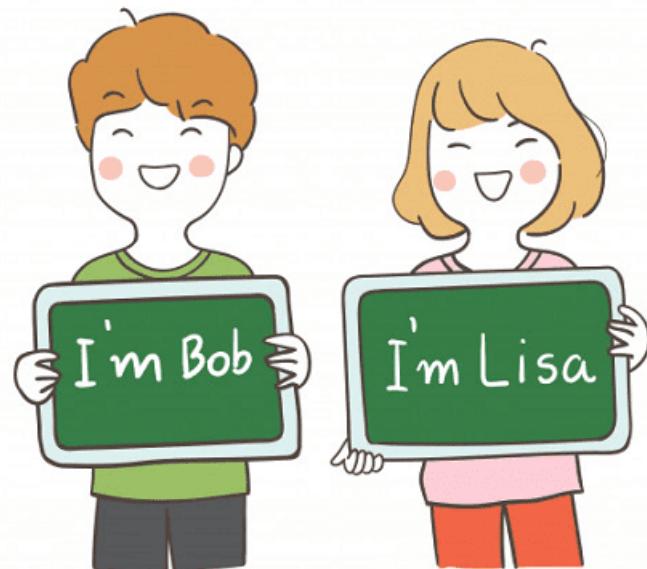


# Compilation

0368-3133

## Lecture 1a: Introduction



# Administrata



# Staff

- Lecturer: Noam Rinetzky
  - maon@cs.tau.ac.il
  - Check Point 342
  - Office hours: set by mail
- T.A.: Noa Schiller

# Lectures & Recitations

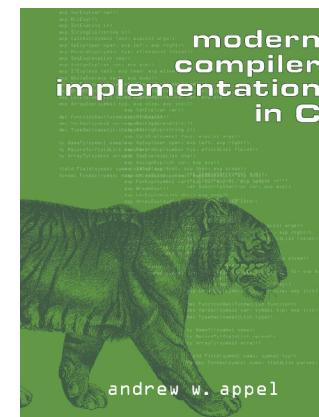
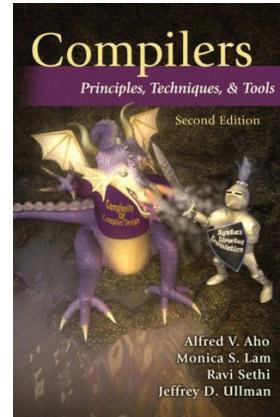
- Lecture (1): Tuesday 9-12 → “Theory”
- Recitations (3): Thursday 16-19 → “Practice”
- Physical participation is highly recommended!
  - Limited online participation (no Q/A)\*
  - Videos will be available



\*Indirectly via chats → friends usually works

# Textbooks

- Modern Compiler Design
  - Grune et al.
- Compilers: Principles, Techniques and Tools
  - Aho et al.
- Modern Compiler Implementation in Java/C/ML
  - Andrew W. Appel



# Practical Compiler Project

- Build a compiler  $L$  (*an OOP language*)  $\Rightarrow$  **MIPS**
  - In Java
  - Groups of 3
  - Managed by the T.A.



# Grades

- Compiler project 40%
  - Ex 1: 3% (lexical analysis)
  - Ex 2: 4% (syntax analysis)
  - Ex 3: 8% (semantic analysis)
  - Ex 4: 10% (intermediate representation + data-flow analysis)
  - Ex 5: 15% (code generation)
- Final exam 60% - must pass!



# Workload—A Heads Up

- Compiler project
  - Ex 1: 3%
  - Ex 2: 4%
  - Ex 3: 8%  
(semantic analysis)
  - Ex 4: 10%
  - Ex 5: 15%  
(code generation)
- Final exam 6%

The compiler project may be the most significant programming task in your studies

$\text{Grade} \sim \log(|\text{Work}|)$

You will eat your own dog food!



# Students in Reserve Duty (Group B/C)

- Can submit only some of the exercises
  - Exercise percentage moves to exam
- Can join 3-person groups
  - i.e., groups of 4



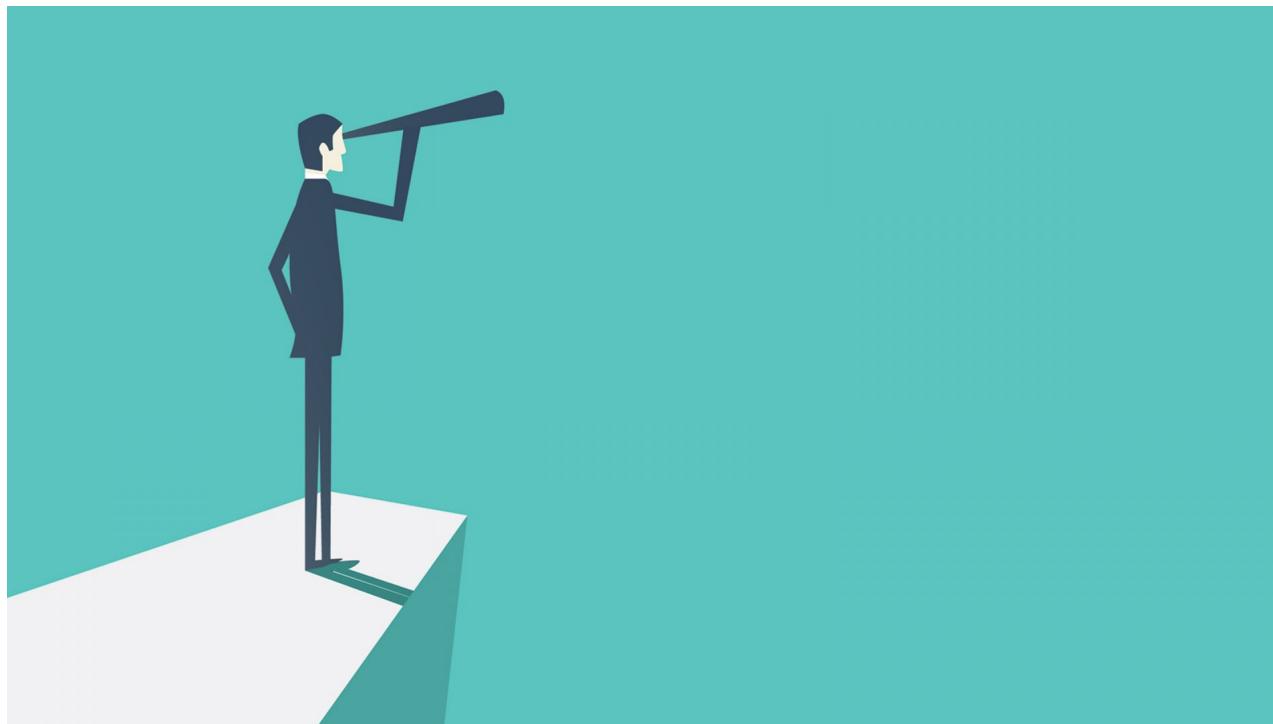
# Students in Reserve Duty (Group C)

- Two options for alternative evaluation:
  - Home exam\* (עובר בינהר) (צון רגיל)
  - Home exam\* + do Ex 5 alone



\*Home exam = Class exam in 7 days

# Course Overview



Modern Compiler Design. Chapter 1

# Course Goals

- What is a compiler
- How does it work
- (Reusable) techniques & tools
- Programming languages implementation
  - Procedure calls, polymorphism, runtime systems (e.g., memory management), ...
- Execution environments
  - Assemblers, linkers, loaders, OS



*To compilers,  
and beyond*

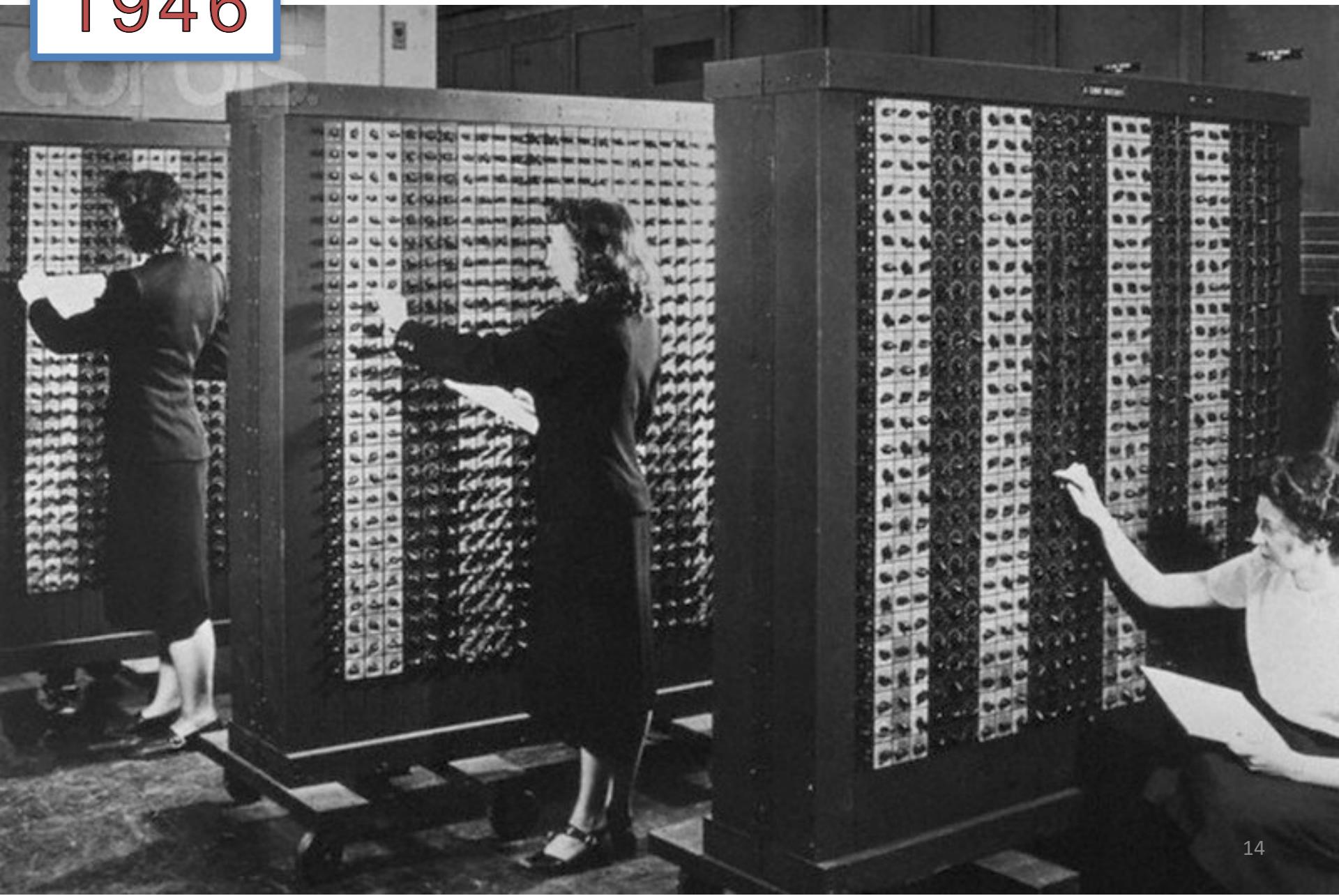
# What is a Compiler?



