# CS 321: Languages and Compiler Design I

Mark P Jones, Portland State University

Winter 2014

Week 1: Basics of Compiler Structure

1

---

# Why?

2

---

Ideas:    • Play a game    High Level
          • Display an image
          • Search a database
          • Visit a web page
          • etc ...

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

Machines:    • Read a value from memory    Low Level
             • Add two numbers
             • Compare two numbers
             • Write a value to memory
             • etc ...

3

---

Ideas:    • Play a game    High Level
          • Display an image
          • Search a database
          • Visit a web page
          • etc ...

express

          • Evaluate an expression
          • Execute a computation multiple times
Languages:    • Call a function
          • Save a result in a variable
          • ...

translate

Machines:    • Read a value from memory    Low Level
             • Add two numbers
             • Compare two numbers
             • Write a value to memory
             • etc ...

4

---

Ideas:    High Level

express

          Admiral Grace
          Hopper
          (Photo: via Wikipedia)

Languages:

translate    Could we program a
             computer to do this?

             human ingenuity
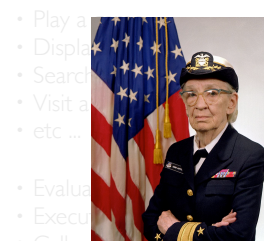             required

Machines:    Low Level

5

---

Ideas:    High Level

express

          Admiral Grace
          Hopper
          (Photo: via Wikipedia)

Languages:
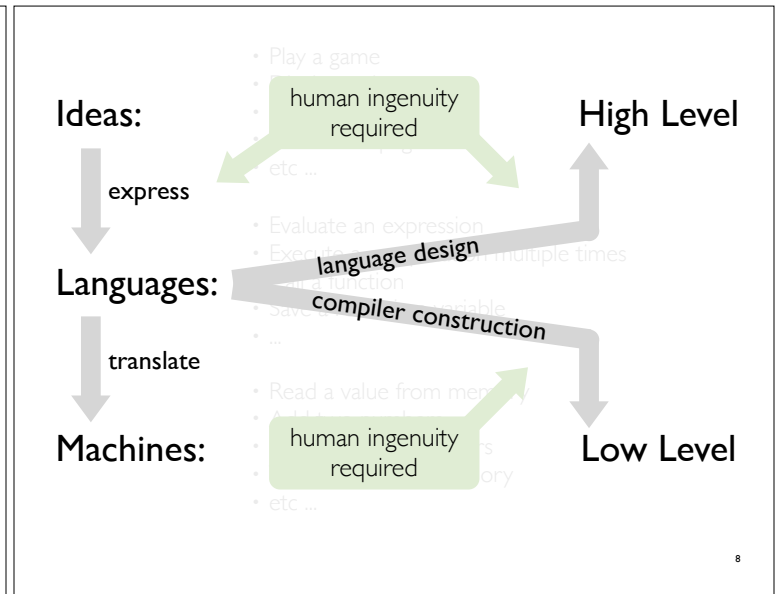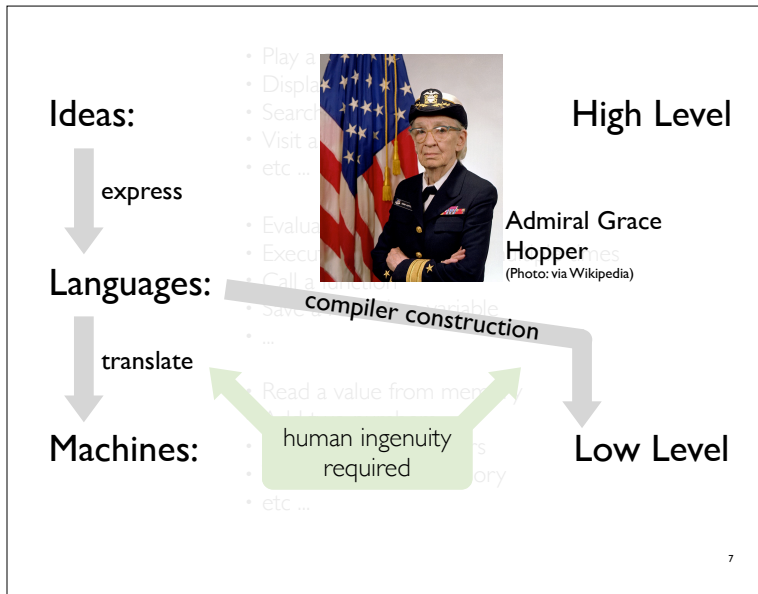
translate    Could we program a
             computer to do this?

