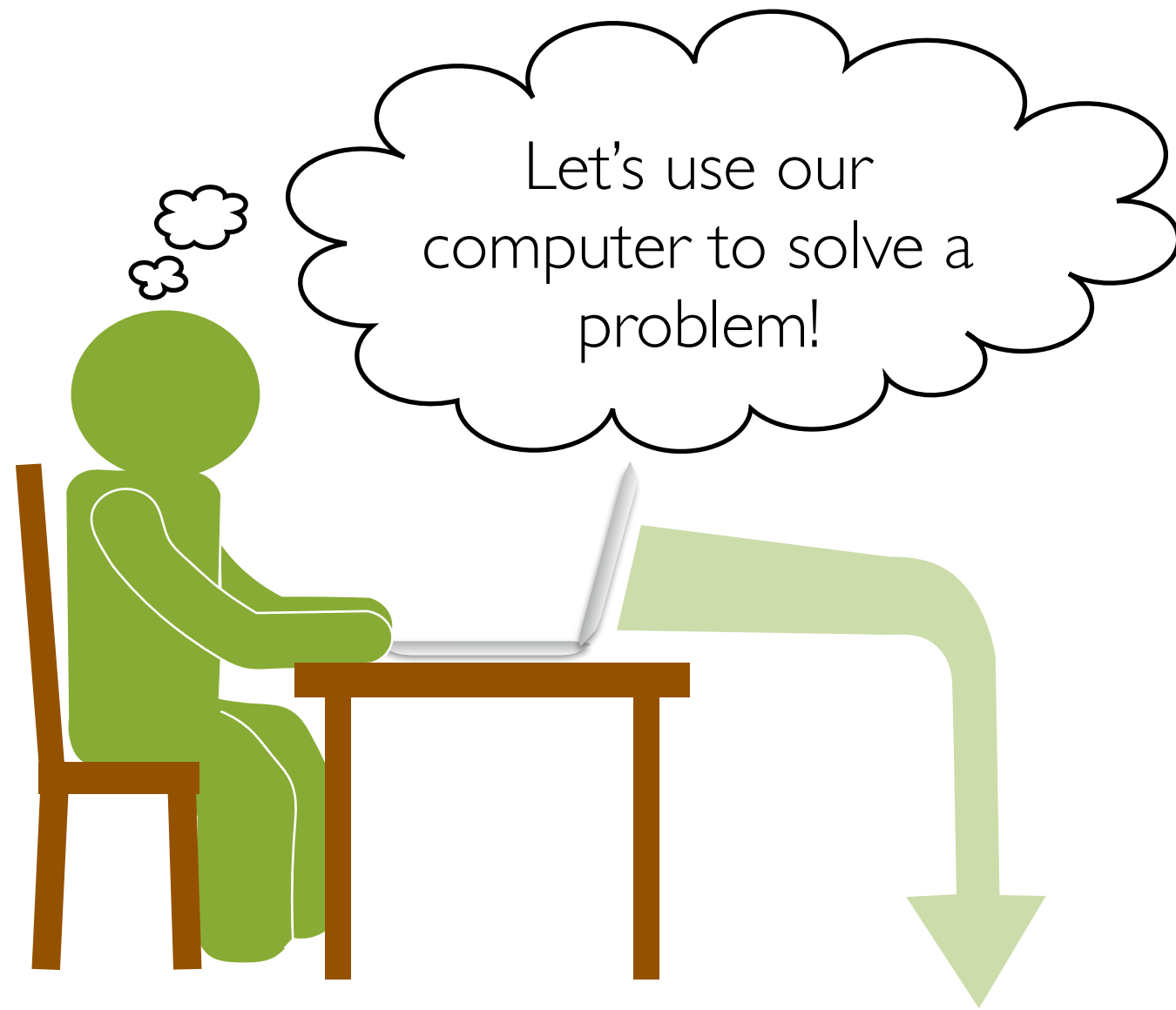# CS 320: Principles of Programming Languages

Mark P Jones & Andrew Tolmach, Portland State University

Winter 2019

Week 2: Programs as Data
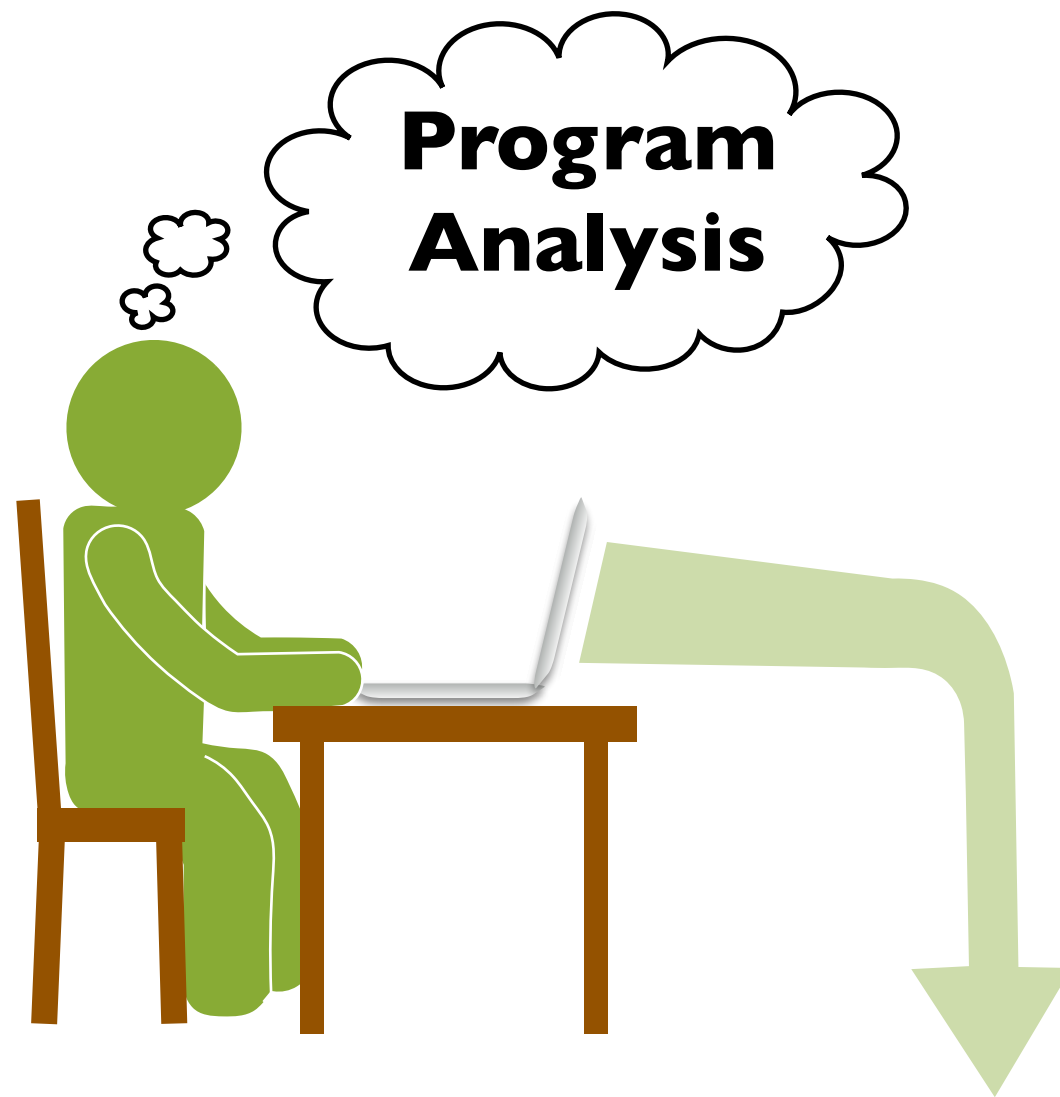
# The computer scientist at work…

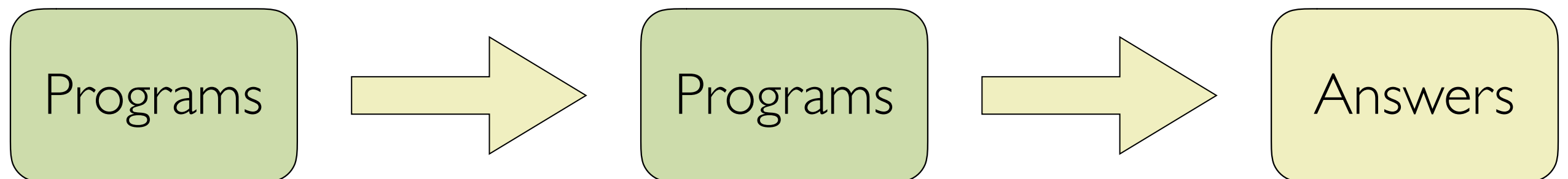But how do I make programs run?

Data → Programs → Answers

3

Program Translation

- Compilers
- Code formatters
- Code update tools
- Macro processors
- Optimizers
- Partial evaluators
- Instrumentation
- Code editors
- …

Programs → Programs → Programs

# How **do** we make high-level programs run on low-level hardware?

Q

# What makes a language "high-level"?

- Complex expressions (arithmetic, logical,...)