1946

# ENIAC

Electronic Numerical Integrator  
and Computer



1952

# “The Education of a Computer” (Grace Hopper)

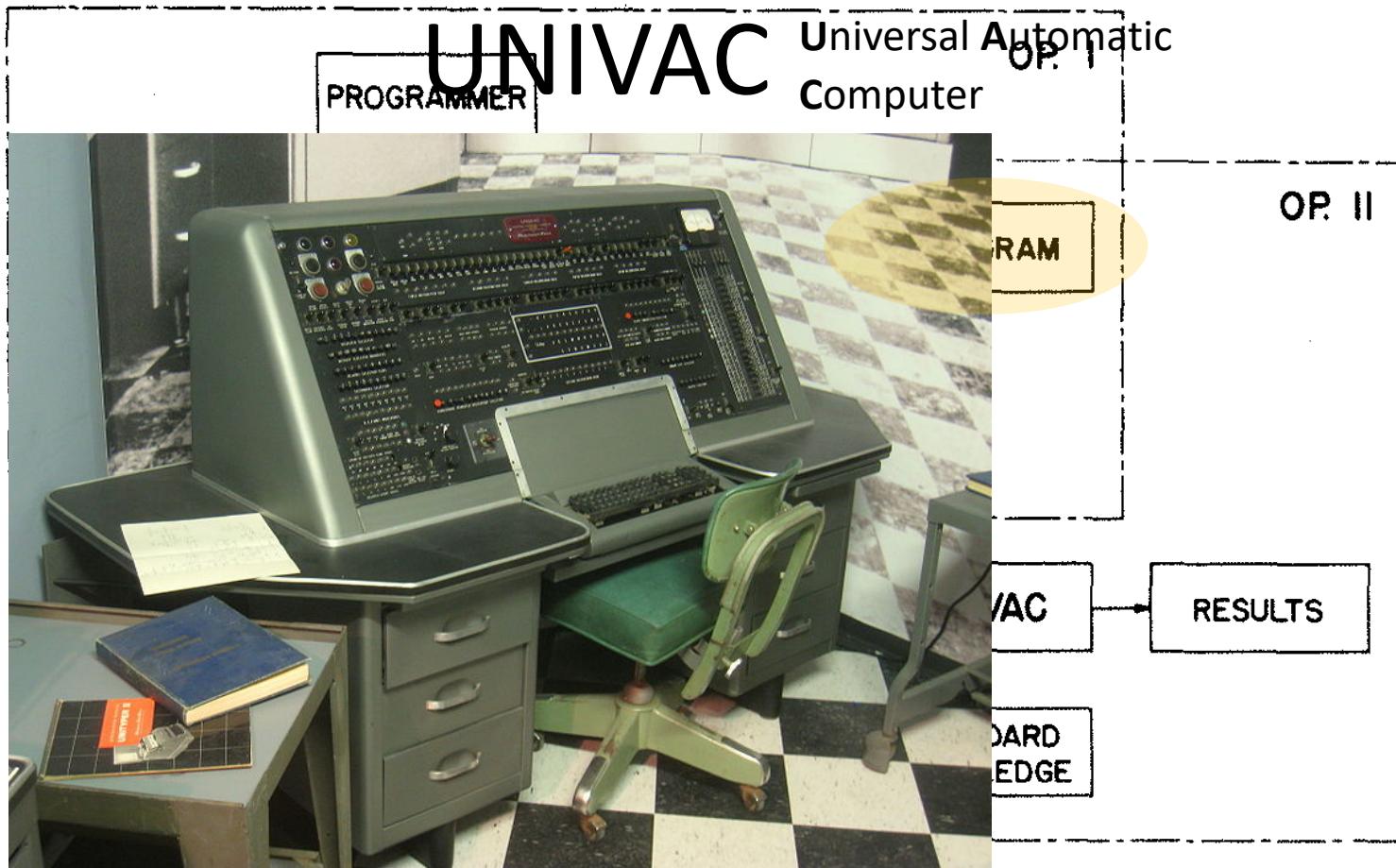


Fig. 4 - SOLUTION OF A PROBLEM

1952

# “The Education of a Computer” (Grace Hopper)

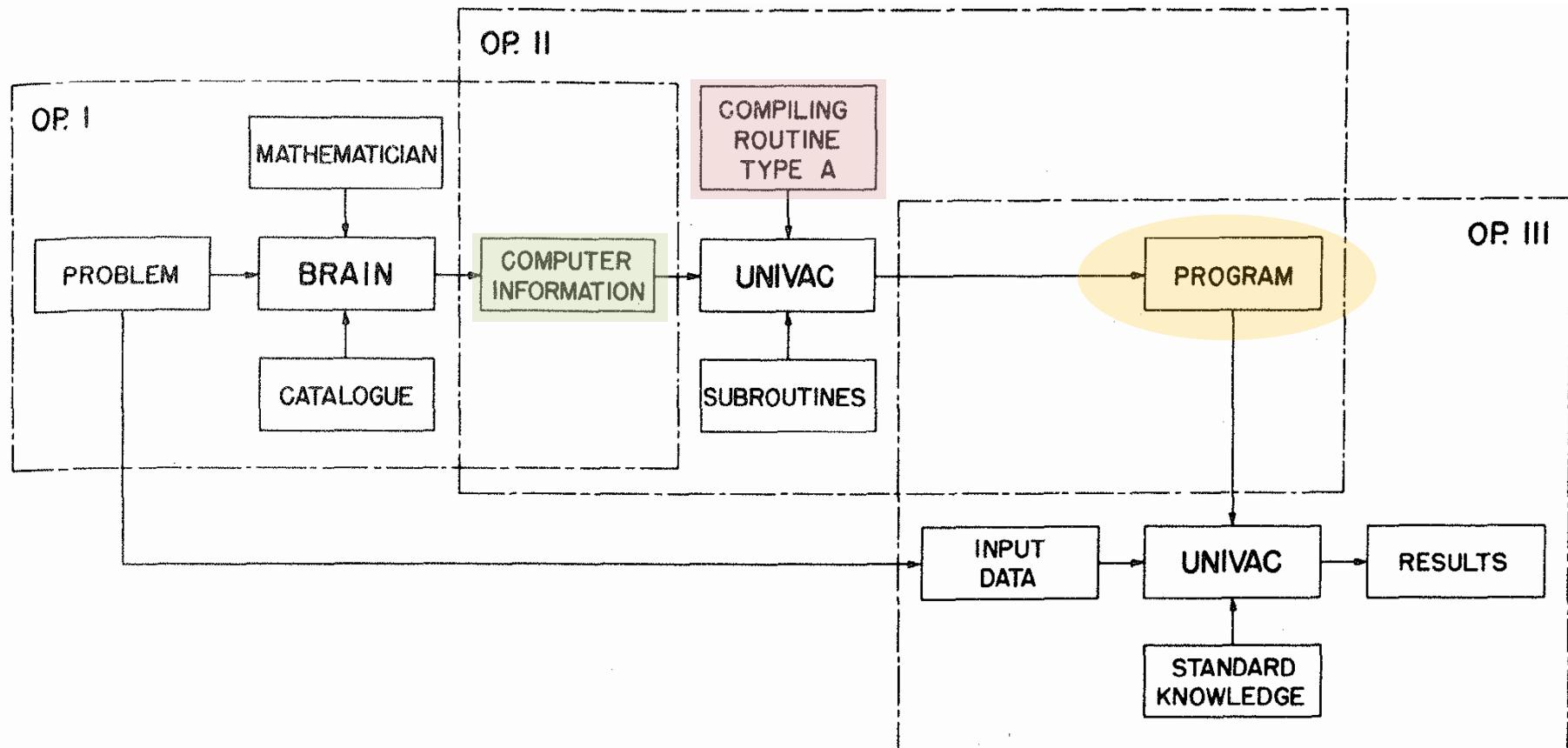
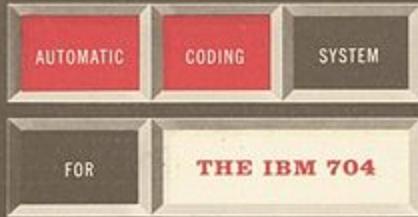


Fig. 5 - COMPILE ROUTINES AND SUBROUTINES

1957

PROGRAMMER'S REFERENCE MANUAL

# Fortran



C STATEMENT NUMBER	FOR COMMENT	COMPUTER IDENTIFICATION	FORTRAN STATEMENT		IDENTIFICATION
			1	2	
C		X	PROGRAM FOR FINDING THE LARGEST VALUE ATTAINED BY A SET OF NUMBERS		
C			DIMENSION A(999)		
			FREQUENCY 30(2,1,10), 5(100)		
			READ 1, N, (A(I), I=1,N)		
			FORMAT (13/(1.2F6.2))		
1			BIGA = A(1)		
5			DO 20 I = 2,N		
30			IF (BIGA-A(I)) 10,20,20		
10			BIGA = A(I)		
20			CONTINUE		
			PRINT 2, N, BIGA		
2			FORMAT (22H1THE LARGEST OF THESE 13, 12H NUMBERS IS F7.2)		
			STOP 77777		

John Backus and team at IBM  
(Fortran = Formula Translator)

The first complete compiler



# High Level Programming Languages

- **Imperative** Algol, PL1, Fortran, Pascal, Ada, Modula, C
  - Closely related to “von Neumann” Computers
- **Object-oriented** Simula, Smalltalk, Modula3, C++, Java, C#, Python
  - Data abstraction and ‘evolutionary’ form of program development
    - Class as implementation of an abstract data type (data+code)
    - Objects instances of a class
    - Inheritance + generics
- **Functional** Lisp, Scheme, ML, Miranda, Hope, Haskel, OCaml, F#
  - First class functions
- **Logic Programming** Prolog

2020

# What is a Compiler?

“A compiler is **computer software** that **translates** computer code written in one programming language (the **source language**) into another language (the **target language**). ...primarily to a lower level language (e.g. assembly language, object code, or machine code) to create an **executable program**.”



-- *Wikipedia*

# What is a Compiler?

source language

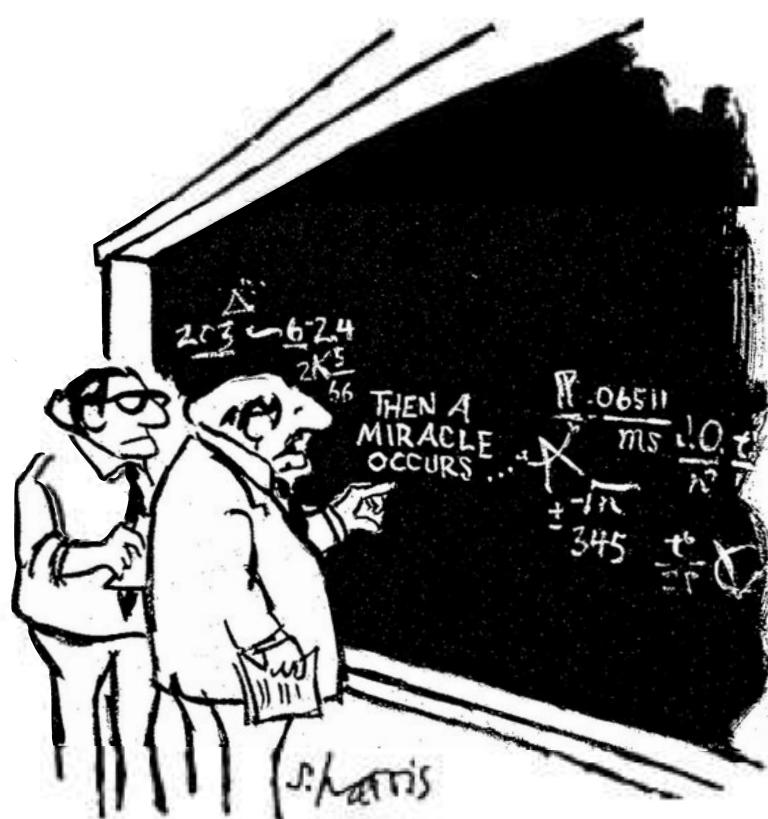
C  
C++  
Pascal  
Java

Perl      txt  
JavaScript  
Python  
Ruby      text  
**Source**

Lisp  
Scheme  
ML  
OCaml

Prolog

Postscript  
TeX



target language

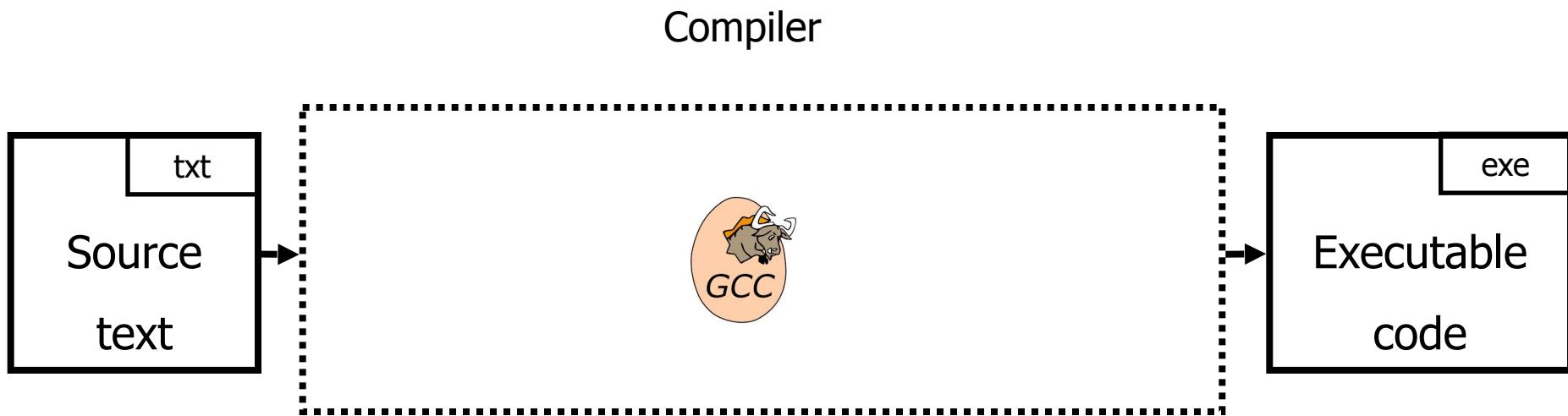
IA32  
IA64  
ARM  
SPARC

exe  
Java Bytecode  
**Executable**  
C      code  
C++  
Pascal  
Java

JavaScript

PDF  
Bit Map  
...

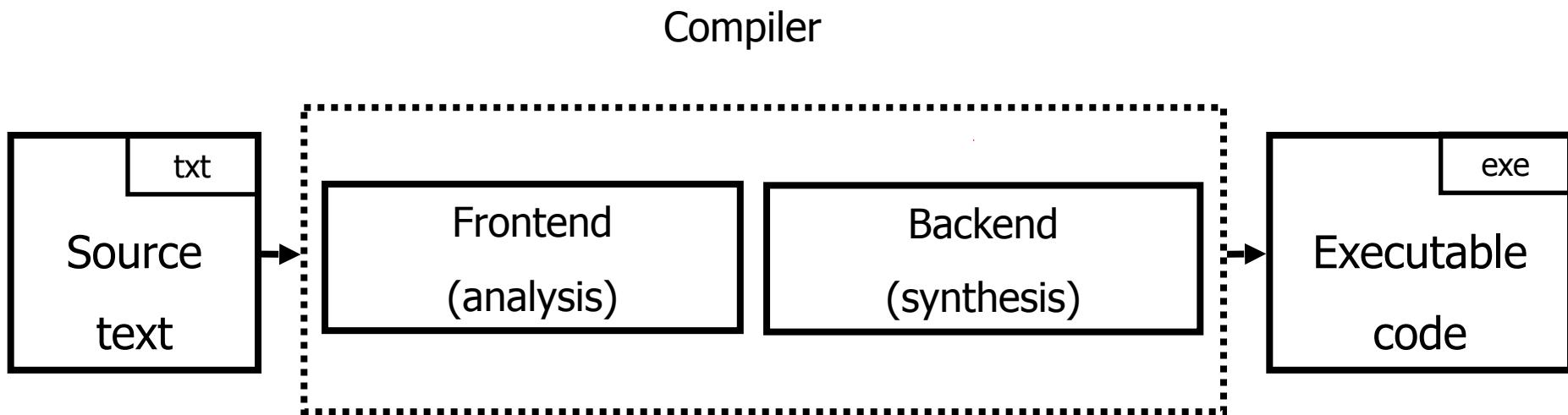
# What is a Compiler?



```
int a, b;  
a = 2;  
b = a*2 + 1;
```

```
MOV R1,2  
SAL R1  
INC R1  
MOV R2,R1
```

# Anatomy of a (2-Stage) Compiler



```
int a, b;  
a = 2;  
b = a*2 + 1;
```

Why?

```
MOV R1, 2  
SAL R1  
INC R1  
MOV R2, R1
```

```

var
  x, y : integer;
begin
  x := 2;
  y := x*2 + 1;
end.

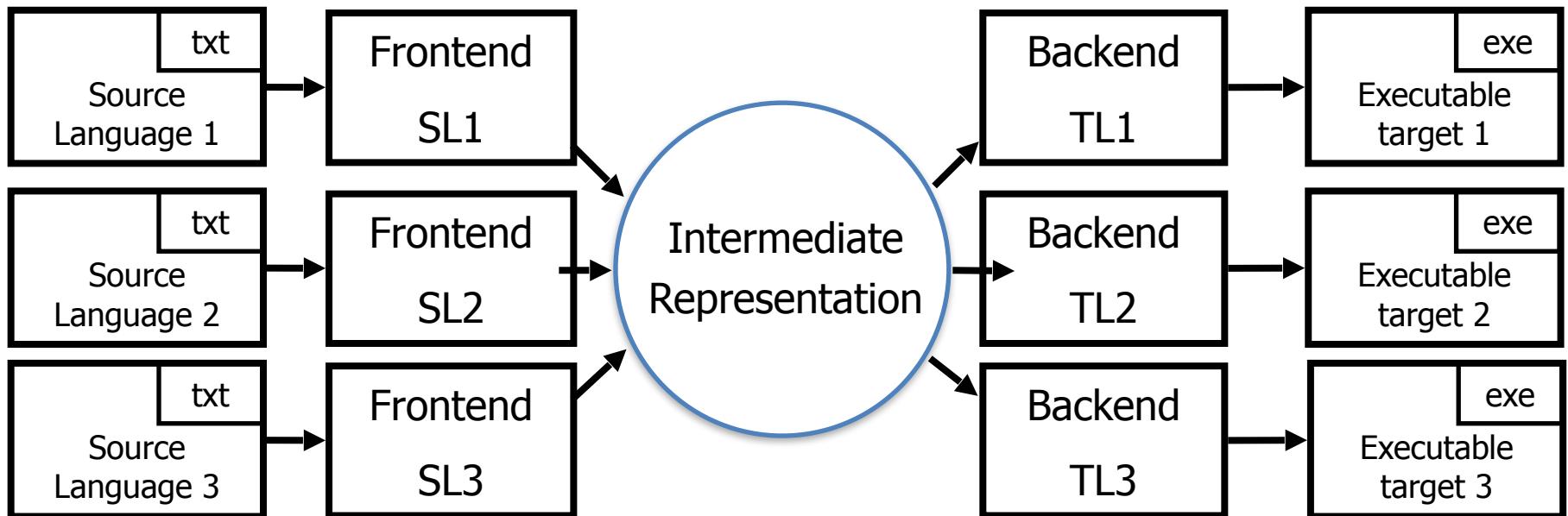
```

