# Compiler Construction

## Lecture 16



## Bootstrapping

Jeremy Yallop, Lent 2025

jeremy.yallop@cl.cam.ac.uk

*Chapter 13 of*
**Basics of Compiler Design**
Torben Ægidius Mogensen
http://hjemmesider.diku.dk/~torbenm/Basics/

# Notation

A program

An interpreter

A machine

f

L

L2

L1

L

Computes function f
written in language L

Interprets language L2
written in language L1

Executes code
in language L

A compiler



Translates language A into language B
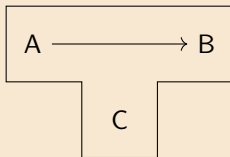Written in language C

# Examples

To execute a program

To execute an interpreter

To execute a compiler

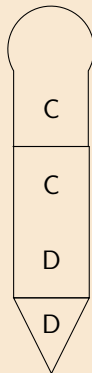

we run it on a machine

we run it on a machine

we run it on a machine

Run a program
written in language C

C

C

on an interpreter for C
written in language D

D

on a D machine

D

(Note: the languages must match)
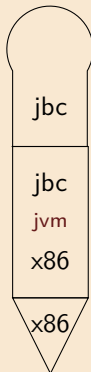
Run a program
written in Java byte code

on an interpreter for Java byte code
written in x86 code

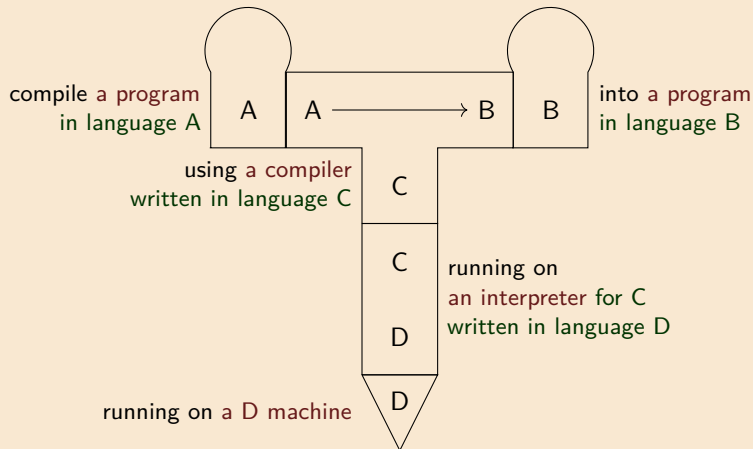on a x86 machine

jbc

jbc
jvm
x86

x86

compile a program
in language A

A   A ⟶ B   B

into a program
in language B

using a compiler
written in language C

C

C

running on
an interpreter for C
written in language D

D

running on a D machine   D

compile a program
in Java — Java | Java $\xrightarrow{\text{javac}}$ jbc | jbc — into a program
in Java byte code

using the javac compiler
written in Java byte code — jbc

jbc
jvm
x86 — running on the JVM
written in x86 code

running on an x86 machine — x86

Thanks to David Greaves for the example

# Compiling compilers

The OCaml compiler
is written in OCaml



Puzzle: how was the compiler compiled?
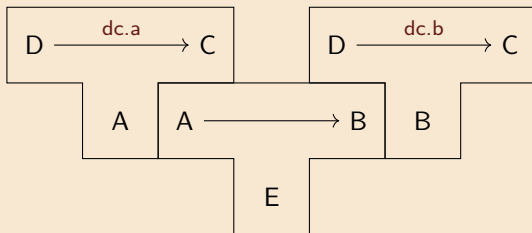
Compilers can be translated, just like any other program:



a compiler from D to C
in language A

a compiler from D to C
in language B

compile programs from A to B

**We have**:
a compiler from ML to arm
that runs on arm

ML → arm

arm

**We want**:
a compiler from ML to x86
that runs on x86

ML → x86

x86

**We have**:
a compiler from ML to arm
that runs on arm

ML → arm

arm

**We want**:
a compiler from ML to x86
that runs on x86

ML → x86

x86

ML ⟶ x86

ML

**1**. write an ML-to-x86 compiler in ML

**We have**:
a compiler from ML to arm
that runs on arm

ML → arm

arm

**We want**:
a compiler from ML to x86
that runs on x86

ML → x86

x86

**1**. write an ML-to-x86 compiler in ML

**2**. compile the compiler for arm

ML ⟶ x86

ML | ML ⟶ arm | arm

arm

arm

ML ⟶ x86

# Porting a compiler to a new platform

Notation

Examples

Compiling
compilers

●●●●

Full
bootstrap

Trusting
trust

**We have**:
a compiler from ML to arm
that runs on arm

ML → arm

arm

**We want**:
a compiler from ML to x86
that runs on x86

ML → x86

x86

1. write an ML-to-x86 compiler in ML

2. compile the compiler for arm

3. run the compiler on arm to compile itself

ML ——→ x86    ML ——→ x86

ML    ML ——→ x86    x86

arm

arm

# Full bootstrap

Previous example: *half bootstrap* (needs existing compiler for the language).

