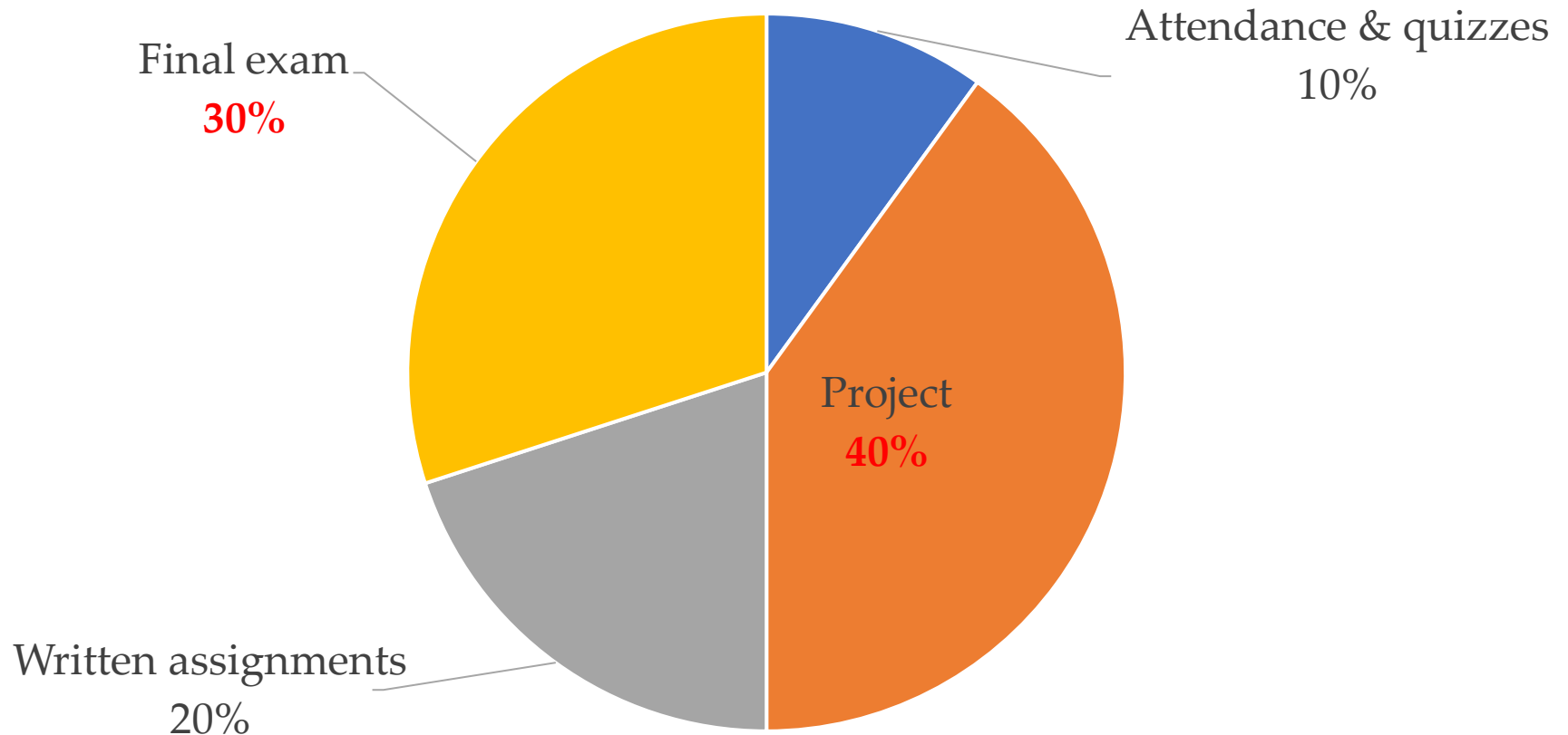
TP314 /4.v2 第8排A面

~60 ¥ on 京东

Course Materials

- Lecture/lab notes (and reference), written assignments, sample answers, project notes, and etc. will be available on SAKAI
 - <https://sakai.sustech.edu.cn/portal/site/3f8c8da1-3427-4770-8572-4be1fd483cbe>
- Other lab/project materials (e.g., code, test cases) will be available on GitHub
 - <https://github.com/sqlab-sustech/CS323-2022F>

Marking Scheme



Academic Integrity

From Spring 2022, the plagiarism policy applied by the Computer Science and Engineering department is the following: ↵

↵

*** If an undergraduate assignment is found to be plagiarized, the first time the score of the assignment will be 0.**↵

*** The second time the score of the course will be 0.**↵

*** If a student does not sign the Assignment Declaration Form or cheats in the course, including regular assignments, midterms, final exams, etc., in addition to the grade penalty, the student will not be allowed to enroll in the two CS majors through 1+3, and cannot receive any recommendation for postgraduate admission exam exemption and all other academic awards.**↵

↵

As it may be difficult when two assignments are identical or nearly identical who actually wrote it, the policy will apply to BOTH students, unless one confesses having copied without the knowledge of the other. ↵