# Modularity

```

iconst 2
iconst 1
ishl
istore_2
iinc 2, 1

```



```

int x, y;
x = 2;
y = x*2 + 1;

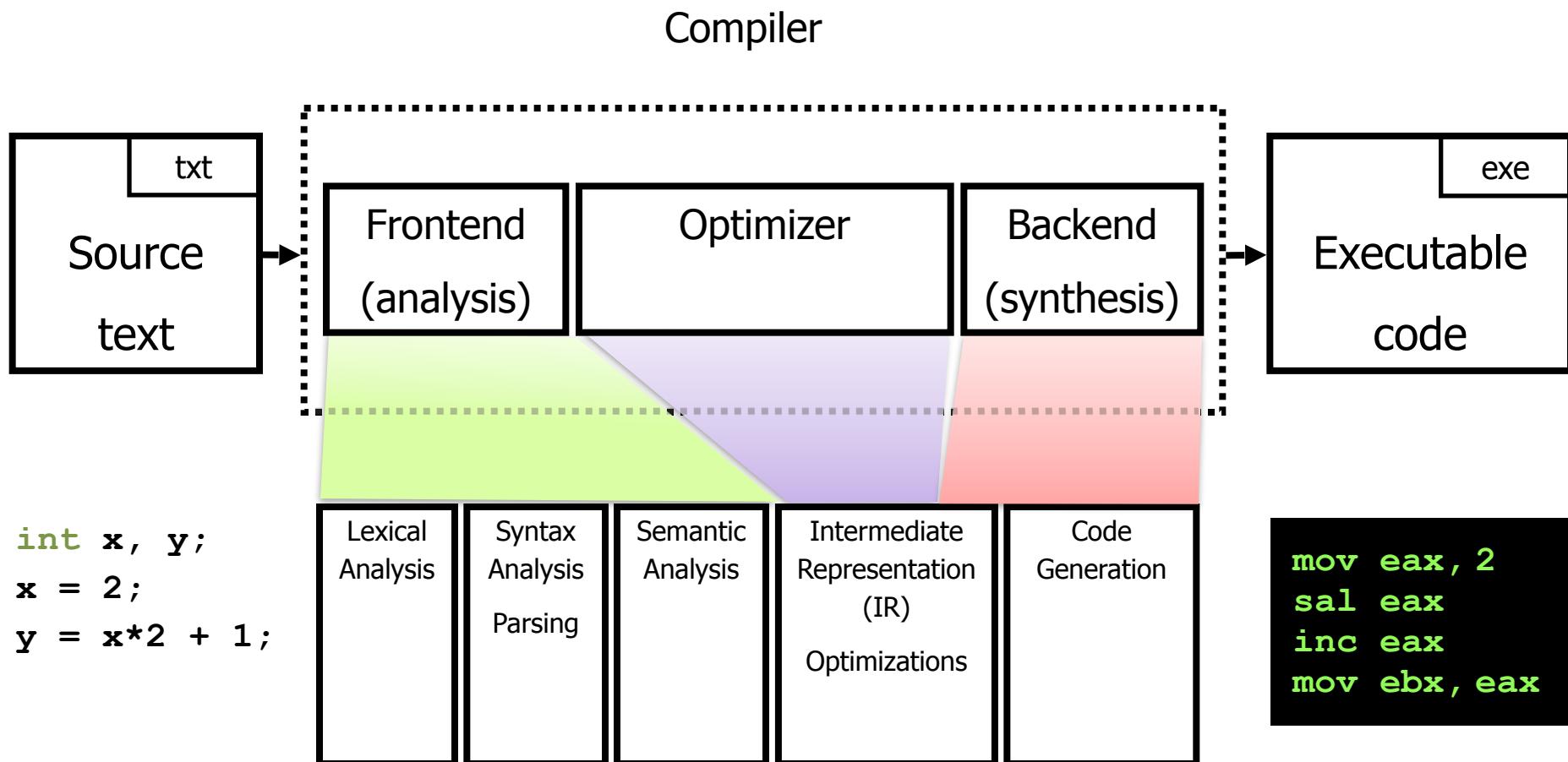
```

```

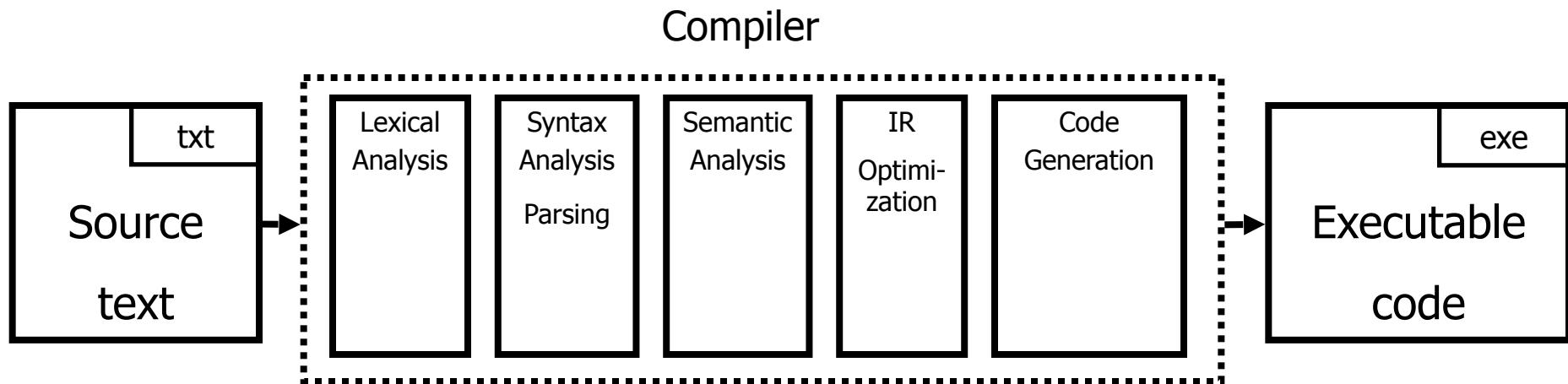
mov eax, 2
sal eax
inc eax
mov ebx, eax

```

# Anatomy of a (3-Stage) Compiler



# Journey inside a compiler



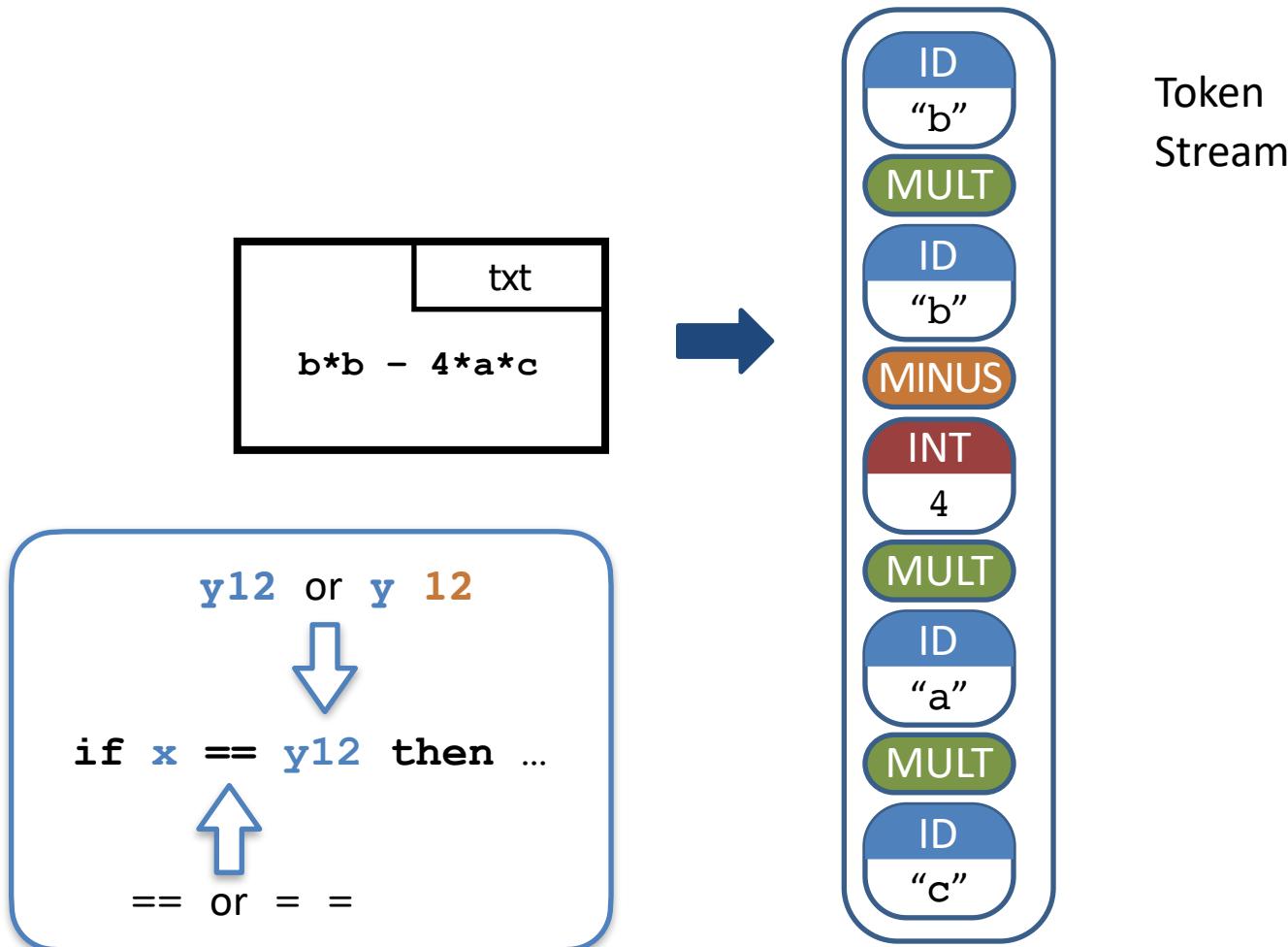
# Journey inside a compiler

txt

$b^2 - 4ac$

Recognizing  
words

# Journey inside a compiler



Lexical  
Analysis

Syntax  
Analysis

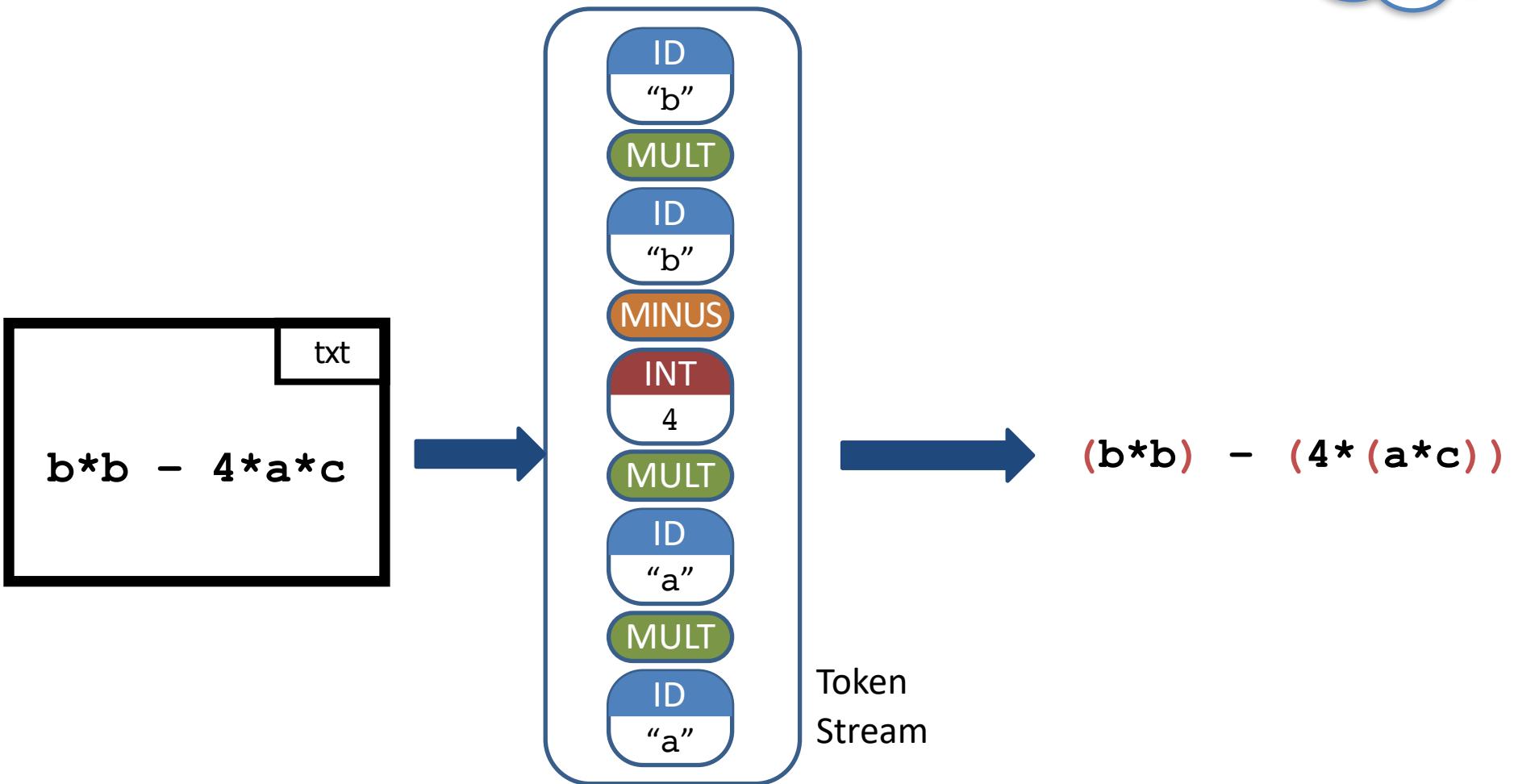
Semantic  
Analysis

IR  
Opt.

Code  
Gen.

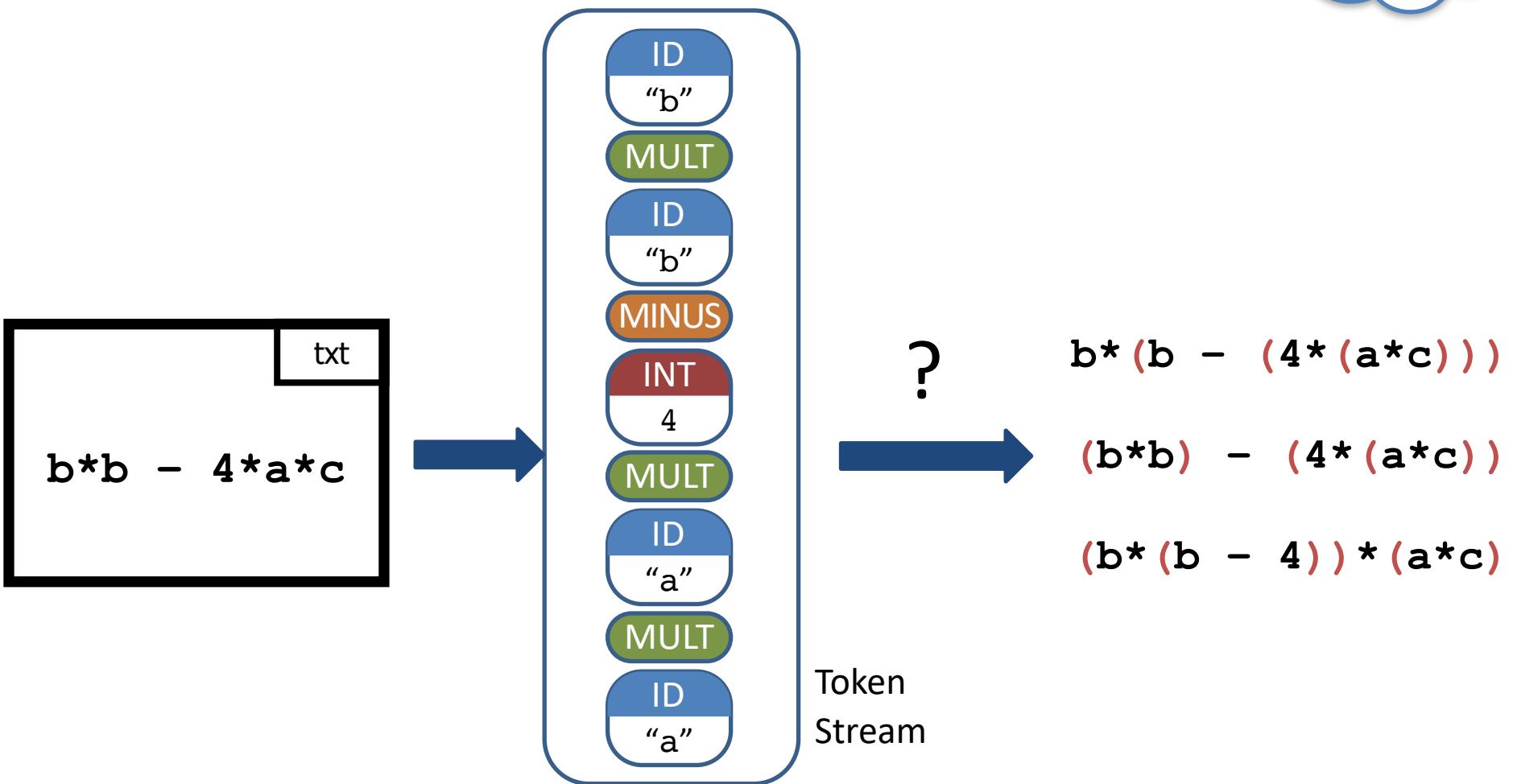
Understanding  
Sentence  
Structure

# Journey inside a compiler



Understanding  
Sentence  
Structure

# Journey inside a compiler



# Journey inside a compiler

Grammar

```

ID      "b"
MULT
ID      "b"
MINUS
INT     4
MULT
ID      "a"
MULT
ID      "a"
  