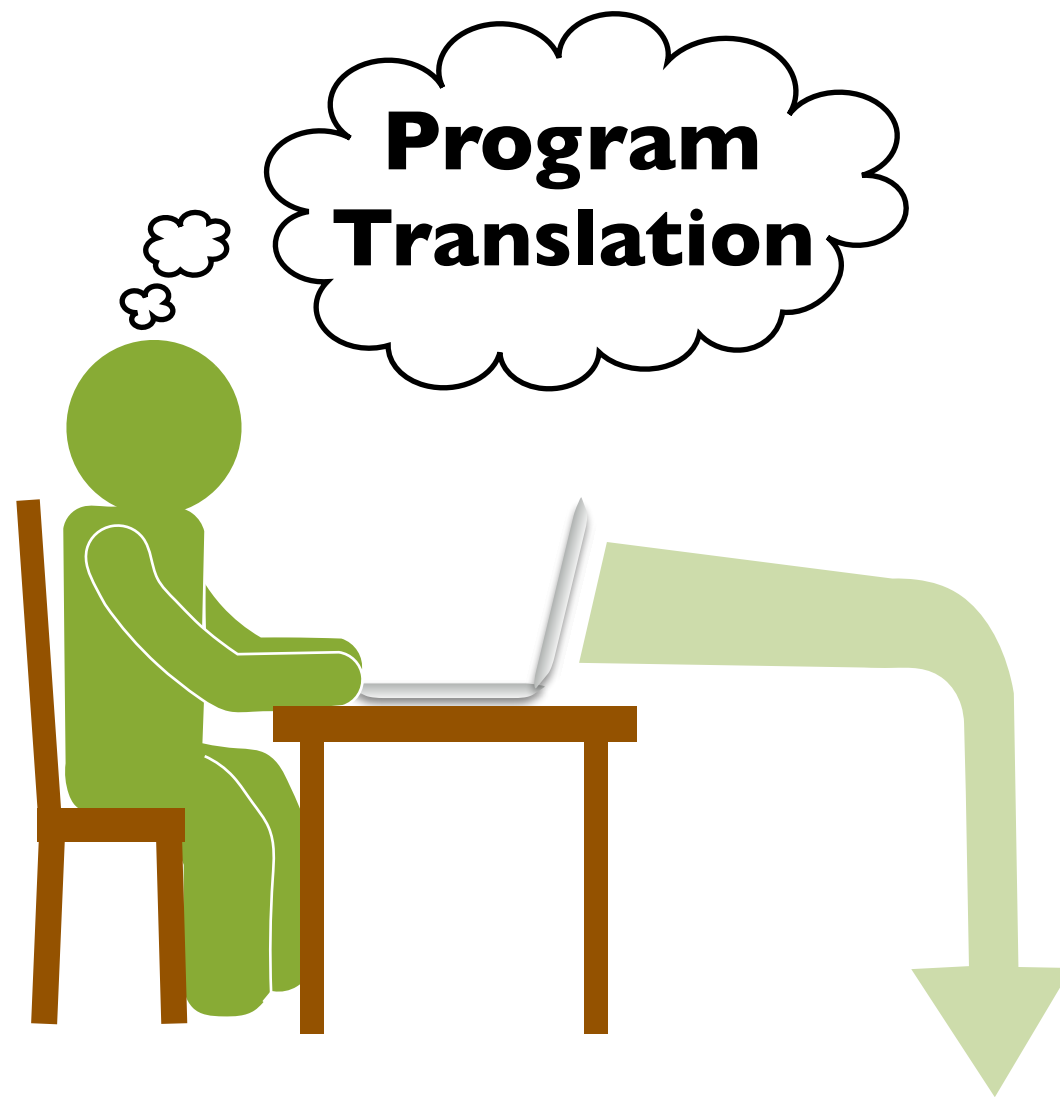- Structured control (loops, conditionals, cases,...)

- Composite types (arrays, records, ...)

- Type declarations and type checking

- Multiple data storage classes (global/local/heap/GC?)

- Procedures/functions (private scope, closures,…)

- Non-local control (exceptions, threads,...)

- Data abstraction (ADTs, modules, objects...)

# What does hardware give us?

- Low-level machine instructions

- Control flow based on labels and conditional branches

- Explicit locations (e.g. registers) for values and intermediate results of computations

- Flat memory model

- Explicit memory management (e.g., stacks for procedure local data)

High-level language

How can we bridge the gap?

Low-level machine

# Interpreters and compilers

# Interpreters and compilers

In conventional English:

- **interpreter**: somebody that translates from one language to another.
  - Example: "I need an interpreter when I'm in Japan"

- **compiler**: somebody who collects, gathers, assembles, or organizes information or things.
  - Latin root: compilare, "plunder or plagiarize"

# Interpreters and compilers

According to my dictionary:

- **in•ter•pret•er** (noun) Computing: a program that can analyze and execute a program line by line

- **com•pile** (verb) Computing (of a computer): convert (a program) into a machine-code or lower-level form in which the program can be executed

  Derivatives: **com•pil•er** (noun)

# Interpreters and compilers

In computer science:

- An interpreter <u>executes</u> (or runs) programs
    - An interpreter for a language L might be thought of as a function: $interp_L : L \rightarrow M$, where M is some set of meanings of programs

- A compiler <u>translates</u> programs
    - A compiler from a language L to a language L' might be thought of as a function $comp : L \rightarrow L'$

- By "language", we mean the set of all strings that correspond to valid programs

# Interpreters and compilers

- Interpreters **execute** programs (turning syntax to semantics)

input

source → | interpreter |

output

- Compilers **translate** programs (turning syntax into syntax)

input

source → | compiler | → | program |

output

# "Doing" vs "Thinking about doing"

- Compilers translate programs (turning syntax to syntax)

- Interpreters run programs (turning syntax to semantics)

- Example:

  - Interpreter (Doing something):
    Use your calculator to evaluate (1+2)+(3+4):

    Answer: 10

  - Compiler (Thinking about doing something):
    Tell me what buttons to press to evaluate (1+2)+(3+4):

    Answer:

    | 1 | + | 2 | = | M | 3 | + | 4 | + | MR | = |

# Basic terminology

**source programs**

many possible source languages, from traditional, to application specific languages.

**target programs**

usually another programming language, often the machine language of a particular computer system.
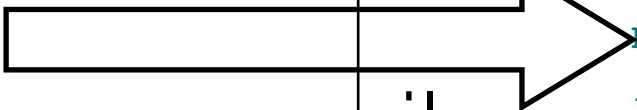
# Example

source program

```
// A simple mini test program

int i = 0;      // initialize
while (i <= 10) {
  print i*i;   // print a square
  i = i + 1;
}
```

target program

compile

execute

```
$ ./squares
0
1
4
9
16
25
36
49
64
81
100
$
```

execute

semantics

```
        .file   "squares.s"
        .comm   _esp0,4
        .globl  _Main_main
_Main_main:
        pushl   %ebp
        movl    %esp,%ebp
        subl    $4,%esp
        movl    $0,%eax
        movl    %eax,-4(%ebp)
        jmp     l1
l0:
        movl    -4(%ebp),%eax
        movl    -4(%ebp),%ebx
        imull   %ebx,%eax
        movl    %esp,_esp0
        subl    $4,%esp
        andl    $0xfffffff0,%esp
        movl    %eax,(%esp)
        call    _print
        movl    _esp0,%esp
        movl    $1,%eax
        movl    -4(%ebp),%ebx
        addl    %ebx,%eax
        movl    %eax,-4(%ebp)
l1:
        movl    $10,%eax
        movl    -4(%ebp),%ebx
        cmpl    %eax,%ebx
        jle     l0
        movl    %ebp,%esp
        popl    %ebp
        ret
```

18

# Compiler correctness

- A compiler should produce valid output for any valid input
- The output should have the same semantics as the input

$$\text{source} \xrightarrow{\text{comp}_{L \to L'}} \text{target}$$

$$\text{source} \xrightarrow{\text{interp}_L} \text{result}_s$$

$$\text{target} \xrightarrow{\text{interp}_{L'}} \text{result}_t$$

$$\text{result}_s \xleftrightarrow{=} \text{result}_t$$

In symbols:     $\forall p.\ \text{interp}_L (p) = \text{interp}_{L'} (\text{comp}_{L \to L'} (p))$

# Desirable properties of a compiler

- Performance:

    - Of compiled code: time, space, power, …

    - Of the compiler: time, space, …

- Diagnostics:

    - High quality error messages and warnings to permit early and accurate diagnosis and resolution of programming mistakes

# Desirable properties, continued

- Support for large programming projects, including:

  - Separate compilation, reducing the amount of recompilation that is needed when part of a program is changed

  - Use of libraries, enabling effective software reuse

- Convenient development environment:

  - Supports program development with an IDE or a range of useful tools, for example: profiling, debugging, cross-referencing, browsing, project management (e.g., make)

# Compiler examples

Compilers show up in many different forms:

- Translating programs in high-level languages like C, C++, Java, etc… to executable machine code

- Just in time compilers: translating byte code to machine code at runtime

- Rendering an HTML web page in a browser window

- Printing a document on a Postscript printer

- Generating audio speech from written text

- Translating from English to Spanish/French/…

- …

# Interpreter characteristics

Common (but not universal) characteristics:

