



**CerCo**

INFORMATION AND COMMUNICATION  
TECHNOLOGIES  
(ICT)  
PROGRAMME

Project FP7-ICT-2009-C-243881 CerCo

**Report n. D2.1**  
**Compiler design and intermediate languages**

Version 1.0

Main Authors:

Roberto M. Amadio, Nicolas Ayache, Yann Régis-Gianas,  
Kayvan Memarian, Ronan Saillard

Project Acronym: CerCo

Project full title: Certified Complexity

Proposal/Contract no.: FP7-ICT-2009-C-243881 CerCo

**Abstract** We discuss the problem of building a compiler which can *lift* in a provably correct way informations on the execution cost of the object code to cost annotations on the source code. To this end, we need a clear and flexible picture of: (i) the meaning of cost annotations, (ii) the method to prove them sound and precise, and (iii) the way such proofs can be composed. We propose two approaches to these three questions which we name *direct* and *labelling*. As a first step, we examine their application to a toy compiler. For this simple framework, we provide a completely formal development which has been partly checked with the **Coq** proof assistant. This formal study suggests that the labelling approach, unlike the direct one, has good compositionality and scalability properties. In order to provide further evidence for this claim, we report our successful experience in implementing and testing the labelling approach on top of a prototype compiler written in **ocaml** for (a large fragment of) the **C** language.

# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Meaning of cost annotations . . . . .	5
1.2	Soundness and precision of cost annotations . . . . .	6
1.3	Compositionality . . . . .	6
1.4	Direct approach to cost annotations . . . . .	6
1.5	Labelling approach to cost annotations . . . . .	7
1.6	A toy compiler . . . . .	8
1.7	A C compiler . . . . .	8
1.8	Organisation . . . . .	9
<b>2</b>	<b>A toy compiler</b>	<b>9</b>
2.1	Imp: language and semantics . . . . .	9
2.2	Big-step semantics . . . . .	9
2.3	Small-step semantics . . . . .	10
2.4	Vm: language and semantics . . . . .	11
2.5	Compilation from Imp to Vm . . . . .	12
2.6	Soundness of compilation for the big-step semantics . . . . .	12
2.7	Soundness of compilation for the small-step semantics . . . . .	12
2.8	Compiled code is well-formed . . . . .	13
2.9	Mips: language and semantics . . . . .	13
2.10	Compilation from Vm to Mips . . . . .	14
<b>3</b>	<b>Direct approach for the toy compiler</b>	<b>14</b>
3.1	Mips and Vm cost annotations . . . . .	15
3.2	Imp cost annotation . . . . .	16
3.3	Composition . . . . .	17
3.4	Coq development . . . . .	18
3.5	Limitations of the direct approach . . . . .	18
<b>4</b>	<b>Labelling approach for the toy compiler</b>	<b>18</b>
4.1	Labelled Imp . . . . .	18
4.2	Labelled Vm . . . . .	19
4.3	Labelled Mips . . . . .	20
4.4	Labellings and instrumentations . . . . .	20
4.5	Sound and precise labellings . . . . .	21
<b>5</b>	<b>A C compiler</b>	<b>23</b>
5.1	Clight . . . . .	23
5.2	Cminor . . . . .	23
5.3	RTLAbs . . . . .	23
5.4	RTL . . . . .	27
5.5	ERTL . . . . .	27
5.6	LTL . . . . .	29
5.7	LIN . . . . .	30
5.8	Mips . . . . .	30
<b>6</b>	<b>Labelling approach for the C compiler</b>	<b>31</b>
6.1	Labelled Clight and labelled Cminor . . . . .	31
6.2	Labels in RTLAbs and the back-end languages . . . . .	32
6.3	Labelling of the source language . . . . .	32
6.3.1	Sequential instructions . . . . .	32
6.3.2	Ternary expressions . . . . .	32
6.3.3	Conditionals . . . . .	33
6.3.4	Loops . . . . .	33
6.3.5	Program Labels and Gotos . . . . .	34
6.3.6	Function calls . . . . .	34

6.4	Verifications on the object code . . . . .	35
6.5	Building the cost annotation . . . . .	35
6.6	Testing . . . . .	36
<b>7</b>	<b>Conclusion and future work</b>	<b>36</b>
<b>A</b>	<b>Assessment of the deliverable within the <i>CerCo</i> project</b>	<b>38</b>
<b>B</b>	<b>Proofs</b>	<b>40</b>
B.1	Notation . . . . .	40
B.2	Proof of proposition 4 . . . . .	40
B.3	Proof of proposition 8 . . . . .	41
B.4	Proof of proposition 10 . . . . .	41
B.5	Proof of proposition 11 . . . . .	42
B.6	Proof of proposition 13 . . . . .	42
B.7	Proof of proposition 14 . . . . .	42
B.8	Proof of proposition 16 . . . . .	42
B.9	Proof of proposition 17 . . . . .	42
B.10	Proof of proposition 20 . . . . .	43
B.11	Proof of proposition 22 . . . . .	43
B.12	Proof of proposition 23 . . . . .	43
B.13	Proof of proposition 25 . . . . .	43

## 1 Introduction

The formal description and certification of software components is reaching a certain level of maturity with impressing case studies ranging from compilers to kernels of operating systems. A well-documented example is the proof of functional correctness of a moderately optimising compiler from a large subset of the C language to a typical assembly language of the kind used in embedded systems [8].

In the framework of the *Certified Complexity (CerCo)* project [2], we aim to refine this line of work by focusing on the issue of the *execution cost* of the compiled code. Specifically, we aim to build a formally verified C compiler that given a source program produces automatically a functionally equivalent object code plus an annotation of the source code which is a sound and precise description of the execution cost of the object code.

We target in particular the kind of C programs produced for embedded applications; these programs are eventually compiled to binaries executable on specific processors. The current state of the art in commercial products such as Scade [3, 6] is that the *reaction time* of the program is estimated by means of abstract interpretation methods (such as those developed by AbsInt [1, 5]) that operate on the binaries. These methods rely on a specific knowledge of the architecture of the processor and may require explicit annotations of the binaries to determine the number of times a loop is iterated (see, *e.g.*, [13] for a survey of the state of the art).

In this context, our aim is to produce a functionally correct compiler which can *lift* in a provably correct way the pieces of information on the execution cost of the binary code to cost annotations on the source C code. Eventually, we plan to manipulate the cost annotations with automatic tools such as Frama – C [4].

In order to carry on our project, we need a clear and flexible picture of: (i) the meaning of cost annotations, (ii) the method to prove them sound and precise, and (iii) the way such proofs can be composed. Our purpose here is to propose two methodologies addressing these three questions and to consider their concrete application to a simple toy compiler and to a moderately optimising C compiler.

### 1.1 Meaning of cost annotations

The execution cost of the source programs we are interested in depends on their control structure. Typically, the source programs are composed of mutually recursive procedures and loops and their execution cost depends, up to some multiplicative constant, on the number of times procedure calls and loop iterations are performed.

Producing a *cost annotation* of a source program amounts to:

- enrich the program with a collection of *global cost variables* to measure resource consumption (time, stack size, heap size, ...)
- inject suitable code at some critical points (procedures, loops, ...) to keep track of the execution cost.

Thus producing a cost-annotation of a source program  $S$  amounts to build an *annotated program*  $An(S)$  which behaves as  $S$  while self-monitoring its execution cost. In particular, if we do *not* observe the cost variables then we expect the annotated program  $An(S)$  to be functionally equivalent to  $S$ . Notice that in the proposed approach an annotated program is

a program in the source language. Therefore the meaning of the cost annotations is automatically defined by the semantics of the source language and tools developed to reason on the source programs can be directly applied to the annotated programs too.

## 1.2 Soundness and precision of cost annotations

Suppose we have a functionally correct compiler  $\mathcal{C}$  that associates with a program  $S$  in the source language a program  $\mathcal{C}(S)$  in the object language. Further suppose we have some obvious way of defining the execution cost of an object code. For instance, we have a good estimate of the number of cycles needed for the execution of each instruction of the object code. Now the annotation of the source program  $An(S)$  is *sound* if its prediction of the execution cost is an upper bound for the ‘real’ execution cost. Moreover, we say that the annotation is *precise* if the *difference* between the predicted and real execution costs is bounded by a constant which depends on the program.

## 1.3 Compositionality

In order to master the complexity of the compilation process (and its verification), the compilation function  $\mathcal{C}$  must be regarded as the result of the composition of a certain number of program transformations  $\mathcal{C} = \mathcal{C}_k \circ \dots \circ \mathcal{C}_1$ . When building a system of cost annotations on top of an existing compiler a certain number of problems arise. First, the estimated cost of executing a piece of source code is determined only at the *end* of the compilation process. Thus while we are used to define the compilation functions  $\mathcal{C}_i$  in increasing order (from left to right), the annotation function  $An$  is the result of a progressive abstraction from the object to the source code (from right to left). Second, we must be able to foresee in the source language the looping and branching points of the object code. Missing a loop may lead to unsound cost annotations while missing a branching point may lead to rough cost predictions. This means that we must have a rather good idea of the way the source code will eventually be compiled to object code. Third, the definition of the annotation of the source code depends heavily on contextual information. For instance, the cost of the compiled code associated with a simple expression such as  $x + 1$  will depend on the place in the memory hierarchy where the variable  $x$  is allocated.

## 1.4 Direct approach to cost annotations

A first ‘direct’ approach to the problem of building cost annotations is summarised by the following diagram.

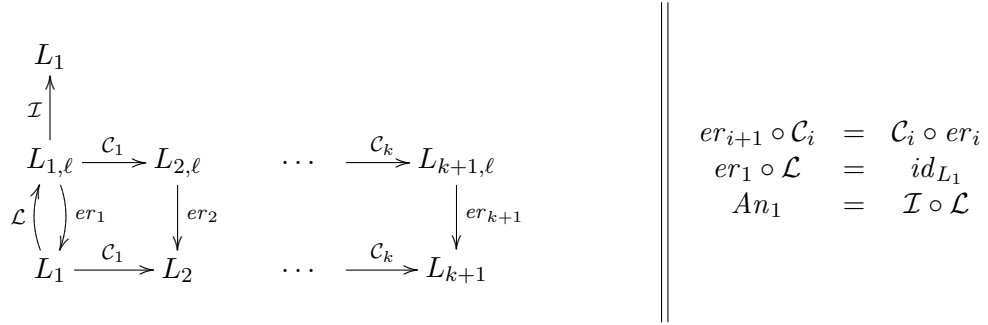
$$\begin{array}{ccccc}
 L_1 & \xrightarrow{\mathcal{C}_1} & L_2 & \dots & \xrightarrow{\mathcal{C}_k} & L_{k+1} \\
 \downarrow An_1 & & \downarrow An_1 & & & \downarrow An_{k+1} \\
 L_1 & & L_2 & & & L_{k+1}
 \end{array}$$

With respect to our previous discussion,  $L_1$  is the source language with the related annotation function  $An_1$  while  $L_{k+1}$  is the object language with a related annotation  $An_{k+1}$ . This annotation of the object code is supposed to be truly straightforward and it is taken as an ‘axiomatic’ definition of the ‘real’ execution cost of the program. The languages  $L_i$ , for  $2 \leq i \leq k$ , are intermediate languages which are also equipped with increasingly ‘realistic’

annotation functions. Suppose we denote with  $S$  the source program, with  $\mathcal{C}(S)$  the compiled program, and that we write  $(P, s) \Downarrow s'$  to mean that the (source or object) program  $P$  in the state  $s$  terminates successfully in the state  $s'$ . The soundness proof of the compilation function guarantees that if  $(S, s) \Downarrow s'$  then  $(\mathcal{C}(S), s) \Downarrow s'$ . In the direct approach, the *proof of soundness* of the cost annotations amounts to lift the proof of functional equivalence of the source program and the object code to a proof of ‘quasi-equivalence’ of the respective instrumented codes. Suppose we write  $s[c/cost]$  for a state that associates  $c$  with the cost variable  $cost$ . Then what we want to show is that whenever  $(An_1(S), s[c/cost]) \Downarrow s'[c'/cost]$  we have that  $(An_{k+1}(\mathcal{C}(S)), s[d/cost]) \Downarrow s'[d'/cost]$  and  $|d' - d| \leq |c' - c| + k$ . This means that the increment in the annotated source program bounds up to an additive constant the increment in the annotated object program. We will also say that the cost annotation is precise if we can also prove that  $|c' - c| \leq |d' - d|$ , *i.e.*, the ‘estimated’ cost is not too far away from the ‘real’ one. We will see that while in theory one can build sound and precise annotation functions, in practice definitions and proofs become unwieldy.

### 1.5 Labelling approach to cost annotations

The ‘labelling’ approach to the problem of building cost annotations is summarized in the following diagram.



For each language  $L_i$  considered in the compilation process, we define an extended *labelled* language  $L_{i,\ell}$  and an extended operational semantics. The labels are used to mark certain points of the control. The semantics makes sure that whenever we cross a labelled control point a labelled and observable transition is produced.

For each labelled language there is an obvious function  $er_i$  erasing all labels and producing a program in the corresponding unlabelled language. The compilation functions  $C_i$  are extended from the unlabelled to the labelled language so that they enjoy commutation with the erasure functions. Moreover, we lift the soundness properties of the compilation functions from the unlabelled to the labelled languages and transition systems.

A *labelling*  $\mathcal{L}$  of the source language  $L_1$  is just a function such that  $er_{L_1} \circ \mathcal{L}$  is the identity function. An *instrumentation*  $\mathcal{I}$  of the source labelled language  $L_{1,\ell}$  is a function replacing the labels with suitable increments of, say, a fresh *cost* variable. Then an *annotation*  $An_1$  of the source program can be derived simply as the composition of the labelling and the instrumentation functions:  $An_1 = \mathcal{I} \circ \mathcal{L}$ .

As for the direct approach, suppose  $s$  is some adequate representation of the state of a program. Let  $S$  be a source program and suppose that its annotation satisfies the following property:

$$(An_1(S), s[c/cost]) \Downarrow s'[c + \delta/cost] \quad (1)$$

where  $\delta$  is some non-negative number. Then the definition of the instrumentation and the fact that the soundness proofs of the compilation functions have been lifted to the labelled languages allows to conclude that

$$(\mathcal{C}(L(S)), s[c/cost]) \Downarrow (s'[c/cost], \lambda) \quad (2)$$

where  $\mathcal{C} = \mathcal{C}_k \circ \dots \circ \mathcal{C}_1$  and  $\lambda$  is a sequence (or a multi-set) of labels whose ‘cost’ corresponds to the number  $\delta$  produced by the annotated program. Then the commutation properties of erasure and compilation functions allows to conclude that the *erasure* of the compiled labelled code  $er_{k+1}(\mathcal{C}(L(S)))$  is actually equal to the compiled code  $\mathcal{C}(S)$  we are interested in. Given this, the following question arises:

Under which conditions the sequence  $\lambda$ , *i.e.*, the increment  $\delta$ , is a sound and possibly precise description of the execution cost of the object code?