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

Machines:    Low Level

6

## Slide 7

Ideas:

express

Languages:

translate

Machines:

High Level

compiler construction

human ingenuity required

Low Level

Admiral Grace Hopper
(Photo: via Wikipedia)

## Slide 8

Ideas:

express

Languages:

translate

Machines:

High Level

human ingenuity required

language design

compiler construction

human ingenuity required

Low Level

## A Big Picture Vision for CS 321/322:

• Build foundations in language design & compiler construction

• Empower developers to:

  • Express their ideas more directly

  • Execute their designs on a computer

• The long term goal is to build better tools that:

  • open programming to more people and more applications

  • increase programmer productivity

  • enhance software quality (functionality, reliability, security, performance, power, ...)
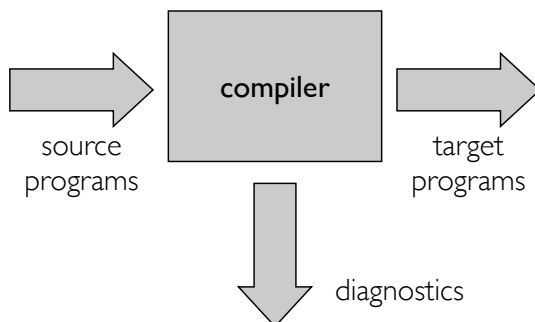
## Basics of Compiler Structure

## What is a compiler?

Compilers are translators:

compiler

source programs

target programs

diagnostics

## Compilers, compilers, compilers, ...

Compilers show up in many different forms:

- Translating programs in high-level languages like C, Fortran, Ada, 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/…

- …

The tools, techniques, and concepts that we cover are more general than you might have thought!

## Compiler inputs and outputs

### 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.

### error diagnostics
essential for serious program development.

13

## Example

source program

```
// A simple mini test program

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

compile

run

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

How does this work?

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
```

run

semantics

14

## Critical properties of a compiler

Always:

- **Correctness**: the compiler should produce valid output for any valid input program, and the output should have the same semantics as the input

Almost Always:

- **Performance** of compiled code: time, space, ...
- **Performance** of compiler: time, space, ...
- **Diagnostics**: to permit early and accurate diagnosis and detection of programming errors

15

## Other desirable features, in practice

- Support for large programming projects, including:
  - **Separate compilation**, reducing the amount of recompilation that is needed when part of a program is changed
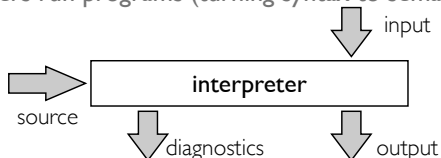  - Use of libraries, enabling effective **software reuse**

- 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)
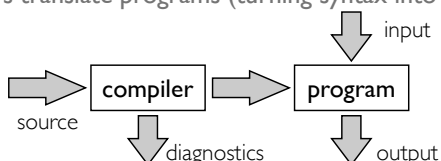
16

## Interpreters vs compilers

- Interpreters run programs (turning syntax to semantics)

input

| interpreter |

source

diagnostics       output

- Compilers translate programs (turning syntax into syntax)

input

| compiler | → | program |

source

diagnostics       output

17

## "Doing" vs "Thinking about doing"

- Compilers translate programs (turning syntax to syntax)
- Interpreters run programs (turning syntax to semantics)
- Example:
  - Use your calculator to evaluate (1+2)+(3+4):
    Answer: 10
  - Tell me what buttons to press to evaluate (1+2)+(3+4):
    Answer: | 1 | + | 2 | = | M | 3 | + | 4 | + | MR | = |

- We'll mostly focus on compilers, but will also use many of the same tools to work with interpreters.

18

## Language vs implementation

- Be very careful to distinguish between languages and their implementations

- It doesn't make much sense to talk about a "slow language"; speed is a property of the implementation, not the language.