- More emphasis on interactive use:
  - Use of a read-eval-print loop (REPL)
  - Examples: language implementations designed for educational or prototyping applications

- Less emphasis on performance:
  - Interpretive overhead that could be eliminated by compilation
  - Performance of scripting code, for example, is less of an issue if the computations that are being scripted are significantly more expensive

# Interpreter characteristics, continued

- Portability:

    - An interpreter is often more easily ported to multiple platforms than a compiler because it does not depend on the details of a particular target language

- Experimental platforms:

    - Specifying programming language semantics

    - More flexible language designs; some features are easier to implement in an interpreter than in a compiler
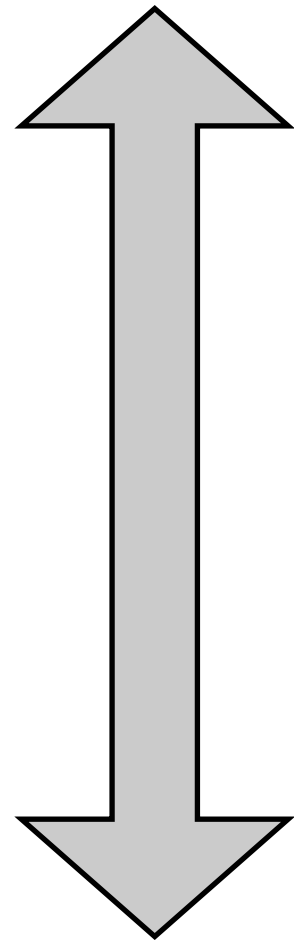
# Interpreter examples

- Programming languages:
  - Scripting languages: PHP, python, ruby, perl, bash, Javascript, ...
  - Educational languages: BASIC, Logo, ...
  - Declarative languages: Lisp, Scheme, ML, Haskell, Prolog, ...
  - Virtual machines: Java, Scala, C#, VB, Pascal (P-Code)
- Document description languages:
  - Postscript, HTML, ...

# Interpreters and Machines

- A virtual machine is one important kind of interpreter
  - Executes programs written in a virtual (i.e. software-defined) instruction set
  - Example: Java Virtual Machine (JVM) executes (interprets) a language of byte codes
- There is no fundamental difference between this and a high-level language interpreter: both execute programs in software
- A CPU executes (machine) programs in hardware:
  - So it is a kind of interpreter too!
  - Faster, but harder to change

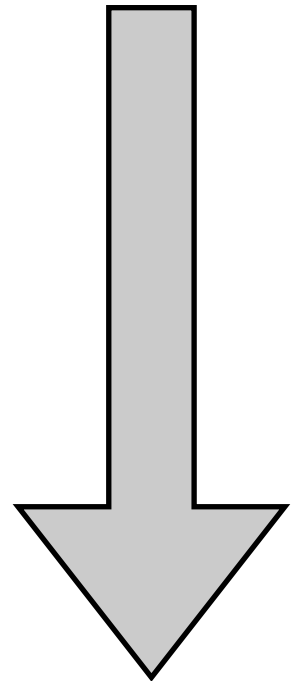# Another look at this question…

High-level language

How can we bridge the gap?

Low-level machine

# We can compile…

High-level language

Compiler

Low-level language

Low-level machine

# We can interpret…
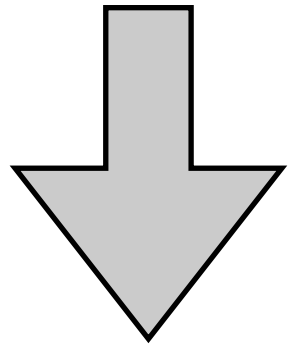
| High-level language |
|---|

| High-level machine |
|---|

↑ Interpreter

| Low-level machine |
|---|

# We can do both…

| High-level language |
|---|

⬇ Compiler

**Mid-level machine/language**

⬆ Interpreter

| Low-level machine |
|---|

# Run-time systems

- Even with a completely compiled approach, we usually need a fixed library of code available at run time, e.g. for:

  - Interfacing to the OS, e.g. to do IO

  - Managing memory, e.g. via garbage collection

  - Managing exception handlers

- This run-time system code is effectively like a (small) virtual machine layer on top of the real hardware and OS process abstraction

- Moral: Every real system involves some elements of interpretation

# Language vs implementation

- Be very careful to distinguish between languages and their implementations

- C is a widely used language

- Haskell is an expressive language

- Java is a well-defined language

- Python is a slow language (NO: speed is a property of an implementation, not a language)

- C++ is a compiled language: (NO: "compiled" describes a property of an implementation, not a language)

Q

# Goals for Compiler Construction

# What is a compiler?

Compilers are translators:

compiler

source
programs

target
programs

diagnostics

# Why translation is needed

- We like to write programs at a higher-level than the machine can execute directly

  - Spreadsheet:     `sum [A1:A3]`

  - Java:     `a[1] + a[2] + a[3]`

  - Machine language:
    ```
    movl $0, %eax
    addl 4(a), %eax
    addl 8(a), %eax
    addl 12(a), %eax
    ```

- High-level languages let us describe what is to be done without worrying about all the details

- In machine languages, every step must be carefully spelled out

**Ideas:**

- Search a database
- Send a message
- Create a song
- Play a game
- etc ...

**High Level**

How do we turn **high level ideas** in to running programs on **low level machines**?

**Machines:**

- Read a value from memory
- Add two numbers
- Compare two numbers
- Write a value to memory
- etc ...

**Low Level**

**Ideas:**

**High Level**

- Search a database
- Send a message
- Create a song
- Play a game
- etc ...

express

**Languages:**

- Evaluate an expression
- Execute a computation multiple times
- Call a function
- Save a result in a variable
- ...

translate

**Machines:**

**Low Level**

- Read a value from memory
- Add two numbers
- Compare two numbers
- Write a value to memory
- etc ...

Ideas:

express

Languages:

translate

Machines:

High Level



Admiral Grace
Hopper (1906-1992)
(Photo: via Wikipedia)

Could we program a
computer to do this?

human ingenuity
required

Low Level

**Ideas:**

express ↓

**Languages:**

translate ↓

**Machines:**

High Level

- Search a database
- Send a
- Create
- Play a
- etc ...

- Evalua
- Execut
- Call a function
- Save a result in a variable
- ...
- Read a value from memory

Admiral Grace
Hopper (1906-1992)
(Photo: via Wikipedia)

Could we program a
computer to do this?

Yes! The A-0 system for
UNIVAC 1 (1951-52):
the first **compiler**

Low Level

39

Ideas:

express

Languages:

translate

Machines:

High Level

Admiral Grace
Hopper (1906-1992)
(Photo: via Wikipedia)

compiler construction

human ingenuity
required

Low Level

Ideas:

Languages:

Machines:

express

translate

human ingenuity required

High Level

language design

compiler construction

human ingenuity required

Low Level

# Languages and tools matter

- Language designs empower developers to:

  - Express their ideas more directly

  - Execute their designs on a computer

- Better tools (compilers, interpreters, etc.) will:

  - open programming to more people and more applications

  - increase programmer productivity

  - enhance software quality (functionality, reliability, security, performance, power, ...)

# Basics of Compiler Structure

# How does a compiler work?

source program

```
// A simple mini test program

int i = 0;      // initialize
while (i <= 10) {
  print i*i;  // print a square
  i = i + 1;
}
```

target
program

```
        .file   "squares.s"
        .comm   _esp0,4
        .globl  _Main_main
_Main_main:
        pushl   %ebp
        movl    %esp,%ebp
        subl    $4,%esp
        movl    $0,%eax
        movl    %eax,-4(%ebp)
        jmp     l1
l0:
        movl    -4(%ebp),%eax
        movl    -4(%ebp),%ebx
        imull   %ebx,%eax
        movl    %esp,_esp0
        subl    $4,%esp
        andl    $0xfffffff0,%esp
        movl    %eax,(%esp)
        call    _print
        movl    _esp0,%esp
        movl    $1,%eax
        movl    -4(%ebp),%ebx
        addl    %ebx,%eax
        movl    %eax,-4(%ebp)
l1:
        movl    $10,%eax
        movl    -4(%ebp),%ebx
        cmpl    %eax,%ebx
        jle     l0
        movl    %ebp,%esp
        popl    %ebp
        ret
```

compile

We need to describe this process in a way that is scalable, precise, mechanical/algorithmic, ...

# What is this?

# False

Dark pixels on a light background

A collection of lines/strokes

A sequence of characters

A single word ("token")

An expression

A boolean expression

A truth value

One thing can be seen in many different ways

We can break a complex process into multiple (hopefully simpler) steps

# "Compiling" English

- ## The symbols must be valid:
  hdk fΩfdh ksdßs dfsjf dslkjé

  ✗  source input

- ## The words must be valid:
  banana jubmod food funning

  ✗  lexical analysis

- ## The text must use correct grammar:  ✗  parser
  my walking up left tree dog

- ## Now we have preliminary abstract syntax:  ✓
  This sentence is a complete.

  ready for "analysis"

# "Compiling" English

- The phrase must make sense    ✗

    This sentence is not true.