```

$E \rightarrow E \text{ PLUS } T$

$E \rightarrow E \text{ MINUS } T$

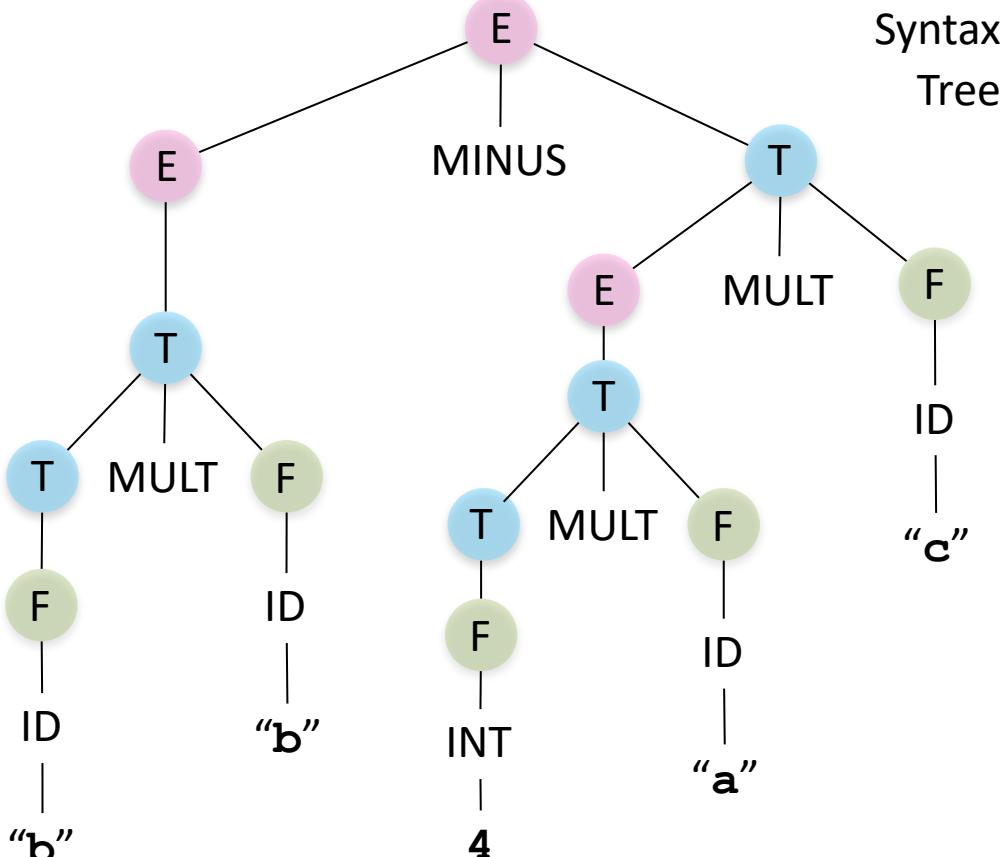
$T \rightarrow T \text{ MULT } F$

$T \rightarrow T \text{ DIV } F$

$F \rightarrow \text{ID}$

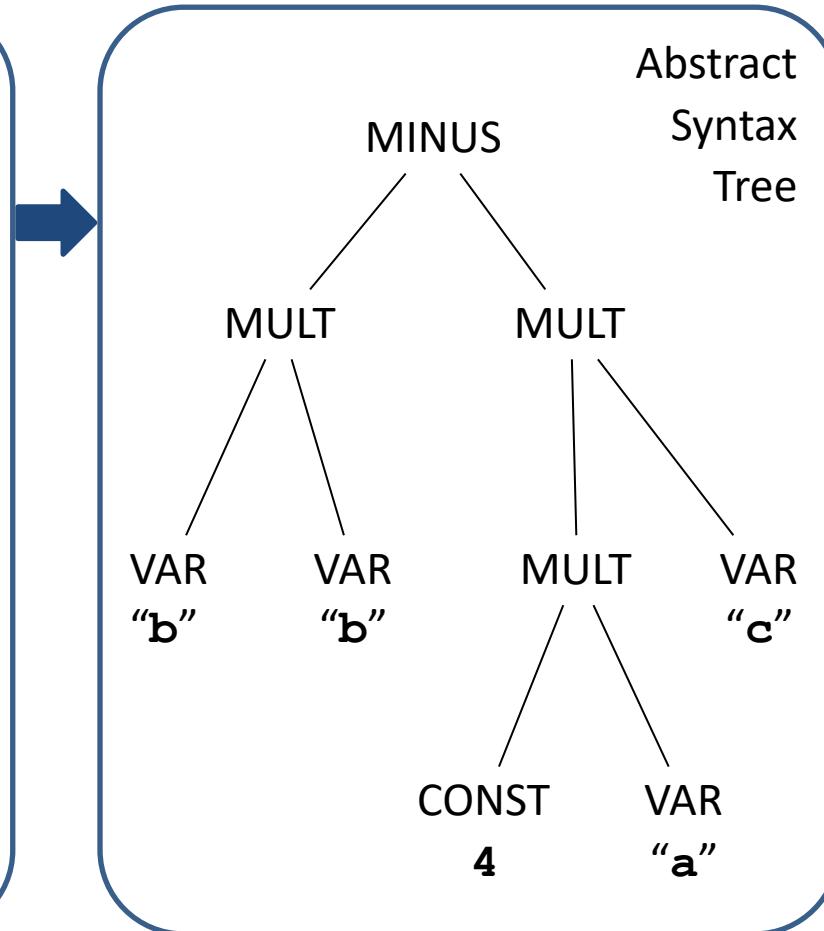
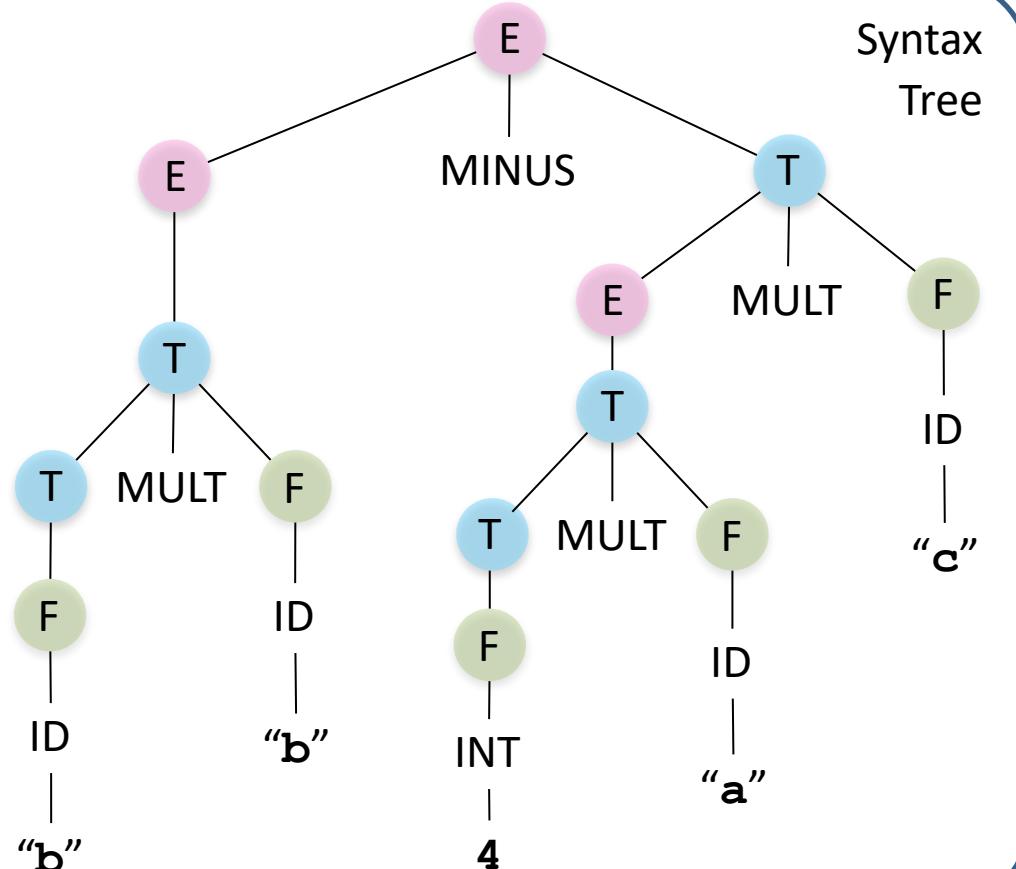
$F \rightarrow \text{INT}$

Token Stream



# Journey inside a compiler

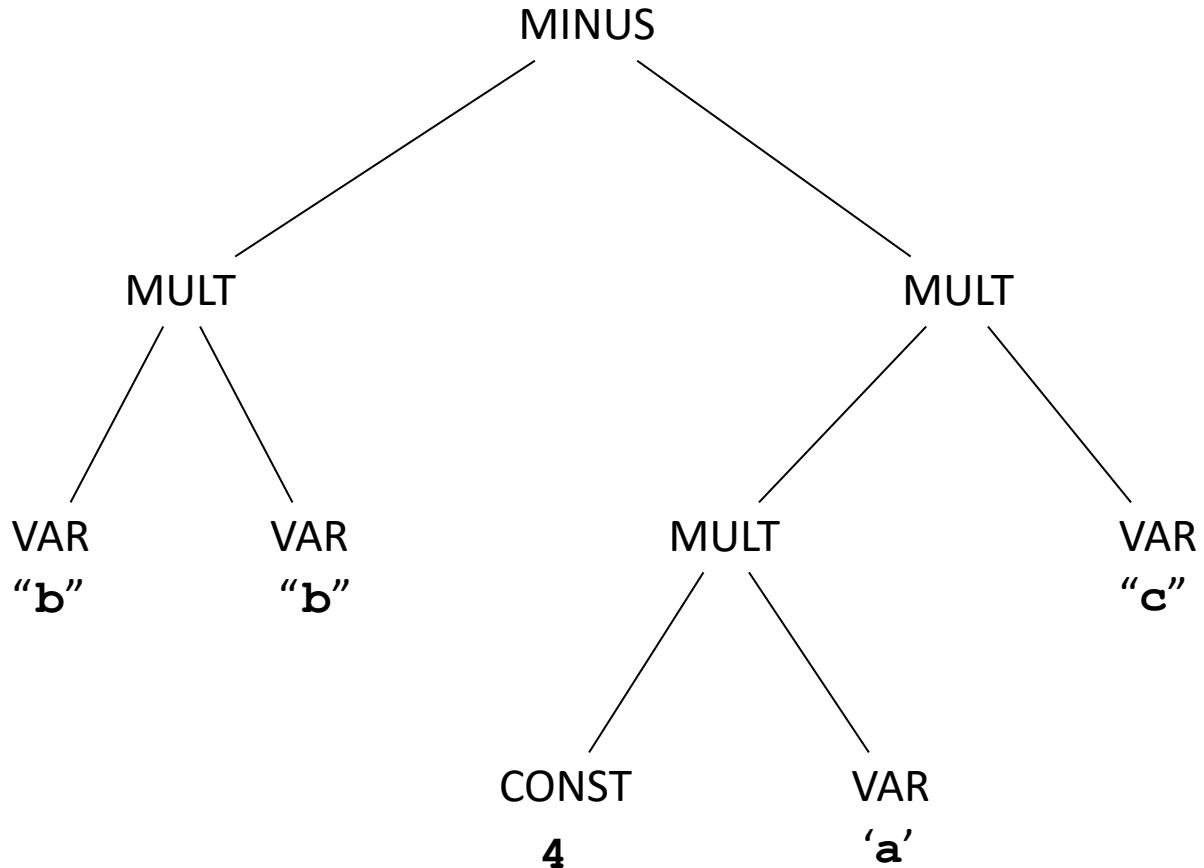
Extracting  
essence



# Journey inside a compiler

Understanding  
“meaning”