南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

计算机科学与工程系

Department of Computer Science and Engineering

本科生作业承诺书

本人_____（学号_____）本学期已选修计算机科学与工程系_____课程。本人已阅读并了解《南方科技大学计算机科学与工程系本科生作业抄袭学术不端行为的认定标准及处理办法》制度中关于禁止本科生作业抄袭的相关规定，并承诺自觉遵守其规定。

承诺人：

年 月 日

Course Content

Introduction to Compilers (引论)	★ ☆ ☆
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)	★ ★ ☆

★ indicates difficulty level, the more the harder

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

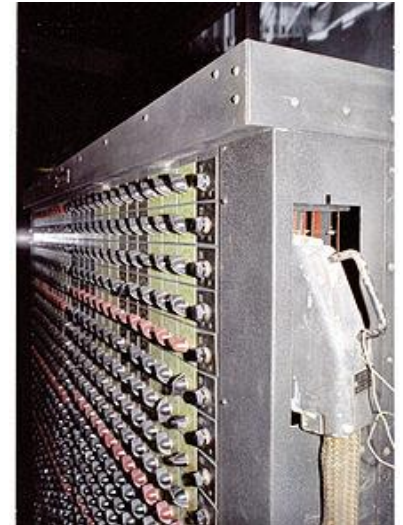
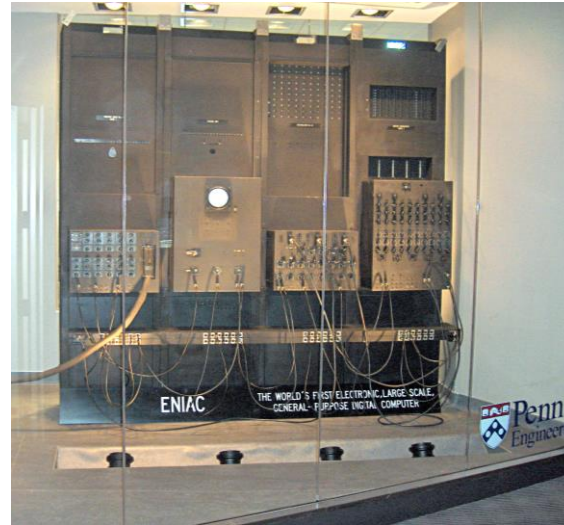
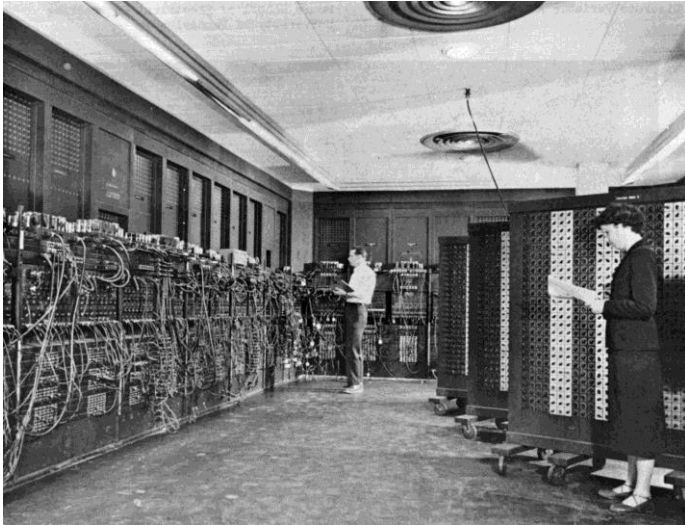
Programming Languages

- 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

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?



```
0000100100101110011001100110100101101100011001010000100
1001000100110110001100101011000110111010001110101011100
1001100101001100010010111001100011001000100000101001100
1110110001101100011001100100101111101100011011011110110
1101011100000110100101101100011001010110010000101110001
1101000001010001011100111001101100101011000110111010001
1010010110111101101110000010010010001000101110011101000
1100101011110000111010000100010000010100000100100101110
0110000101101100011010010110011101101110001000000011010
0000010100000100100101110011001110110110001101111011000
1001100001011011000010000001101101011000010110100101101
1100000101000001001001011100111010001111001011100000110
0101000010010010000001101101011000010110100101101110...
```

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 **tedious** and **slow**
 - Programs are **not understandable** by human beings
 - Programs are **error-prone** and **hard to debug**
- High-level programming languages appeared in the second half of the 1950s
 - **Fortran**: for scientific computation
 - **Cobol**: for business data processing
 - **Lisp**: 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)