To answer this question, we observe that the object code we are interested in is some kind of assembly code and its control flow can be easily represented as a control flow graph. The fact that we have to prove the soundness of the compilation functions means that we have plenty of pieces of information on the way the control flows in the compiled code, in particular as far as procedure calls and returns are concerned. These pieces of information allow to build a rather accurate representation of the control flow of the compiled code at run time.

The idea is then to perform two simple checks on the control flow graph. The first check is to verify that all loops go through a labelled node. If this is the case then we can associate a finite cost with every label and prove that the cost annotations are sound. The second check amounts to verify that all paths starting from a label have the same cost. If this check is successful then we can conclude that the cost annotations are precise.

## 1.6 A toy compiler

As a first case study for the two approaches to cost annotations we have sketched, we introduce a *toy compiler* which is summarised by the following diagram.

$$\text{Imp} \xrightarrow{\mathcal{C}} \text{Vm} \xrightarrow{\mathcal{C}'} \text{Mips}$$

The three languages considered can be shortly described as follows: **Imp** is a very simple imperative language with pure expressions, branching and looping commands, **Vm** is an assembly-like language enriched with a stack, and **Mips** is a Mips-like assembly language with registers and main memory. The first compilation function  $\mathcal{C}$  relies on the stack of the **Vm** language to implement expression evaluation while the second compilation function  $\mathcal{C}'$  allocates (statically) the base of the stack in the registers and the rest in main memory. This is of course a naive strategy but it suffices to expose some of the problems that arise in defining a compositional approach (cf. section 1.3).

## 1.7 A C compiler

As a second, more complex, case study we consider a **C** compiler we have built in **ocaml** whose structure is summarised by the following diagram:

$$\begin{array}{ccccccc} \text{C} & \rightarrow & \text{Clight} & \rightarrow & \text{Cminor} & \rightarrow & \text{RTLabs} & \text{(front end)} \\ & & & & & & \downarrow & \\ \text{Mips} & \leftarrow & \text{LIN} & \leftarrow & \text{LTL} & \leftarrow & \text{ERTL} & \leftarrow & \text{RTL} & \text{(back-end)} \end{array}$$



The structure follows rather closely the one of the **CompCert** compiler. Notable differences are that some compilation steps are fusioned, that the front-end goes till **RTLabs** (rather than **Cminor**) and that we target the **Mips** assembly language (rather than **PowerPc**). These differences are contingent to the way we built the compiler. The compilation from **C** to **Clight** relies on the **CIL** front-end [11]. The one from **Clight** to **RTL** has been programmed from scratch and it is partly based on the **Coq** definitions available in the **CompCert** compiler. Finally, the back-end from **RTL** to **Mips** is based on a compiler developed in **ocaml** for pedagogical purposes [12]. The main optimisations it performs are common subexpression elimination, liveness analysis and register allocation, and graph compression.

## 1.8 Organisation

The rest of the paper is organised as follows. Section 2 describes the 3 languages and the 2 compilation steps of the toy compiler. Section 3 describes the application of the direct approach to the toy compiler and points out its limitations. Section 4 describes the application of the labelling approach to the toy compiler. Section 5 provides some details on the structure of the **C** compiler we have implemented. Section 6 reports our experience in implementing and testing the labelling approach on the **C** compiler. Section 7 summarizes our contribution and outlines some perspectives for future work. Section B sketches the proofs that have not been mechanically checked in **Coq**.

## 2 A toy compiler

We formalise the toy compiler introduced in section 1.6.

### 2.1 Imp: language and semantics

The syntax of the **Imp** language is described below. This is a rather standard imperative language with while loops and if-then-else.

$id$	$::= x \mid y \mid \dots$	(identifiers)
$n$	$::= 0 \mid -1 \mid +1 \mid \dots$	(integers)
$v$	$::= n \mid \text{true} \mid \text{false}$	(values)
$e$	$::= id \mid n \mid e + e$	(numerical expressions)
$b$	$::= e < e$	(boolean conditions)
$S$	$::= \text{skip} \mid id := e \mid S; S \mid \text{if } b \text{ then } S \text{ else } S \mid \text{while } b \text{ do } S$	(commands)
$P$	$::= \text{prog } S$	(programs)

Let  $s$  be a total function from identifiers to integers representing the **state**. If  $s$  is a state,  $x$  an identifier, and  $n$  an integer then  $s[n/x]$  is the ‘updated’ state such that  $s[n/x](x) = n$  and  $s[n/x](y) = s(y)$  if  $x \neq y$ .

### 2.2 Big-step semantics

The big-step operational semantics of **Imp** programs is defined by the following judgements:

$$(e, s) \Downarrow v \quad (b, s) \Downarrow v \quad (S, s) \Downarrow s' \quad (P, s) \Downarrow s'$$

$$\begin{array}{c}
\frac{}{(v, s) \Downarrow v} \quad \frac{}{(x, s) \Downarrow s(x)} \quad \frac{(e, s) \Downarrow v \quad (e', s) \Downarrow v'}{(e + e', s) \Downarrow (v +_{\mathbf{Z}} v')} \quad \frac{(e, s) \Downarrow v \quad (e', s) \Downarrow v'}{(e <_{\mathbf{Z}} e', s) \Downarrow (v <_{\mathbf{Z}} v')} \\
\frac{}{(\text{skip}, s) \Downarrow s} \quad \frac{(e, s) \Downarrow v}{(x := e, s) \Downarrow s[v/x]} \quad \frac{(S_1, s) \Downarrow s' \quad (S_2, s') \Downarrow s''}{(S_1; S_2, s) \Downarrow s''} \\
\frac{(b, s) \Downarrow \text{true} \quad (S, s) \Downarrow s'}{(\text{if } b \text{ then } S \text{ else } S', s) \Downarrow s'} \quad \frac{(b, s) \Downarrow \text{false} \quad (S', s) \Downarrow s'}{(\text{if } b \text{ then } S \text{ else } S', s) \Downarrow s'} \\
\frac{(b, s) \Downarrow \text{false}}{(\text{while } b \text{ do } S, s) \Downarrow s} \quad \frac{(b, s) \Downarrow \text{true} \quad (S; \text{while } b \text{ do } S, s) \Downarrow s'}{(\text{while } b \text{ do } S, s) \Downarrow s'} \\
\frac{(S, s) \Downarrow s'}{(\text{prog } S, s) \Downarrow s'}
\end{array}$$

Table 1: Big-step operational semantics of `Imp`

$$\begin{array}{lcl}
(x := e, K, s) & \rightarrow & (\text{skip}, K, s[v/x]) \quad \text{if } (e, s) \Downarrow v \\
(S; S', K, s) & \rightarrow & (S, S' \cdot K, s) \\
(\text{if } b \text{ then } S \text{ else } S', K, s) & \rightarrow & \begin{cases} (S, K, s) & \text{if } (b, s) \Downarrow \text{true} \\ (S', K, s) & \text{if } (b, s) \Downarrow \text{false} \end{cases} \\
(\text{while } b \text{ do } S, K, s) & \rightarrow & \begin{cases} (S, (\text{while } b \text{ do } S) \cdot K, s) & \text{if } (b, s) \Downarrow \text{true} \\ (\text{skip}, K, s) & \text{if } (b, s) \Downarrow \text{false} \end{cases} \\
(\text{skip}, S \cdot K, s) & \rightarrow & (S, K, s)
\end{array}$$

Table 2: Small-step operational semantics of `Imp` commands

and it is described in table 1. We assume that the addition  $v +_{\mathbf{Z}} v'$  is only defined if  $v$  and  $v'$  are integers and the comparison  $v <_{\mathbf{Z}} v'$  produces a value `true` or `false` only if  $v$  and  $v'$  are integers.

### 2.3 Small-step semantics

We also introduce an alternative small-step semantics of the `Imp` language. A *continuation*  $K$  is a list of commands which terminates with a special symbol `halt`.

$$K ::= \text{halt} \mid S \cdot K$$

Table 2 defines a small-step semantics of `Imp` commands whose basic judgement has the shape:

$$(S, K, s) \rightarrow (S', K', s').$$

We define the semantics of a program `prog S` as the semantics of the command  $S$  with continuation `halt`. We derive a big step semantics from the small step one as follows:

$$(S, s) \Downarrow s' \quad \text{if } (S, \text{halt}, s) \rightarrow \cdots \rightarrow (\text{skip}, \text{halt}, s').$$

Rule	$C[i] =$
$C \vdash (i, \sigma, s) \rightarrow (i + 1, n \cdot \sigma, s)$	<code>cnst(<math>n</math>)</code>
$C \vdash (i, \sigma, s) \rightarrow (i + 1, s(x) \cdot \sigma, s)$	<code>var(<math>x</math>)</code>
$C \vdash (i, n \cdot \sigma, s) \rightarrow (i + 1, \sigma, s[n/x])$	<code>setvar(<math>x</math>)</code>
$C \vdash (i, n \cdot n' \cdot \sigma, s) \rightarrow (i + 1, (n +_Z n') \cdot \sigma, s)$	<code>add</code>
$C \vdash (i, \sigma, s) \rightarrow (i + k + 1, \sigma, s)$	<code>branch(<math>k</math>)</code>
$C \vdash (i, n \cdot n' \cdot \sigma, s) \rightarrow (i + 1, \sigma, s)$	<code>bge(<math>k</math>)</code> and $n <_Z n'$
$C \vdash (i, n \cdot n' \cdot \sigma, s) \rightarrow (i + k + 1, \sigma, s)$	<code>bge(<math>k</math>)</code> and $n \geq_Z n'$

Table 3: Operational semantics Vm programs

## 2.4 Vm: language and semantics

Following [9], we define a virtual machine Vm and its programming language. The machine includes the following elements: (1) A fixed code  $C$  (a possibly empty sequence of instructions), (2) A program counter  $pc$ , (3) A store  $s$  (as for the source program), (4) A stack of integers  $\sigma$ .

We will rely on the following instructions with the associated informal semantics:

<code>cnst(<math>n</math>)</code>	push on the stack
<code>var(<math>x</math>)</code>	push value $x$
<code>setvar(<math>x</math>)</code>	pop value and assign it to $x$
<code>add</code>	pop 2 values and push their sum
<code>branch(<math>k</math>)</code>	jump with offset $k$
<code>bge(<math>k</math>)</code>	pop 2 values and jump if greater or equal with offset $k$
<code>halt</code>	stop computation

In the branching instructions,  $k$  is an integer that has to be added to the current program counter in order to determine the following instruction to be executed. Given a sequence  $C$ , we denote with  $|C|$  its length and with  $C[i]$  its  $i^{\text{th}}$  element (the leftmost element being the  $0^{\text{th}}$  element). The operational semantics of the instructions is formalised by rules of the shape:

$$C \vdash (i, \sigma, s) \rightarrow (j, \sigma', s')$$

and it is fully described in table 3. Notice that `Imp` and `Vm` semantics share the same notion of store. We write, e.g.,  $n \cdot \sigma$  to stress that the top element of the stack exists and is  $n$ . We will also write  $(C, s) \Downarrow s'$  if  $C \vdash (0, \epsilon, s) \xrightarrow{*} (i, \epsilon, s')$  and  $C[i] = \text{halt}$ .

Code coming from the compilation of `Imp` programs has specific properties that are used in the following compilation step when values on the stack are allocated either in registers or in main memory. In particular, it turns out that for every instruction of the compiled code it is possible to predict statically the *height of the stack* whenever the instruction is executed. We now proceed to define a simple notion of *well-formed* code and show that it enjoys this property. In the following section, we will define the compilation function from `Imp` to `Vm` and show that it produces well-formed code.

**Definition 1** *We say that a sequence of instructions  $C$  is well formed if there is a function  $h : \{0, \dots, |C|\} \rightarrow \mathbf{N}$  which satisfies the conditions listed in table 4 for  $0 \leq i \leq |C| - 1$ . In this case we write  $C : h$ .*

The conditions defining the predicate  $C : h$  are strong enough to entail that  $h$  correctly predicts the stack height and to guarantee the uniqueness of  $h$  up to the initial condition.

**Proposition 2** (1) *If  $C : h$ ,  $C \vdash (i, \sigma, s) \xrightarrow{*} (j, \sigma', s')$ , and  $h(i) = |\sigma|$  then  $h(j) = |\sigma'|$ .*  
(2) *If  $C : h$ ,  $C : h'$  and  $h(0) = h'(0)$  then  $h = h'$ .*

$C[i] =$	Conditions for $C : h$
$\text{cnst}(n)$ or $\text{var}(x)$	$h(i+1) = h(i) + 1$
$\text{add}$	$h(i) \geq 2, \quad h(i+1) = h(i) - 1$
$\text{setvar}(x)$	$h(i) = 1, \quad h(i+1) = 0$
$\text{branch}(k)$	$0 \leq i+k+1 \leq  C , \quad h(i) = h(i+1) = h(i+k+1) = 0$
$\text{bge}(k)$	$0 \leq i+k+1 \leq  C , \quad h(i) = 2, \quad h(i+1) = h(i+k+1) = 0$
$\text{halt}$	$i =  C  - 1, \quad h(i) = h(i+1) = 0$

Table 4: Conditions for well-formed code

$$\begin{aligned}
C(x) = \text{var}(x) \quad C(n) = \text{cnst}(n) \quad C(e + e') = C(e) \cdot C(e') \cdot \text{add} \\
C(e < e', k) = C(e) \cdot C(e') \cdot \text{bge}(k) \\
C(x := e) = C(e) \cdot \text{setvar}(x) \quad C(S; S') = C(S) \cdot C(S') \\
C(\text{if } b \text{ then } S \text{ else } S') = C(b, k) \cdot C(S) \cdot (\text{branch}(k')) \cdot C(S') \\
\text{where: } k = \text{sz}(S) + 1, \quad k' = \text{sz}(S') \\
C(\text{while } b \text{ do } S) = C(b, k) \cdot C(S) \cdot \text{branch}(k') \\
\text{where: } k = \text{sz}(S) + 1, \quad k' = -(\text{sz}(b) + \text{sz}(S) + 1) \\
C(\text{prog } S) = C(S) \cdot \text{halt}
\end{aligned}$$

Table 5: Compilation from Imp to Vm

## 2.5 Compilation from Imp to Vm

In table 5, we define compilation functions  $\mathcal{C}$  from Imp to Vm which operate on expressions, boolean conditions, statements, and programs. We write  $\text{sz}(e)$ ,  $\text{sz}(b)$ ,  $\text{sz}(S)$  for the number of instructions the compilation function associates with the expression  $e$ , the boolean condition  $b$ , and the statement  $S$ , respectively.

## 2.6 Soundness of compilation for the big-step semantics

We follow [9] for the proof of correctness of the compilation function with respect to the big-step semantics (see also [10] for a much older reference).