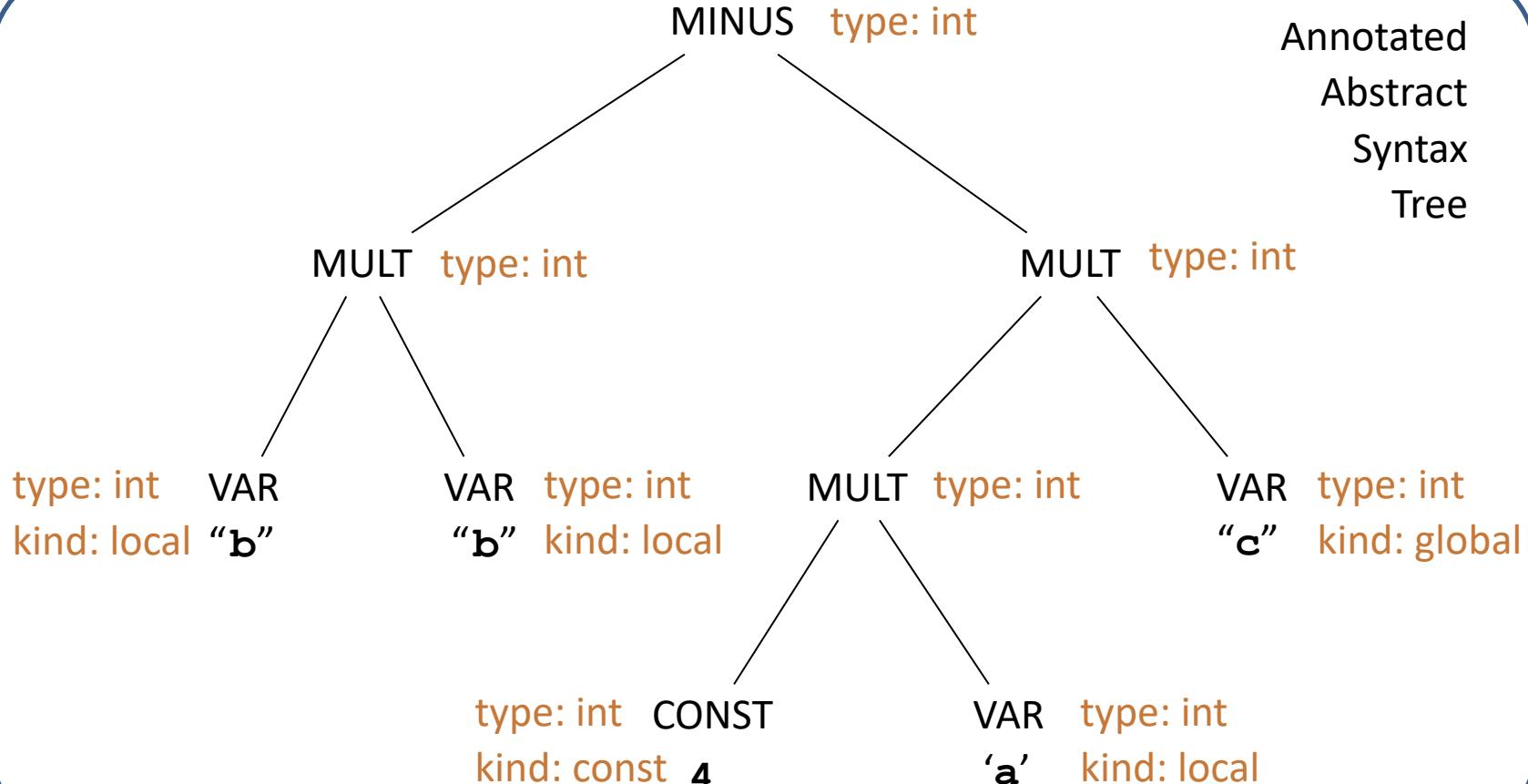
Abstract  
Syntax  
Tree



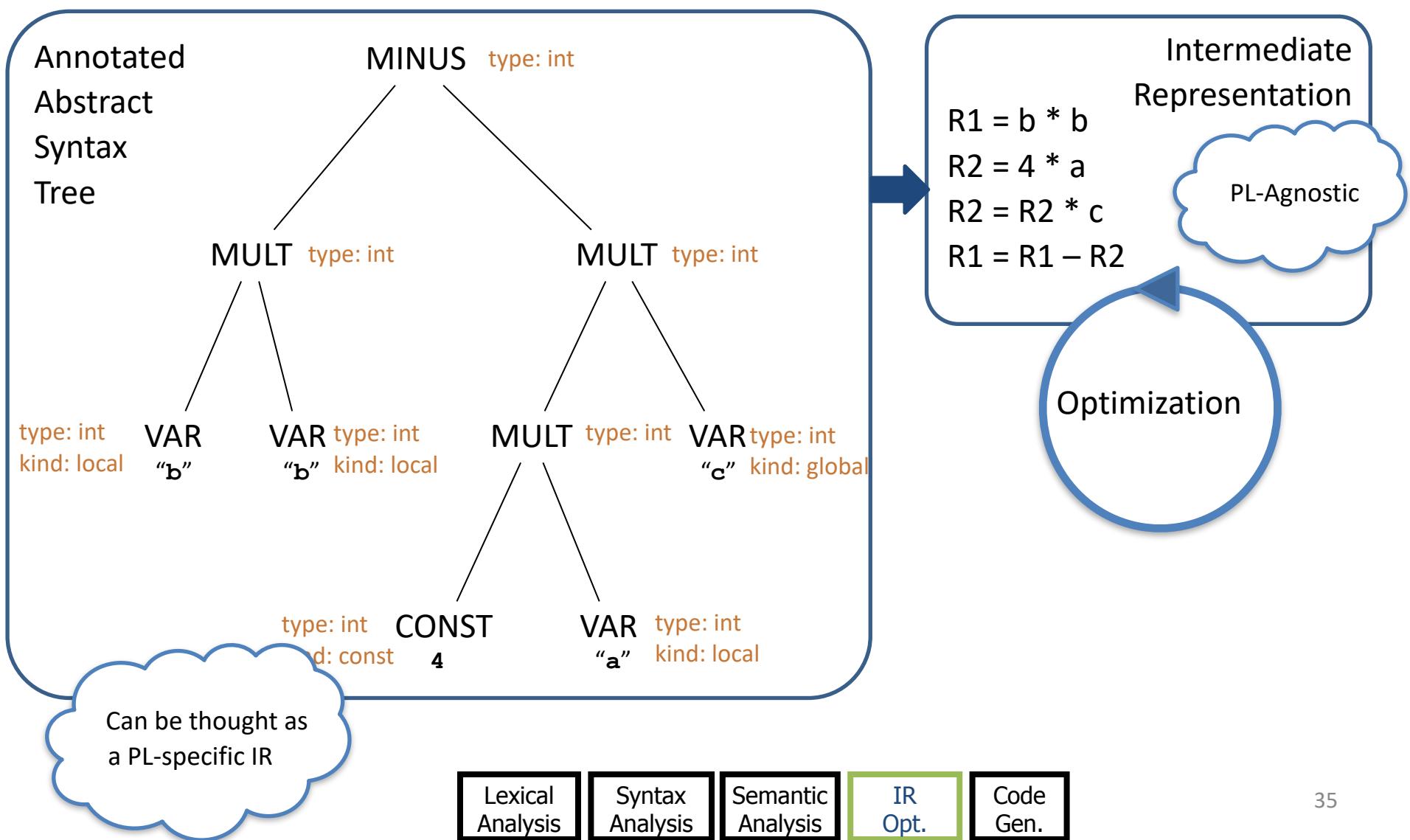
# Journey inside a compiler

Understanding  
“meaning”

Annotated  
Abstract  
Syntax  
Tree



# Journey inside a compiler



# IR Optimizations

- Many optimization heuristics
  - Constant propagation
    - Leverage compile-time information to save work at runtime (pre-computation)

```
width = 10;  
length = 20;  
z = x * width * length;
```

Actually,  
done on IR

---

```
width = 10;  
length = 20;  
z = x * 200;
```



# IR Optimizations

- Many optimization heuristics
  - Constant propagation
    - Leverage compile-time information to save work at runtime (pre-computation)
  - Dead code elimination

```
x = y;  
z = t;  
x = y + z;
```

Actually,  
done on IR

---

```
z = t;  
x = y + z;
```

# IR Optimizations

- Many optimization heuristics
  - Constant propagation
    - Leverage compile-time information to save work at runtime (pre-computation)
  - Dead code elimination
  - Loop optimizations: **hoisting**, unrolling, ...

```
for (int i = 0; i < 100; ++i) {  
    array[i] = x + y;  
}
```

Actually,  
done on IR

---

```
int t = x + y;  
for (int i = 0; i < 100; ++i) {  
    array[i] = t;  
}
```

# IR Optimizations

- Many optimization heuristics
  - Constant propagation
    - Leverage compile-time information to save work at runtime (pre-computation)
  - Dead code elimination
  - Loop optimizations: hoisting, **unrolling**, ...

```
for (int i = 0; i < 80; ++i) {  
    delete array[i];  
}
```

Actually,  
done on IR

---

```
for (int i = 0; i < 80; i += 4) {  
    delete array[i];  
    delete array[i+1];  
    delete array[i+2];  
    delete array[i+3];  
}
```

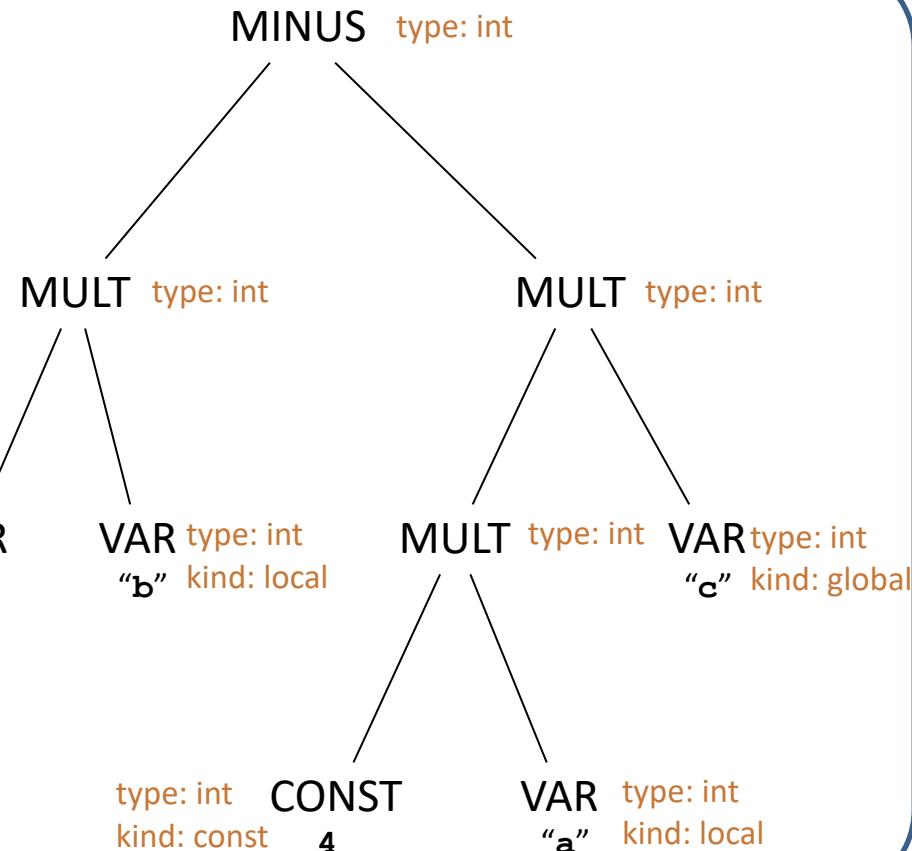
# IR Optimizations

- Many optimization heuristics
  - Constant propagation
    - Leverage compile-time information to save work at runtime (pre-computation)
  - Dead code elimination
  - Loop optimizations: hoisting, unrolling, ...
- “Optimal code” is out of reach
  - Undecidable/NP-complete problems
  - Use approximation and/or heuristics
  - Must preserve correctness, should (mainly) improve code
- The majority of compilation time is spent in optimization



# Journey inside a compiler

Annotated  
Abstract  
Syntax  
Tree



Intermediate  
Representation

```
R1 = b * b  
R2 = 4 * a  
R2 = R2 * c  
R1 = R1 - R2
```

Assembly  
Code

```
mov eax,[esp+16]  
mul eax,[esp+16]  
mov ebx,[esp+8]  
sal ebx,2  
mul ebx,[sp+24]  
sub eax,ebx
```

# Machine code generation

- Assigning variables to memory locations
- Register allocation
  - Optimal register assignment is NP-Complete
  - In practice, known heuristics perform well
- Instruction selection
  - Convert IR to actual machine instructions
  - Machine code optimizations: peephole optimization

R2 = 4 \* a

IR

mov ebx,[esp+8]

mul ebx, 4

Assembly

mov ebx,[esp+16]

sal ebx, 2

Assembly

- Modern architecture challenges
  - Multicores
  - Memory hierarchies
  - SIMD instructions

(not covered)

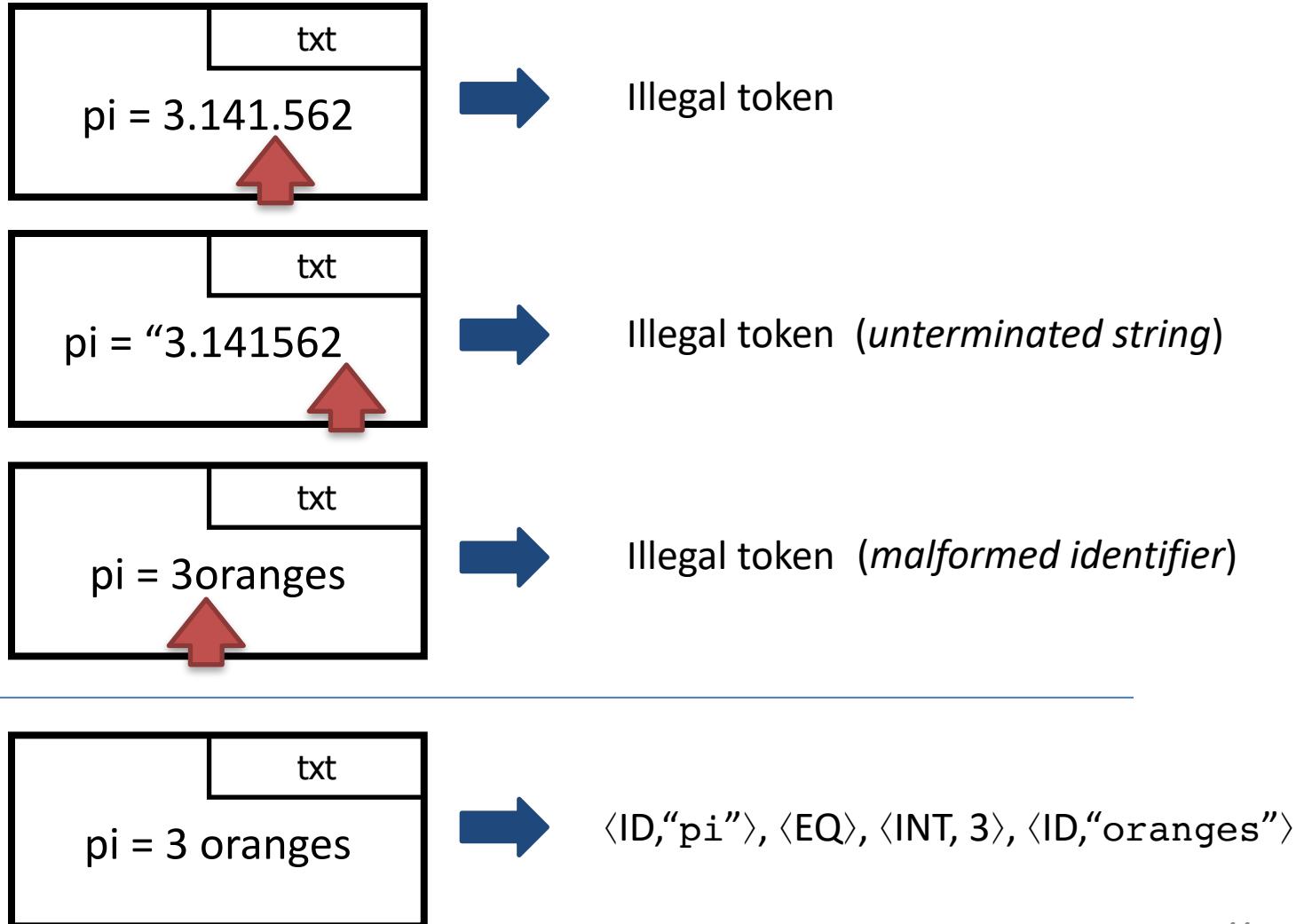
# Error Checking

*Check on every stage.  
Report as soon as possible.*

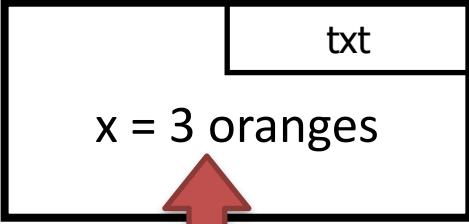
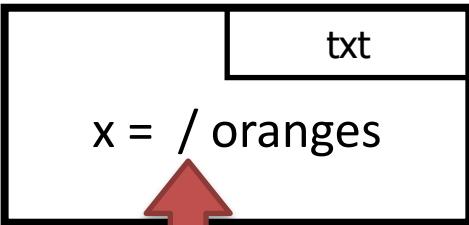
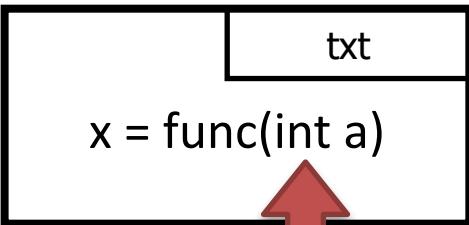
- Lexical analysis: illegal tokens
- Syntax analysis: illegal syntax
- Semantic analysis: incompatible types, undefined variables, ...
- Even at runtime: division by zero, array bounds,...
  - (Compiler generates the error checking code)
- Every phase tries to recover and proceed with compilation



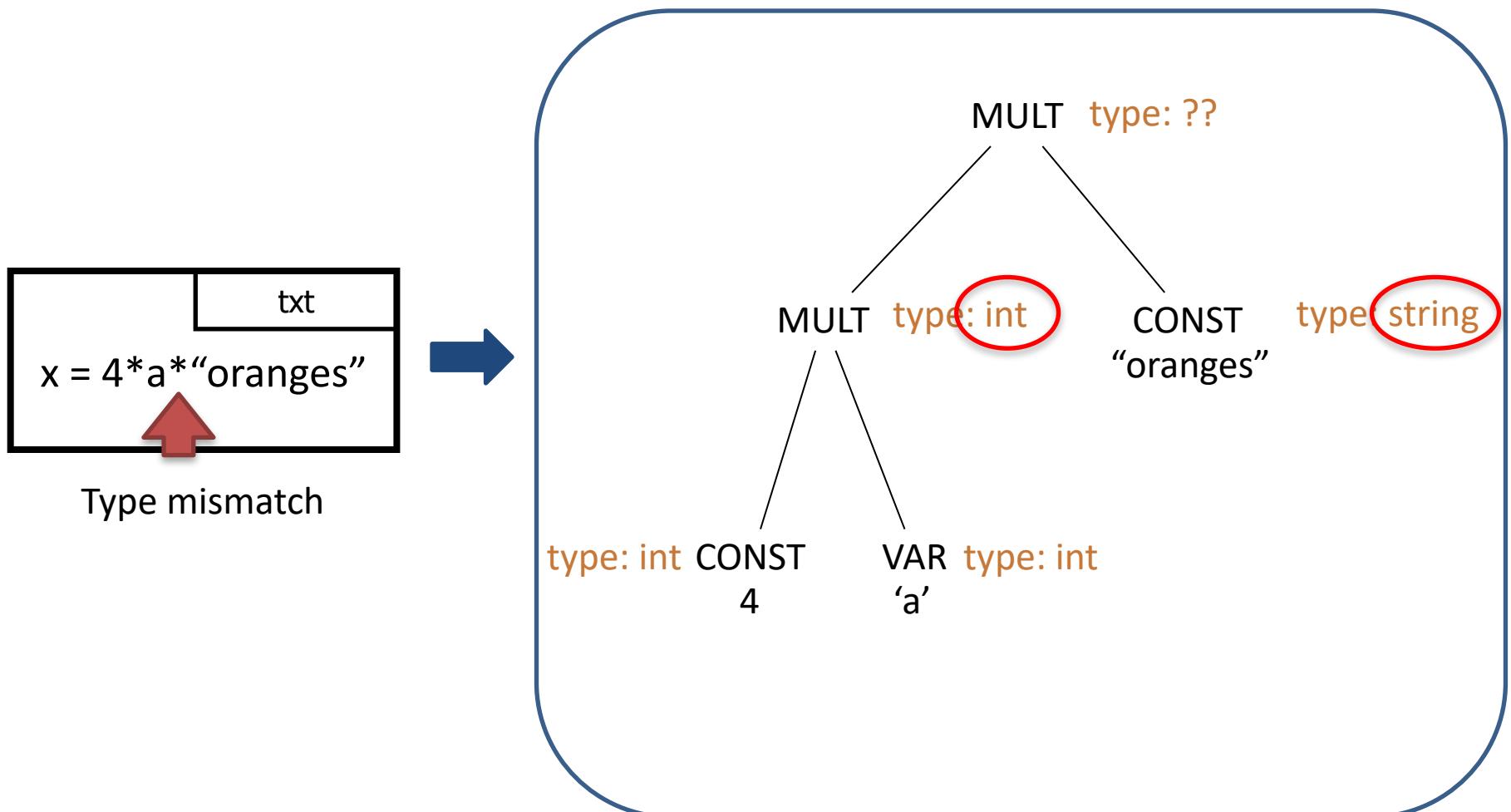
# Lexical Errors



# Syntax Errors

-   Missing operator
-   Wrong number of arguments  
to operator “/”
-   A declaration is not expected  
inside a call