- The phrase must not be ambiguous    ✗    static analysis

    Close the window.  My old friend.

- The sentence must fit in context    ✗

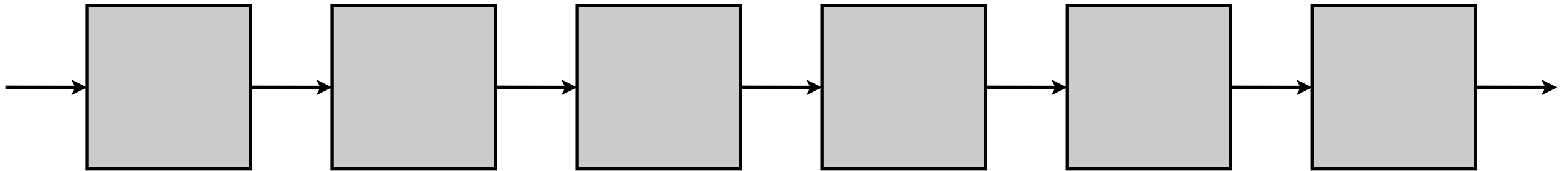    The next song is about geography.

- Finally, we have valid abstract syntax! ✓

    Languages are very interesting.
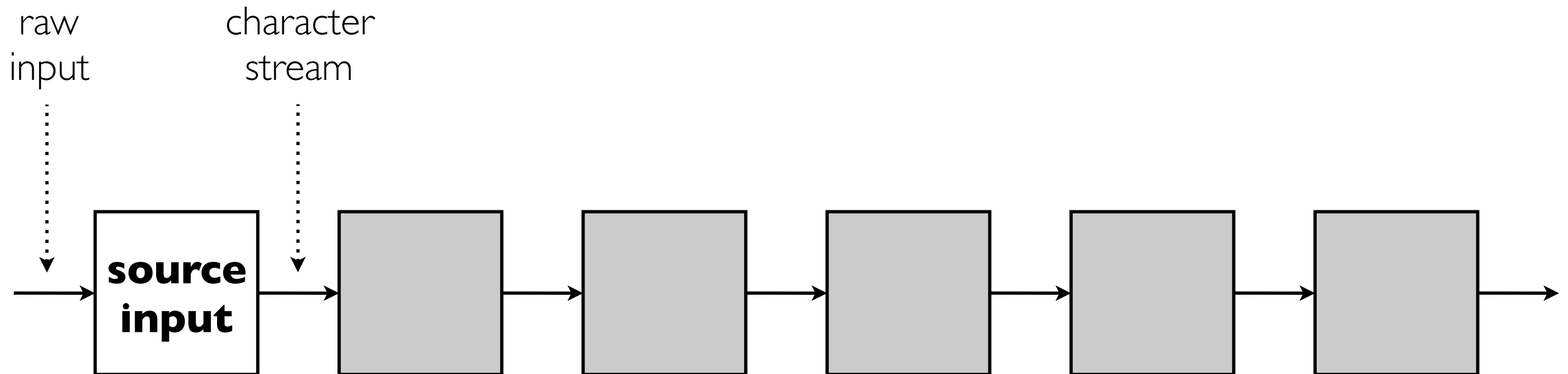
ready for "code generation"

# The compiler pipeline

- Traditionally, the task of compilation is broken down into several steps, or compilation <u>phases</u>:

# Source input      (not a standard term)

raw
input

character
stream

```
→ | source input | → [ ] → [ ] → [ ] → [ ] → [ ] →
```
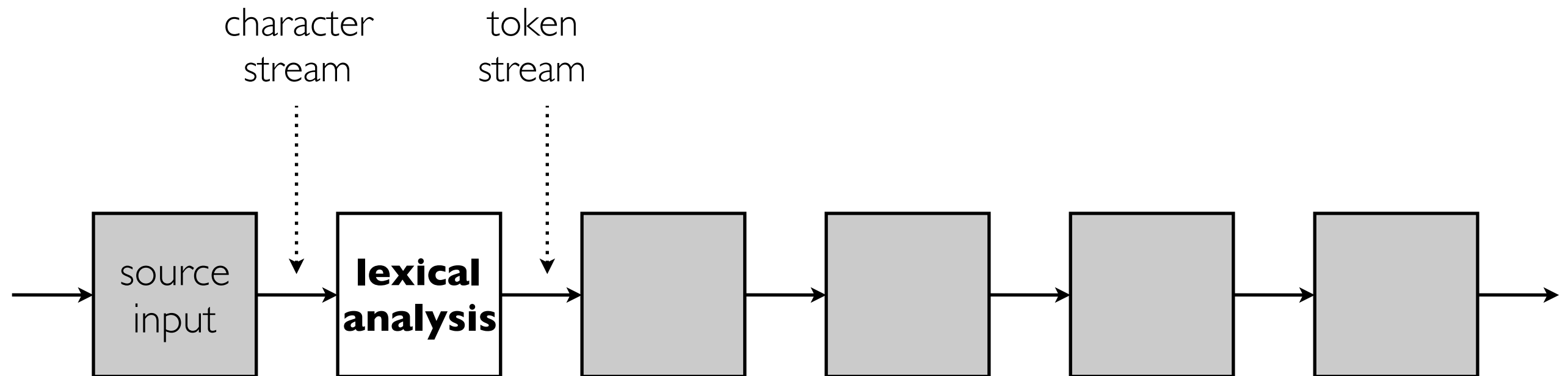
- Turn data from a raw input source into a sequence of characters or lines

  Data might come from a disk, memory, a keyboard, a network, a thumb drive, ...
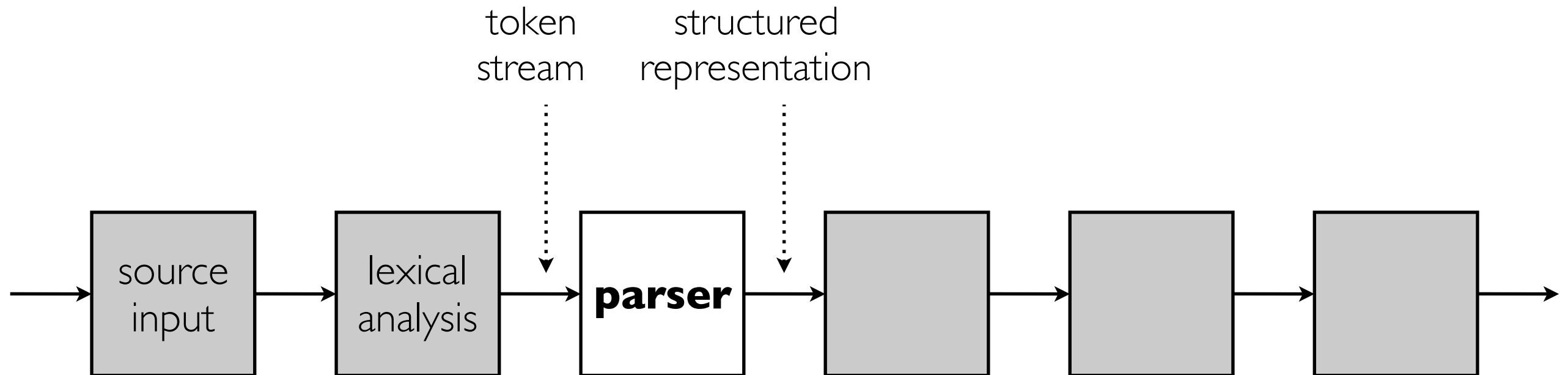
  The operating system usually takes care of most of this ...

# Lexical analysis
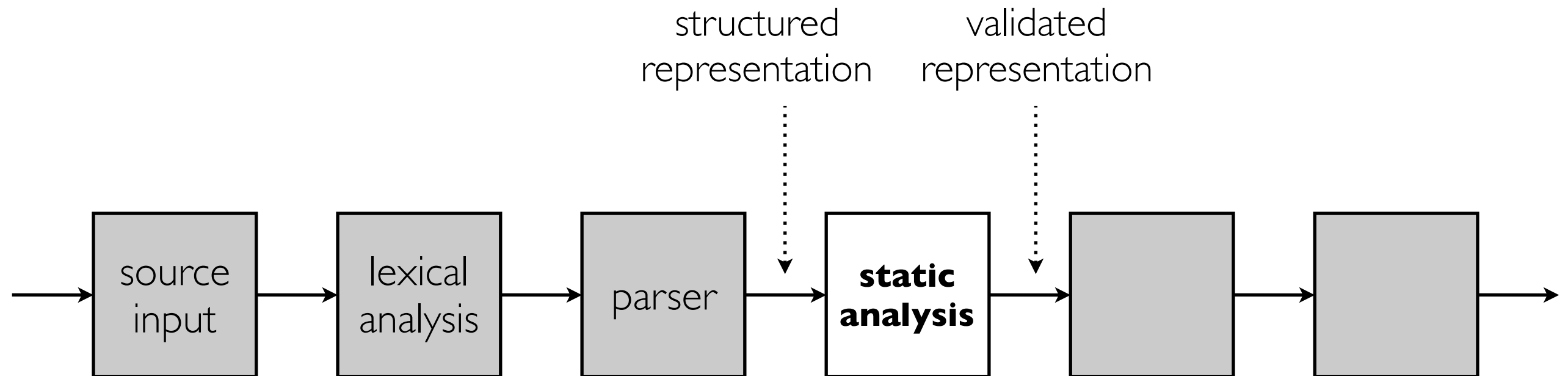


- Convert the input stream of characters into a stream of tokens

- For example, the keyword `for` is treated as a single token, and not as three separate characters