**Proposition 3** *The following properties hold:*

- (1) *If  $(e, s) \Downarrow v$  then  $C \cdot \mathcal{C}(e) \cdot C' \vdash (i, \sigma, s) \xrightarrow{*} (j, v \cdot \sigma, s)$  where  $i = |C|$  and  $j = |C \cdot \mathcal{C}(e)|$ .*
- (2) *If  $(b, s) \Downarrow \text{true}$  then  $C \cdot \mathcal{C}(b, k) \cdot C' \vdash (i, \sigma, s) \xrightarrow{*} (j+k, \sigma, s)$  where  $i = |C|$  and  $j = |C \cdot \mathcal{C}(b, k)|$ .*
- (3) *If  $(b, s) \Downarrow \text{false}$  then  $C \cdot \mathcal{C}(b, k) \cdot C' \vdash (i, \sigma, s) \xrightarrow{*} (j, \sigma, s)$  where  $i = |C|$  and  $j = |C \cdot \mathcal{C}(b, k)|$ .*
- (4) *If  $(S, s) \Downarrow s'$  then  $C \cdot \mathcal{C}(S) \cdot C' \vdash (i, \sigma, s) \xrightarrow{*} (j, \sigma, s')$  where  $i = |C|$  and  $j = |C \cdot \mathcal{C}(e)|$ .*

## 2.7 Soundness of compilation for the small-step semantics

We prove soundness with respect to the small step semantics too. To this end, given a Vm code  $C$ , we define an ‘accessibility relation’  $\overset{\mathcal{C}}{\rightsquigarrow}$  as the least binary relation on  $\{0, \dots, |C| - 1\}$  such that:

$$\frac{}{i \overset{C}{\rightsquigarrow} i} \quad \frac{C[i] = \text{branch}(k) \quad (i+k+1) \overset{C}{\rightsquigarrow} j}{i \overset{C}{\rightsquigarrow} j}$$

We also introduce a ternary relation  $R(C, i, K)$  which relates a  $\mathbf{Vm}$  code  $C$ , a number  $i \in \{0, \dots, |C| - 1\}$  and a continuation  $K$ . The intuition is that relative to the code  $C$ , the instruction  $i$  can be regarded as having continuation  $K$ . Formally, the relation  $R$  is defined as the least one that satisfies the following conditions.

$$\frac{i \overset{C}{\rightsquigarrow} j \quad C[j] = \text{halt}}{R(C, i, \text{halt})} \quad \frac{i \overset{C}{\rightsquigarrow} i' \quad C = C_1 \cdot C(S) \cdot C_2 \quad i' = |C_1| \quad j = |C_1 \cdot C(S)| \quad R(C, j, K)}{R(C, i, S \cdot K)} .$$

We can then state the correctness of the compilation function as follows.

**Proposition 4** *If  $(S, K, s) \rightarrow (S', K', s')$  and  $R(C, i, S \cdot K)$  then  $C \vdash (i, \sigma, s) \xrightarrow{*} (j, \sigma, s')$  and  $R(C, j, S' \cdot K')$ .*

## 2.8 Compiled code is well-formed

As announced, we can prove that the result of the compilation is a well-formed code.

**Proposition 5** *For any expression  $e$ , statement  $S$ , and program  $P$  the following holds:*

- (1) *For any  $n \in \mathbf{N}$  there is a unique  $h$  such that  $\mathcal{C}(e) : h$ ,  $h(0) = n$ , and  $h(|\mathcal{C}(e)|) = h(0) + 1$ .*
- (2) *For any  $S$ , there is a unique  $h$  such that  $\mathcal{C}(S) : h$ ,  $h(0) = 0$ , and  $h(|\mathcal{C}(e)|) = 0$ .*
- (3) *There is a unique  $h$  such that  $\mathcal{C}(P) : h$ .*

## 2.9 Mips: language and semantics

We consider a Mips-like machine [7] which includes the following elements: (1) a fixed code  $M$  (a sequence of instructions), (2) a program counter  $pc$ , (3) a finite set of registers including the registers  $A$ ,  $B$ , and  $R_0, \dots, R_{b-1}$ , and (4) an (infinite) main memory which maps locations to integers.

We denote with  $R, R', \dots$  registers, with  $l, l', \dots$  locations and with  $m, m', \dots$  memories which are total functions from registers and locations to (unbounded) integers. We will rely on the following instructions with the associated informal semantics:

loadi $R, n$	store value $n$ in the register $R$
load $R, l$	store contents of location $l$ in the register $R$
store $R, l$	store contents of register $R$ in the location $l$
add $R, R', R''$	add contents of $R', R''$ and store it in $R$
branch $k$	jump with offset $k$
bge $R, R', k$	jump with offset $k$ if contents $R$ greater or equal than contents $R'$
halt	stop computation

We denote with  $M$  a list of instructions. The operational semantics is formalised in table 6 by rules of the shape:

$$M \vdash (i, m) \rightarrow (j, m')$$

where  $M$  is a list of Mips instructions,  $i, j$  are natural numbers and  $m, m'$  are memories. We write  $(M, m) \Downarrow m'$  if  $M \vdash (0, m) \xrightarrow{*} (j, m')$  and  $M[j] = \text{halt}$ .

Rule	$M[i] =$
$M \vdash (i, m) \rightarrow (i + 1, m[n/R])$	loadi $R, n$
$M \vdash (i, m) \rightarrow (i + 1, m[m(l)/R])$	load $R, l$
$M \vdash (i, m) \rightarrow (i + 1, m[m(R)/l])$	store $R, l$
$M \vdash (i, m) \rightarrow (i + 1, m[m(R') + m(R'')/R])$	add $R, R', R''$
$M \vdash (i, m) \rightarrow (i + k + 1, m)$	branch $k$
$M \vdash (i, m) \rightarrow (i + 1, m)$	bge $R, R', k$ and $m(R) <_{\mathbf{Z}} m(R')$
$M \vdash (i, m) \rightarrow (i + k + 1, m)$	bge $R, R', k$ and $m(R) \geq_{\mathbf{Z}} m(R')$

Table 6: Operational semantics Mips programs

## 2.10 Compilation from Vm to Mips

In order to compile Vm programs to Mips programs we make the following hypotheses. (1) For every Vm program variable  $x$  we reserve an address  $l_x$ , (2) For every natural number  $h \geq b$ , we reserve an address  $l_h$  (the addresses  $l_x, l_h, \dots$  are all distinct), (3) We store the first  $b$  elements of the stack  $\sigma$  in the registers  $R_0, \dots, R_{b-1}$  and the remaining (if any) at the addresses  $l_b, l_{b+1}, \dots$

We say that the memory  $m$  represents the stack  $\sigma$  and the store  $s$ , and write  $m \parallel -\sigma, s$ , if the following conditions are satisfied: (1)  $s(x) = m(l_x)$ , and (2) if  $0 \leq i < |\sigma|$  then

$$\sigma[i] = \begin{cases} m(R_i) & \text{if } i < b \\ m(l_i) & \text{if } i \geq b \end{cases}$$

The compilation function  $\mathcal{C}'$  from Vm to Mips is described in table 7. It operates on a well-formed Vm code  $C$  whose last instruction is halt. Hence, by proposition 5(3), there is a unique  $h$  such that  $C : h$ . We denote with  $\mathcal{C}'(C)$  the concatenation  $\mathcal{C}'(0, C) \cdots \mathcal{C}'(|C| - 1, C)$ . Given a well formed Vm code  $C$  with  $i < |C|$  we denote with  $p(i, C)$  the position of the first instruction in  $\mathcal{C}'(C)$  which corresponds to the compilation of the instruction with position  $i$  in  $C$ . This is defined as:<sup>1</sup>

$$p(i, C) = \sum_{0 \leq j < i} d(j, C), \quad (3)$$

where the function  $d(i, C)$  is defined as follows:

$$d(i, C) = |\mathcal{C}'(i, C)|. \quad (4)$$

Hence  $d(i, C)$  is the number of Mips instructions associated with the  $i^{\text{th}}$  instruction of the (well-formed)  $C$  code.

The functional correctness of the compilation function can then be stated as follows.

**Proposition 6** *Let  $C : h$  be a well formed code. If  $C \vdash (i, \sigma, s) \rightarrow (j, \sigma', s')$  with  $h(i) = |\sigma|$  and  $m \parallel -\sigma, s$  then  $\mathcal{C}'(C) \vdash (p(i, C), m) \xrightarrow{*} (p(j, C), m')$  and  $m' \parallel -\sigma', s'$ .*

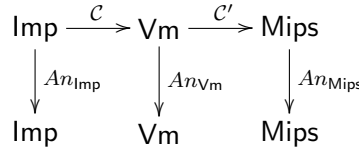
## 3 Direct approach for the toy compiler

We apply the direct approach discussed in section 1.4 to the toy compiler which results in the following diagram:

<sup>1</sup>There is an obvious circularity in this definition that can be easily eliminated by defining first the function  $d$  following the case analysis in table 7, then the function  $p$  as in (3), and finally the function  $\mathcal{C}'$  as in table 7.

$C[i] =$	$C'(i, C) =$
$\text{cnst}(n)$	$\begin{cases} (\text{loadi } R_h, n) & \text{if } h = h(i) < b \\ (\text{loadi } A, n) \cdot (\text{store } A, l_h) & \text{otherwise} \end{cases}$
$\text{var}(x)$	$\begin{cases} (\text{load } R_h, l_x) & \text{if } h = h(i) < b \\ (\text{load } A, l_x) \cdot (\text{store } A, l_h) & \text{otherwise} \end{cases}$
$\text{add}$	$\begin{cases} (\text{add } R_{h-2}, R_{h-2}, R_{h-1}) & \text{if } h = h(i) < (b-1) \\ (\text{load } A, l_{h-1}) \cdot (\text{add } R_{h-2}, R_{h-2}, A) & \text{if } h = h(i) = (b-1) \\ (\text{load } A, l_{h-1}) \cdot (\text{load } B, l_{h-2}) & \text{if } h = h(i) > (b-1) \\ (\text{add } A, B, A) \cdot (\text{store } A, l_{h-2}) & \end{cases}$
$\text{setvar}(x)$	$\begin{cases} (\text{store } R_{h-1} l_x) & \text{if } h = h(i) < b \\ (\text{load } A, l_{h-1}) \cdot (\text{store } A, l_x) & \text{if } h = h(i) \geq b \end{cases}$
$\text{branch}(k)$	$(\text{branch } k') \quad \text{if } k' = p(i+k+1, C) - p(i+1, C)$
$\text{bge}(k)$	$\begin{cases} (\text{bge } R_{h-2}, R_{h-1}, k') & \text{if } h = h(i) < (b-1) \\ (\text{load } A, l_{h-1}) \cdot (\text{bge } R_{h-2}, A, k') & \text{if } h = h(i) = (b-1) \\ (\text{load } A, l_{h-2}) \cdot (\text{load } B, l_{h-1}) \cdot (\text{bge } A, B, k') & \text{if } h = h(i) > (b-1) \text{ and} \\ & k' = p(i+k+1, C) - p(i+1, C) \end{cases}$
$\text{halt}$	$\text{halt}$

Table 7: Compilation from Vm to Mips



### 3.1 Mips and Vm cost annotations

The definition of the cost annotation  $An_{\text{Mips}}(M)$  for a Mips code  $M$  goes as follows assuming that all the Mips instructions have cost 1.

1. We select fresh locations  $l_{\text{cost}}, l_A, l_B$  not used in  $M$ .
2. Before each instruction in  $M$  we insert the following list of 8 instructions whose effect is to increase by 1 the contents of location  $l_{\text{cost}}$ :

$$\begin{aligned}
& (\text{store } A, l_A) \cdot (\text{store } B, l_B) \cdot (\text{loadi } A, 1) \cdot (\text{load } B, l_{\text{cost}}) \cdot \\
& (\text{add } A, A, B) \cdot (\text{store } A, l_{\text{cost}}) \cdot (\text{load } A, l_A) \cdot (\text{load } B, l_B) \cdot
\end{aligned}$$

3. We update the offset of the branching instructions according to the map  $k \mapsto (9 \cdot k)$ .

The definition of the cost annotation  $An_{\text{Vm}}(C)$  for a well-formed Vm code  $C$  such that  $C : h$  goes as follows.

1. We select a fresh *cost* variable not used in  $C$ .

$C[i] =$	cnst( $n$ ) or var( $x$ )	add	setvar( $x$ )	branch( $k$ )	bge( $k$ )	halt
$\kappa(C[i]) =$	2	4	2	1	3	1

Table 8: Upper bounds on the cost of Vm instructions

2. Before the instruction  $i^{\text{th}}$  in  $C$ , we insert the following list of 4 instructions whose effect is to increase the *cost* variable by the maximum number  $\kappa(C[i])$  of instructions associated with  $C[i]$  (as specified in table 8):

$$(\text{cnst}(\kappa(C[i]))) \cdot (\text{var}(\text{cost})) \cdot (\text{add}) \cdot (\text{setvar}(\text{cost})) .$$

3. We update the offset of the branching instructions according to the map  $k \mapsto (5 \cdot k)$ .

To summarise, the situation is that the instruction  $i$  of the Vm code  $C$  corresponds to: the instruction  $5 \cdot i + 4$  of the annotated Vm code  $An_{\text{Vm}}(C)$ , the instruction  $p(i, C)$  of the Mips code  $C'(C)$ , and the instruction  $9 \cdot p(i, C) + 8$  of the instrumented Mips code  $An_{\text{Mips}}(C'(C))$ .

The following lemma describes the effect of the injected sequences of instructions.

**Lemma 7** *The following properties hold:*

- (1) *If  $M$  is a Mips instruction then  $C_1 \cdot An_{\text{Mips}}(M) \cdot C_2 \vdash (|C_1|, m[c/l_{\text{cost}}]) \xrightarrow{*} (|C_1| + 8, m[c + 1/l_{\text{cost}}, m(A)/l_A, m(B)/l_B])$ .*
- (2) *If  $C$  is a Vm instruction then  $C_1 \cdot An_{\text{Vm}}(C) \cdot C_2 \vdash (|C_1|, \sigma, s[c/cost]) \xrightarrow{*} (|C_1| + 4, \sigma, s[c + \kappa(C[i])/cost])$ , where  $i = |C_1|$ .*

The simulation proposition 6 is extended to the annotated codes as follows where we write  $\rightarrow^k$  for the relation obtained by composing  $k$  times the reduction relation  $\rightarrow$ .

**Proposition 8** *Let  $C : h$  be a well-formed code. If  $An_{\text{Vm}}(C) \vdash (5 \cdot i, \sigma, s) \rightarrow^5 (5 \cdot j, \sigma', s')$ ,  $m \Vdash \sigma, s$ , and  $h(i) = |\sigma|$  then  $An_{\text{Mips}}(C'(C)) \vdash (9 \cdot p(i, C), m) \xrightarrow{*} (9 \cdot p(j, C), m')$ , with  $m'[s'(cost)/l_{\text{cost}}] \Vdash \sigma', s'$  and*

$$m'(l_{\text{cost}}) - m(l_{\text{cost}}) \leq s'(cost) - s(cost) .$$

### 3.2 Imp cost annotation

The definition of the cost annotation  $An_{\text{Imp}}(P)$  for an Imp program  $P$  is defined in table 9 and it relies on an auxiliary function  $\kappa$  which provides an upper bound on the cost of executing the various operators of the language. For the sake of simplicity, this annotation introduces two main *approximations*:

- It over-approximates the cost of an `if_then_else` by taking the maximum cost of the cost of the two branches.
- It always considers the worst case where the data in the stack resides in the main memory.