# Semantic Errors: Type Checking



# Runtime Errors

```
txt  
x = det(singular_matrix)  
  
y = 100 / x
```



Division by zero

```
txt  
a = new int[9] // a[0..8]  
  
b = a[answer]
```



Array index out of bounds: 42 > 8

# Code Generated to Capture Runtime Errors

```
txt  
x = det(singular_matrix)  
if (x==0) error()  
y = 100 / x
```



Division by zero

```
txt  
a = new int[9] // a[0..8]  
if (answer > 8) error()  
b = a[answer]
```



Array index out of bounds: 42 > 8

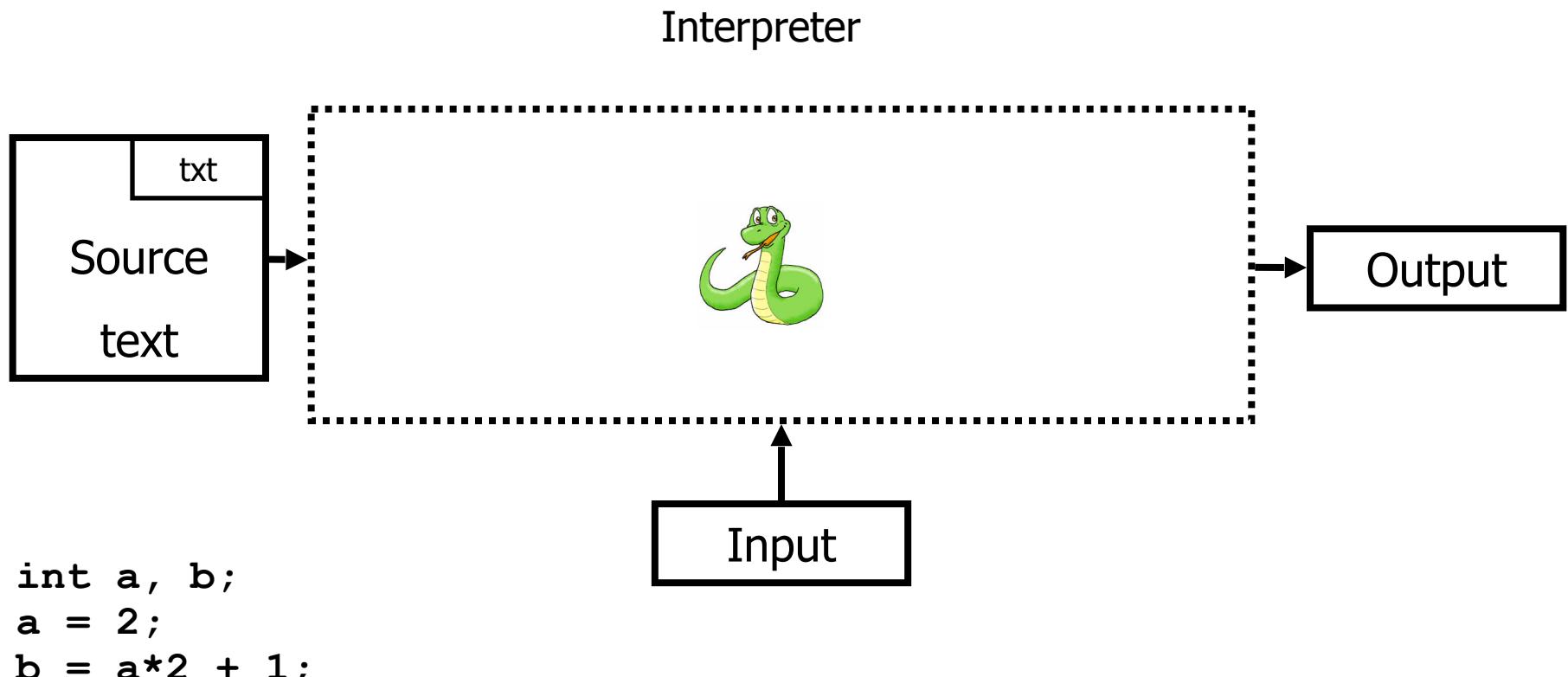
# Shortcuts

- Avoid generating machine code
  - Use assembler
- Generate C code

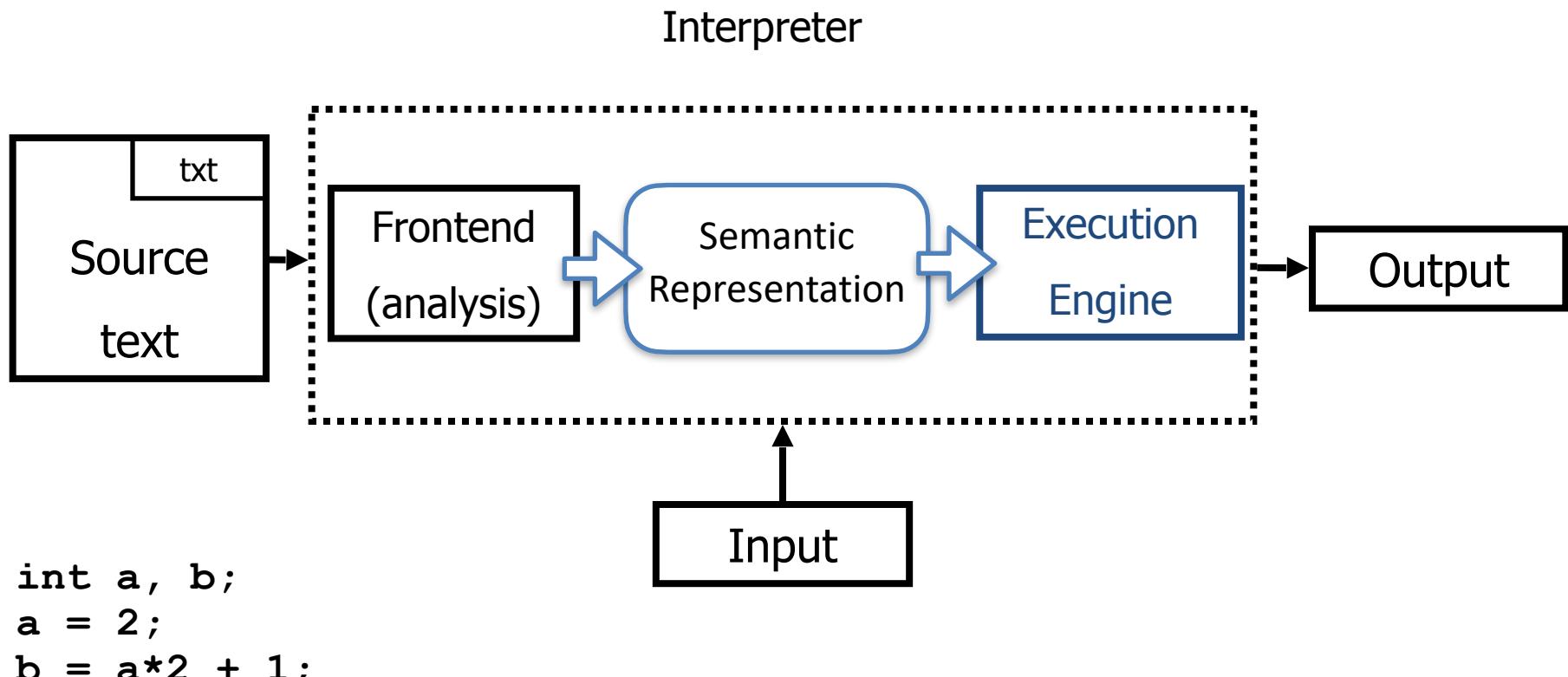
# Compilers and Interpreters



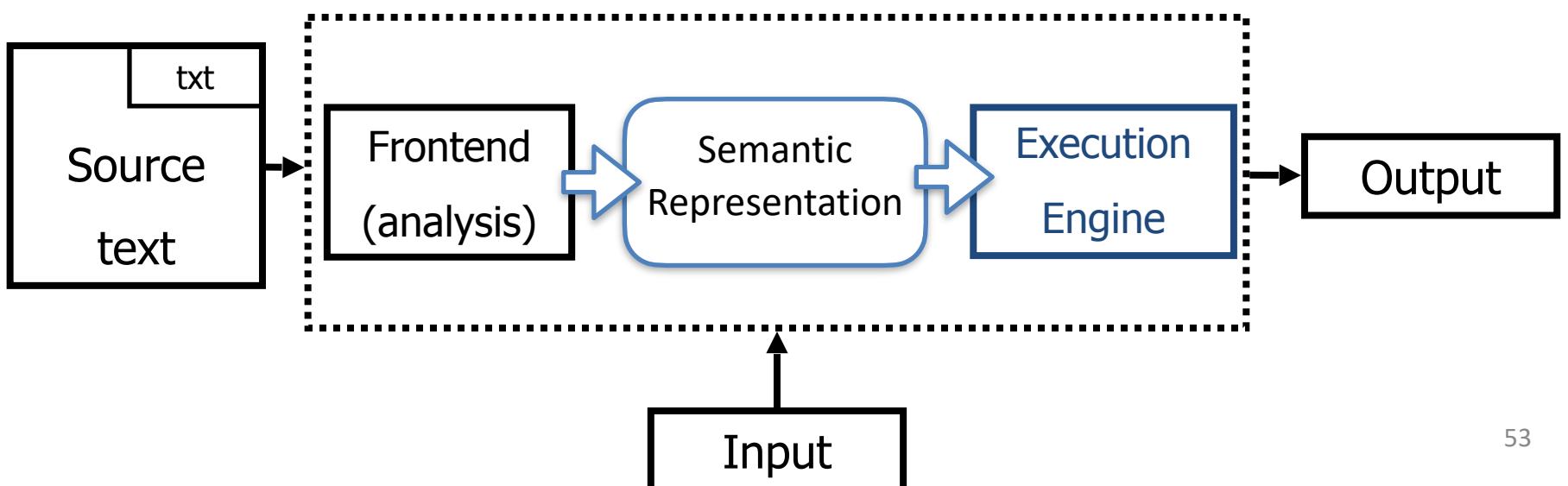
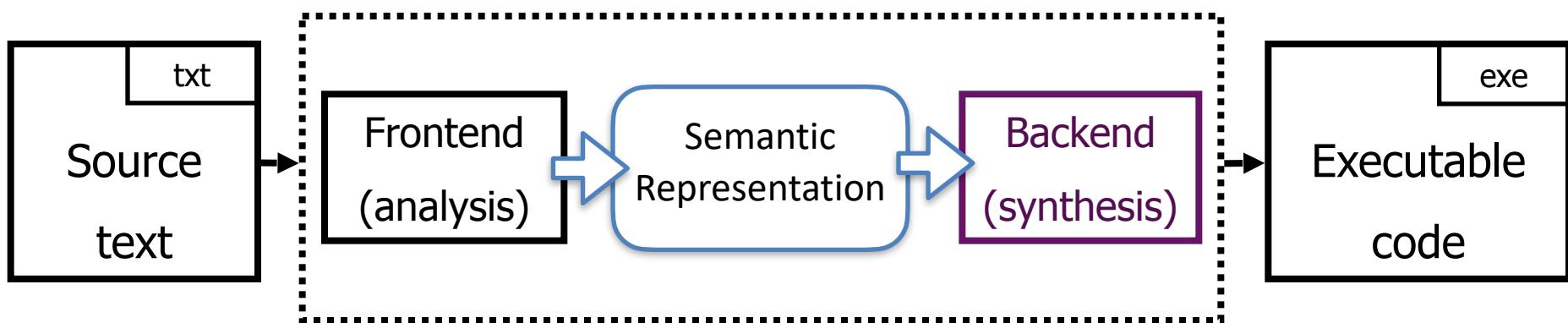
# What is a Interpreter?