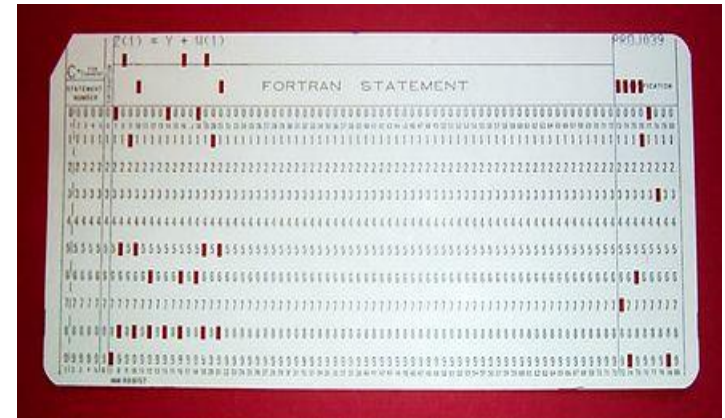
IBM 704 mainframe

Fortran: The 1st High-Level Language

- The 1st Fortran (**F**ormula **T**ranslation) 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. 19, TIOBE Index Aug 2022)

```
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



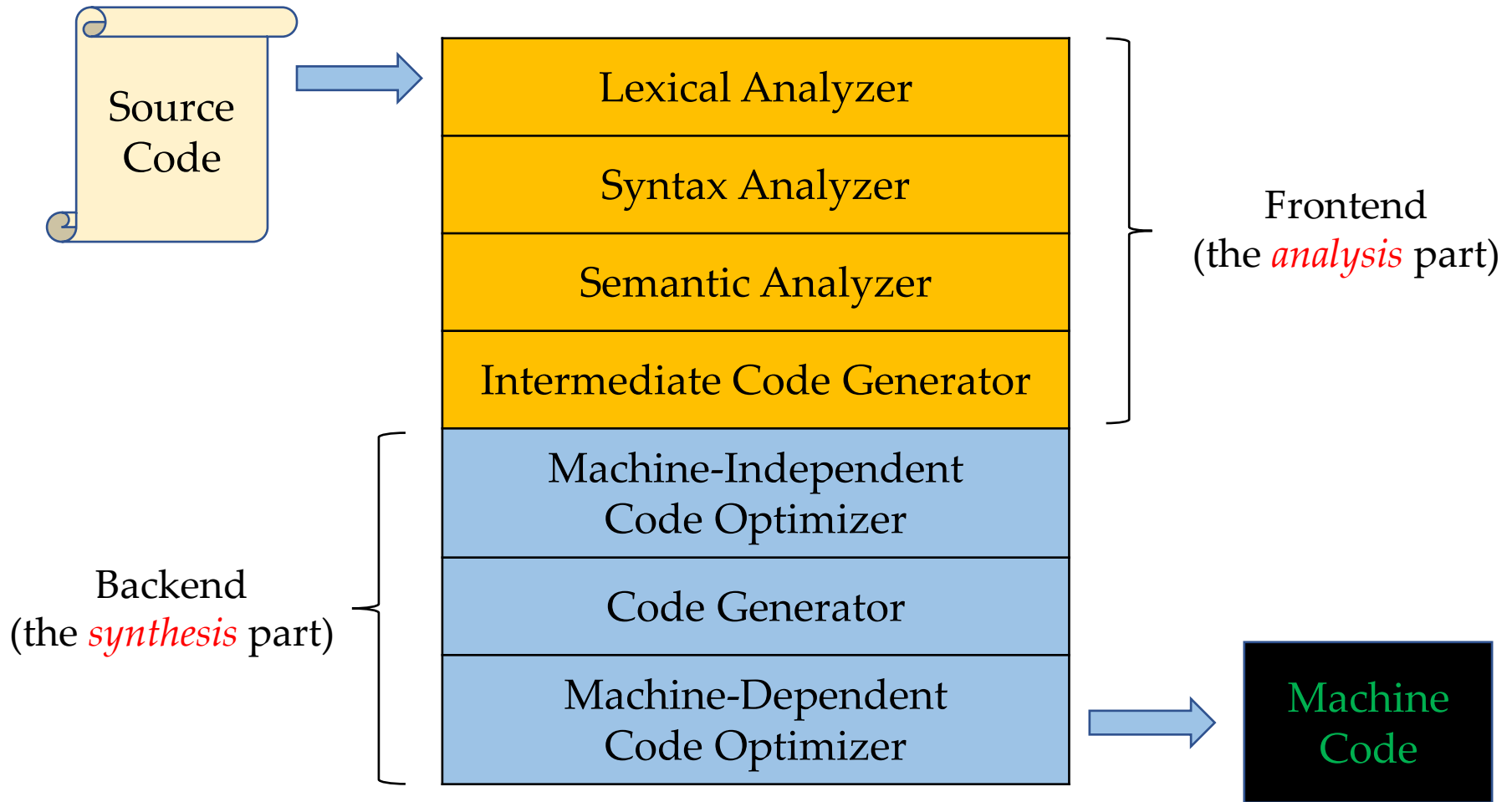
Fortran code on a punch card

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

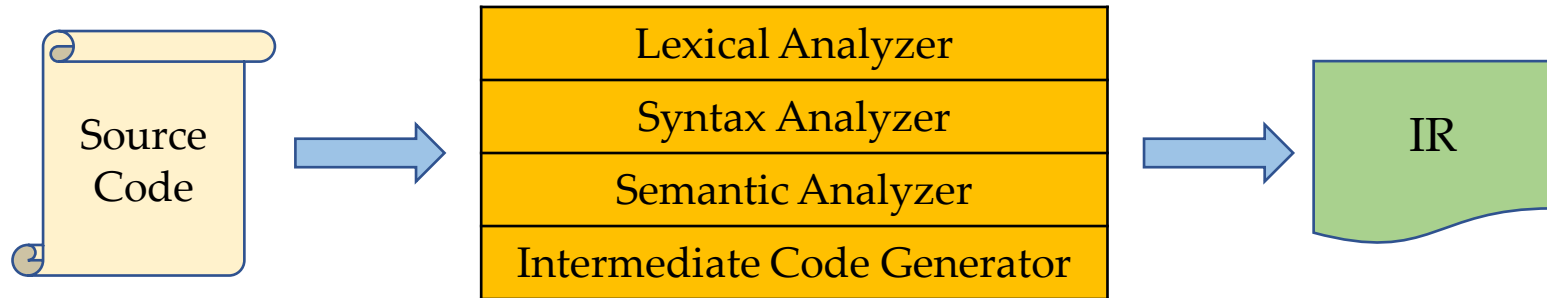
Outline

- Course Information
- Why Study Compilers?
- The Evolution of Programming Languages
- Compiler Structure and Phases

The Structure of a Compiler