- "lexical":

  "of or relating to the words or vocabulary of a language"

# Parser

token stream    structured representation

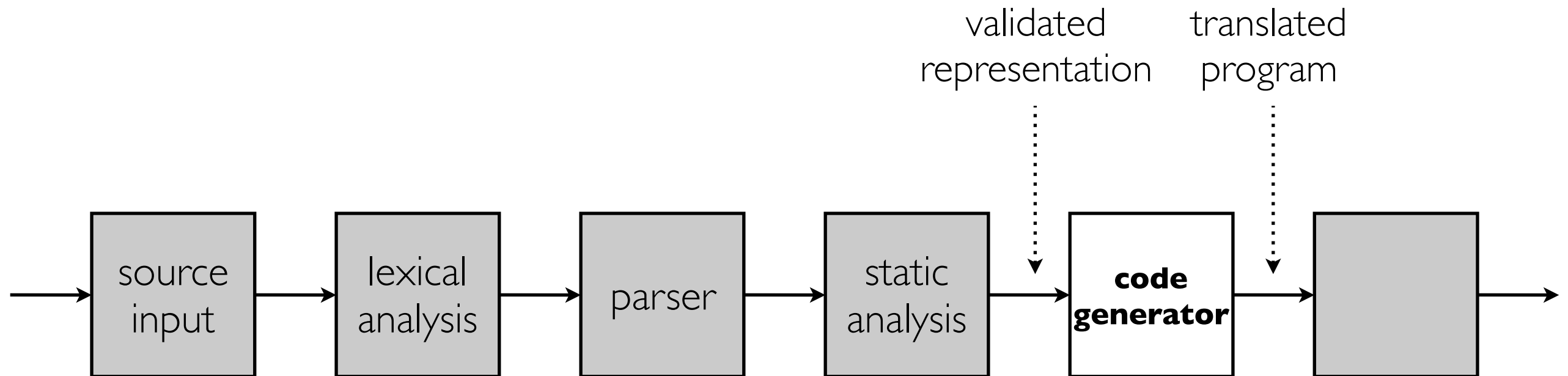source input → lexical analysis → **parser** → [ ] → [ ] → [ ]

- Build data structures that capture the underlying structure (abstract syntax) of the input program

- Determines whether inputs are grammatically well-formed (and reports a syntax error when they are not)

# Static analysis
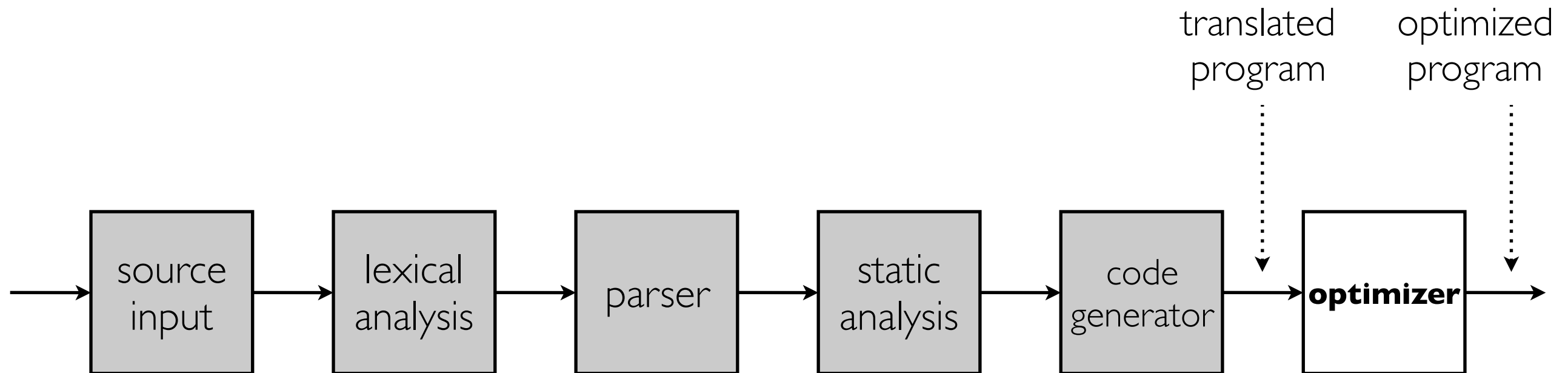
structured
representation

validated
representation

| source input | lexical analysis | parser | **static analysis** | | |

- Check that the program is reasonable:

  no references to unbound variables

  no type inconsistencies

  etc...

# Code generation

validated representation          translated program

| source input | lexical analysis | parser | static analysis | **code generator** | |

- Generate an appropriate sequence of machine instructions as output

- Different strategies are needed for different target machines

# Optimization

translated program → optimizer → optimized program

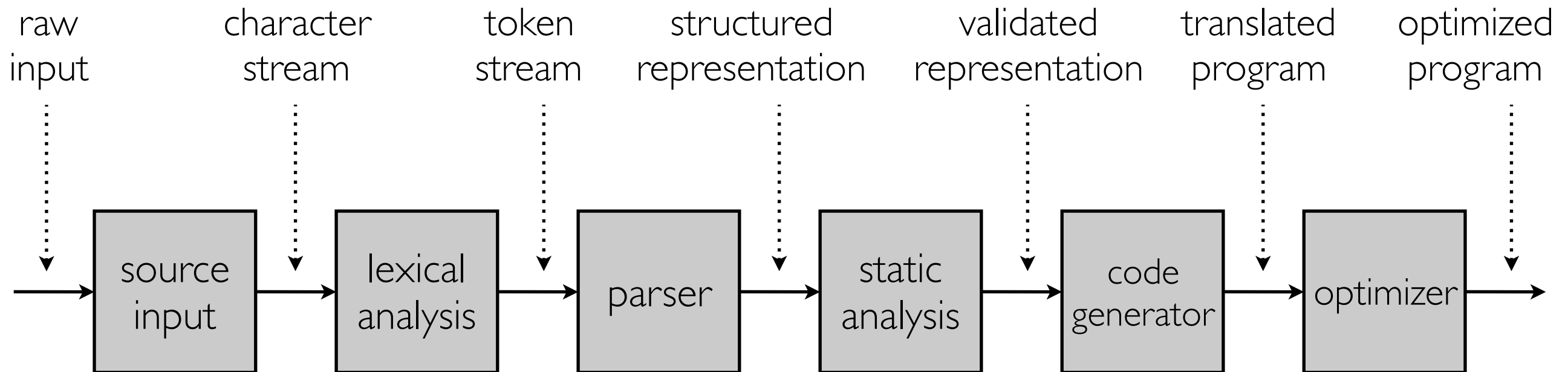source input → lexical analysis → parser → static analysis → code generator → **optimizer** →

- Look for opportunities to improve the quality of the output code:

  There may be conflicting ways to "improve" a given program; the choice depends on the context/the user's priorities

  Producing genuinely "optimal" code is theoretically impossible; "improved" is as good as it gets!

# The full pipeline

| raw input | character stream | token stream | structured representation | validated representation | translated program | optimized program |
|---|---|---|---|---|---|---|

source input → lexical analysis → parser → static analysis → code generator → optimizer

- There are many variations on this approach that you'll see in practical compilers:

    extra phases (e.g., preprocessing)

    iterated phases (e.g., multiple optimization passes)

    additional data may be passed between phases

Q

# Snapshots from a "mini" compiler pipeline

# Snapshots from a "mini" compiler pipeline

- In this week's labs, we'll trace the results of passing the following program through a compiler for a language called "mini"
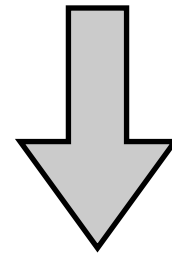
- A sample mini program:

```
// A simple mini test program

int i = 0;     // initialize
while (i <= 10) {
  print i*i;  // print a square
  i = i + 1;
}
```

- The goal here is just to get a sense of how compiler phases work together in practice; you don't need to understand all of the fine details

# Source input (as numbers)
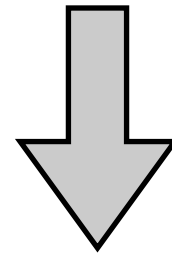
```
// A simple mini test program

int i = 0;     // initialize
while (i <= 10) {
  print i*i;  // print a square
  i = i + 1;
}
```

|47|47|32|65|32|115|105|109|112|108|101|32|77|105|110|105|
32|116|101|115|116|32|112|114|111|103|114|97|109|10|10|105|
110|116|32|105|32|61|32|48|59|32|32|32|32|47|47|32|105|110|
105|116|105|97|108|105|122|101|10|119|104|105|108|101|32|40|
105|32|60|61|32|49|48|41|32|123|10|32|32|112|114|105|110|
116|32|105|42|105|59|32|32|47|47|32|112|114|105|110|116|32|
97|32|115|113|117|97|114|101|10|32|32|105|32|61|32|105|32|
43|32|49|59|10|125|10|