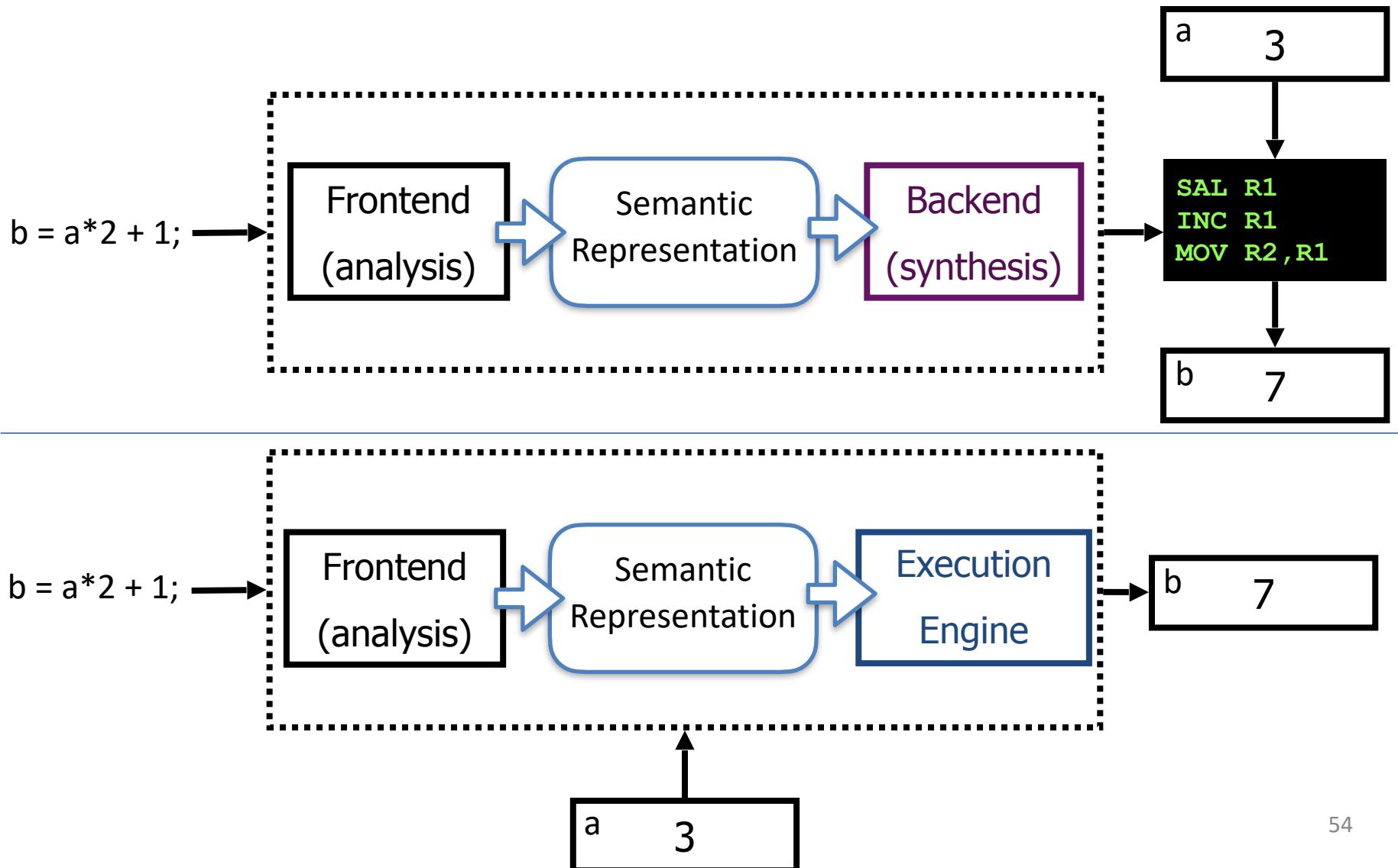
# Anatomy of an Interpreter



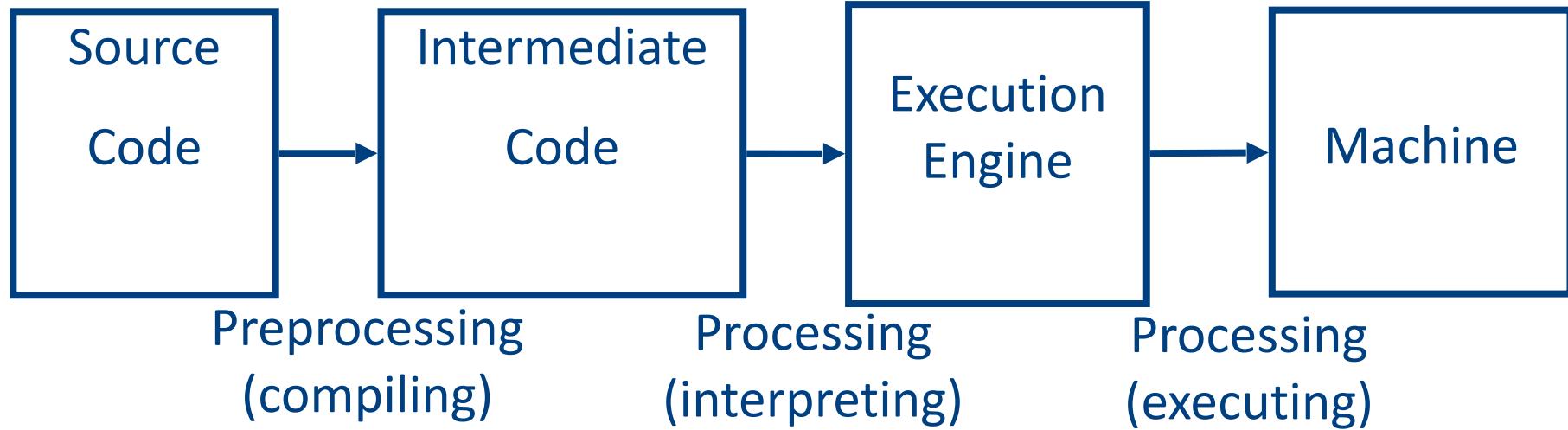
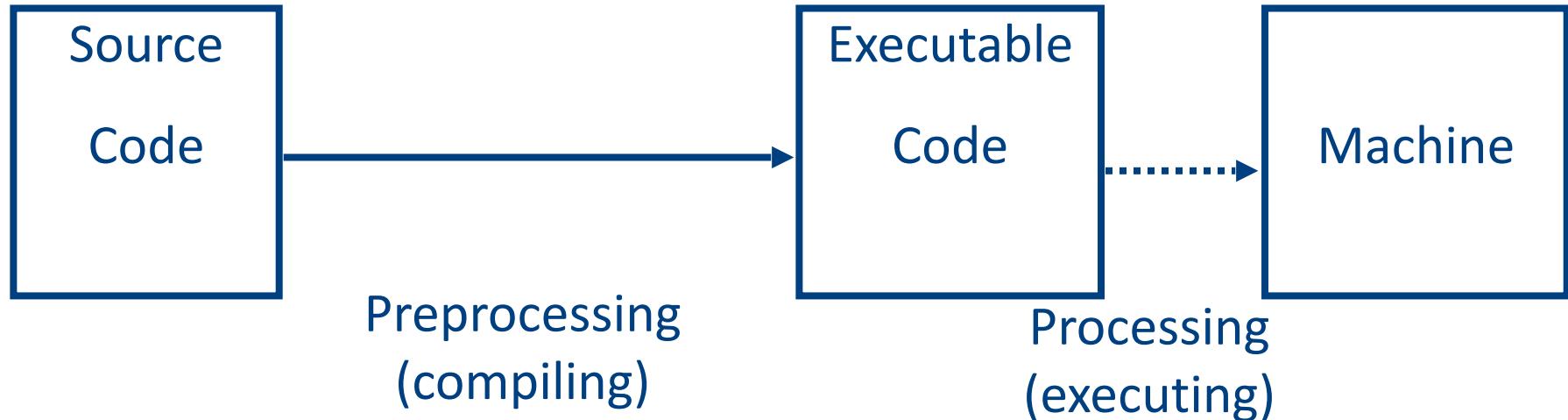
# Compiler vs. Interpreter



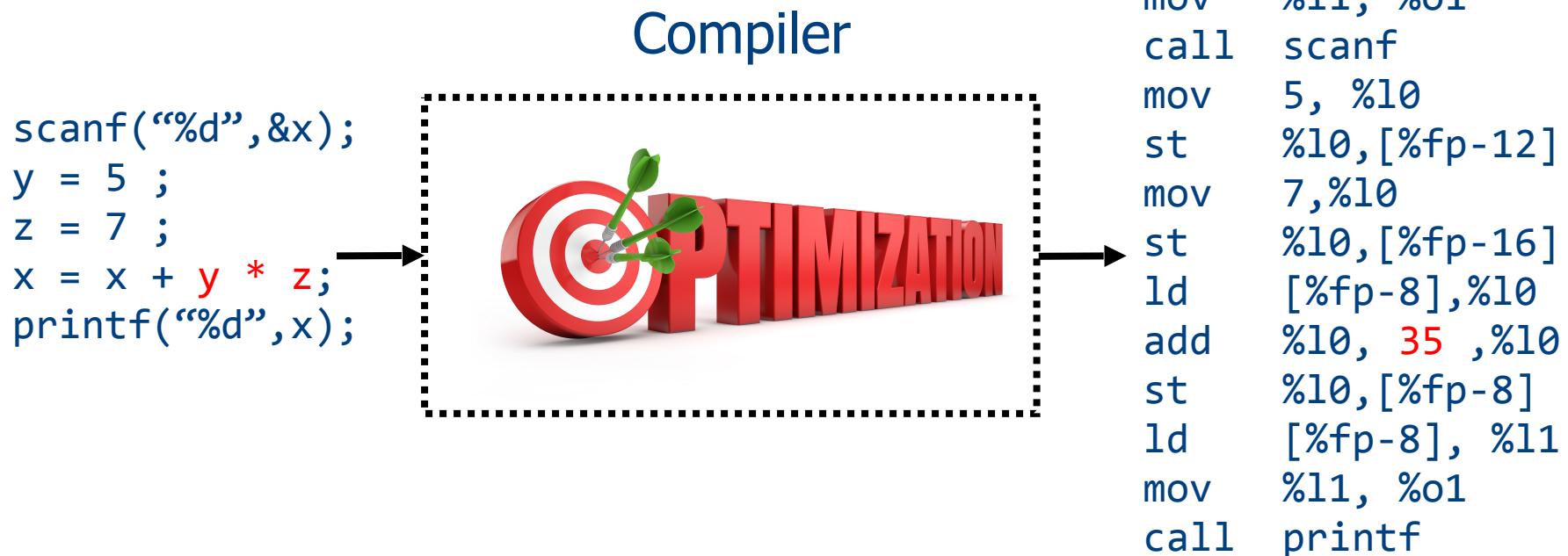
# Compiler vs. Interpreter



# Compiler vs. Interpreter



# Compiled programs are usually more efficient than interpreted



# Compilers report input-independent **possible** errors



- Input-program

```
scanf("%d", &y);
if (y < 0)
    x = 5;

...
if (y <= 0)
    z = 1 / (x + 5 * y);
```

- Compiler-Output

- “line 88: x may be used before set”

# Interpreters report input-specific **definite** errors



- Input-program

```
scanf("%d", &y);
if (y < 0)
    x = 5;
...
if (y <= 0)
    z = 1 / (x + 5 * y);
```

- Input data

- y = -1 : division by zero
- y = 0 : variable x is not defined
- y = 1 : ✓



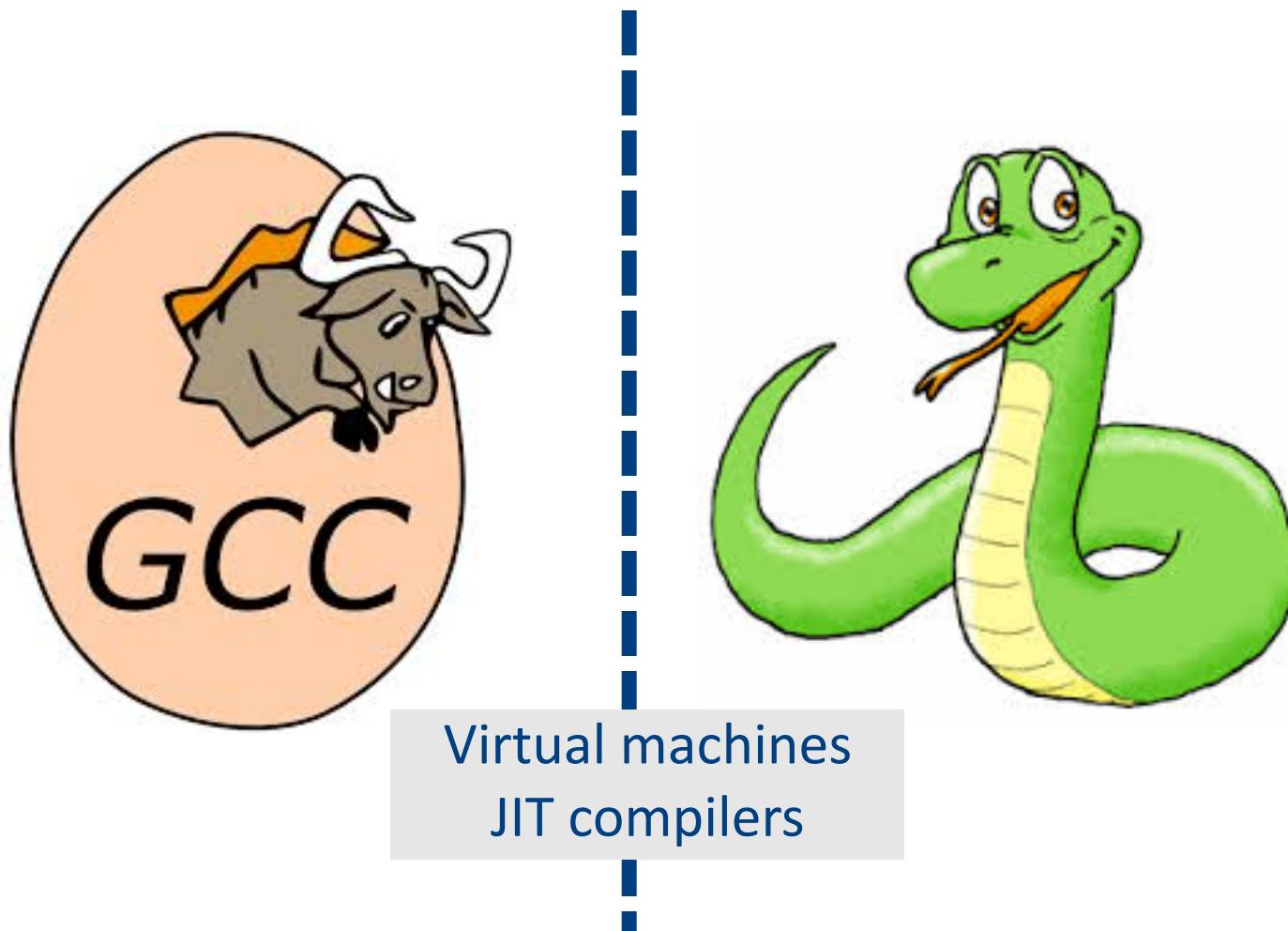
# Compiler

- Programs written in high-level languages
  - Process programs (written in high-level languages)
  - Share functionality
  - Transform code
  - Report potential errors
  - More efficient code
- 
- Run code
  - Report input-specific (definite) errors
  - More responsive
  - Conceptually simpler
    - “Define the language”
  - Easier to port

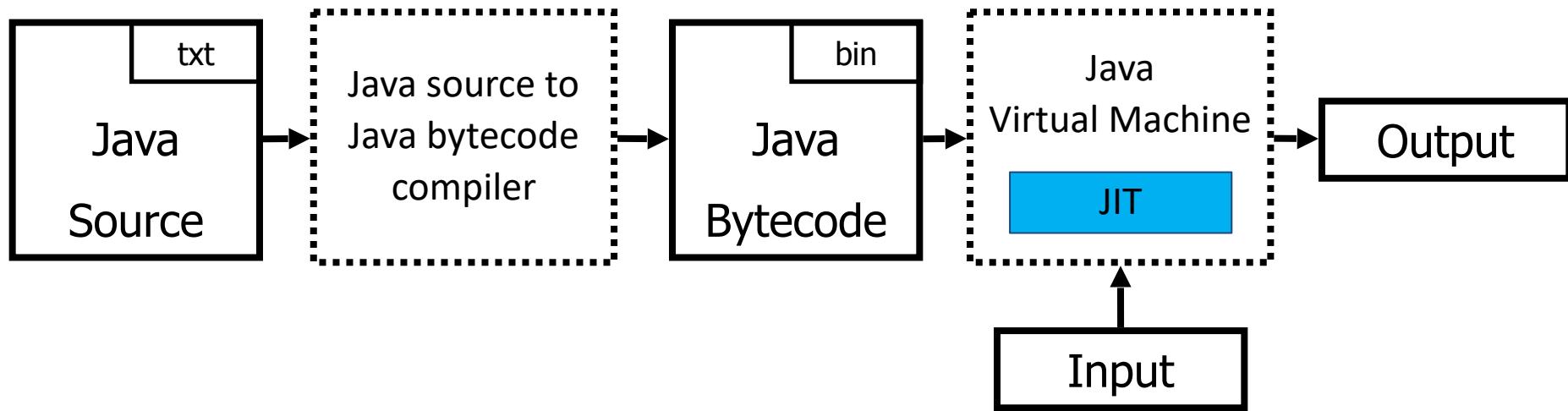
# Interpreter



# “Compiled Lang.” vs. “Interpreted Lang.”



# Just-in-time Compiler (Java HotSpot VM)



**Just-in-time (JIT) compilation:** bytecode interpreter (in the JVM) compiles **hot** program fragments to native machine code during interpretation to avoid expensive re-interpretation