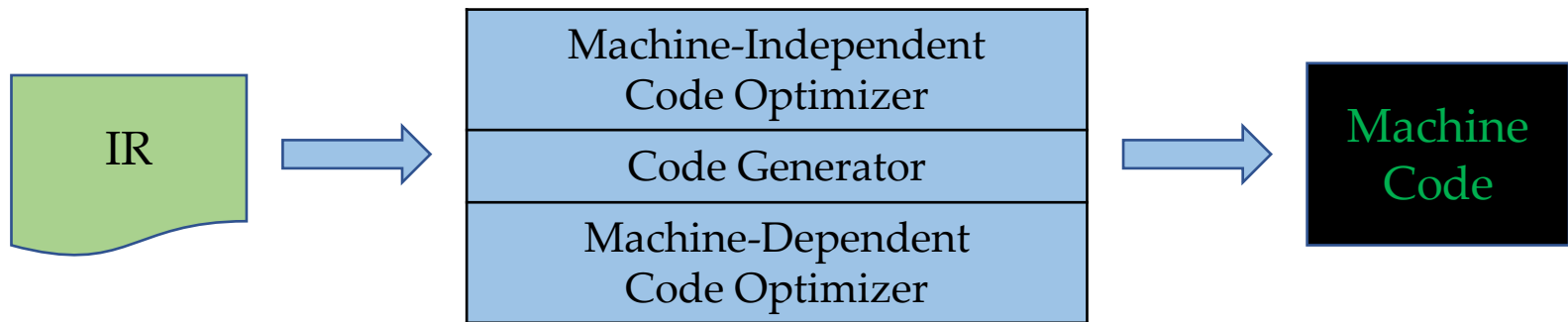


The Frontend (前端) of a Compiler



- 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



- 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)



- 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.

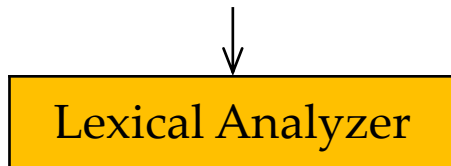
Lexemes vs. Tokens

- 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



<id, 1>

<=>

<id, 2>

<+>

<id, 3>

<*>

<60>

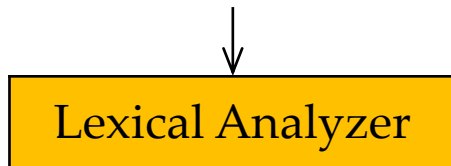
SYMBOL TABLE

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`



`<id, 1>`

`<=>`

`<id, 2>`

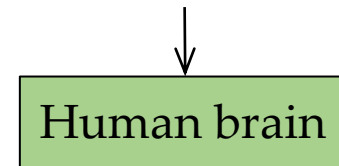