# Source input (as characters)

```
|47|47|32|65|32|115|105|109|112|108|101|32|109|105|110|105|
32|116|101|115|116|32|112|114|111|103|114|97|109|10|10|105|
110|116|32|105|32|61|32|48|59|32|32|32|32|47|47|32|105|110|
105|116|105|97|108|105|122|101|10|119|104|105|108|101|32|40|
105|32|60|61|32|49|48|41|32|123|10|32|32|112|114|105|110|
116|32|105|42|105|59|32|32|47|47|32|112|114|105|110|116|32|
97|32|115|113|117|97|114|101|10|32|32|105|32|61|32|105|32|
43|32|49|59|10|125|10|
```

```
|/|/| |A| |s|i|m|p|l|e| |m|i|n|i| |t|e|s|t| |p|r|o|g|r|a|m|\n
|\n
|i|n|t| |i| |=| |0|;| | | | |/|/| |i|n|i|t|i|a|l|i|z|e|\n
|w|h|i|l|e| |(|i| |<|=| |1|0|)| |{|\n
| | |p|r|i|n|t| |i|*|i|;| | |/|/| |p|r|i|n|t| |a| |s|q|u|a|r|e|\n
| | |i| |=| |i| |+| |1|;|\n
|}|\n
|\n
```

# Lexical analysis

```
|/|/|  |A|  |s|i|m|p|l|e|  |m|i|n|i|  |t|e|s|t|  |p|r|o|g|r|a|m|\n
|\n
|i|n|t|  |i|  |=|  |0|;|  |  |  |  |/|/|  |i|n|i|t|i|a|l|i|z|e|\n
|w|h|i|l|e|  |(|i|  |<|=|  |1|0|)|  |{|\n
|  |  |p|r|i|n|t|  |i|*|i|;|  |  |/|/|  |p|r|i|n|t|  |a|  |s|q|u|a|r|e|\n
|  |  |i|  |=|  |i|  |+|  |1|;|\n
|}|\n
|\n
```



```
| INT | ID(i) | = | INTLIT(0) | Semicolon ";" | WHILE
| Open parenthesis "(" | ID(i) | <= | INTLIT(10)
| Close parenthesis ")" | Open brace "{" | PRINT | ID(i)
| * | ID(i) | Semicolon ";" | ID(i) | = | ID(i) | +
| INTLIT(1) | Semicolon ";" | Close brace "}" |
```
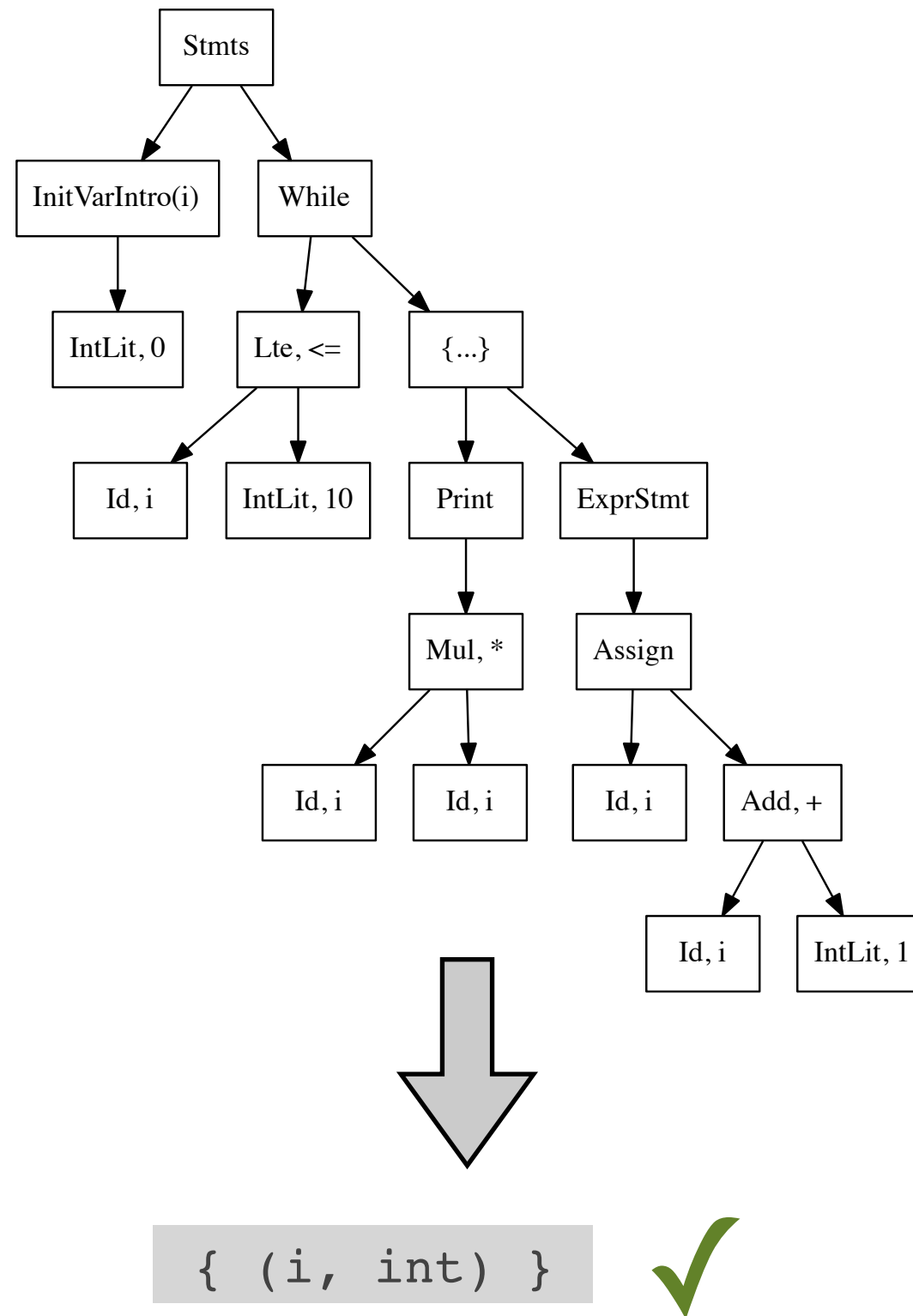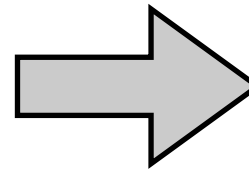
# Parsing

```
| INT | ID(i) | = | INTLIT(0) | Semicolon ";" | WHILE
| Open parenthesis "(" | ID(i) | <= | INTLIT(10)
| Close parenthesis ")" | Open brace "{" | PRINT | ID(i)
| * | ID(i) | Semicolon ";" | ID(i) | = | ID(i) | +
| INTLIT(1) | Semicolon ";" | Close brace "}" |
```

# Static analysis

# Code generation



```
            .file      "squares.s"
            .comm      _esp0,4
            .globl     _Main_main
_Main_main:
            pushl      %ebp
            movl       %esp,%ebp
            subl       $4,%esp
            movl       $0,%eax
            movl       %eax,-4(%ebp)
            jmp        l1
l0:
            movl       -4(%ebp),%eax
            movl       -4(%ebp),%ebx
            imull      %ebx,%eax
            movl       %esp,_esp0
            subl       $4,%esp
            andl       $0xfffffff0,%esp
            movl       %eax,(%esp)
            call       _print
            movl       _esp0,%esp
            movl       $1,%eax
            movl       -4(%ebp),%ebx
            addl       %ebx,%eax
            movl       %eax,-4(%ebp)
l1:
            movl       $10,%eax
            movl       -4(%ebp),%ebx
            cmpl       %eax,%ebx
            jle        l0
            movl       %ebp,%esp
            popl       %ebp
            ret
```

