

Certifying cost annotations in compilers

Nicolas Ayache

Post-doctorate at PPS — Paris 7

18th november 2010

The CerCo project

3 years European project (FP7)

Laboratories

- ▶ University of Bologna — Claudio Sacerdoti Cohen
- ▶ University of Edinburgh — Randy Pollack
- ▶ **University of Paris Diderot (PPS)**

Roberto Amadio (site leader)

Nicolas Ayache (post-doc)

Yann Régis-Gianas

Ronan Saillard (Ph.D. student)

State of the art: Worst Case Execution Time

AbsInt (Abstract Interpretation)

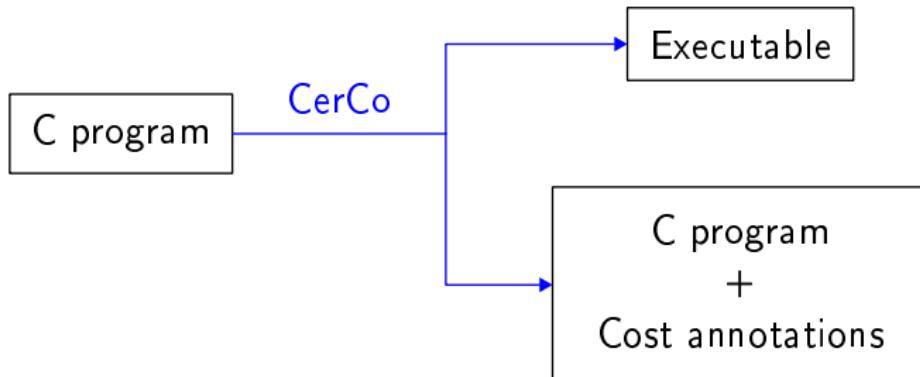
- ✓ Concrete time
- ✗ Binary analysis
- ✗ User interaction (loop iteration)
- ✗ Not formally proven

Linear Logic

- ✓ Formal framework
- ✗ Complexity classes
- ✗ Verbose

Goal

Formally sound, cost annotating compiler



✓ Overapproximated concrete complexity

CerCo's approach: problematic

What is the cost of evaluating $\text{tab}[i]+1$?

Depends on:

- ▶ The variable final locations
- ▶ The way memory accesses are compiled
- ▶ The way operations are compiled

⇒ Depends on the compilation process

CerCo's approach: solution

Bad solution

Consider worst case scenarios

- ✗ Too imprecise
- ✗ Optimizations lost
- ✗ Not modular

CerCo's approach: solution

Bad solution

Consider worst case scenarios

- ✗ Too imprecise
- ✗ Optimizations lost
- ✗ Not modular

CerCo's solution

Symbolic cost update: label



CerCo's approach: common considerations

Operational semantics

- ▶ *Without labels*: $\mathcal{P}(\text{Prog} \times \text{State} \times \text{State})$
- ▶ *With labels*: $\mathcal{P}(\text{Prog} \times \text{State} \times \text{label trace} \times \text{State})$

CerCo's approach: common considerations

Operational semantics

- ▶ *Without labels*: $\mathcal{P}(\text{Prog} \times \text{State} \times \text{State})$
- ▶ *With labels*: $\mathcal{P}(\text{Prog} \times \text{State} \times \text{label trace} \times \text{State})$

Annotation semantics

✓ Annotation = Instruction of the source language

- ▶ Explicits cost annotations
- ▶ Analysis tools for cost synthesis

CerCo's approach: common considerations

Operational semantics

- ▶ *Without labels*: $\mathcal{P}(\text{Prog} \times \text{State} \times \text{State})$
- ▶ *With labels*: $\mathcal{P}(\text{Prog} \times \text{State} \times \text{label trace} \times \text{State})$

Annotation semantics

✓ Annotation = Instruction of the source language

- ▶ Explicits cost annotations
- ▶ Analysis tools for cost synthesis

Demo

Outline

1 Toy compiler

- Languages
- Compilation
- Labelling
- Annotation

2 Realistic C compiler

- Architecture
- Labelling
- Experiments

3 Future work

Outline

1 Toy compiler

- Languages
- Compilation
- Labelling
- Annotation

2 Realistic C compiler

- Architecture
- Labelling
- Experiments

3 Future work

Overview

Imp → **VM** → **ASM**

- ▶ **Imp**: simple while language
- ▶ **VM**: virtual machine (stack operations)
- ▶ **ASM**: assembly

Overview

Imp → **VM** → **ASM**

- ▶ **Imp**: simple while language
- ▶ **VM**: virtual machine (stack operations)
- ▶ **ASM**: assembly

In the following...