New example: *full bootstrap* (no existing ML compiler for the language)

**We want**:
a compiler from **XL** to arm
that runs on arm

XL ⟶ arm

arm

**We have**:
a compiler from ML to arm
that runs on arm

ML ⟶ arm

arm

**We have**:
a compiler from ML to arm
that runs on arm

ML → arm
arm

**We want**:
a compiler from **XL** to arm
that runs on arm

XL → arm
arm

**We have**:
a compiler from ML to arm
that runs on arm

| ML → arm |
| arm |

**We want**:
a compiler from **XL** to arm
that runs on arm

| XL → arm |
| arm |

**1**. write a quick-and-dirty (QAD)
**XL-to-ML** compiler in ML

| XL $\xrightarrow{\text{qad}}$ ML |
| ML |

**We have**:
a compiler from ML to arm
that runs on arm

ML → arm

arm

**We want**:
a compiler from **XL** to arm
that runs on arm

XL → arm

arm

**1**. write a quick-and-dirty (QAD)
**XL-to-ML** compiler in ML

**2**. compile the QAD compiler for arm

XL $\xrightarrow{\text{qad}}$ ML

ML  ML ⟶ arm  arm

XL $\xrightarrow{\text{qad}}$ ML

arm

arm

**We have**:
a compiler from ML to arm
that runs on arm

| ML → arm |
|---|
| arm |

**We want**:
a compiler from **XL** to arm
that runs on arm

| XL → arm |
|---|
| arm |

**1**. write a quick-and-dirty (QAD)
**XL-to-ML** compiler in ML

**2**. compile the QAD compiler for arm

**3**. Write a real **XL-to-arm** compiler in **XL**

$$XL \xrightarrow{\text{real}} arm$$

| XL |
|---|

**We have**:
a compiler from ML to arm
that runs on arm

| ML → arm |
| arm |

**We want**:
a compiler from **XL** to arm
that runs on arm

| XL → arm |
| arm |

**1**. write a quick-and-dirty (QAD)
**XL-to-ML** compiler in ML

**2**. compile the QAD compiler for arm

**3**. Write a real **XL-to-arm** compiler in **XL**

**4**. Use the QAD compiler to compile
the real compiler to ML

**We have**:
a compiler from ML to arm
that runs on arm

$ML \rightarrow arm$

arm

**We want**:
a compiler from **XL** to arm
that runs on arm

$XL \rightarrow arm$

arm

**1**. write a quick-and-dirty (QAD)
**XL-to-ML** compiler in ML

**2**. compile the QAD compiler for arm

**3**. Write a real **XL-to-arm** compiler in **XL**

**4**. Use the QAD compiler to compile
the real compiler to ML

**5**. Compile the resulting ML program to arm

$XL \xrightarrow{\text{real}} arm$

ML  |  ML $\longrightarrow$ arm  |  arm

arm

arm

$XL \xrightarrow{\text{real}} arm$

**We have**:
a compiler from ML to arm
that runs on arm

$$ML \rightarrow arm$$
arm

**We want**:
a compiler from **XL** to arm
that runs on arm
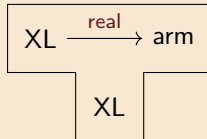
$$XL \rightarrow arm$$
arm

**1**. write a quick-and-dirty (QAD)
**XL-to-ML** compiler in ML

**2**. compile the QAD compiler for arm

**3**. Write a real **XL-to-arm** compiler in **XL**
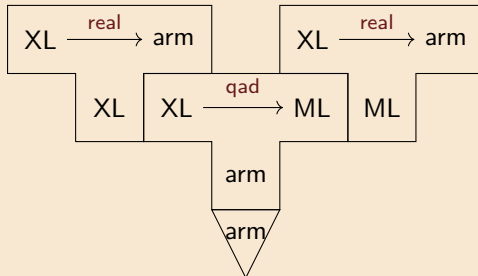
**4**. Use the QAD compiler to compile
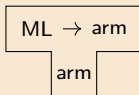the real compiler to ML

**5**. Compile the resulting ML program to arm

**6**. Use the generated compiler to compile itself

$$XL \xrightarrow{\text{real}} arm$$
XL    $$XL \xrightarrow{\text{real}} arm$$    arm
arm
arm

$$XL \xrightarrow{\text{real}} arm$$

The *speed* of the quick-and-dirty compiler does not matter much
(We could even use a quick-and-dirty interpreter instead)

We don't need to give the quick-and-dirty compiler to users

Once the real compiler works,
we can discard the quick-and-dirty compiler altogether

# Trusting trust

"The cutest program I ever wrote"
– Ken Thompson

*(Reflections on Trusting Trust)*

**Aim**: modify a compiler to compromise `login`

**Warm up**: teach a compiler about vertical tabs

C compilers have code to interpret escape sequences like `\n` in `"Hello, world\n`:

```
...
c = next ();
if (c != '\\') return c;
c = next ();
if (c == '\\') return '\\';
if (c == 'n') return '\n';
...
```

**Q**: how can we add support for vertical tabs `\v`?

(Assume the C compiler is bootstrapped.)