`<+>`

`<id, 3>`

`<*>`

`<60>`

`SUSTech is a great university`



`<noun, "SUSTech">`

`<verb, "is">`

`<article, "a">`

`<adjective, "great">`

`<noun, "university">`

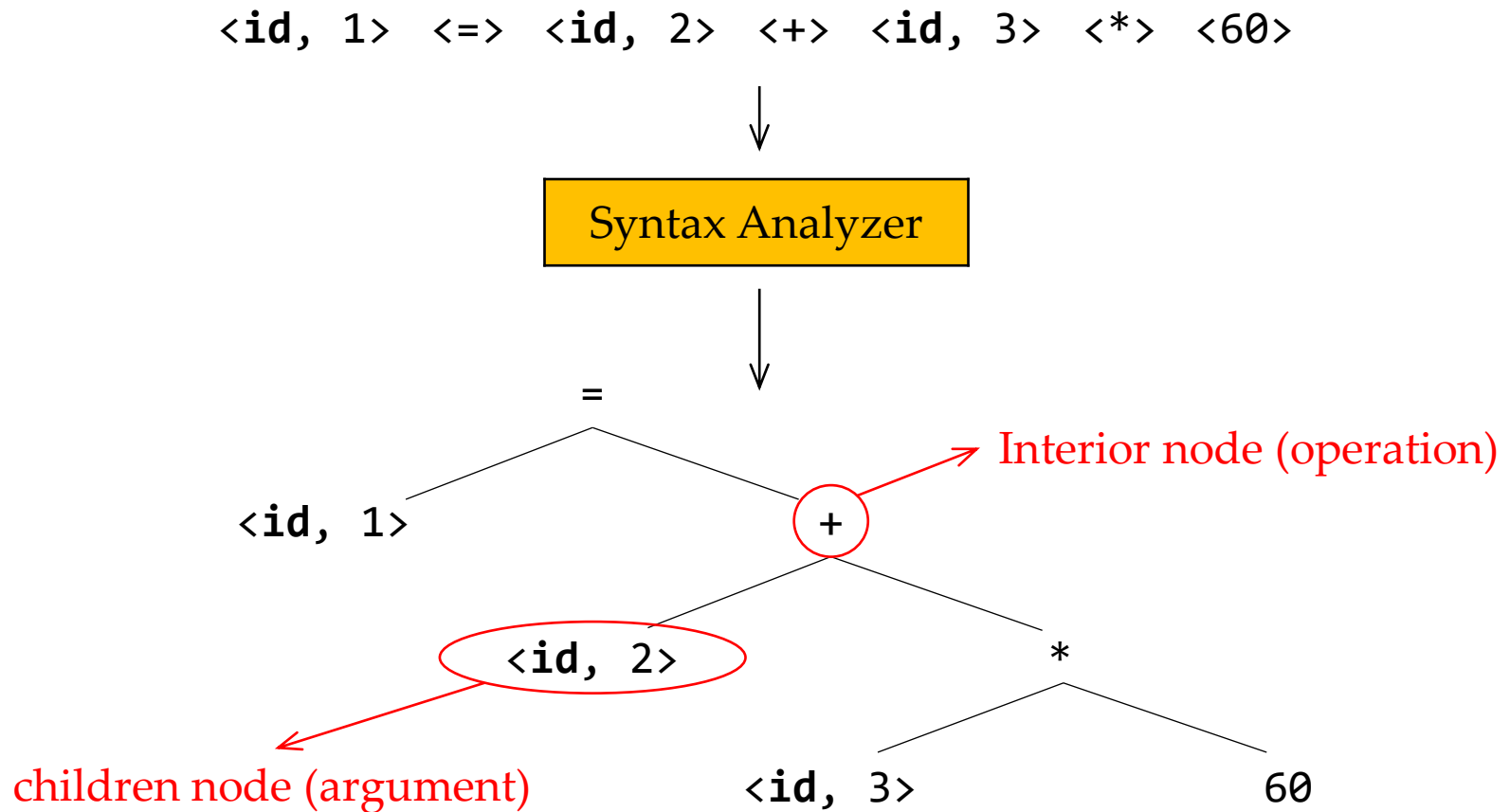
Example adapted from Aiken's notes (Stanford CS143)

Syntax Analysis (Parsing, 语法分析)



- 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

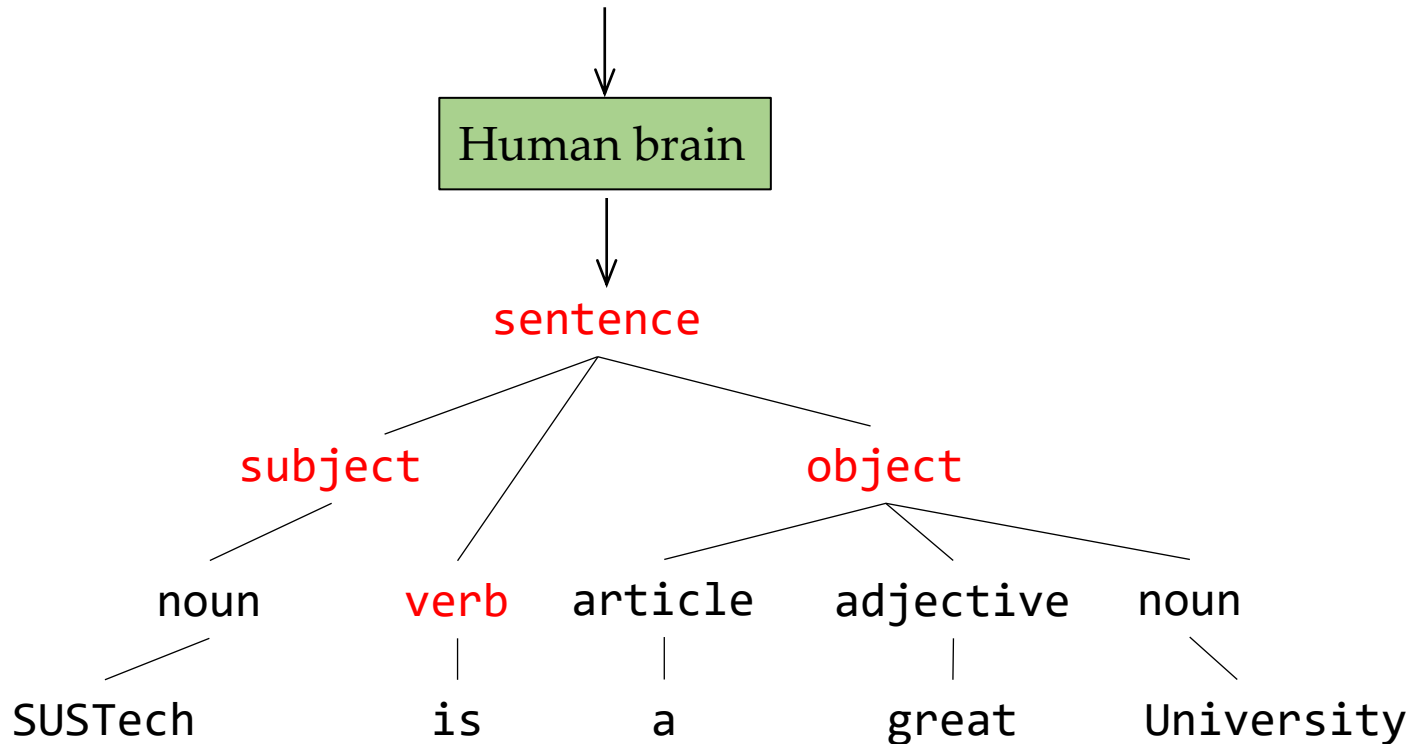
Syntax Analysis (Example)



Syntax Analysis in English

<article, “SUSTech”> <verb, “is”> <article, “a”>

<adjective, “great”> <noun, “university”>



Semantic Analysis (语义分析)