- It doesn't make much sense to talk about a "compiled language"; again, "compiled" is a detail of the implementation, not the language

---

## 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;
}
```

compile

target program

```
        .file   "squares.s"
        .comm   _sqr0,4
        .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    %eax,_sqr0
        subl    $4,%esp
        andl    $0xfffffff0,%esp
        movl    %eax,(%esp)
        call    _print
        movl    _sqr0,%eax
        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
```

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

---

## What is this?

# 41+1

Black pixels on a white background

A sequence of characters

A sequence of "tokens"

An expression

A valid expression

Meaning: 42

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 sentence must use correct grammar: ✗    parser

    my walking up left tree dog

---

## "Compiling" English

- The sentence must make sense ✗

    This sentence is not true

- The sentence must not be ambiguous ✗

    Fruit flies like a banana

- The sentence must fit in context ✗

    This lecture is about geography

static analysis

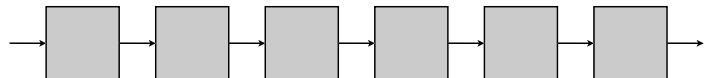- Finally, we are ready to translate! ✓

    Compilers are very interesting

ready for "code generation"
(i.e., for CS 322)
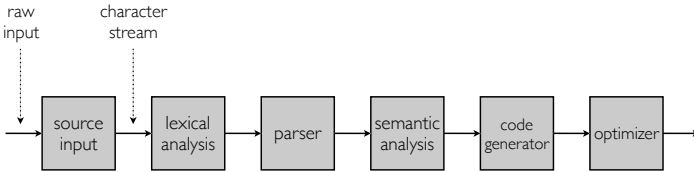
---

## The compiler pipeline

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

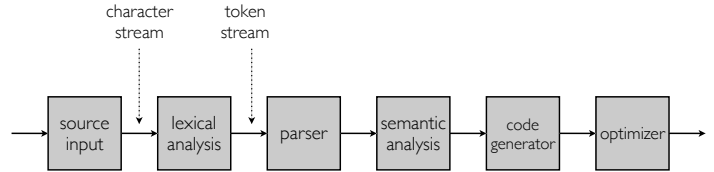## Source input — (not a standard term)



- 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, ...

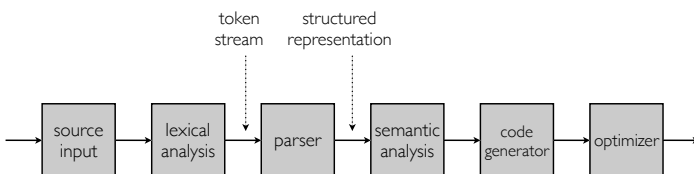  The operating system usually takes care of most of this ...

## 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



- 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)

## Semantic analysis



- Check that the program is reasonable:

  no references to unbound variables

  no type inconsistencies

  etc...

## Code generation



- Generate an appropriate sequence of machine instructions as output

- Different strategies are needed for different target machines

## Optimization



- 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 → semantic 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

## Phases and passes

- A <u>phase</u> is a *logical* stage in a compiler pipeline

- A <u>pass</u> is a *physical* traversal over the representation of a program

- Several phases may be combined in one pass

- Passes may be run in sequence or in parallel

- Some languages are specifically designed so that they can be implemented in a single pass

## Snapshots from a "mini" compiler pipeline

## Snapshots from a "mini" compiler pipeline

- In this section, we'll trace the results of passing the following program through a compiler for a language called "mini":

```
// 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

- We'll see more about the "mini" compiler in this week's lab session

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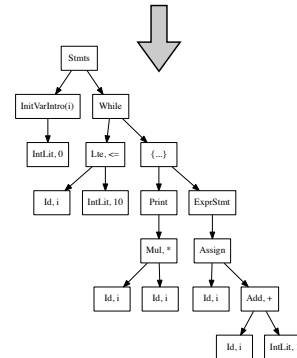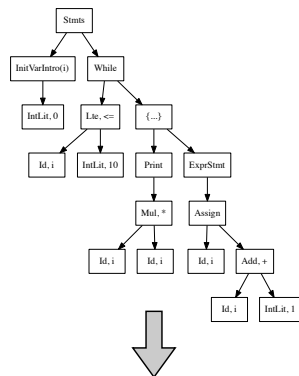
| INT | ID(i) | = | INTLIT(0) | Semicolon ";" | WHILE
| Open parenthesis "(" | ID(i) | <= | INTLIT(10)
| Close parenthesis ")" | Open brace "{" | PRINT | ID(i)
| * | ID(i) | Semicolon ";" | ID(i) | = | ID(i) | +
| INTLIT(1) | Semicolon ";" | Close brace "}" |