We will discuss in section 3.5 to what extent these approximations can be removed.

Next we formulate the soundness of the Imp annotation with respect to the Vm annotation along the pattern of proposition 3.



$An_{\text{Imp}}(\text{prog } S)$	$= \text{cost} := \text{cost} + \kappa(S); An_{\text{Imp}}(S)$
$An_{\text{Imp}}(\text{skip})$	$= \text{skip}$
$An_{\text{Imp}}(x := e)$	$= x := e$
$An_{\text{Imp}}(S; S')$	$= An_{\text{Imp}}(S); An_{\text{Imp}}(S')$
$An_{\text{Imp}}(\text{if } b \text{ then } S \text{ else } S')$	$= (\text{if } b \text{ then } An_{\text{Imp}}(S)$ $\quad \text{else } An_{\text{Imp}}(S'))$
$An_{\text{Imp}}(\text{while } b \text{ do } S)$	$= (\text{while } b \text{ do } \text{cost} := \text{cost} + \kappa(b) + \kappa(S) + 1; An_{\text{Imp}}(S))$
$\kappa(\text{skip}) = 0$	$\kappa(x := e) = \kappa(e) + \kappa(\text{setvar})$
$\kappa(S; S') = \kappa(S) + \kappa(S')$	$\kappa(\text{if } b \text{ then } S \text{ else } S') = \kappa(b) + \max(\kappa(S) + \kappa(\text{branch}), \kappa(S'))$
$\kappa(\text{while } b \text{ do } S) = \kappa(b)$	
$\kappa(e < e') = \kappa(e) + \kappa(e') + \kappa(\text{bge})$	$\kappa(e + e') = \kappa(e) + \kappa(e') + \kappa(\text{add})$
$\kappa(n) = \kappa(\text{cnst})$	$\kappa(x) = \kappa(\text{var})$

Table 9: Annotation for Imp programs

**Proposition 9** *The following properties hold.*

(1) *If  $(e, s) \Downarrow v$  then*

$$C \cdot An_{\text{Vm}}(\mathcal{C}(e)) \cdot C' \vdash (i, \sigma, s[c/\text{cost}]) \xrightarrow{*} (j, v \cdot \sigma, s[c + \kappa(e)/\text{cost}])$$

where  $i = |C|$ ,  $j = |C \cdot An_{\text{Vm}}(\mathcal{C}(e))|$ .

(2) *If  $(b, s) \Downarrow \text{true}$  then*

$$C \cdot An_{\text{Vm}}(\mathcal{C}(b, k)) \cdot C' \vdash (i, \sigma, s[c/\text{cost}]) \xrightarrow{*} (5 \cdot k + j, \sigma, s[c + \kappa(b)/\text{cost}])$$

where  $i = |C|$ ,  $j = |C \cdot An_{\text{Vm}}(\mathcal{C}(b, k))|$ .

(3) *If  $(b, s) \Downarrow \text{false}$  then*

$$C \cdot An_{\text{Vm}}(\mathcal{C}(b, k)) \cdot C' \vdash (i, \sigma, s[c/\text{cost}]) \xrightarrow{*} (j, \sigma, s[c + \kappa(b)/\text{cost}])$$

where  $i = |C|$ ,  $j = |C \cdot An_{\text{Vm}}(\mathcal{C}(b, k))|$ .

(4) *If  $(An_{\text{Imp}}(S), s[c/\text{cost}]) \Downarrow s'[c'/\text{cost}]$  then*

$$C \cdot An_{\text{Vm}}(\mathcal{C}(S)) \cdot C' \vdash (i, \sigma, s[d/\text{cost}]) \xrightarrow{*} (j, \sigma, s'[d'/\text{cost}])$$

where  $i = |C|$ ,  $j = |C \cdot An_{\text{Vm}}(\mathcal{C}(S))|$ , and  $(d' - d) \leq (c' - c) + \kappa(S)$ .

### 3.3 Composition

The soundness of the cost annotations can be composed so as to obtain the soundness of the cost annotations of the source program relatively to the one of the object code.

**Proposition 10** *If  $(An_{\text{Imp}}(P), s[0/\text{cost}]) \Downarrow s'[c'/\text{cost}]$  and  $m \Vdash \epsilon, s[0/l_{\text{cost}}]$  then*

$$(An_{\text{Mips}}(\mathcal{C}'(\mathcal{C}(P))), m) \Downarrow m'$$

where  $m'(l_{\text{cost}}) \leq c'$  and  $m'[c'/l_{\text{cost}}] \Vdash \epsilon, s'[c'/\text{cost}]$ .

### 3.4 Coq development

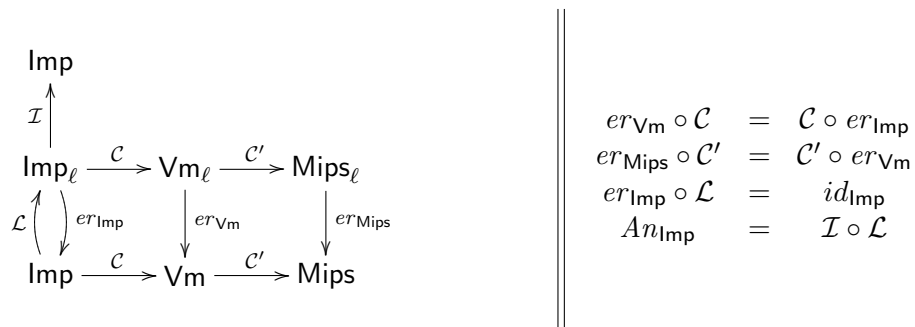
We have formalised and mechanically checked in COQ the application of the direct approach to the toy compiler (but for propositions 8 and 10 for which we provide a ‘paper proof’ in the appendix). By current standards, this is a small size development including 1000 lines of specifications and 4000 lines of proofs. Still there are a couple of points that deserve to be mentioned. First, we did not find a way of re-using the soundness proof of the compiler in a modular way. As a matter of fact, the soundness proof of the annotations is intertwined with the soundness proof of the compiler. Second, the manipulation of the cost variables in the annotated programs entails a significant increase in the size of the proofs. In particular, the soundness proof for the compilation function  $\mathcal{C}$  from `Imp` to `Vm` available in [9] is roughly 7 times smaller than the soundness proof of the annotation function of the `Imp` code relatively to the one of the `Vm` code.

### 3.5 Limitations of the direct approach

As mentioned in section 3.2, the annotation function for the source language introduces some over-approximation thus *failing to be precise*. On one hand, it is easy to modify the definitions in table 9 so that they compute the cost of each branch of an `if_then_else` separately rather than taking the maximum cost of the two branches. On the other hand, it is rather difficult to refine the annotation function so that it accounts for the memory hierarchy in the `Mips` machine; one needs to pass to the function  $\kappa$  an ‘hidden parameter’ which corresponds to the stack height. This process of pushing hidden parameters into the definition of the annotation is error prone and it seems unlikely to work in practice for a realistic compiler. We programmed sound (but not precise) cost annotations for the C compiler introduced in section 1.7 and found that the approach is difficult to test because an over-approximation of the cost at some point may easily compensate an under-approximation somewhere else. By contrast, in the labelling approach introduced in section 1.5, we manipulate costs at an abstract level as labels and produce numerical values only at the very end of the construction.

## 4 Labelling approach for the toy compiler

We apply the labelling approach introduced in section 1.5 to the toy compiler which results in the following diagram.



### 4.1 Labelled Imp

We extend the syntax so that statements and boolean conditions in while loops can be labelled.

$$\begin{array}{ll}
lb ::= \ell : b & \text{(labelled boolean conditions)} \\
S ::= \dots \mid \ell : S \mid \text{while } lb \text{ do } S & \text{(labelled commands)}
\end{array}$$

For instance  $\ell : (\text{while } \ell' : (n < x) \text{ do } \ell : S)$  is a labelled command. Also notice that the definition allows labels in commands to be nested as in  $\ell : (\ell' : (x := e))$ , though this feature is not really used. The evaluation predicate  $\Downarrow$  on labelled boolean conditions is extended as follows:

$$\frac{(b, s) \Downarrow v}{(\ell : b, s) \Downarrow (v, \ell)} \quad (5)$$

So a labelled boolean condition evaluates into a pair composed of a boolean value and a label. The small step semantics of statements defined in table 2 is extended as follows.

$$\begin{array}{ll}
(\ell : S, K, s) & \xrightarrow{\ell} (S, K, s) \\
(\text{while } lb \text{ do } S, K, s) & \xrightarrow{\ell} \begin{cases} (S, (\text{while } lb \text{ do } S); K, s) & \text{if } (lb, s) \Downarrow (\text{true}, \ell) \\ (\text{skip}, K, s) & \text{if } (lb, s) \Downarrow (\text{false}, \ell) \end{cases}
\end{array}$$

We denote with  $\lambda, \lambda', \dots$  finite sequences of labels. In particular, we denote with  $\epsilon$  the empty sequence and identify an unlabelled transition with a transition labelled with  $\epsilon$ . Then the small step reduction relation we have defined on statements becomes a *labelled transition system*. There is an obvious *erasure* function  $er_{\text{Imp}}$  from the labelled language to the unlabelled one which is the identity on expressions, removes labels from labelled boolean conditions and is the identity on unlabelled ones, and traverses commands removing all labels. We derive a *labelled* big-step semantics as follows:

$$(S, s) \Downarrow (s', \lambda) \quad \text{if } (S, \text{halt}, s) \xrightarrow{\lambda_1} \dots \xrightarrow{\lambda_n} (\text{skip}, \text{halt}, s') \text{ and } \lambda = \lambda_1 \dots \lambda_n .$$

## 4.2 Labelled Vm

We introduce a new instruction  $\text{nop}(\ell)$  whose semantics is defined as follows:

$$C \vdash (i, \sigma, s) \xrightarrow{\ell} (i + 1, \sigma, s) \quad \text{if } C[i] = \text{nop}(\ell) .$$

The erasure function  $er_{\text{Vm}}$  amounts to remove from a Vm code  $C$  all the  $\text{nop}(\ell)$  instructions and recompute jumps accordingly. Specifically, let  $n(C, i, j)$  be the number of  $\text{nop}$  instructions in the interval  $[i, j]$ . Then, assuming  $C[i] = \text{branch}(k)$  we replace the offset  $k$  with an offset  $k'$  determined as follows:

$$k' = \begin{cases} k - n(C, i, i + k) & \text{if } k \geq 0 \\ k + n(C, i + 1 + k, i) & \text{if } k < 0 \end{cases}$$

The compilation function  $\mathcal{C}$  is extended to  $\text{Imp}_\ell$  by defining:

$$\mathcal{C}(\ell : b, k) = (\text{nop}(\ell)) \cdot \mathcal{C}(b, k) \quad \mathcal{C}(\ell : S) = (\text{nop}(\ell)) \cdot \mathcal{C}(S) .$$

**Proposition 11** *For all commands  $S$  in  $\text{Imp}_\ell$  we have that:*

- (1)  $er_{\text{Vm}}(\mathcal{C}(S)) = \mathcal{C}(er_{\text{Imp}}(S))$ .
- (2) *If  $(S, s) \Downarrow (s', \lambda)$  then  $(\mathcal{C}(S), s) \Downarrow (s', \lambda)$ .*

**Remark 12** *In the current formulation, a sequence of transitions  $\lambda$  in the source code must be simulated by the same sequence of transitions in the object code. However, in the actual computation of the costs, the order of the labels occurring in the sequence is immaterial. Therefore one may consider a more relaxed notion of simulation where  $\lambda$  is a multi-set of labels.*

### 4.3 Labelled Mips

The labelled extension of **Mips** is similar to the one of **Vm**. We add an instruction **nop**  $\ell$  whose semantics is defined as follows:

$$M \vdash (i, m) \xrightarrow{\ell} (i + 1, m) \quad \text{if } M[i] = (\text{nop } \ell) .$$

The *erasure function*  $er_{\text{Mips}}$  is also similar to the one of **Vm** as it amounts to remove from a **Mips** code all the **(nop**  $\ell$ ) instructions and recompute jumps accordingly. The compilation function  $\mathcal{C}'$  is extended to  $\text{Vm}_\ell$  by simply translating **nop**( $\ell$ ) as **(nop**  $\ell$ ):

$$\mathcal{C}'(i, C) = (\text{nop } \ell) \quad \text{if } C[i] = \text{nop}(\ell)$$

The evaluation predicate for labelled **Mips** is defined as  $(M, m) \Downarrow (m', \lambda)$  if  $M \vdash (0, m) \xrightarrow{\lambda_1} \dots \xrightarrow{\lambda_n} (j, m')$ ,  $\lambda = \lambda_1 \cdots \lambda_n$  and  $M[j] = \text{halt}$ . The following proposition relates  $\text{Vm}_\ell$  code and its compilation and it is similar to proposition 11.

**Proposition 13** *Let  $C$  be a  $\text{Vm}_\ell$  code. Then:*

- (1)  $er_{\text{Mips}}(\mathcal{C}'(C)) = \mathcal{C}'(er_{\text{Vm}}(C))$ .
- (2) *If  $(C, s) \Downarrow (s', \lambda)$  and  $m \Vdash -\epsilon, s$  then  $(\mathcal{C}'(C), m) \Downarrow (m', \lambda)$  and  $m' \Vdash -\epsilon, s'$ .*

### 4.4 Labellings and instrumentations

Assuming a function  $\kappa$  which associates an integer number with labels and a distinct variable *cost* which does not occur in the program  $P$  under consideration, we abbreviate with  $inc(\ell)$  the assignment  $cost := cost + \kappa(\ell)$ . Then we define the instrumentation  $\mathcal{I}$  (relative to  $\kappa$  and *cost*) as follows.

$$\begin{aligned} \mathcal{I}(\ell : S) &= inc(\ell); \mathcal{I}(S) \\ \mathcal{I}(\text{while } \ell : b \text{ do } S) &= inc(\ell); \text{while } b \text{ do } (\mathcal{I}(S); inc(\ell)) \end{aligned}$$

The function  $\mathcal{I}$  just distributes over the other operators of the language. We extend the function  $\kappa$  on labels to sequences of labels by defining  $\kappa(\ell_1, \dots, \ell_n) = \kappa(\ell_1) + \dots + \kappa(\ell_n)$ .

The instrumented **Imp** program relates to the labelled one has follows.