- ▶ How to label the languages?
- ▶ How to adapt the proofs?
- ▶ How to keep the proofs modular?

Syntax

Imp

$$e ::= id \mid n \in \mathbb{N} \mid e + e \quad b ::= e < e$$
$$S ::= \text{skip} \mid id := e \mid S ; S \mid \text{if } b \text{ then } S \text{ else } S \mid \text{while } b \text{ do } S$$
$$P ::= \text{prog } S$$

VM

$$\mathcal{I} ::= \text{cnst}(n) \mid \text{var}(id) \mid \text{setvar}(id) \mid \text{add} \mid \text{branch}(k) \mid \text{bge}(k) \mid \text{halt}$$
$$P ::= \mathcal{I} \text{ list}$$

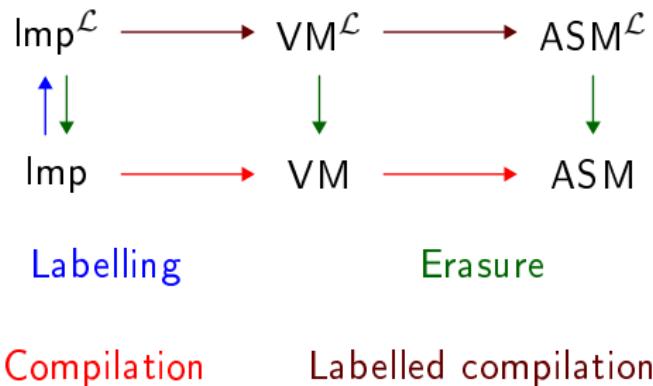
ASM

$$\begin{aligned} \mathcal{I} &::= \text{loadi } R, n \mid \text{load } R, \text{addr} \mid \text{store } R, \text{addr} \mid \text{add } R, R, R \\ &\quad \mid \text{branch } k \mid \text{bge } R, R, k \mid \text{halt} \end{aligned}$$
$$P ::= \mathcal{I} \text{ list}$$

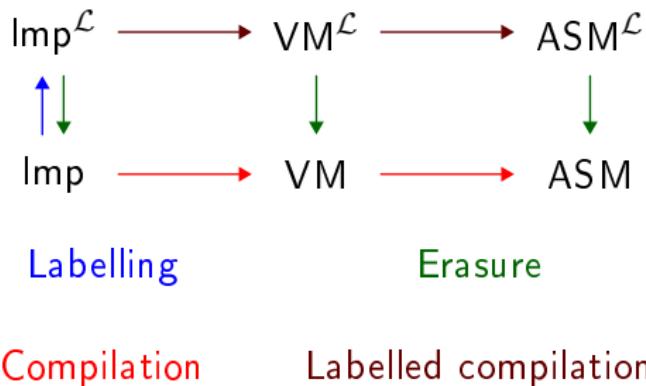
Labelled syntax

 $\text{Imp}^{\mathcal{L}}$
$$\begin{aligned} e &::= id \mid n \in \mathbb{N} \mid e + e \quad b ::= e < e \\ S &::= \text{skip} \mid id := e \mid S ; S \mid \text{if } b \text{ then } S \text{ else } S \mid \text{while } b \text{ do } S \mid \color{red}{\ell : S} \\ P &::= \text{prog } S \end{aligned}$$
 $\text{VM}^{\mathcal{L}}$
$$\begin{aligned} \mathcal{I} &::= \text{cnst}(n) \mid \text{var}(id) \mid \text{setvar}(id) \mid \text{add} \mid \text{branch}(k) \mid \text{bge}(k) \mid \text{halt} \mid \color{red}{\text{emit}(\ell)} \\ P &::= \mathcal{I} \text{ list} \end{aligned}$$
 $\text{ASM}^{\mathcal{L}}$
$$\begin{aligned} \mathcal{I} &::= \text{loadi } R, n \mid \text{load } R, \text{addr} \mid \text{store } R, \text{addr} \mid \text{add } R, R, R \\ &\quad \mid \text{branch } k \mid \text{bge } R, R, k \mid \text{halt} \mid \color{red}{\text{emit } \ell} \\ P &::= \mathcal{I} \text{ list} \end{aligned}$$