# Assembly

```
        .file   "squares.s"
        .comm   _esp0,4
        .globl  _Main_main
_Main_main:
        pushl   %ebp
        movl    %esp,%ebp
        subl    $4,%esp
        movl    $0,%eax
        movl    %eax,-4(%ebp)
        jmp     l1
l0:
        movl    -4(%ebp),%eax
        movl    -4(%ebp),%ebx
        imull   %ebx,%eax
        movl    %esp,_esp0
        subl    $4,%esp
        andl    $0xfffffff0,%esp
        movl    %eax,(%esp)
        call    _print
        movl    _esp0,%esp
        movl    $1,%eax
        movl    -4(%ebp),%ebx
        addl    %ebx,%eax
        movl    %eax,-4(%ebp)
l1:
        movl    $10,%eax
        movl    -4(%ebp),%ebx
        cmpl    %eax,%ebx
        jle     l0
        movl    %ebp,%esp
        popl    %ebp
        ret
```
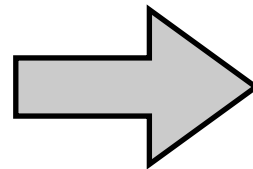
```
$ od -A x -t x1 squares.o
0000000    ce  fa  ed  fe  07  00  00  00  03  00  00  00  01  00  00  00
0000010    03  00  00  00  e4  00  00  00  00  00  00  00  01  00  00  00
0000020    7c  00  00  00  00  00  00  00  00  00  00  00  00  00  00  00
0000030    00  00  00  00  00  00  00  00  87  00  00  00  01  00  00  00
0000040    87  00  00  00  07  00  00  00  07  00  00  00  01  00  00  00
0000050    00  00  00  00  5f  5f  74  65  78  74  00  00  00  00  00  00
0000060    00  00  00  00  5f  5f  54  45  58  54  00  00  00  00  00  00
0000070    00  00  00  00  00  00  00  00  87  00  00  00  01  00  00  00
0000080    00  00  00  00  88  01  00  00  08  00  00  00  04  00  80
0000090    00  00  00  00  00  00  00  00  02  00  00  00  18  00  00  00
00000a0    c8  01  00  00  05  00  00  00  04  02  00  00  20  00  00  00
00000b0    0b  00  00  00  50  00  00  00  00  00  00  00  02  00  00  00
00000c0    02  00  00  00  01  00  00  00  03  00  00  00  02  00  00  00
00000d0    00  00  00  00  00  00  00  00  00  00  00  00  00  00  00  00
*
0000100    55  89  e5  83  ec  08  b8  00  00  00  00  89  45  fc  b8  00
0000110    00  00  00  89  45  f8  e9  3b  00  00  00  8b  45  fc  8b  5d
0000120    fc  0f  af  c3  89  25  00  00  00  00  83  ec  04  83  e4  f0
0000130    89  04  24  e8  c8  ff  ff  ff  8b  25  00  00  00  00  8b  45
0000140    fc  8b  5d  f8  01  d8  89  45  f8  b8  01  00  00  00  8b  5d
0000150    fc  01  d8  89  45  fc  b8  0a  00  00  00  8b  5d  fc  39  c3
0000160    0f  8e  b5  ff  ff  ff  8b  45  f8  89  25  00  00  00  00  83
0000170    ec  04  83  e4  f0  89  04  24  e8  83  ff  ff  ff  8b  25  00
0000180    00  00  00  89  ec  5d  c3  00  7f  00  00  00  03  00  00  0c
0000190    79  00  00  00  04  00  00  0d  6b  00  00  00  03  00  00  0c
00001a0    62  00  00  00  01  00  00  05  3a  00  00  00  03  00  00  0c
00001b0    34  00  00  00  04  00  00  0d  26  00  00  00  03  00  00  0c
00001c0    17  00  00  00  01  00  00  05  19  00  00  00  0e  01  00  00
00001d0    56  00  00  00  1c  00  00  0e  01  00  00  1b  00  00  00
00001e0    07  00  00  00  0f  01  00  00  00  00  00  00  01  00  00  00
00001f0    01  00  00  00  04  00  00  00  12  00  00  00  01  00  00  00
0000200    00  00  00  00  00  5f  65  73  70  30  00  5f  4d  61  69  6e
0000210    5f  6d  61  69  6e  00  5f  70  72  69  6e  74  00  6c  31  00
0000220    6c  30  00  00
0000224
```
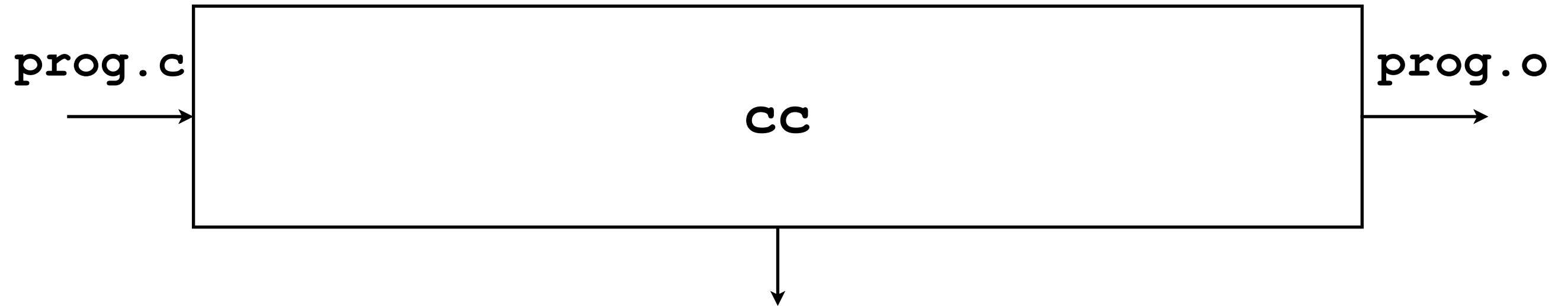
64

# Modularity in compiler design

# Modularity

- Modularity is all about building large systems from collections of smaller components

- Modular implementations can be easier to write, test, debug, understand, and maintain than monolithic implementations

- For example:
  - Components can be developed independently
  - Some components can be reused in other contexts
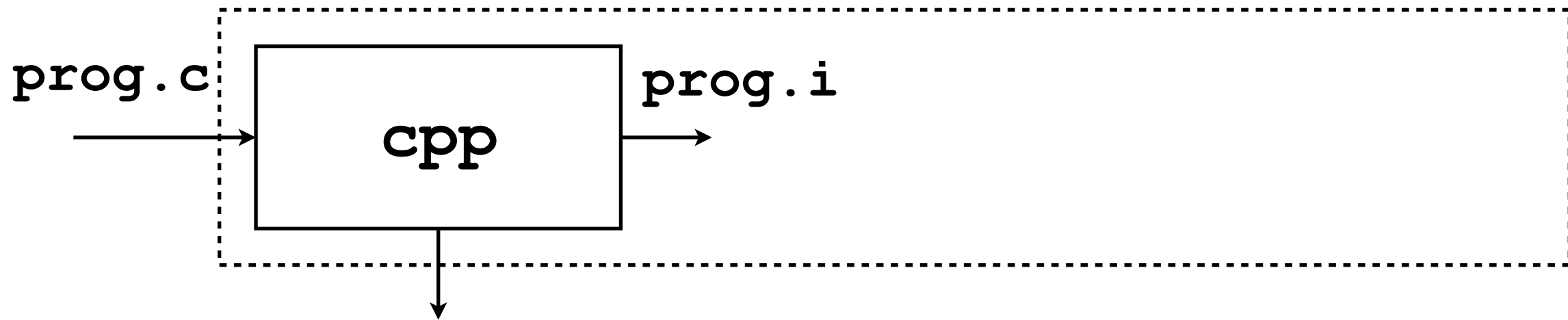  - Some components may even be useful as standalone tools

# Combining compilers

- The classic Unix C compiler, `cc`, is implemented by a pipeline of compilers:

`prog.c`  →  [ **cc** ]  →  `prog.o`

# Combining compilers

- The classic Unix C compiler, `cc`, is implemented by a pipeline of compilers:

```
prog.c ──▶ ┌─────────┐ prog.i
           │   cpp   │ ──▶
           └────┬────┘
                │
                ▼
```

**cpp:**  the C preprocessor, expands the use of macros and compiler directives in the source program

# Combining compilers

- The classic Unix C compiler, **cc**, is implemented by a pipeline of compilers:

```
prog.c ──→ [ cpp ] ──→ [ cc1 ] ──→ prog.S
```

**cc1**: the main C compiler, which translates C code to the assembly language for a particular machine

# Combining compilers

- The classic Unix C compiler, **cc**, is implemented by a pipeline of compilers:

```
prog.c                                                          prog.o
  ──►  [ cpp ]  ──►  [ cc1 ]  ──►  [ as ]  ──►
          │            │            │
          ▼            ▼            ▼
```
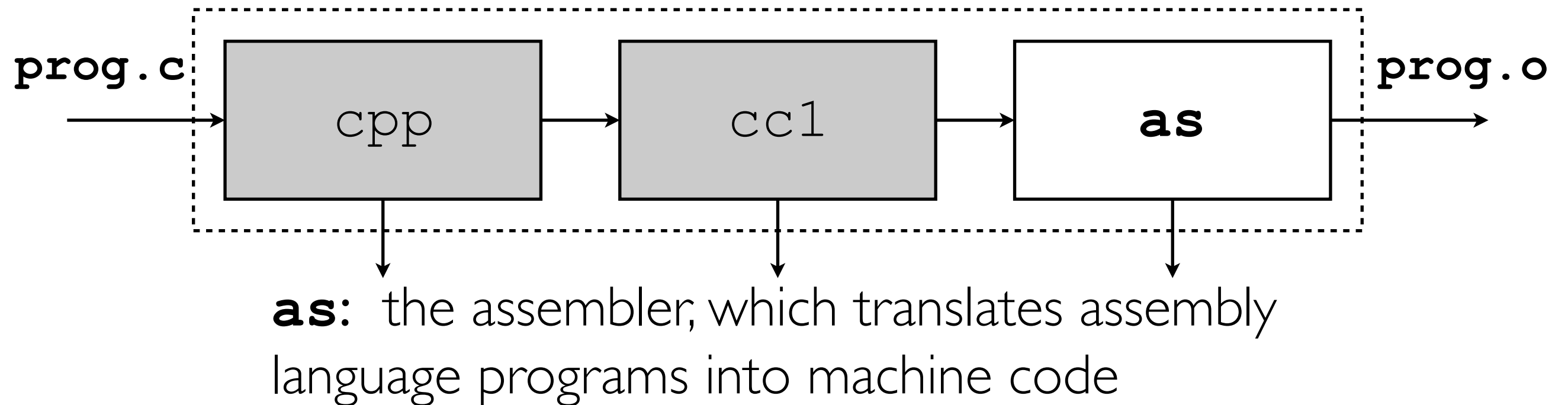
**as**: the assembler, which translates assembly language programs into machine code