**Proposition 14** *Let  $S$  be an  $\text{Imp}_\ell$  command. If  $(\mathcal{I}(S), s[c/cost]) \Downarrow s'[c+\delta/cost]$  then  $\exists \lambda \kappa(\lambda) = \delta$  and  $(S, s[c/cost]) \Downarrow (s'[c/cost], \lambda)$ .*

**Definition 15** *A labelling is a function  $\mathcal{L}$  from an unlabelled language to the corresponding labelled one such that  $er_{\text{Imp}} \circ \mathcal{L}$  is the identity function on the **Imp** language.*

**Proposition 16** *For any labelling function  $\mathcal{L}$ , and **Imp** program  $P$ , the following holds:*

$$er_{\text{Mips}}(\mathcal{C}'(\mathcal{C}(\mathcal{L}(P)))) = \mathcal{C}'(\mathcal{C}(P)) . \tag{6}$$

**Proposition 17** *Given a function  $\kappa$  for the labels and a labelling function  $\mathcal{L}$ , for all programs  $P$  of the source language if  $(\mathcal{I}(\mathcal{L}(P)), s[c/cost]) \Downarrow s'[c + \delta/cost]$  and  $m \Vdash -s[c/cost]$  then  $(\mathcal{C}'(\mathcal{C}(\mathcal{L}(P))), m) \Downarrow (m', \lambda)$ ,  $m' \Vdash -s'[c/cost]$  and  $\kappa(\lambda) = \delta$ .*

#### 4.5 Sound and precise labellings

With any  $\text{Mips}_\ell$  code  $M$  we can associate a directed and rooted (control flow) graph whose nodes are the instruction positions  $\{0, \dots, |M| - 1\}$ , whose root is the node 0, and whose directed edges correspond to the possible transitions between instructions. We say that a node is labelled if it corresponds to an instruction  $\text{nop } \ell$ .

**Definition 18** *A simple path in a  $\text{Mips}_\ell$  code  $M$  is a directed finite path in the graph associated with  $M$  where the first node is labelled, the last node is the predecessor of either a labelled node or a leaf, and all the other nodes are unlabelled.*

**Definition 19** *A  $\text{Mips}_\ell$  code  $M$  is soundly labelled if in the associated graph the root node 0 is labelled and there are no loops that do not go through a labelled node.*

In a soundly labelled graph there are finitely many simple paths. Thus, given a soundly labelled  $\text{Mips}$  code  $M$ , we can associate with every label  $\ell$  a number  $\kappa(\ell)$  which is the maximum (estimated) cost of executing a simple path whose first node is labelled with  $\ell$ . We stress that in the following we assume that the cost of a simple path is proportional to the number of  $\text{Mips}$  instructions that are crossed in the path.

**Proposition 20** *If  $M$  is soundly labelled and  $(M, m) \Downarrow (m', \lambda)$  then the cost of the computation is bounded by  $\kappa(\lambda)$ .*

Thus for a soundly labelled  $\text{Mips}$  code the sequence of labels associated with a computation is a significant information on the execution cost.

**Definition 21** *We say that a soundly labelled code is precise if for every label  $\ell$  in the code, the simple paths starting from a node labelled with  $\ell$  have the same cost.*

In particular, a code is precise if we can associate at most one simple path with every label.

**Proposition 22** *If  $M$  is precisely labelled and  $(M, m) \Downarrow (m', \lambda)$  then the cost of the computation is  $\kappa(\lambda)$ .*

The next point we have to check is that there are labelling functions (of the source code) such that the compilation function does produce sound and possibly precise labelled  $\text{Mips}$  code. To discuss this point, we introduce in table 10 two labelling functions  $\mathcal{L}_s$  and  $\mathcal{L}_p$  for the  $\text{Imp}$  language where the operator *new* is meant to return fresh labels.

**Proposition 23** *For all  $\text{Imp}$  programs  $P$ :*

- (1)  $\mathcal{C}'(\mathcal{C}(\mathcal{L}_s(P)))$  is a soundly labelled  $\text{Mips}$  code.
- (2)  $\mathcal{C}'(\mathcal{C}(\mathcal{L}_p(P)))$  is a soundly and precisely labelled  $\text{Mips}$  code.

$\mathcal{L}_s(\text{prog } S)$	$= \text{prog } \ell : \mathcal{L}_s(S)$	
$\mathcal{L}_s(\text{skip})$	$= \text{skip}$	
$\mathcal{L}_s(x := e)$	$= x := e$	
$\mathcal{L}_s(S; S')$	$= \mathcal{L}_s(S); \mathcal{L}_s(S')$	
$\mathcal{L}_s(\text{if } b \text{ then } S_1 \text{ else } S_2)$	$= \text{if } b \text{ then } \mathcal{L}_s(S_1) \text{ else } \mathcal{L}_s(S_2)$	
$\mathcal{L}_s(\text{while } b \text{ do } S)$	$= \text{while } b \text{ do } \ell : \mathcal{L}_s(S)$	
$\mathcal{L}_p(\text{prog } S)$	$= \text{let } \ell = \text{new in } \text{prog } \ell : \mathcal{L}_p(S)$	
$\mathcal{L}_p(\text{skip})$	$= \text{let } \ell = \text{new in } \ell : (\text{skip})$	
$\mathcal{L}_p(x := e)$	$= \text{let } \ell = \text{new in } \ell : (x := e)$	
$\mathcal{L}_p(S; S')$	$= \mathcal{L}_p(S); \mathcal{L}_p(S')$	
$\mathcal{L}_p(\text{if } b \text{ then } S_1 \text{ else } S_2)$	$= \text{let } \ell = \text{new in } \text{if } \ell : b \text{ then } \mathcal{L}_p(S_1) \text{ else } \mathcal{L}_p(S_2)$	
$\mathcal{L}_p(\text{while } b \text{ do } S)$	$= \text{let } \ell = \text{new in } \text{while } \ell : b \text{ do } \mathcal{L}_p(S)$	

Table 10: Two labellings for the Imp language

For an example of command which is not soundly labelled consider  $\ell : \text{while } 0 < x \text{ do } x := x + 1$  which when compiled produces a loop that does not go through any label. On the other hand, for an example of a program which is not precisely labelled consider  $(\text{while } \ell : (0 < x) \text{ do } x := x + 1)$ . In the compiled code, we find two simple paths associated with the label  $\ell$  whose cost will be quite different in general.

Once a sound and possibly precise labelling  $\mathcal{L}$  has been designed, we can determine the cost of each label and define an instrumentation  $\mathcal{I}$  whose composition with  $\mathcal{L}$  will produce the desired cost annotation.

**Definition 24** *Given a labelling function  $\mathcal{L}$  for the source language Imp and a program  $P$  in the Imp language, we define an annotation for the source program as follows:*

$$An_{\text{Imp}}(P) = \mathcal{I}(L(P)) .$$

**Proposition 25** *If  $P$  is a program and  $\mathcal{C}'(\mathcal{C}(\mathcal{L}(P)))$  is a sound (sound and precise) labelling then  $(An_{\text{Imp}}(P), s[c/\text{cost}]) \Downarrow s'[c + \delta/\text{cost}]$  and  $m \Vdash s[c/\text{cost}]$  entails that  $(\mathcal{C}'(\mathcal{C}(P)), m) \Downarrow m'$ ,  $m' \Vdash s'[c/\text{cost}]$  and the cost of the execution is bound (is exactly)  $\delta$ .*

To summarise, producing sound and precise labellings is mainly a matter of designing the labelled source language so that the labelling is sufficiently *fine grained*. For instance, in the toy compiler, the fact that boolean conditions are labelled is instrumental to precision while the labelling of expressions turns out to be unnecessary.

Besides soundness and precision, a third criteria to evaluate labellings is that they do not introduce too many unnecessary labels. We call this property *economy*. There are two reasons for this requirement. On one hand we would like to minimise the number of labels so that the source program is not cluttered by too many cost annotations and on the other hand we would like to maximise the length of the simple paths because in a non-trivial processor the longer the sequence of instructions we consider the more accurate is the estimation of their execution cost (on a long sequence certain costs are amortized). In practice, it seems that one can produce first a sound and possibly precise labelling and then apply heuristics to eliminate unnecessary labels.

## 5 A C compiler

This section gives an informal overview of the compiler, in particular it highlights the main features of the intermediate languages, the purpose of the compilation steps, and the optimizations.

### 5.1 Clight

Clight is a large subset of the C language that we adopt as the source language of our compiler. It features most of the types and operators of C. It includes pointer arithmetic, pointers to functions, and `struct` and `union` types, as well as all C control structures. The main difference with the C language is that Clight expressions are side-effect free, which means that side-effect operators (`=`, `+=`, `++`, ...) and function calls within expressions are not supported. Given a C program, we rely on the CIL tool [11] to deal with the idiosyncrasy of C concrete syntax and to produce an equivalent program in Clight abstract syntax. We refer to the CompCert project [8] for a formal definition of the Clight language. Here we just recall in figure 5.1 its syntax which is classically structured in expressions, statements, functions, and whole programs. In order to limit the implementation effort, our current compiler for Clight does *not* cover the operators relating to the floating point type `float`. So, in a nutshell, the fragment of C we have implemented is Clight without floating point.

### 5.2 Cminor

Cminor is a simple, low-level imperative language, comparable to a stripped-down, typeless variant of C. Again we refer to the CompCert project for its formal definition and we just recall in figure 5.2 its syntax which as for Clight is structured in expressions, statements, functions, and whole programs.

**Translation of Clight to Cminor.** As in Cminor stack operations are made explicit, one has to know which variables are stored in the stack. This information is produced by a static analysis that determines the variables whose address may be ‘taken’. Also space is reserved for local arrays and structures. In a second step, the proper compilation is performed: it consists mainly in translating Clight control structures to the basic ones available in Cminor.

### 5.3 RTLabs

RTLabs is the last architecture independent language in the compilation process. It is a rather straightforward *abstraction* of the *architecture-dependent* RTL intermediate language available in the CompCert project and it is intended to factorize some work common to the various target assembly languages (e.g. optimizations) and thus to make retargeting of the compiler a simpler matter.

We stress that in RTLabs the structure of Cminor expressions is lost and that this may have a negative impact on the following instruction selection step. Still, the subtleties of instruction selection seem rather orthogonal to our goals and we deem the possibility of retargeting easily the compiler more important than the efficiency of the generated code.

Expressions:	$a ::=$ <ul style="list-style-type: none"> <li><math>id</math></li> <li><math>  n</math></li> <li><math>  \text{sizeof}(\tau)</math></li> <li><math>  op_1 a</math></li> <li><math>  a op_2 a</math></li> <li><math>  *a</math></li> <li><math>  a.id</math></li> <li><math>  \&amp;a</math></li> <li><math>  (\tau)a</math></li> <li><math>  a?a : a</math></li> </ul>	<ul style="list-style-type: none"> <li>variable identifier</li> <li>integer constant</li> <li>size of a type</li> <li>unary arithmetic operation</li> <li>binary arithmetic operation</li> <li>pointer dereferencing</li> <li>field access</li> <li>taking the address of</li> <li>type cast</li> <li>conditional expression</li> </ul>
Statements:	$s ::=$ <ul style="list-style-type: none"> <li><b>skip</b></li> <li><math>  a = a</math></li> <li><math>  a = a(a^*)</math></li> <li><math>  a(a^*)</math></li> <li><math>  s; s</math></li> <li><b>if</b> <math>a</math> <b>then</b> <math>s</math> <b>else</b> <math>s</math></li> <li><b>switch</b> <math>a</math> <math>sw</math></li> <li><b>while</b> <math>a</math> <b>do</b> <math>s</math></li> <li><b>do</b> <math>s</math> <b>while</b> <math>a</math></li> <li><b>for</b>(<math>s, a, s</math>) <math>s</math></li> <li><b>break</b></li> <li><b>continue</b></li> <li><b>return</b> <math>a^?</math></li> <li><b>goto</b> <math>lbl</math></li> <li><math>lbl : s</math></li> </ul>	<ul style="list-style-type: none"> <li>empty statement</li> <li>assignment</li> <li>function call</li> <li>procedure call</li> <li>sequence</li> <li>conditional</li> <li>multi-way branch</li> <li>“while” loop</li> <li>“do” loop</li> <li>“for” loop</li> <li>exit from current loop</li> <li>next iteration of the current loop</li> <li>return from current function</li> <li>branching</li> <li>labelled statement</li> </ul>
Switch cases:	$sw ::=$ <ul style="list-style-type: none"> <li><b>default</b> : <math>s</math></li> <li><math>  \text{case } n : s; sw</math></li> </ul>	<ul style="list-style-type: none"> <li>default case</li> <li>labelled case</li> </ul>
Variable declarations:	$dcl ::=$ <ul style="list-style-type: none"> <li><math>(\tau \ id)^*</math></li> </ul>	<ul style="list-style-type: none"> <li>type and name</li> </ul>
Functions:	$Fd ::=$ <ul style="list-style-type: none"> <li><math>\tau \ id(dcl)\{dcl; s\}</math></li> <li><math>  \text{extern } \tau \ id(dcl)</math></li> </ul>	<ul style="list-style-type: none"> <li>internal function</li> <li>external function</li> </ul>
Programs:	$P ::=$ <ul style="list-style-type: none"> <li><math>dcl; Fd^*; \text{main} = id</math></li> </ul>	<ul style="list-style-type: none"> <li>global variables, functions, entry point</li> </ul>

Figure 1: Syntax of the Clight language



Signatures:	$sig ::= sig \vec{int} \ (int void)$	arguments and result
Expressions:	$a ::=$ <ul style="list-style-type: none"> <li><math>id</math></li> <li><math>  n</math></li> <li><math>  \text{addrsymbol}(id)</math></li> <li><math>  \text{addrstack}(\delta)</math></li> <li><math>  op_1 a</math></li> <li><math>  op_2 a a</math></li> <li><math>  \kappa[a]</math></li> <li><math>  a?a : a</math></li> </ul>	<ul style="list-style-type: none"> <li>local variable</li> <li>integer constant</li> <li>address of global symbol</li> <li>address within stack data</li> <li>unary arithmetic operation</li> <li>binary arithmetic operation</li> <li>memory read</li> <li>conditional expression</li> </ul>
Statements:	$s ::=$ <ul style="list-style-type: none"> <li><b>skip</b></li> <li><math>  id = a</math></li> <li><math>  \kappa[a] = a</math></li> <li><math>  id^? = a(\vec{a}) : sig</math></li> <li><math>  \text{tailcall } a(\vec{a}) : sig</math></li> <li><math>  \text{return}(a^?)</math></li> <li><math>  s; s</math></li> <li><math>  \text{if } a \text{ then } s \text{ else } s</math></li> <li><math>  \text{loop } s</math></li> <li><math>  \text{block } s</math></li> <li><math>  \text{exit } n</math></li> <li><math>  \text{switch } a \text{ tbl}</math></li> <li><math>  lbl : s</math></li> <li><math>  \text{goto } lbl</math></li> </ul>	<ul style="list-style-type: none"> <li>empty statement</li> <li>assignment</li> <li>memory write</li> <li>function call</li> <li>function tail call</li> <li>function return</li> <li>sequence</li> <li>conditional</li> <li>infinite loop</li> <li>block delimiting <b>exit</b> constructs</li> <li>terminate the <math>(n + 1)^{th}</math> enclosing block</li> <li>multi-way test and exit</li> <li>labelled statement</li> <li>jump to a label</li> </ul>
Switch tables:	$tbl ::=$ <ul style="list-style-type: none"> <li><b>default:exit</b>(<math>n</math>)</li> <li><math>  \text{case } i: \text{exit}(n);tbl</math></li> </ul>	
Functions:	$Fd ::=$ <ul style="list-style-type: none"> <li><b>internal</b> <math>sig \vec{id} \vec{id} n s</math></li> <li><math>  \text{external } id \ sig</math></li> </ul>	<ul style="list-style-type: none"> <li>internal function: signature, parameters, local variables, stack size and body</li> <li>external function</li> </ul>
Programs:	$P ::= \text{prog } (id = data)^* (id = Fd)^* id$	global variables, functions and entry point