Notation

Examples

Compiling
compilers

Full
bootstrap

Trusting
trust

**Step 1**: hard-code the ASCII code for \v in the compiler source:

```
c = next ();
if (c == '\\') return '\\';
if (c == 'n') return '\n';
if (c == 'v') return 11;
...
```
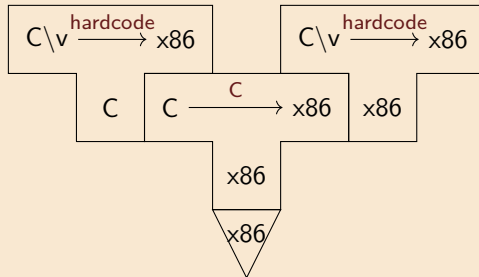
Recompile the compiler source using the installed C compiler:



Now we have a C compiler that supports \v in C programs. Install it.

# The compiler has learnt about \v

Notation

Examples

Compiling
compilers

Full
bootstrap

Trusting
trust

**Step 2**: modify the compiler source again to remove the hardcoded constant:

```
c = next();
if (c == '\\') return '\\';
if (c == 'n') return '\n';
if (c == 'v') return '\v';
...
```
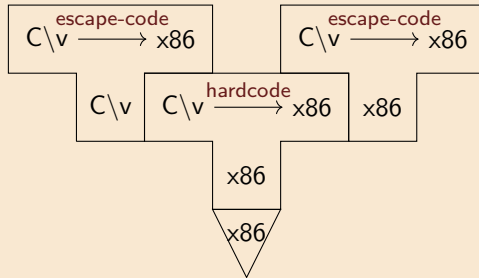
Recompile the modified source using the freshly installed C compiler:



The C compiler has learnt to translate \v (but there's no record in the source!)

Notation

Examples

Compiling
compilers

Full
bootstrap

Trusting
trust

**Plan**: repeat the process to compromise the login command.

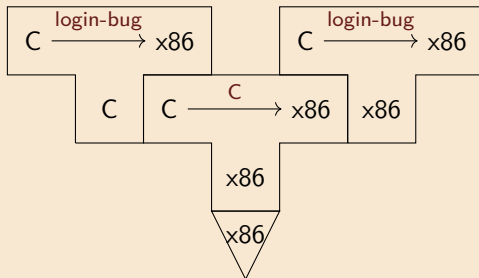**Step 1**: update the C compiler's code to detect login.c and insert a bug:

```c
void compile(const char *program) {
  if (matches(program, "< login code >")) {
      compile("< code for backdoor >");
  }
  ...
```
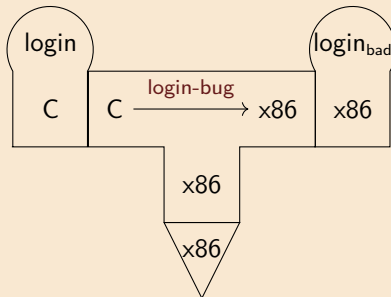
Compile and install the new C compiler:

Now the compiler will miscompile login:



**Problem**: people will easily spot the bug in the compiler source.
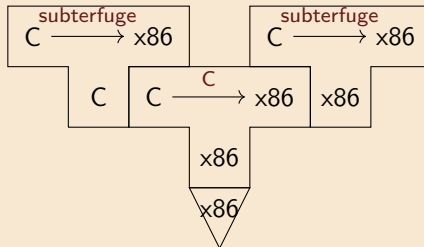
**Step 2**: update the C compiler code to detect compiler.c and insert a 2nd bug:

```c
void compile(const char *program) {
  if (matches(program, "< login code >")) {
    compile("< code for backdoor >");
  }
  if (matches(program, "< compiler code >")) {
    compile("< code for miscompilation >");
  }
```

Compile and install the new C compiler:



**Finally**: remove the bugs from the compiler source.

**The compiler has learnt to insert backdoors**

Notation

Examples

Compiling
compilers

Full
bootstrap

Trusting
trust

The compiler will still miscompile
login:

The compiler will now also miscompile
the compiler:



The system is **compromised**, with **no trace** in the login or compiler source.
We need to *debootstrap* to recover an uncompromised compiler.