Simulation



Simulation

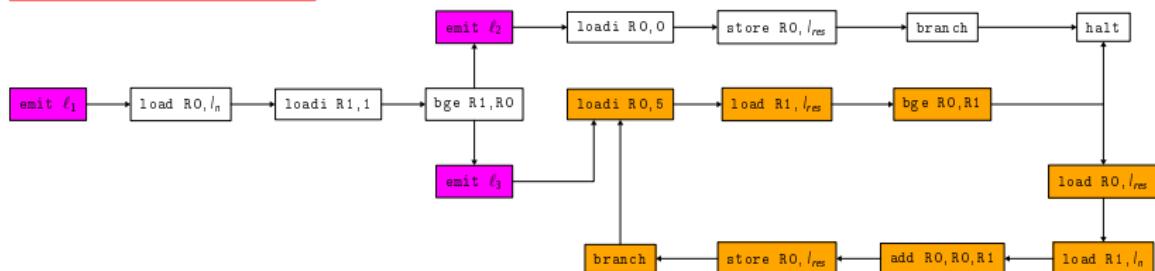


Theorem: diagram commutativity

$$\left. \begin{array}{l}
 \mathcal{E}_{\text{Imp}} \circ \mathcal{L}_{\text{Imp}} = \text{Id}_{\text{Imp}} \\
 \mathcal{C}_{\text{Imp}} \circ \mathcal{E}_{\text{Imp}} = \mathcal{E}_{\text{VM}} \circ \mathcal{C}_{\text{Imp}}^{\mathcal{L}} \\
 \mathcal{C}_{\text{VM}} \circ \mathcal{E}_{\text{VM}} = \mathcal{E}_{\text{ASM}} \circ \mathcal{C}_{\text{VM}}^{\mathcal{L}}
 \end{array} \right\} \Rightarrow \mathcal{E}_{\text{ASM}} \circ \mathcal{C}^{\mathcal{L}} \circ \mathcal{L}_{\text{Imp}} = \mathcal{C}$$

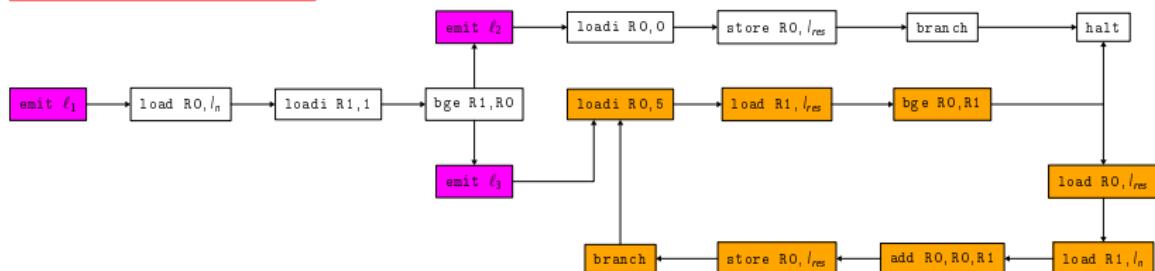
Example: soundness

```
prog
 $\ell_1$ :
if n < 1 then
   $\ell_2$ : res := 0
else
   $\ell_3$ :
    while res < 5 do
      res := res + n
```



Example: soundness

```
prog
 $\ell_1$ :
if n < 1 then
   $\ell_2$ : res := 0
else
   $\ell_3$ :
    while res < 5 do
      res := res + n
```