Figure 2: Syntax of the Cminor language

$return\_type ::= int \mid void$	$signature ::= (int \rightarrow)^* return\_type$	
$memq ::= int8s \mid int8u \mid int16s \mid int16u \mid int32$	$fun\_ref ::= fun\_name \mid psd\_reg$	
$instruction ::=$	$skip \rightarrow node$ $psd\_reg := op(psd\_reg^*) \rightarrow node$ $psd\_reg := \&var\_name \rightarrow node$ $psd\_reg := \&locals[n] \rightarrow node$ $psd\_reg := fun\_name \rightarrow node$ $psd\_reg := memq(psd\_reg[psd\_reg]) \rightarrow node$ $memq(psd\_reg[psd\_reg]) := psd\_reg \rightarrow node$ $psd\_reg := fun\_ref(psd\_reg^*) : signature \rightarrow node$ $fun\_ref(psd\_reg^*) : signature$ $test\ op(psd\_reg^*) \rightarrow node, node$ $return\ psd\_reg?$	(no instruction) (operation) (address of a global) (address of a local) (address of a function) (memory load) (memory store) (function call) (function tail call) (branch) (return)
	$fun\_def ::= fun\_name(psd\_reg^*) : signature$	
	$result : psd\_reg?$	
	$locals : psd\_reg^*$	
	$stack : n$	
	$entry : node$	
	$exit : node$	
	$(node : instruction)^*$	
$init\_datum ::= reserve(n) \mid int8(n) \mid int16(n) \mid int32(n)$	$init\_data ::= init\_datum^+$	
$global\_decl ::= var\ var\_name\{init\_data\}$	$fun\_decl ::= extern\ fun\_name(signature) \mid fun\_def$	
	$program ::= global\_decl^*$	
	$fun\_decl^*$	

Table 11: Syntax of the RTLabs language

**Syntax.** In RTLabs, programs are represented as *control flow graphs* (CFGs for short). We associate with the nodes of the graphs instructions reflecting the Cminor commands. As usual, commands that change the control flow of the program (e.g. loops, conditionals) are translated by inserting suitable branching instructions in the CFG. The syntax of the language is depicted in table 11. Local variables are now represented by *pseudo registers* that are available in unbounded number. The grammar rule *op* that is not detailed in table 11 defines usual arithmetic and boolean operations (+, xor,  $\leq$ , etc.) as well as constants and conversions between sized integers.

**Translation of Cminor to RTLabs.** Translating Cminor programs to RTLabs programs mainly consists in transforming Cminor commands in CFGs. Most commands are sequential and have a rather straightforward linear translation. A conditional is translated in a branch instruction; a loop is translated using a back edge in the CFG.

<i>size</i> ::= Byte   HalfWord   Word	<i>fun_ref</i> ::= <i>fun_name</i>   <i>psd_reg</i>
<i>instruction</i> ::=	
skip → <i>node</i>	(no instruction)
<i>psd_reg</i> := <i>n</i> → <i>node</i>	(constant)
<i>psd_reg</i> := <i>unop</i> ( <i>psd_reg</i> ) → <i>node</i>	(unary operation)
<i>psd_reg</i> := <i>binop</i> ( <i>psd_reg</i> , <i>psd_reg</i> ) → <i>node</i>	(binary operation)
<i>psd_reg</i> := & <i>globals</i> [ <i>n</i> ] → <i>node</i>	(address of a global)
<i>psd_reg</i> := & <i>locals</i> [ <i>n</i> ] → <i>node</i>	(address of a local)
<i>psd_reg</i> := <i>fun_name</i> → <i>node</i>	(address of a function)
<i>psd_reg</i> := <i>size</i> ( <i>psd_reg</i> [ <i>n</i> ]) → <i>node</i>	(memory load)
<i>size</i> ( <i>psd_reg</i> [ <i>n</i> ]) := <i>psd_reg</i> → <i>node</i>	(memory store)
<i>psd_reg</i> := <i>fun_ref</i> ( <i>psd_reg</i> *) → <i>node</i>	(function call)
<i>fun_ref</i> ( <i>psd_reg</i> *)	(function tail call)
test <i>uncon</i> ( <i>psd_reg</i> ) → <i>node</i> , <i>node</i>	(branch unary condition)
test <i>bincon</i> ( <i>psd_reg</i> , <i>psd_reg</i> ) → <i>node</i> , <i>node</i>	(branch binary condition)
return <i>psd_reg</i> ?	(return)
<i>fun_def</i> ::=	<i>program</i> ::=
<i>fun_name</i> ( <i>psd_reg</i> *)	<i>globals</i> : <i>n</i>
<i>result</i> : <i>psd_reg</i> ?	<i>fun_def</i> *
<i>locals</i> : <i>psd_reg</i> *	
<i>stack</i> : <i>n</i>	
<i>entry</i> : <i>node</i>	
<i>exit</i> : <i>node</i>	
( <i>node</i> : <i>instruction</i> )*	

Table 12: Syntax of the RTL language

## 5.4 RTL

As in RTLabs, the structure of RTL programs is based on CFGs. RTL is the first architecture-dependant intermediate language of our compiler which, in its current version, targets the Mips assembly language.

**Syntax.** RTL is very close to RTLabs. It is based on CFGs and explicit the Mips instructions corresponding to the RTLabs instructions. Type information disappears: everything is represented using 32 bits integers. Moreover, each global of the program is associated with an offset. The syntax of the language can be found in table 12. The grammar rules *unop*, *binop*, *uncon*, and *bincon*, respectively, represent the sets of unary operations, binary operations, unary conditions and binary conditions of the Mips language.

**Translation of RTLabs to RTL.** This translation is mostly straightforward. A RTLabs instruction is often directly translated to a corresponding Mips instruction. There are a few exceptions: some RTLabs instructions are expanded in two or more Mips instructions. When the translation of a RTLabs instruction requires more than a few simple Mips instruction, it is translated into a call to a function defined in the preamble of the compilation result.

## 5.5 ERTL

As in RTL, the structure of ERTL programs is based on CFGs. ERTL explicit the calling conventions of the Mips assembly language.

$size ::=$	Byte   HalfWord   Word	$fun\_ref ::=$	$fun\_name$   $psd\_reg$
$instruction ::=$	skip $\rightarrow node$		(no instruction)
	NewFrame $\rightarrow node$		(frame creation)
	DelFrame $\rightarrow node$		(frame deletion)
	$psd\_reg := stack[slot, n] \rightarrow node$		(stack load)
	$stack[slot, n] := psd\_reg \rightarrow node$		(stack store)
	$hdw\_reg := psd\_reg \rightarrow node$		(pseudo to hardware)
	$psd\_reg := hdw\_reg \rightarrow node$		(hardware to pseudo)
	$psd\_reg := n \rightarrow node$		(constant)
	$psd\_reg := unop(psd\_reg) \rightarrow node$		(unary operation)
	$psd\_reg := binop(psd\_reg, psd\_reg) \rightarrow node$		(binary operation)
	$psd\_reg := fun\_name \rightarrow node$		(address of a function)
	$psd\_reg := size(psd\_reg[n]) \rightarrow node$		(memory load)
	$size(psd\_reg[n]) := psd\_reg \rightarrow node$		(memory store)
	$fun\_ref(n) \rightarrow node$		(function call)
	$fun\_ref(n)$		(function tail call)
	test uncon( $psd\_reg$ ) $\rightarrow node, node$		(branch unary condition)
	test bincon( $psd\_reg, psd\_reg$ ) $\rightarrow node, node$		(branch binary condition)
	return $b$		(return)
$fun\_def ::=$	$fun\_name(n)$	$program ::=$	globals : $n$
	locals : $psd\_reg^*$		$fun\_def^*$
	stack : $n$		
	entry : $node$		
	( $node : instruction$ )*		

Table 13: Syntax of the ERTL language

**Syntax.** The syntax of the language is given in table 13. The main difference between RTL and ERTL is the use of hardware registers. Parameters are passed in specific hardware registers; if there are too many parameters, the remaining are stored in the stack. Other conventionally specific hardware registers are used: a register that holds the result of a function, a register that holds the base address of the globals, a register that holds the address of the top of the stack, and some registers that need to be saved when entering a function and whose values are restored when leaving a function. Following these conventions, function calls do not list their parameters anymore; they only mention their number. Two new instructions appear to allocate and deallocate on the stack some space needed by a function to execute. Along with these two instructions come two instructions to fetch or assign a value in the parameter sections of the stack; these instructions cannot yet be translated using regular load and store instructions because we do not know the final size of the stack area of each function. At last, the return instruction has a boolean argument that tells whether the result of the function may later be used or not (this is exploited for optimizations).

**Translation of RTL to ERTL.** The work consists in expliciting the conventions previously mentioned. These conventions appear when entering, calling and leaving a function, and when referencing a global variable or the address of a local variable.

**Optimizations.** A *liveness analysis* is performed on ERTL to replace unused instructions by a skip. An instruction is tagged as unused when it performs an assignment on a register that will not be read afterwards. Also, the result of the liveness analysis is exploited by a *register*

<i>size</i> ::= Byte   HalfWord   Word	<i>fun_ref</i> ::= <i>fun_name</i>   <i>hdw_reg</i>
<i>instruction</i> ::=	
skip → <i>node</i>	(no instruction)
NewFrame → <i>node</i>	(frame creation)
DelFrame → <i>node</i>	(frame deletion)
<i>hdw_reg</i> := <i>n</i> → <i>node</i>	(constant)
<i>hdw_reg</i> := <i>unop</i> ( <i>hdw_reg</i> ) → <i>node</i>	(unary operation)
<i>hdw_reg</i> := <i>binop</i> ( <i>hdw_reg</i> , <i>hdw_reg</i> ) → <i>node</i>	(binary operation)
<i>hdw_reg</i> := <i>fun_name</i> → <i>node</i>	(address of a function)
<i>hdw_reg</i> := <i>size</i> ( <i>hdw_reg</i> [ <i>n</i> ]) → <i>node</i>	(memory load)
<i>size</i> ( <i>hdw_reg</i> [ <i>n</i> ]) := <i>hdw_reg</i> → <i>node</i>	(memory store)
<i>fun_ref</i> () → <i>node</i>	(function call)
<i>fun_ref</i> ()	(function tail call)
test <i>uncon</i> ( <i>hdw_reg</i> ) → <i>node</i> , <i>node</i>	(branch unary condition)
test <i>bincon</i> ( <i>hdw_reg</i> , <i>hdw_reg</i> ) → <i>node</i> , <i>node</i>	(branch binary condition)
return	(return)
<i>fun_def</i> ::= <i>fun_name</i> ( <i>n</i> )	<i>program</i> ::= <b>globals</b> : <i>n</i>
<i>locals</i> : <i>n</i>	<i>fun_def</i> *
<i>stack</i> : <i>n</i>	
<i>entry</i> : <i>node</i>	
( <i>node</i> : <i>instruction</i> )*	

Table 14: Syntax of the LTL language

*allocation* algorithm whose result is to efficiently associate a physical location (a hardware register or an address in the stack) to each pseudo register of the program.

## 5.6 LTL

As in ERTL, the structure of LTL programs is based on CFGs. Pseudo registers are not used anymore; instead, they are replaced by physical locations (a hardware register or an address in the stack).

**Syntax.** Except for a few exceptions, the instructions of the language are those of ERTL with hardware registers replacing pseudo registers. Calling and returning conventions were explicitated in ERTL; thus, function calls and returns do not need parameters in LTL. The syntax is defined in table 14.

**Translation of ERTL to LTL.** The translation relies on the results of the liveness analysis and of the register allocation. Unused instructions are eliminated and each pseudo register is replaced by a physical location. In LTL, the size of the stack frame of a function is known; instructions intended to load or store values in the stack are translated using regular load and store instructions.

**Optimizations.** A *graph compression* algorithm removes empty instructions generated by previous compilation passes and by the liveness analysis.

<i>size</i> ::= Byte   HalfWord   Word	<i>fun_ref</i> ::= <i>fun_name</i>   <i>hdw_reg</i>
<i>instruction</i> ::=	
NewFrame	(frame creation)
DelFrame	(frame deletion)
<i>hdw_reg</i> := <i>n</i>	(constant)
<i>hdw_reg</i> := <i>unop</i> ( <i>hdw_reg</i> )	(unary operation)
<i>hdw_reg</i> := <i>binop</i> ( <i>hdw_reg</i> , <i>hdw_reg</i> )	(binary operation)
<i>hdw_reg</i> := <i>fun_name</i>	(address of a function)
<i>hdw_reg</i> := <i>size</i> ( <i>hdw_reg</i> [ <i>n</i> ])	(memory load)
<i>size</i> ( <i>hdw_reg</i> [ <i>n</i> ]) := <i>hdw_reg</i>	(memory store)
call <i>fun_ref</i>	(function call)
tailcall <i>fun_ref</i>	(function tail call)
<i>uncon</i> ( <i>hdw_reg</i> ) → <i>node</i>	(branch unary condition)
<i>bincon</i> ( <i>hdw_reg</i> , <i>hdw_reg</i> ) → <i>node</i>	(branch binary condition)
<i>mips_label</i> :	(Mips label)
goto <i>mips_label</i>	(goto)
return	(return)
<i>fun_def</i> ::= <i>fun_name</i> ( <i>n</i> )	<i>program</i> ::= <i>globals</i> : <i>n</i>
<i>locals</i> : <i>n</i>	<i>fun_def</i> *
<i>instruction</i> *	

Table 15: Syntax of the LIN language

## 5.7 LIN

In LIN, the structure of a program is no longer based on CFGs. Every function is represented as a sequence of instructions.

**Syntax.** The instructions of LIN are very close to those of LTL. *Program labels*, *gotos* and branch instructions handle the changes in the control flow. The syntax of LIN programs is shown in table 15.

**Translation of LTL to LIN.** This translation amounts to transform in an efficient way the graph structure of functions into a linear structure of sequential instructions.

## 5.8 Mips

Mips is a rather simple assembly language. As for other assembly languages, a program in Mips is a sequence of instructions. The Mips code produced by the compilation of a Clight program starts with a preamble in which some useful and non-primitive functions are predefined (e.g. conversion from 8 bits unsigned integers to 32 bits integers). The subset of the Mips assembly language that the compilation produces is defined in table 16.

**Translation of LIN to Mips.** This final translation is simple enough. Stack allocation and deallocation are explicit and the function definitions are sequentialized.