- The semantic analyzer uses the syntax tree and the information in the symbol table to check the source program for **semantic consistency** with the language definition
- Also gathers **type information** 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)

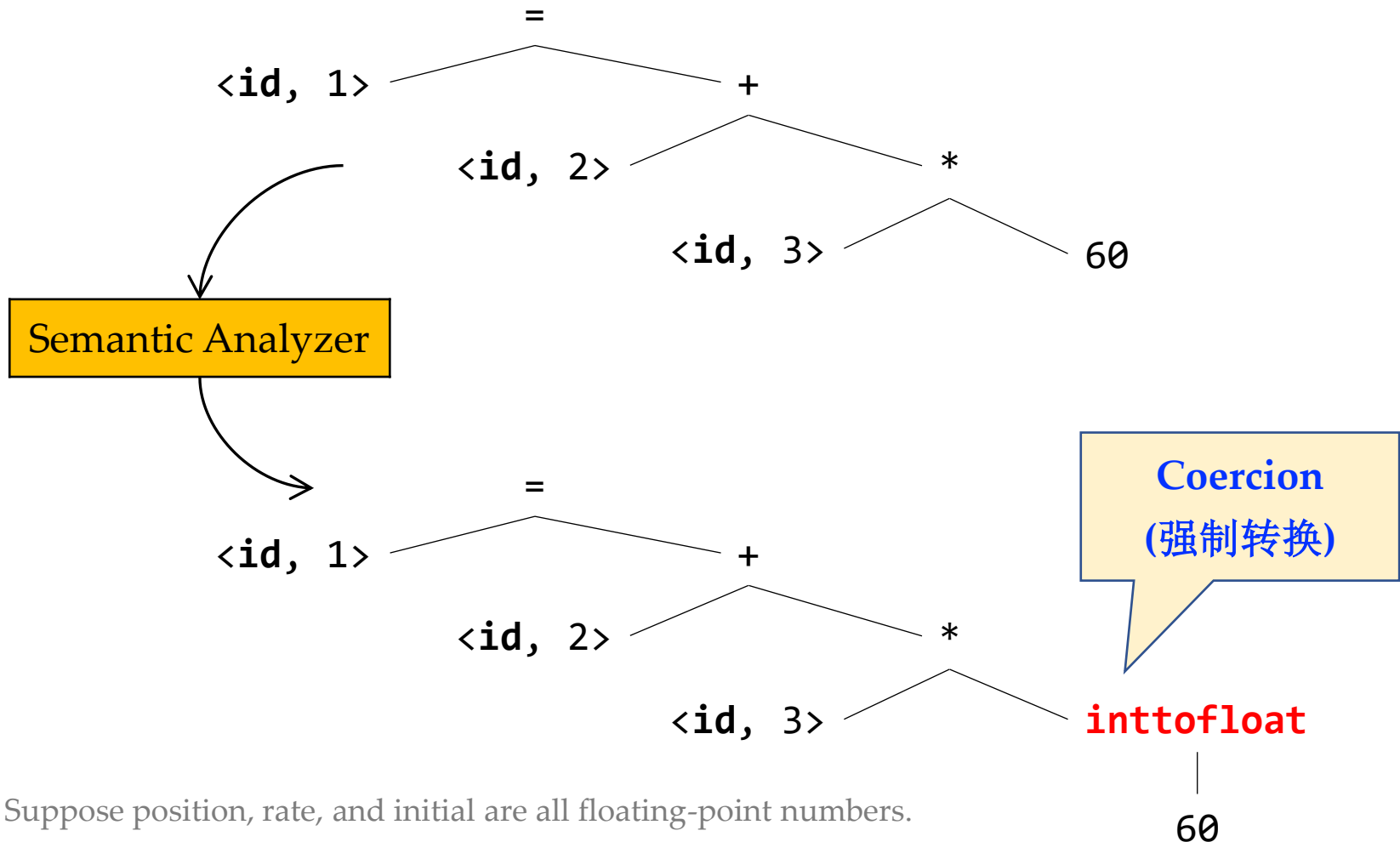
- An important part of semantic analysis is type checking
- Compilers check that each operator has matching operands (of correct types)

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)

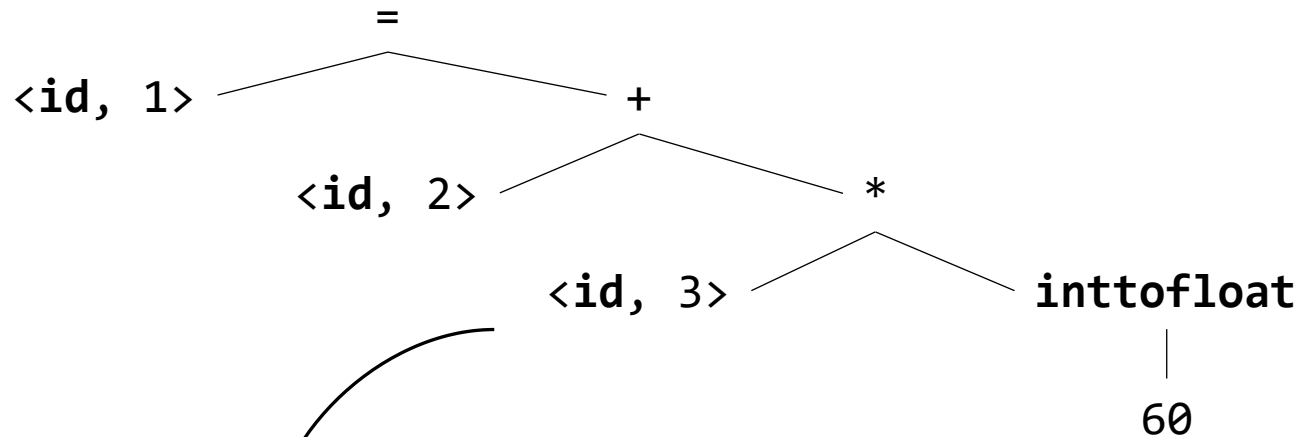


Intermediate Code Generation (中间代码生成)