# Advantages of modularity

- Some components (e.g., `as`) are useful in their own right

- Some components can be reused (e.g., replace `cc1` to build a C++ compiler)

- Some components (e.g., `cpp`) are machine independent, so they do not need to be rewritten for each new machine

- Modular implementations can be easier to write, test, debug, understand, and maintain

# Disadvantages of modularity?

- Performance

    It takes extra time to write out the data produced at the end of each stage

    It takes extra time to read it back in at the beginning of the next stage

    Later stages may need to repeat calculations from earlier stages if the information that they need is not included in the output of those earlier stages
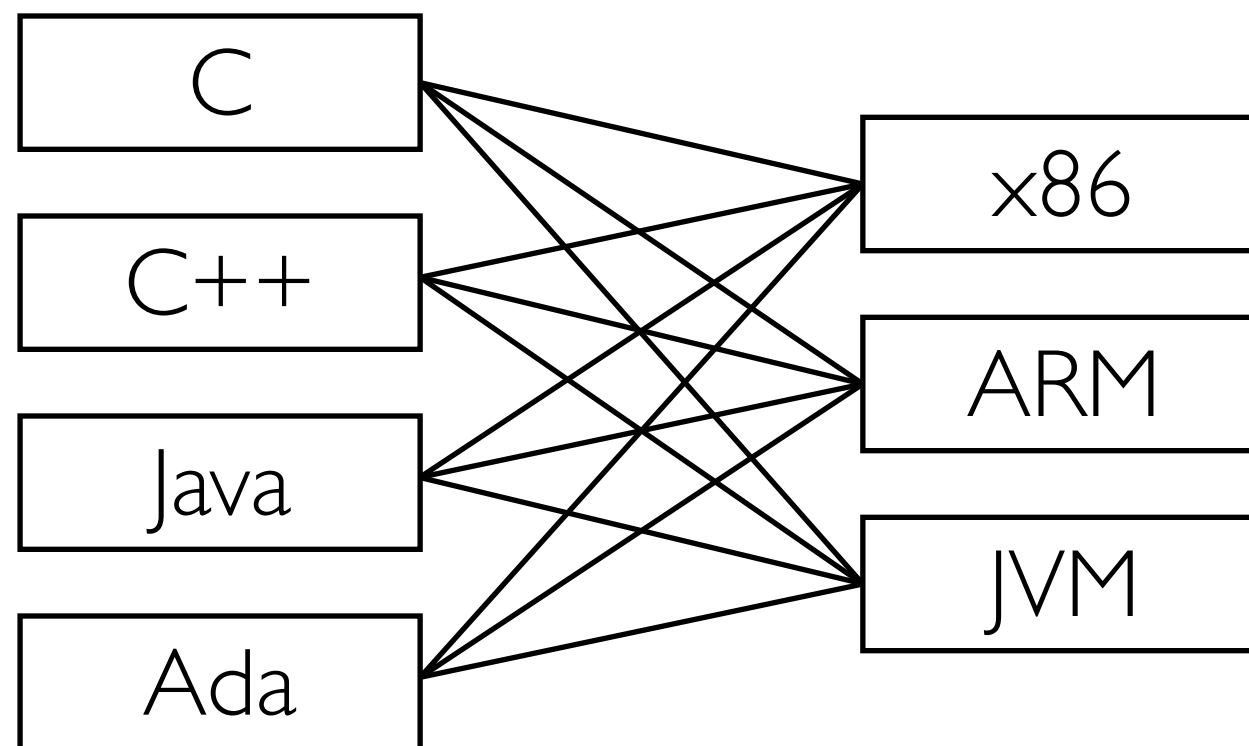
- **But modern machines and disks are pretty fast, and compilers are often complex, so modularity usually wins!**

# General building blocks

- A **front end** reads source programs (e.g., flat text files) and captures the corresponding abstract syntax in a collection of data structures (e.g., trees, graphs, arrays, …)

- A **middle end** analyzes and manipulates the abstract syntax data structures of a program

- A **back end** generates output(e.g., a flat, binary executable file) from the abstract syntax data structures of a program

- Substantial parts of these components can be shared by multiple tools
  - Example: the ghc (compiler) and ghci (interpreter) for Haskell use the same front and middle end components
  - Example: the g++ compiler for C++ and gcc compiler for C use the same middle and back end components
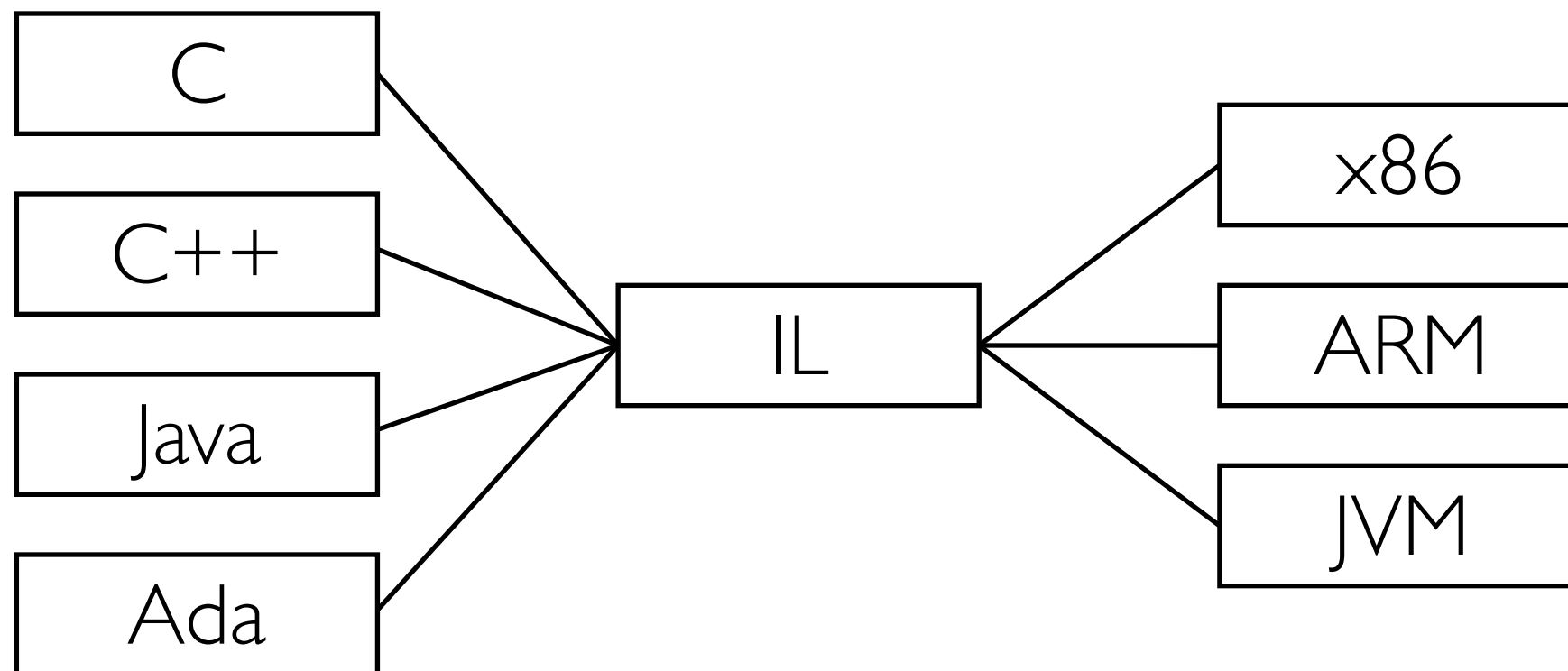
# Multiple languages and targets

- Suppose that we want to write compilers for n different languages, with m different target platforms.

| C | | x86 |
| C++ | | ARM |
| Java | | JVM |
| Ada | | |

- That's n x m different compilers!

# An intermediate language

- Alternatively: design a general purpose, shared "intermediate language":

```
  C  ─┐
       ├──┐           ┌── x86
 C++  ─┤   ├─── IL ───┤── ARM
       │   │           └── JVM
 Java ─┤   │
       │   │
  Ada ─┘   ┘
```

- Now we only have n front ends and m back ends to write!

- The biggest challenge is to find an intermediate language that is general enough to accommodate a wide range of languages and machine types

# Summary

- **Basic principles**

  programs as data

- **Interpreters and compilers**

  correctness means preserving semantics

- **The compiler pipeline / "phase structure"**

  source input, lexical analysis, parsing, static analysis, code generation, optimization

- **Modularity**

  Techniques for simplifying compiler construction tasks