<i>load</i> ::= lb   lhw   lw	<i>store</i> ::= sb   shw   sw	<i>fun_ref</i> ::= <i>fun_name</i>   <i>hdw_reg</i>
<i>instruction</i> ::=	nop	(empty instruction)
	li <i>hdw_reg</i> , <i>n</i>	(constant)
	unop <i>hdw_reg</i> , <i>hdw_reg</i>	(unary operation)
	binop <i>hdw_reg</i> , <i>hdw_reg</i> , <i>hdw_reg</i>	(binary operation)
	la <i>hdw_reg</i> , <i>fun_name</i>	(address of a function)
	load <i>hdw_reg</i> , <i>n</i> ( <i>hdw_reg</i> )	(memory load)
	store <i>hdw_reg</i> , <i>n</i> ( <i>hdw_reg</i> )	(memory store)
	call <i>fun_ref</i>	(function call)
	uncon <i>hdw_reg</i> , <i>node</i>	(branch unary condition)
	bincon <i>hdw_reg</i> , <i>hdw_reg</i> , <i>node</i>	(branch binary condition)
	<i>mips_label</i> :	(Mips label)
	j <i>mips_label</i>	(goto)
	return	(return)
	<i>program</i> ::=	globals : <i>n</i>
		entry : <i>mips_label</i> *
		<i>instruction</i> *

Table 16: Syntax of the Mips language

1. Label the input Clight program.
2. Compile the labelled Clight program in the labelled world. This produces a labelled Mips code.
3. For each label of the labelled Mips code, compute the cost of the instructions under its scope and generate a *label-cost mapping*. An unlabelled Mips code — the result of the compilation — is obtained by removing the labels from the labelled Mips code.
4. Add a fresh *cost variable* to the labelled Clight program and replace the labels by an increment of this cost variable according to the label-cost mapping. The result is an *annotated* Clight program with no label.

Table 17: Building the annotation of a Clight program in the labelling approach

## 6 Labelling approach for the C compiler

This section informally describes the labelled extensions of the languages in the compilation chain, the way the labels are propagated by the compilation functions, the labelling of the source code, the hypotheses on the control flow of the labelled Mips code and the verification that we perform on it, the way we build the instrumentation, and finally the way the labelling approach has been tested. The process of annotating a Clight program using the labelling approach is summarized in table 17 and is detailed in the following sections.

### 6.1 Labelled Clight and labelled Cminor

Both the Clight and Cminor languages are extended in the same way by labelling both statements and expressions (by comparison, in the toy language lmp we just labelled statements and boolean conditions). The labelling of expressions aims to capture precisely their execution cost. Indeed, Clight and Cminor include expressions such as  $a_1 ? a_2 ; a_3$  whose evaluation cost depends on the boolean value  $a_1$ .

As both languages are extended in the same way, the extended compilation does nothing

more than sending Clight labelled statements and expressions to Cminor labelled statements and expressions.

## 6.2 Labels in RTLAb and the back-end languages

The labelled version of RTLAb and the languages in the back-end language simply consists in adding a new instruction whose semantics is to emit a label without modifying the state. For the CFG based languages (RTLAb to LTL), this new instruction is `emit label → node`. For LIN and Mips, it is `emit label`. The translation of these label instructions is immediate. In Mips, we also rely on a reserved label `begin_function` to pinpoint the beginning of a function code (cf. section 6.3.6).

## 6.3 Labelling of the source language

The goal here is to add labels in the source program that cover every reachable instruction of the program and avoid unlabelled loops; this can be seen as a *soundness* property. Another important point is *precision*, meaning that a label might cover several paths to the next labels only if those paths have equal costs. Several labellings might satisfy the soundness and precision conditions, but from an engineering point of view, a labelling that makes obvious which instruction is under the scope of which label would be better. There is a thin line to find between too many labels — which may obfuscate the code — and too few labels — which makes it harder to see which instruction is under the scope of which label. The balance leans a bit towards the *economy* of labels because the cost of executing an assembly instruction often depends on its context (for instance by the status of the cache memory). We explain our labelling by considering the constructions of Clight and their compilation to Mips.

### 6.3.1 Sequential instructions

A sequence of Clight instructions that compile to sequential Mips code, such as a sequence of assignments, can be handled by a single label which covers the unique execution path. The example below illustrates the labelling of ‘sequential’ Clight instructions.

Clight	$\xrightarrow{\text{Labelling}}$	Labelled Clight	$\xrightarrow{\text{Compilation}}$	Labelled Mips
<pre>i = 0; tab[i] = x; x++;</pre>		<pre>_cost: i = 0; tab[i] = x; x++;</pre>		<pre>emit _cost li \$v0, 4 mul \$v0, \$zero, \$v0 add \$v0, \$a1, \$v0 sw \$a0, 0(\$v0) li \$v0, 1 add \$a0, \$a0, \$v0</pre>

### 6.3.2 Ternary expressions

Most Clight expressions compile to sequential Mips code. There is one exception: *ternary expressions* that introduce a branching in the control flow. Because of the precision condition, we must associate a label with each branch.



Clight	$\xrightarrow{\text{Labelling}}$	Labelled Clight	$\xrightarrow{\text{Compilation}}$	Labelled Mips
<pre>b ? x+1 :   y</pre>		<pre>b ? (_cost1: x+1) :   (_cost2: y)</pre>		<pre>beq \$a0, \$zero, c_false emit _cost1 li \$v0, 1 add \$v0, \$a1, \$v0 j exit c_false: emit _cost2 move \$v0, \$a2 exit:</pre>

**Related cases.** The two Clight boolean operations `&&` and `||` have a lazy semantics: depending on the evaluation of the first argument, the second one might be evaluated or not. There is an obvious translation to ternary expressions. For instance, the expression `x && y` is translated into the expression `x?(y?1:0):0`. Our compiler performs this translation *before* computing the labelling.

### 6.3.3 Conditionals

Conditionals are another way to introduce a branching. As for ternary expressions, the labelling of a conditional consists in adding a starting label to the labelling of each branch.

Clight	$\xrightarrow{\text{Labelling}}$	Labelled Clight	$\xrightarrow{\text{Compilation}}$	Labelled Mips
<pre>if (b) {   x = 1;   ... } else {   x = 2;   ... }</pre>		<pre>if (b) {   _cost1:   x = 1;   ... } else {   _cost2:   x = 2;   ... }</pre>		<pre>beq \$a0, \$zero, c_false emit _cost1 li \$v0, 1 ... j exit c_false: emit _cost2 li \$v0, 2 ... exit:</pre>

### 6.3.4 Loops

Loops in Clight are guarded by a condition. Following the arguments of the previous cases, we add two labels when encountering a loop construct: one label to start the loop's body, and one label when exiting the loop. This is enough to guarantee that the loop in the compiled code goes through a label.

Clight	$\xrightarrow{\text{Labelling}}$	Labelled Clight	$\xrightarrow{\text{Compilation}}$	Labelled Mips
<pre>while (b) {   i++;   ... } x = i;</pre>		<pre>while (b) {   _cost1:   i++;   ... } _cost2: x = i;</pre>		<pre>loop: beq \$a0, \$zero, exit emit _cost1 li \$v0, 1 add \$a1, \$a1, \$v0 ... j loop exit: emit _cost2 move \$a2, \$a1</pre>

### 6.3.5 Program Labels and Gotos

In Clight, program labels and gotos are intraprocedural. Their only effect on the control flow of the resulting assembly code is to potentially introduce an unguarded loop. This loop must contain at least one cost label in order to satisfy the soundness condition, which we ensure by adding a cost label right after a program label.

Clight	$\xrightarrow{\text{Labelling}}$	Labelled Clight	$\xrightarrow{\text{Compilation}}$	Labelled Mips
lbl: i++; ... goto lbl;		lbl: _cost: i++; ... goto lbl;		lbl: emit _cost li \$v0, 1 add \$a0, \$a0, \$v0 ... j lbl

### 6.3.6 Function calls

Function calls in Mips are performed by indirect jumps, the address of the callee being in a register. In the general case, this address cannot be inferred statically. Even though the destination point of a function call is unknown, when the considered Mips code has been produced by our compiler, we know for a fact that this function ends with a return statement that transfers the control back to the instruction following the function call in the caller. As a result, we treat function calls according to the following principles: (1) the instructions of a function are covered by the labels inside this function, (2) we assume a function call always returns and runs the instruction following the call.

Principle (1) entails in particular that each function must contain at least one label. To ensure this, we simply add a starting label in every function definition. The example below illustrates this point:

Clight	$\xrightarrow{\text{Labelling}}$	Labelled Clight	$\xrightarrow{\text{Compilation}}$	Labelled Mips
void f () { f's body }		void f () { _cost: f's body }		f_start: Frame Creation Initializations emit _cost f's body Frame Deletion return

We notice that some instructions in Mips will be inserted *before* the first label is emitted. These instructions relate to the frame creation and/or variable initializations, and are composed of sequential instructions (no branching). To deal with this issue, we take the convention that the instructions that precede the first label in a function code are actually under the scope of the first label.

Principle (2) is of course an over-approximation of the program behaviour as a function might fail to return because of an infinite loop. In this case, the proposed labelling remains correct: it just assumes that the instructions following the function call will be executed, and takes their cost into consideration. The final computed cost is still an over-approximation of the actual cost.

## 6.4 Verifications on the object code

The labelling previously described has been designed so that the compiled Mips code satisfies the soundness and precision conditions. However, we do not need to prove this, instead we have to devise an algorithm that checks the conditions on the compiled code. The algorithm assumes a correct management of function calls in the compiled code. In particular, when we call a function we always jump to the first instruction of the corresponding code segment and when we return we always jump to an instruction that follows a call. We stress that this is a reasonable hypothesis that is essentially subsumed by the proof that the object code *simulates* the source code.

In our current implementation, we check the soundness and the precision conditions while building at the same time the label-cost mapping. To this end, the algorithm takes the following main steps.

- First, for each function a control flow graph is built.
- For each graph, we check whether there is a unique label that is reachable from the root by a unique path. This unique path corresponds to the instructions generated by the calling conventions as discussed in section 6.3.6. We shift the occurrence of the label to the root of the graph.
- By a strongly connected components algorithm, we check whether every loop in the graphs goes through at least one label.
- We perform a (depth-first) search of the graph. Whenever we reach a labelled node, we perform a second (depth-first) search that stops at labelled nodes and computes an upper on the cost of the occurrence of the label. Of course, when crossing a branching instruction, we take the maximum cost of the branches. When the second search stops we update the current cost of the label-cost mapping (by taking a maximum) and we continue the first search.
- Warning messages are emitted whenever the maximum is taken between two different values as in this case the precision condition may be violated.

## 6.5 Building the cost annotation

Once the label-cost mapping is computed, instrumenting the labelled source code is an easy task. A fresh global variable which we call *cost variable* is added to the source program with the purpose of holding the cost value and it is initialised at the very beginning of the main program. Then, every label is replaced by an increment of the cost variable according to the label-cost mapping. Following this replacement, the cost labels disappear and the result is a Clight program with annotations in the form of assignments.

There is one final problem: labels inside expressions. As we already mentioned, Clight does not allow writing side-effect instructions — such as cost increments — inside expressions. To cope with this restriction, we produce first an instrumented C program — with side-effects in expressions — that we translate back to Clight using CIL. This process is summarized below.

$$\left. \begin{array}{l} \text{Labelled Clight} \\ \text{label-cost mapping} \end{array} \right\} \xrightarrow{\text{Instrumentation}} \text{Instrumented C} \xrightarrow{\text{CIL}} \text{Instrumented Clight}$$

## 6.6 Testing

It is desirable to test the coherence of the labelling from Clight to Mips. To this end, each labelled language comes with an interpreter that produces the trace of the labels encountered during the computation. Then, one naive approach is to test the equality of the traces produced by the program at the different stages of the compilation. Our current implementation passes this kind of tests. For some optimisations that may re-order computations, a weaker condition could be considered which consists in abstracting the traces as *multi-sets* of labels before comparing them.

## 7 Conclusion and future work

We have discussed the problem of building a compiler which can *lift* in a provably correct way pieces of information on the execution cost of the object code to cost annotations on the source code. To this end, we have introduced the so called *direct* and *labelling* approaches and discussed their formal application to a toy compiler. Based on this experience, we have argued that the second approach has better scalability properties. To substantiate this claim, we have reported on our successful experience in implementing and testing the labelling approach on top of a prototype compiler written in ocaml for a large fragment of the C language which can be shortly described as Clight without floating point.

We discuss next a few directions for future work. First, we plan to test the current compiler on the kind of C code produced for embedded applications, a typical example being the C code produced by the compilation of synchronous languages such as Lustre or Esterel. Starting from the annotated C code, we expect to produce automatically meaningful information on, say, the reaction time of a given synchronous program. Second, we plan to port the current compiler to commercial assembly languages. In particular, it would be interesting to target one of the assembly languages covered by the AbsInt tool so as to obtain more realistic estimations of the execution cost of sequences of instructions. Third, we plan to formalise and validate in the *Calculus of Inductive Constructions* the prototype implementation of the labelling approach for the C compiler described in section 5. This requires a major implementation effort which will be carried on in collaboration with our partners of the CerCo project [2]. Fourth, we plan to study the applicability of the labelling approach to other optimisation techniques in the realm of the C compilers technology such as *loop optimisations*, and to other languages which rely on rather distinct compilation technologies such as a language of the ML family.

## References

- [1] AbsInt Angewandte Informatik. <http://www.absint.com/>.
- [2] Certified Complexity (Project description). ICT-2007.8.0 FET Open, Grant 243881. <http://cerco.cs.unibo.it>.
- [3] Esterel Technologies. <http://www.esterel-technologies.com>.
- [4] Frama – C software analysers. <http://frama-c.com/>.
- [5] C. Ferdinand, R. Heckmann, T. Le Sergent, D. Lopes, B. Martin, X. Fornari, and F. Martin. Combining a high-level design tool for safety-critical systems with a tool for WCET analysis of executables. In Proc. of the 4th European Congress on *Embedded Real Time Software (ERTS)*, Toulouse, 2008.
- [6] X. Fornari. Understanding how SCADE suite KCG generates safe C code. White paper, Esterel Technologies, 2010.

- [7] J. Larus. Assemblers, linkers, and the SPIM simulator. Appendix of *Computer Organization and Design: the hardware/software interface*, book by Hennessy and Patterson, 2005.
- [8] X. Leroy. Formal verification of a realistic compiler. *Commun. ACM*, 52(7):107-115, 2009.
- [9] X. Leroy. Mechanized semantics, with applications to program proof and compiler verification (lecture notes and Coq development). *Marktoberdorf summer school*, Germany, 2009.
- [10] J. McCarthy and J. Painter. Correctness of a compiler for arithmetic expressions. In *Mathematical aspects of computer science 1*, volume 19 of Symposia in Applied Mathematics, AMS, 1967.
- [11] G. Necula, S. McPeak, S.P. Rahul, and W. Weimer. CIL: Intermediate Language and Tools for Analysis and Transformation of C Programs. In *Proceedings of Conference on Compiler Construction*, Springer LNCS 2304:213–228, 2002.
- [12] F. Pottier. Compilation (INF 564), École Polytechnique, 2009-2010. <http://www.enseignement.polytechnique.fr/informatique/INF564/>.
- [13] R. Wilhelm et al. The worst-case execution-time problem - overview of methods and survey of tools. *ACM Trans. Embedded Comput. Syst.*, 7(3), 2008.

## A Assessment of the deliverable within the *CerCo* project

Following the extensive experiments described in the previous sections, we now feel confident about the possibility of scaling our approach to a realistic, mildly optimizing, complexity preserving<sup>2</sup> compiler for a target architecture of the kind described in the proposal. While we believe that our approach should scale to a retargetable compiler, in the timeframe of the European Project and according to the project work-plan we only commit to the investigation of a compiler having a single target processor.