- 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



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)



- 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. **t3** is only used for value transmitting

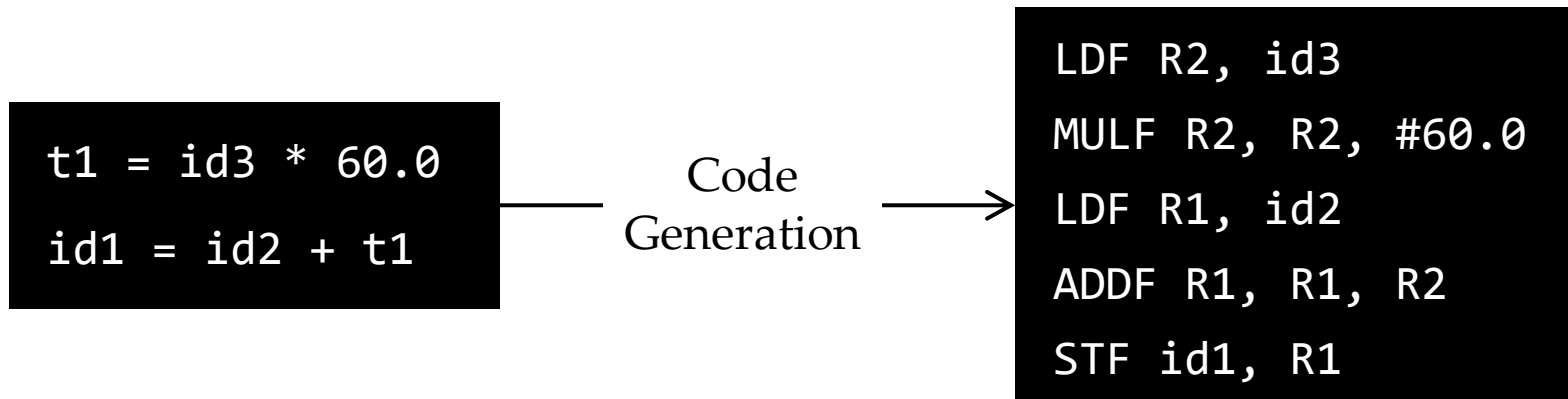
Optimization

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


Code Generation (代码生成)



- Map IR to target language, analogous to human translation
- It is crucial to **allocate register and memory** to hold values

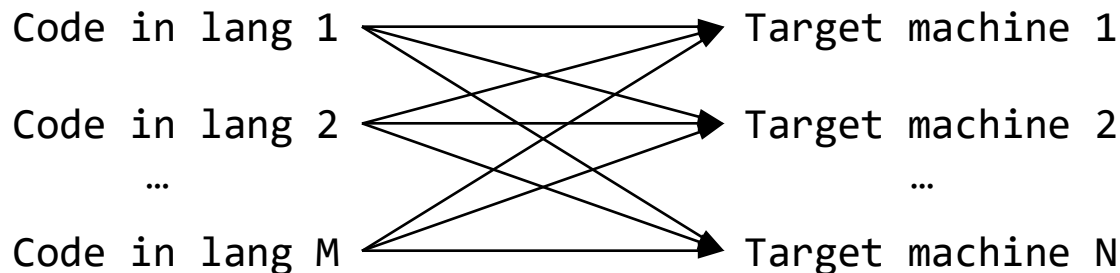


Symbol Table Management

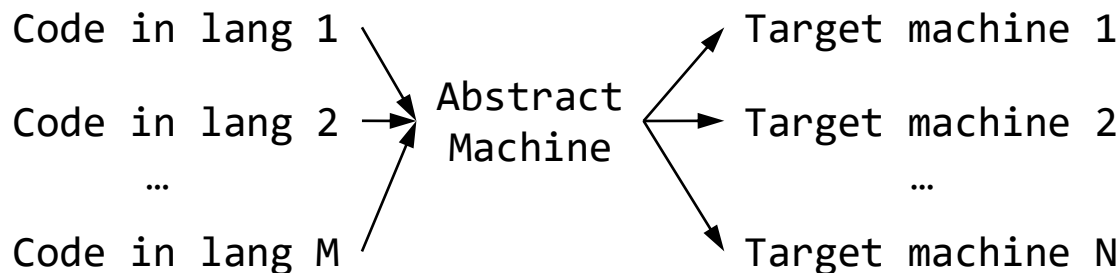