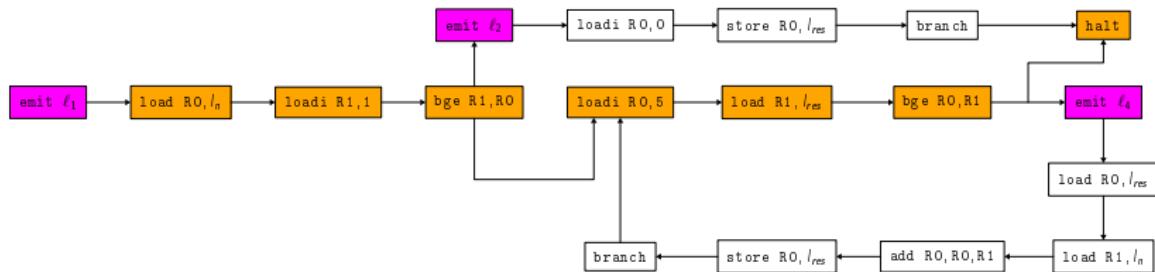


Loop with no label \Rightarrow not constant cost code

Example: precision

```
prog
 $\ell_1$ :
if n < 1 then
   $\ell_2$ : res := 0
else

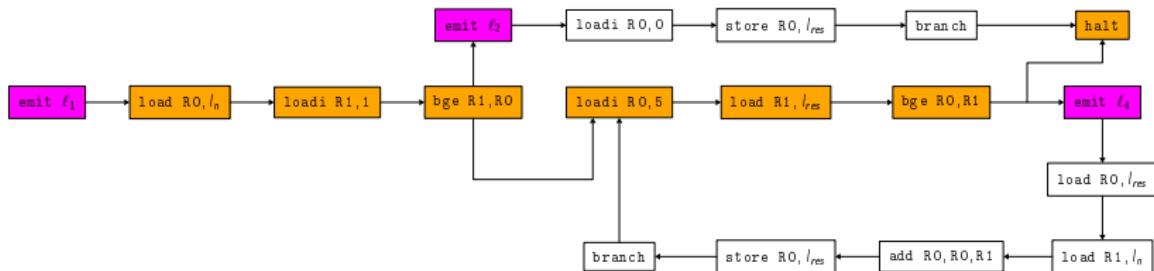
  while res < 5 do
     $\ell_4$ : res := res + n
```



Example: precision

```
prog
ℓ1:
if n < 1 then
  ℓ2: res := 0
else

  while res < 5 do
    ℓ4: res := res + n
```



From emit ℓ_1 : three paths with different costs \Rightarrow imprecision

Labelling criteria

Soundness

Every reachable code is in the scope of a label.
At least one label inside each loop.

Precision

Two different paths to the next labels have the same cost.

Criteria

The labelling must be **sound** and **precise**.
(This can be syntactically checked on the assembly code.)

✓ Nice plus is a reasonable **economy**: not too many labels.

Instrumentation

Cost deduction

Given: $\phi : \text{ASM Instruction list} \rightarrow \mathbb{N}$ loop free
may overapproximate

Deduced: $\kappa : \mathcal{L} \rightarrow \mathbb{N}$

Instrumentation

Cost deduction

Given: $\phi : \text{ASM Instruction list} \rightarrow \mathbb{N}$ $\left(\begin{array}{l} \text{loop free} \\ \text{may overapproximate} \end{array} \right)$

Deduced: $\kappa : \mathcal{L} \rightarrow \mathbb{N}$

Instrumentation

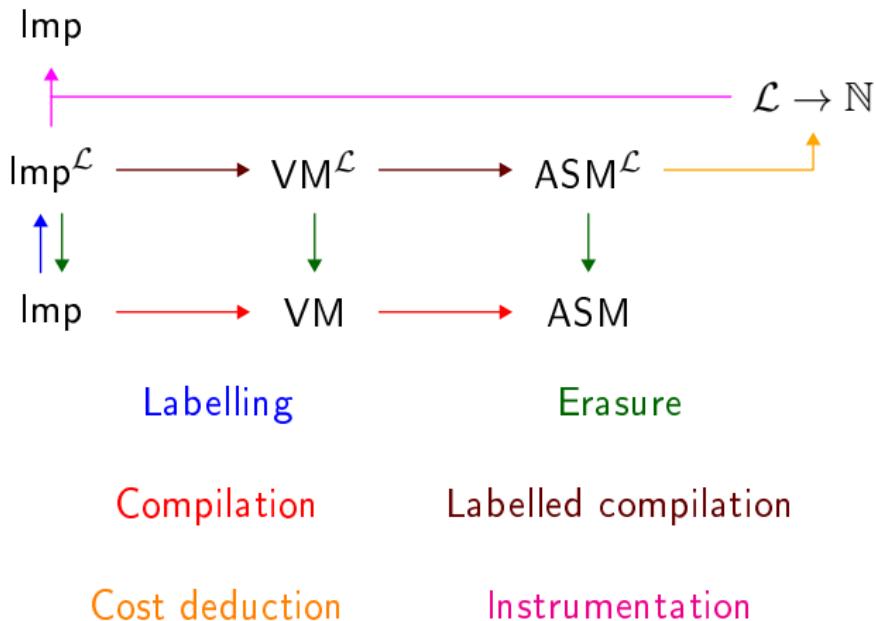
- ▶ Use a **fresh variable**
- ▶ Initialize it to **0**
- ▶ Replace **labels** with **increments** (following κ)

```
prog
  ℓ1:
    while res < 5 do      →
      ℓ2:
        res := res + n
```

```
prog
  _cost := 0;
  _cost := _cost + 2;
  while res < 5 do
    _cost := _cost + 4;
    res := res + n
```