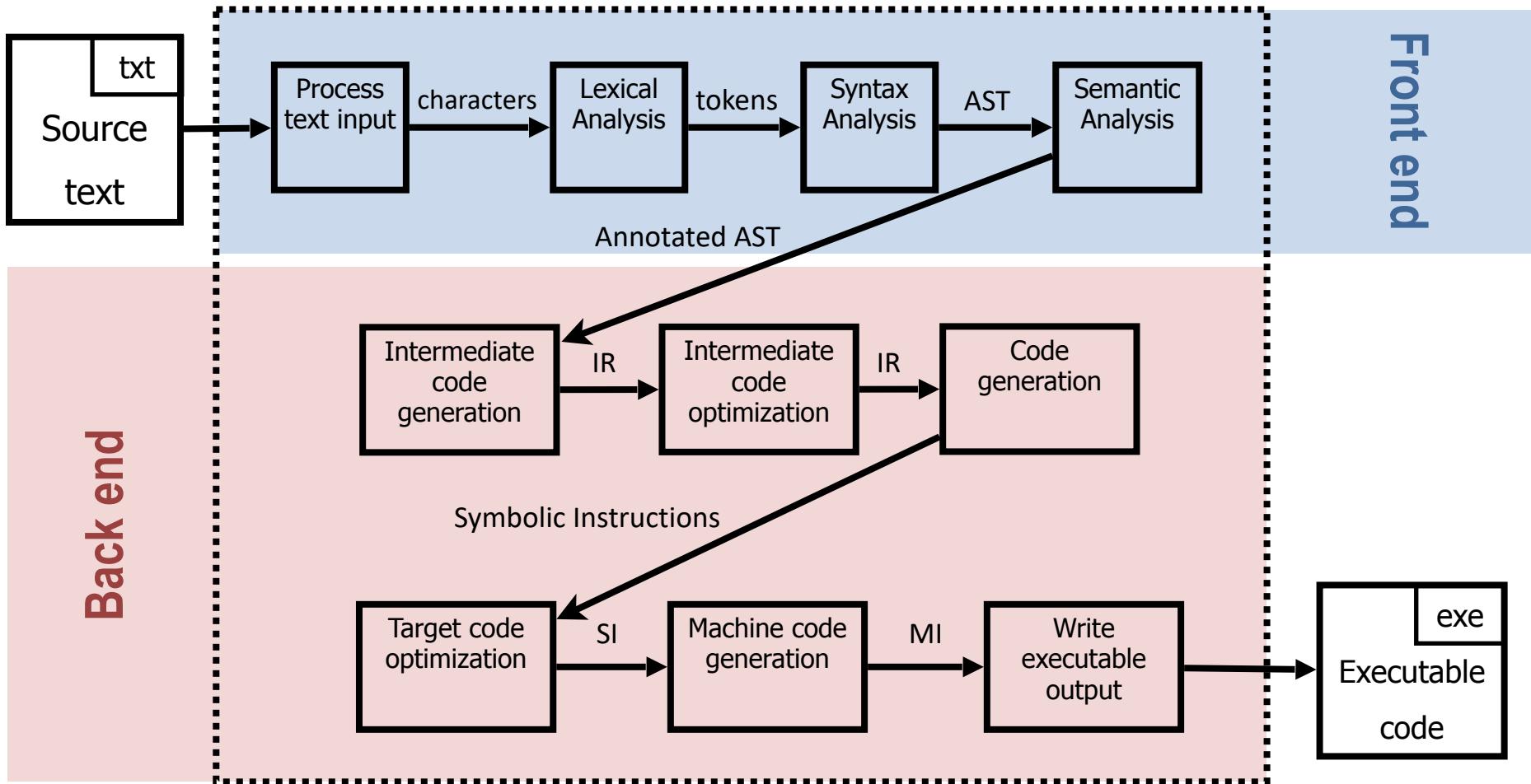
- C1 compiler - a quick compiler (simple optimizations)
- C2 compiler - a slow compiler (aggressive advance optimization for "**very hot**" code fragments)

conclusion

# Summary

- ✓ Compiler = a program that translates code from source language (high level) to target language (low level)
- ✓ Compiler constructed from modular phases
  - Reusable components
  - Different front/back end combinations
- ✓ Compilers play a critical role
  - Bridge programming languages to the machine
  - Many useful techniques and algorithms
  - Many useful tools (e.g., lexer/parser generators)

# The Real Anatomy of a Compiler



# Compiler: Past vs Present

- Past: Lexical and syntactic analysis considered difficult
- Today: Lexical and syntactic analysis “solved.” Bulk effort in analysis and optimizations

# Why should you care?

- Every person in this class will build a “compiler” some day
  - Or wish they knew how to build one...
- But that is not all!

# Why should you care?

Better  
understanding  
of programming  
languages

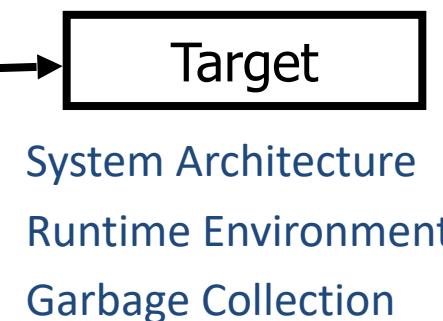


Programming Languages  
Software Engineering

Compiler

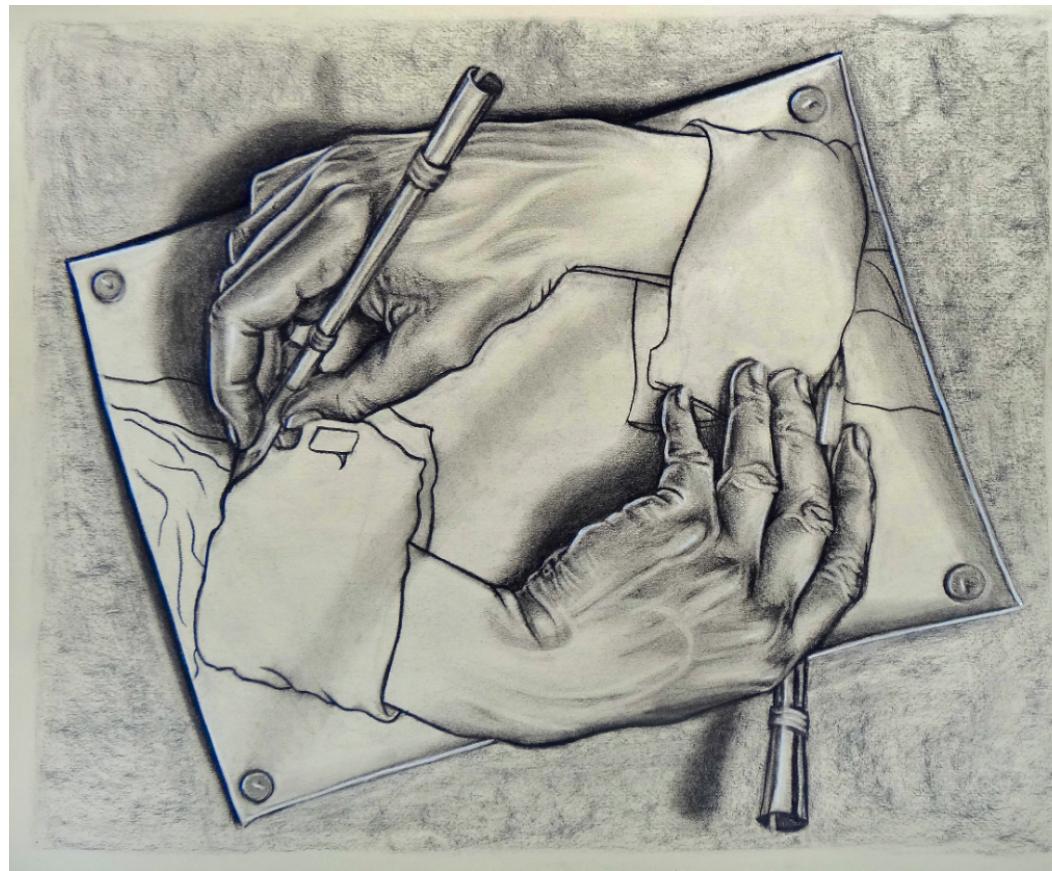
- Useful formalisms
  - Regular expressions
  - Context-free grammars
- Data structures
- Algorithms
- Design concepts
  - Abstractions
  - Modularity
- Program (compiler)-generating tools

Understand  
(some) details  
of target  
architectures

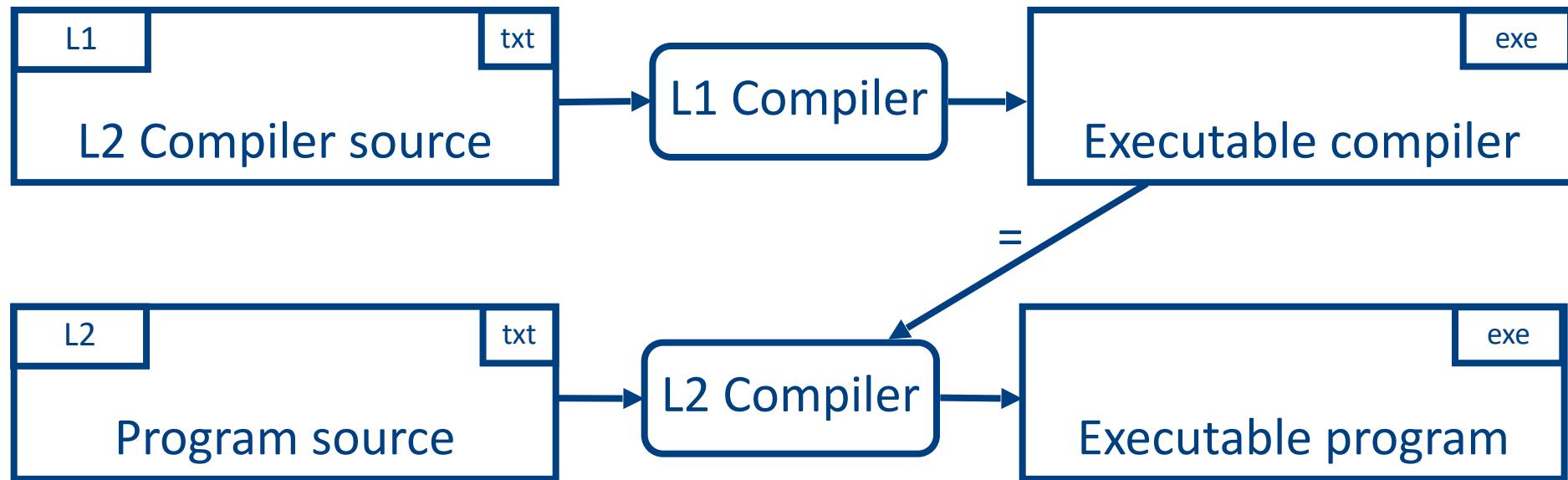


Understand internals of compilers

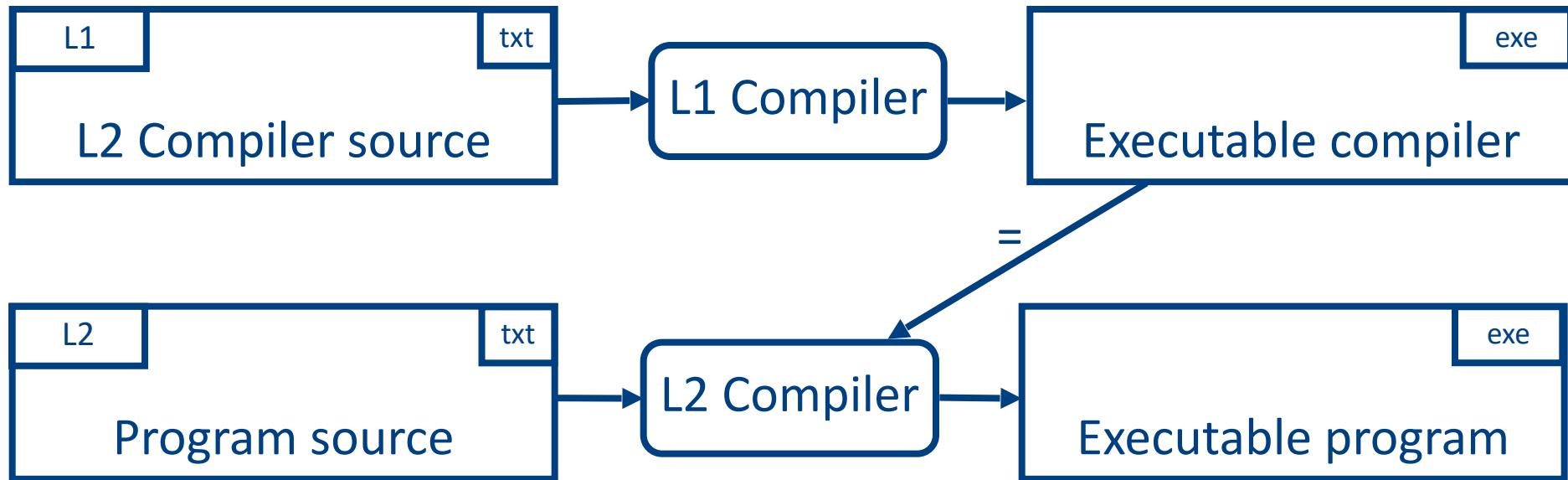
# One More Thing: How to compiler a compiler?



# How to compile a compiler?



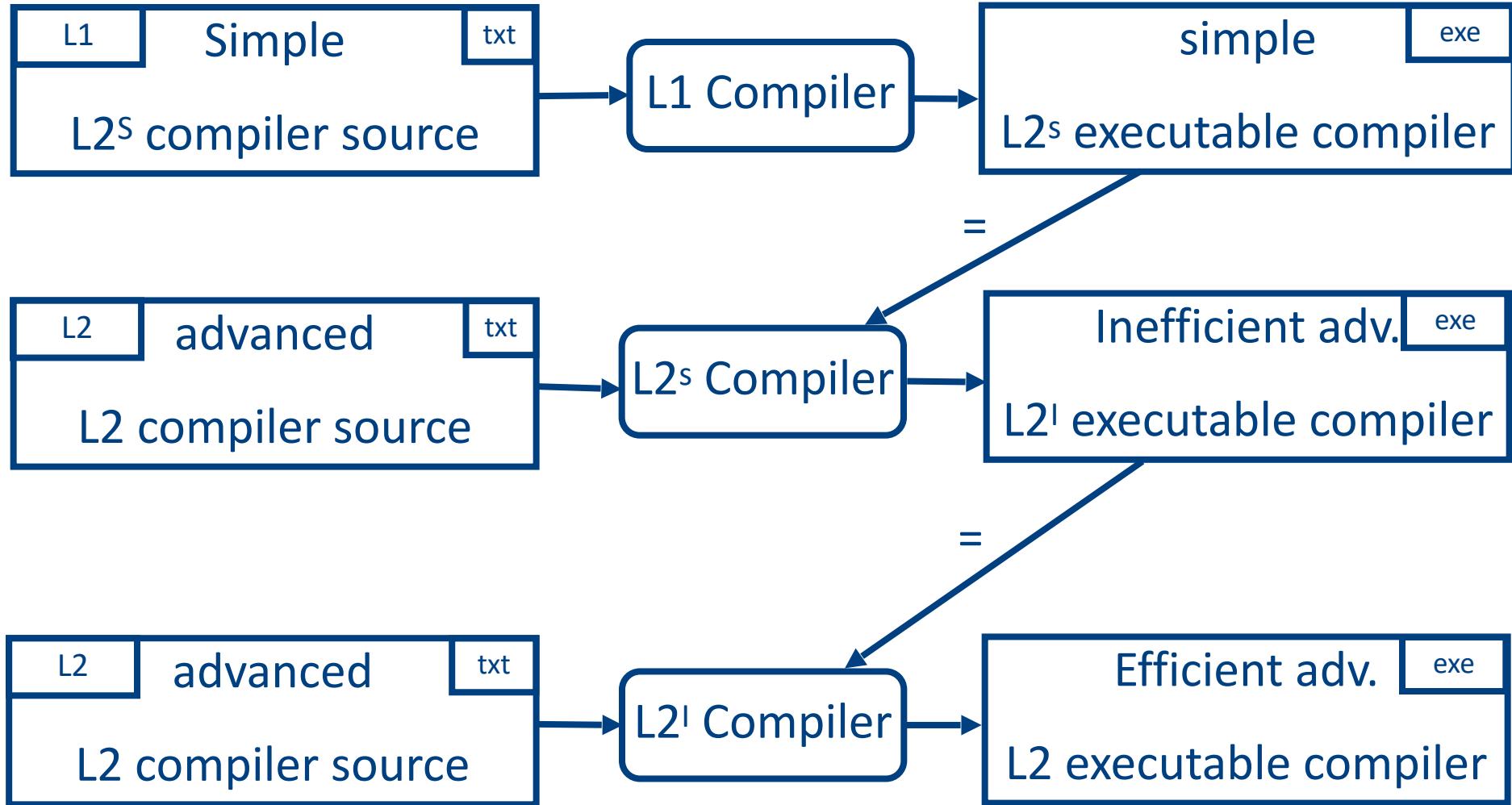
# How to compile a compiler?



Can a compiler for L2 be written in L2?

Can a compiler compile itself?

# Bootstrapping a compiler



# The End