- Performed by the frontend, symbol table is passed along with the intermediate code to the backend
- Record the variable names and various attributes
 - storage allocated, type, scope
- Record the procedure names and various attributes
 - the number and type of arguments
 - the way of passing arguments (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



$M * N$ compilers
without a good IL



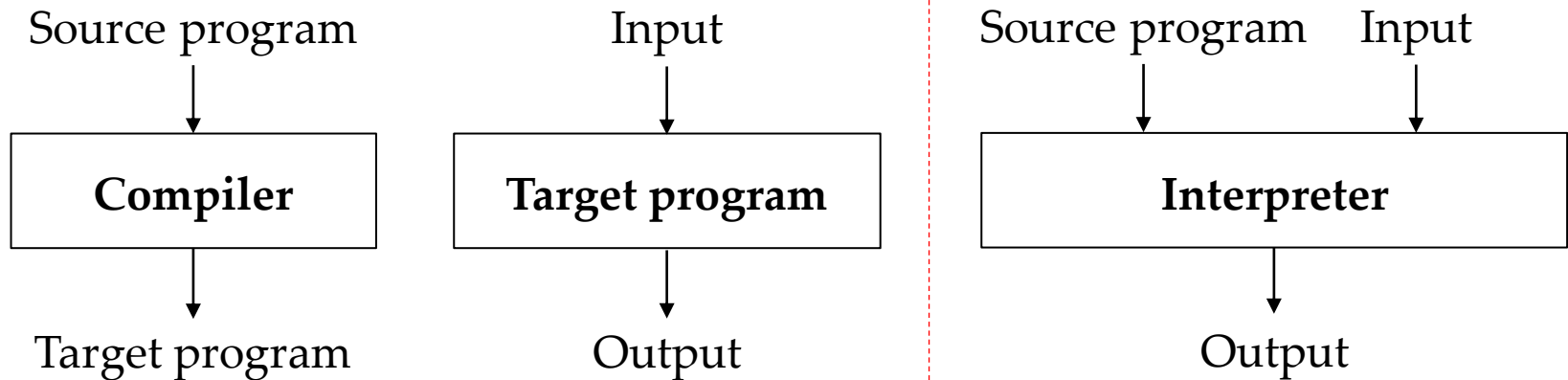
$M + N$ compilers
with a good IL

Compilers vs. Interpreters

1

A **compiler** 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.



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