Annotation

Annotation = Instrumentation \circ Labelling



Outline

1 Toy compiler

- Languages
- Compilation
- Labelling
- Annotation

2 Realistic C compiler

- Architecture
- Labelling
- Experiments

3 Future work

C → C_{light} → C_{minor} → RTL_{abs} → RTL → ERTL → LTL → LIN → MIPS

C to Cminor

C → Clight → Cminor → RTL_{abs} → RTL → ERTL → LTL → LIN → MIPS

- ▶ Inspired from **CompCert**
(GNU GPL and INRIA Non-Commercial)
- ▶ **C to Clight:** CIL
- ▶ **Clight to Cminor:** manual port from Coq to OCaml

RTL_{abs}

C → Clight → Cminor → RTL_{abs} → RTL → ERTL → LTL → LIN → MIPS

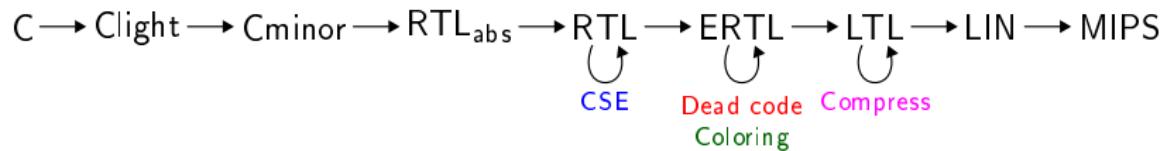
- ▶ Home made
- ▶ Architecture independent
- ▶ Retargetting simplified (inspired from Gimple in GCC)
 - ✓ Some common optimizations
 - ✗ Some optimisations lost

RTL to MIPS

C → Clight → Cminor → RTL_{abs} → **RTL** → ERTL → LTL → LIN → MIPS

- ▶ Adapted from pedagogical Pseudo-Pascal compiler
(François Pottier, Creative Commons)

Optimizations



- ▶ **CSE**: Common Subexpression Elimination
- ▶ **Dead code elimination**
- ▶ **Coloring**: variable locations
- ▶ **Graph compression**

Restrictions

C → **C_{light}** → C_{minor} → RTL_{abs} → **RTL** → ERTL → LTL → LIN → MIPS

C_{light}

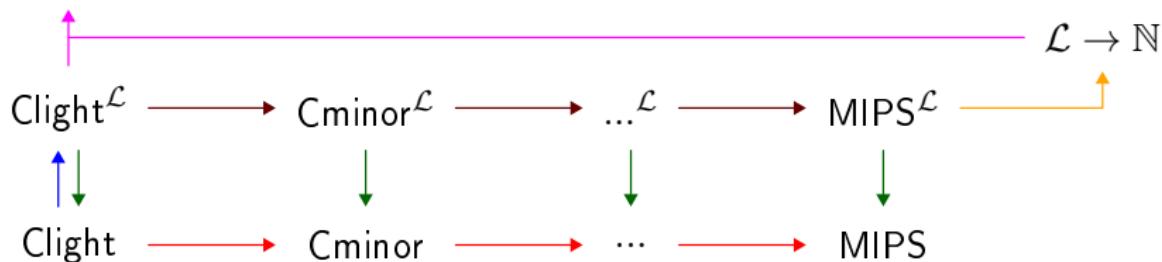
- ✗ long long and long double types
- ✗ longjmp and setjmp instructions
- ✗ Unreasonable forms of switch
- ✗ Unprototyped and variable-arity functions

RTL

- ✗ float

Overview

Clight



Labelling

Compilation

Cost deduction

Erasure

Labelled compilation

Instrumentation

Differences Imp / C

Side effect expressions `y = x++;`

Ternary expressions `x ? y+2*z : z`

Labels and Gotos `lbl: ... goto lbl;`

Function calls `f(&x);`

Differences Imp / C

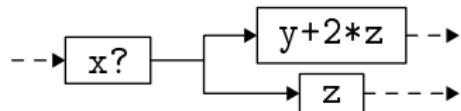
Side effect expressions `y = x++;`
eliminated by CIL → `_tmp = x; x = _tmp+1; y = _tmp;`