## 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



{ (i, int) }

## 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
```
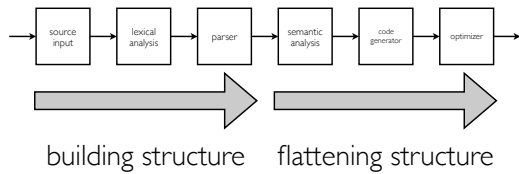
## Reflections

## Reflections: flat vs structured

- We might prefer to think of them as texts, but the source and target programs here are really just "flat" sequences of numbers

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

building structure    flattening structure

---

## Reflections: Syntax and semantics

- Syntax

  Written/spoken/symbolic/physical form; how things are communicated

- Semantics

  What those things mean

---

# Syntax    and    Semantics

(written form)                    (meaning)

$$\text{``}41 + 1\text{''} \Rightarrow 42$$

---

# Syntax    and    Semantics

(written form)                    (meaning)

$$\frac{A \wedge B}{A} \qquad \frac{A \wedge B}{B}$$

$$\frac{A \quad B}{A \wedge B} \qquad \Rightarrow$$

| A | B | A∧B |
|---|---|-----|
| F | F | F |
| F | T | F |
| T | F | F |
| T | T | T |

---

# Syntax    and    Semantics

(written form)                    (meaning)

"The elephant
    sat on a chair"

⇓

---

# Syntax    and    Semantics

(spoken form)                    (meaning)

⇓

'The elephant sat on a chair"

⇓

# Syntax    and    Semantics

(written form)                    (meaning)



49

---

# Syntax    and    Semantics

(written form)                    (meaning)

```
// A simple Mini test program

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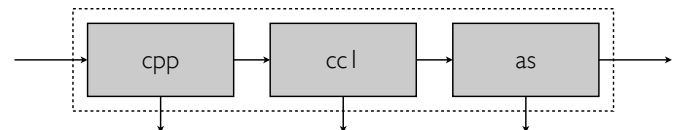
compiler                    execution

50

---

# Modularity in compiler design

51

---

## Combining compilers

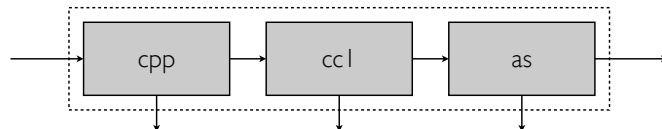• The standard Unix C compiler is structured as a pipeline of compilers:



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

52

---

## Combining compilers

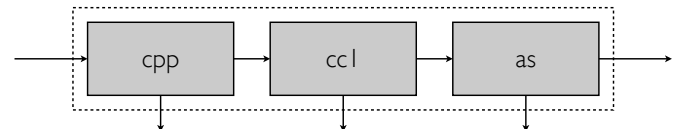• The standard Unix C compiler is structured as a pipeline of compilers:



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

53

---

## Combining compilers

• The standard Unix C compiler is structured as a pipeline of compilers:



**as:** the assembler, which translates assembly language programs into machine code

54

## 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

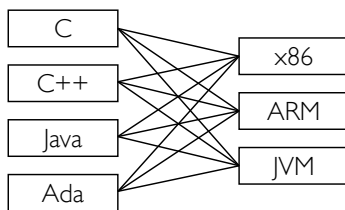- Modular implementations can be easier to understand, test, debug, write, ...

## 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!

## Multiple languages and targets

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



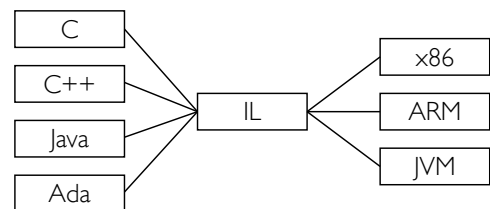- That's n x m different compilers!

## An intermediate language

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



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

## Front ends and back ends

- Front end: those parts of a compiler that depend most heavily on the source language

    Source input, lexical analysis, parsing, static analysis

- Back end: those parts of a compiler that depend most heavily on the target language

    Code generation, optimization, assembly

- 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

    syntax and semantics
    correctness means preserving semantics

- The compiler pipeline

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

- Modularity

    Techniques for simplifying compiler construction tasks