Ternary expressions `x ? y+2*z : z`

Labels and Gotos `lbl: ... goto lbl;`

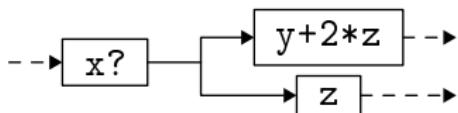
Function calls `f(&x);`

Ternary expressions

 $x? (y+2*z) : z$ 

Branching \Rightarrow 1 label per branch for precision

Ternary expressions

 $x? (y+2*z) : z$ 

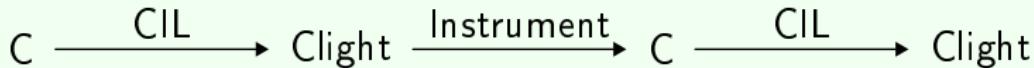
Branching \Rightarrow 1 label per branch for precision

Labelled expressions

 $x? (\ell_1: y+2*z) : (\ell_2: z)$

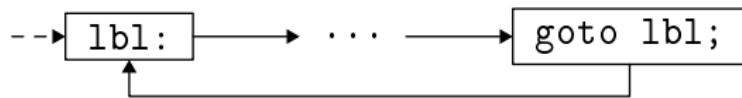
Instrumentation

- 1) Side effects inside expressions
- 2) Elimination by CIL



Labels and Gotos

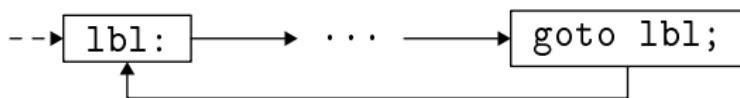
```
lbl: ... goto lbl;
```



Label \Rightarrow potential loop \Rightarrow 1 cost label needed

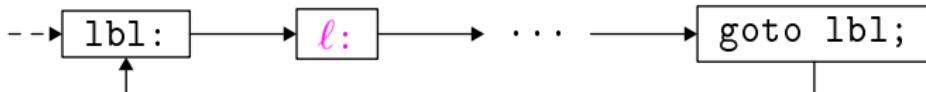
Labels and Gotos

```
lbl: ... goto lbl;
```



Label \Rightarrow potential loop \Rightarrow 1 cost label needed

```
lbl: ℓ: ... goto lbl;
```



Function calls

Function call: sequential instruction

```
x++;           → void f(int* x) {  
f(&x); ←     ...  
y = x; ←     return; }
```

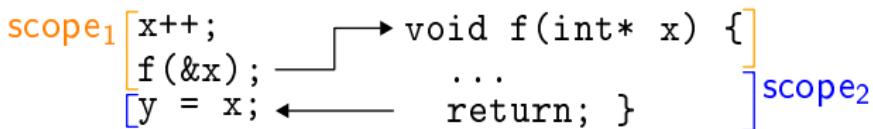
Function calls

Function call: sequential instruction

```
scope1 [x++;  
        f(&x);  
        y = x; ]  
              void f(int* x) {  
                ...  
                return; } ] scope2
```

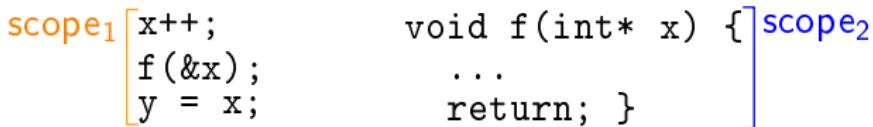
Function calls

Function call: sequential instruction



Function pointer: statically unresolvable destination

✓ Each function handles its cost



Benchmarks

	gcc -O0	acc	gcc -O1
badsort	55.93	34.51	12.96
fib	76.24	34.28	45.68
mat_det	163.42	156.20	54.76
min	12.21	16.25	3.95
quicksort	27.46	17.95	9.41
search	463.19	623.79	155.38

Outline

1 Toy compiler

- Languages
- Compilation
- Labelling
- Annotation

2 Realistic C compiler

- Architecture
- Labelling
- Experiments

3 Future work

- ▶ Formal proofs in Matita
- ▶ Real world example (Lustre)
- ▶ Frama-C plugin (demo)
- ▶ Retarget to 8051 (moduralize the target architecture)
- ▶ Functional languages

CerCo

- ▶ C sound compiler (untrusted for now...)
- ▶ Correct cost annotations (overapproximation)
- ▶ Symbolic cost labels
- ▶ Modular proofs

Thank you for your attention!

Questions?