In particular, we will target the 8051 (also known as 8052 or MCS51) family of processors. The 8051 is an 8 bit CISC microprocessor introduced in the 1980 by Intel, very popular and still manufactured by a host of companies, many European. It is widely used in embedded systems and, thanks to its predictable behaviour and execution cost, it will allow us to compute fully accurate measurement of the actual computational complexity of  $O(1)$  assembly program slices<sup>3</sup>, to be manifested at the C level.

With respect to the test-cases previously described, in particular the C compiler for MIPS, the main difficulty introduced by the 8051 is the non uniform concrete memory model: the processor has different types of memory (on-chip RAM, external RAM, on-chip and/or external ROM) that can be accessed using different access modes and pointer types. Moreover, memory mapped I/O is heavily used in the design of the chip, to the point that all registers (apart from the inaccessible program counter) are seen as memory locations and have their own memory address. To complete the picture, the different memories are split into regions that may overlap, and the same accessing mode can point to different regions (or even to memory mapped registers) according to the value of the pointer. Finally, the amount of memory available is quite limited, with a stack which is at most 80 bytes wide and different speeds and opcode sizes to access different memory areas. For this reasons, compilers that target 8051 processors usually abound in directives to drive the tool in assigning memory locations to values.

Hence, because of the peculiar choice of target processor, in the next six months of the *CerCo* project and partly as an extension of the original work-plan we will:

- Extend the memory model used so far (and taken from CompCert) to accurately describe the different memory types and regions of the 8051, as well as to obey to the `volatile` directive used to map memory mapped regions to program variables. This work will be done as part of Tasks T2.3 and T4.1 and the outcome will be described in Deliverables D2.2 and D4.1 where we will also provide a formalization in Matita of the executable semantics of the extended memory model. In case we decide to adopt the same memory model for all compilation phases, as it is done in CompCert, the memory model extension will also span over Task T3.1, requiring close cooperation between the front-end formalization (lead by Edinburgh) and the back-end formalization (lead by Bologna).
- Devise language extensions (possibly inspired by the variable modifiers of the SDCC compiler) to let the user suggest or force the compiler to put data into particular memory regions. Similarly, we could refine the pointer types of ANSI C into classes of pointer types pointing to particular regions, in order to reflect the difference in sizes of pointers

---

<sup>2</sup>In the sense of the project proposal.

<sup>3</sup>In the sense of the project proposal.

to different regions (8 bits to address internal RAM and memory mapped registers; 8 bits to address bits in the bit memory or in registers using an ad-hoc addressing mode; 16 bits for code memory and external memory; 24 bits for generic pointers). All of these language extensions, if any, must be reflected all over the compilation chain, with modifications to the intermediate languages and/or memory model. This work will be done as part of Tasks T2.3, T3.1 and T4.1 and the outcome will be described in Deliverables D2.2, D3.1, D3.3, D4.1, D4.3 where we will provide executable semantics in Matita of the source, target and intermediate languages.

In order to remain compatible with the project schedule, we will first focus on compiling standard Clight without floating points to the 8051, adding the language extensions in a second moment, within Task T2.3.

- Modify and extend the experimental compiler from Clight (without floating point) to MIPS in order to target the 8051 architecture (Task T2.3, Deliverable D2.2).

But for the modifications required by the selection of the target processor (8051), we plan to rely as much as possible on the architecture and the the intermediate languages described in this deliverable D2.1.

## B Proofs

### B.1 Notation

Let  $\xrightarrow{t}$  be a family of reduction relations where  $t$  ranges over the set of labels and  $\epsilon$ . Then we define:

$$\xrightarrow{t} = \begin{cases} (\xrightarrow{\epsilon})^* & \text{if } t = \epsilon \\ (\xrightarrow{\epsilon})^* \circ \xrightarrow{t} \circ (\xrightarrow{\epsilon})^* & \text{otherwise} \end{cases}$$

where as usual  $R^*$  denote the reflexive and transitive closure of the relation  $R$  and  $\circ$  denotes the composition of relations.

### B.2 Proof of proposition 4

The following properties are useful.

**Lemma 26** (1) *The relation  $\xrightarrow{C}$  is transitive.*

(2) *If  $i \xrightarrow{C} j$  and  $R(C, j, K)$  then  $R(C, i, K)$ .*

The first property can be proven by induction on the definition of  $\xrightarrow{C}$  and the second by induction on the structure of  $K$ . Next we can focus on the proposition. The notation  $C \dot{i} C'$  means that  $i = |C|$ . Suppose that:

$$(S, K, s) \rightarrow (S', K', s') \quad (1) \quad \text{and} \quad R(C, i, S \cdot K) \quad (2).$$

From (2), we know that there exist  $i'$  and  $i''$  such that:

$$i \xrightarrow{C} i' \quad (3), \quad C = C_1 \dot{i}' C(S) \dot{i}'' C_2 \quad (4), \quad \text{and} \quad R(C, i'', K) \quad (5)$$

and from (3) it follows that:

$$C \vdash (i, \sigma, s) \xrightarrow{*} (i', \sigma, s) \quad (3').$$

We are looking for  $j$  such that:

$$C \vdash (i, \sigma, s) \xrightarrow{*} (j, \sigma, s') \quad (6), \quad \text{and} \quad R(C, j, S' \cdot K') \quad (7).$$

We proceed by case analysis on  $S$ . We just detail the case of the conditional command as the the remaining cases have similar proofs. If  $S = \text{if } e_1 < e_2 \text{ then } S_1 \text{ else } S_2$  then (4) is rewritten as follows:

$$C = C_1 \dot{i}' C(e_1) \cdot C(e_2) \cdot \text{bge}(k_1) \dot{a} C(S_1) \dot{b} \text{branch}(k_2) \dot{c} C(S_2) \dot{i}'' C_2$$

where  $c = a + k_1$  and  $i'' = c + k_2$ . We distinguish two cases according to the evaluation of the boolean condition. We describe the case  $(e_1 < e_2) \Downarrow \text{true}$ . We set  $j = a$ .

- The instance of (1) is  $(S, K, s) \rightarrow (S_1, K, s)$ .
- The reduction required in (6) takes the form  $C \vdash (i, \sigma, s) \xrightarrow{*} (i', \sigma, s) \xrightarrow{*} (a, \sigma, s')$ , and it follows from (3'), the fact that  $(e_1 < e_2) \Downarrow \text{true}$ , and proposition 3(2).
- Property (7), follows from lemma 26(2), fact (5), and the following proof tree:

$$\frac{j \xrightarrow{C} j \quad \frac{b \xrightarrow{C} i'' \quad R(C, i'', K)}{R(C, b, K)}}{R(C, j, S_1 \cdot K')} .$$

□



### B.3 Proof of proposition 8

We recall that the compiled code  $\mathcal{C}'(C)$  does not read or write the locations  $l_{cost}$ ,  $l_A$ , and  $l_B$ . Then we note the following properties.

- (a) If  $An_{Vm}(C) \vdash (5 \cdot i, \sigma, s) \rightarrow^5 (5 \cdot j, \sigma', s')$  then  $C \vdash (i, \sigma, s) \rightarrow (j, \sigma', s'[s(cost)/cost])$ .  
 (b) If  $M \vdash (i, m) \rightarrow (j, m')$  then  $An_{Mips}(M) \vdash (9 \cdot i, m) \xrightarrow{*} (9 \cdot j, m'[m(l_{cost}) + 1/l_{cost}, m(A)/l_A, m(B)/l_B])$ .  
 Using (a) and the hypothesis, we derive:

$$C \vdash (i, \sigma, s) \rightarrow (j, \sigma', s'[s(cost)/cost])$$

Then, by proposition 6 and the hypothesis on  $m$  we derive:

$$\mathcal{C}'(C) \vdash (p(i, C), m) \xrightarrow{*} (p(j, C), m') \text{ and } m' \Vdash \sigma', s'[s(cost)/cost] \quad (7)$$

We perform a case analysis on the instruction  $C[i]$  and  $h(i)$  to derive:

$$An_{Mips}(\mathcal{C}'(C)) \vdash (9 \cdot p(i, C), m) \xrightarrow{*} (9 \cdot p(j, C), m'[m(l_{cost}) + d(i, C)/l_{cost}, x/l_A, y/l_B])$$

where  $x$  and  $y$  are respectively the value of  $l_A$  and  $l_B$  in the last intermediate configuration.

**The memory realisation.** The final memory state is:  $m'[m(l_{cost}) + d(i, C)/l_{cost}, x/l_A, y/l_B]$ . Then we have to verify:

$$(m'[m(l_{cost}) + d(i, C)/l_{cost}, x/l_A, y/l_B])[s'(cost)/l_{cost}] \Vdash \sigma', s'.$$

This is equivalent to:

$$m'[s'(cost)/l_{cost}, x/l_A, y/l_B] \Vdash \sigma', s'.$$

Since the realisation predicate does not depend on the locations  $l_A$  and  $l_B$ , we can conclude using the second part of (7).

**The inequation.** To conclude, we verify the inequation:

$$m'[m(l_{cost}) + d(i, C)/l_{cost}, x/l_A, y/l_B](l_{cost}) - m(l_{cost}) \leq s'(cost) - s(cost).$$

By rewriting, we get:  $d(i, C) \leq s'(cost) - s(cost)$ . This follows by case analysis of  $C[i]$  and  $h(i)$  since, by definition,  $An_{Vm}$  uses the worst possible value of  $d(i, C)$ .  $\square$

### B.4 Proof of proposition 10

We have  $P = \mathbf{prog} S$  and by definition,  $An_{Imp}(P) = cost := cost + \kappa(S); An_{Imp}(S)$ . With our simulation hypothesis, we derive  $(An_{Imp}(S), s[\kappa(S)/cost]) \Downarrow s'[c'/cost]$ . Using the proposition 9 with  $C = C' = \epsilon$ ,  $d = 0$  and  $\sigma = \epsilon$ , we have:

$$An_{Vm}(\mathcal{C}(S)) \vdash (0, \epsilon, s[0/cost]) \xrightarrow{*} (|An_{Vm}(\mathcal{C}(P))|, \epsilon, s'[d'/cost]),$$

where  $d' \leq (c' - \kappa(S)) + \kappa(S) = c'$ . By definition,  $\mathcal{C}(P) = \mathcal{C}(S) \cdot \mathbf{halt}$  and  $|An_{Vm}(\mathcal{C}(P))| = 5 \cdot |\mathcal{C}(P)|$ . We can rewrite the simulation:

$$An_{Vm}(\mathcal{C}(P)) \vdash (0, \epsilon, s[0/cost]) \xrightarrow{*} (5 \cdot |\mathcal{C}(P)|, \epsilon, s'[d'/cost]).$$

By proposition 5, there is a decoration  $h$  such that  $\mathcal{C}(P) : h$  and  $h(0) = 0$ . By iterating proposition 8, we derive:

$$An_{Mips}(\mathcal{C}'(\mathcal{C}(P))) \vdash (9 \cdot p(0, \mathcal{C}(P)), m) \xrightarrow{*} (9 \cdot p(|\mathcal{C}(P)|, \mathcal{C}(P)), m'),$$

with  $m'[d'/l_{cost}] \Vdash \epsilon$ ,  $s'[d'/cost]$  and  $m'(l_{cost}) \leq d'$ . We know that the last instruction of  $\mathcal{C}(P)$  is an **halt**, therefore the last instruction of  $\mathcal{C}'(\mathcal{C}(P))$  is also an **halt**. By definition,  $p(|\mathcal{C}(P)|, \mathcal{C}(P))$  is the position after this instruction. That gives us:

$$(An_{\text{Mips}}(\mathcal{C}'(\mathcal{C}(P))), m) \Downarrow m' .$$

And by transitivity we have:  $m'(l_{cost}) \leq d' \leq c'$ .  $\square$

## B.5 Proof of proposition 11

- (1) By induction on the structure of the command  $S$ .
- (2) By iterating the following proposition.

**Proposition 27** *If  $(S, K, s) \xrightarrow{t} (S', K', s')$  and  $R(C, i, S \cdot K)$  with  $t = \ell$  or  $t = \epsilon$  then  $C \vdash (i, \sigma, s) \xrightarrow{t} (j, \sigma, s')$  and  $R(C, j, S' \cdot K')$ .*

This is an extension of proposition 4 and it is proven in the same way with an additional case for labelled commands.  $\square$

## B.6 Proof of proposition 13

- (1) The compilation of the **Vm** instruction  $\text{nop}(\ell)$  is the **Mips** instruction  $(\text{nop } \ell)$ .
- (2) By iterating the following proposition.

**Proposition 28** *Let  $C : h$  be a well formed code. If  $C \vdash (i, \sigma, s) \xrightarrow{t} (j, \sigma', s')$  with  $t = \ell$  or  $t = \epsilon$ ,  $h(i) = |\sigma|$  and  $m \Vdash \sigma, s$  then  $\mathcal{C}'(C) \vdash (p(i, C), m) \xrightarrow{t} (p(j, C), m')$  and  $m' \Vdash \sigma', s'$ .*

This is an extension of proposition 6 and it is proven in the same way with an additional case for the **nop** instruction.  $\square$

## B.7 Proof of proposition 14

In order to carry on the proof, one needs to develop a bit more the properties of the small-step operational semantics. In particular, one needs to show that a continuation such as  $S \cdot (S' \cdot K)$  is ‘observationally equivalent’ to the continuation  $(S; S') \cdot K$ .  $\square$

## B.8 Proof of proposition 16

By diagram chasing using propositions 11(1), 13(1), and the definition 15 of labelling.  $\square$

## B.9 Proof of proposition 17

Suppose that:

$$(\mathcal{I}(\mathcal{L}(P)), s[c/cost]) \Downarrow s'[c + \delta/cost] \text{ and } m \Vdash s[c/cost] .$$

Then, by proposition 14, for some  $\lambda$ :

$$(\mathcal{L}(P), s[c/cost]) \Downarrow (s'[c/cost], \lambda) \text{ and } \kappa(\lambda) = \delta .$$

Finally, by propositions 11(2) and 13(2) :

$$(\mathcal{C}'(\mathcal{C}(\mathcal{L}(P))), m) \Downarrow (m', \lambda) \text{ and } m' \Vdash s'[c/cost] .$$

$\square$

### B.10 Proof of proposition 20

If  $\lambda = \ell_1 \cdots \ell_n$  then the computation is the concatenation of simple paths labelled with  $\ell_1, \dots, \ell_n$ . Since  $\kappa(\ell_i)$  bounds the cost of a simple path labelled with  $\ell_i$ , the cost of the overall computation is bounded by  $\kappa(\lambda) = \kappa(\ell_1) + \cdots \kappa(\ell_n)$ .  $\square$

### B.11 Proof of proposition 22

Same proof as proposition 20, by replacing the word *bounds* by *is exactly* and the words *bounded by* by *exactly*.  $\square$

### B.12 Proof of proposition 23

In both labellings under consideration the root node is labelled. An obvious observation is that only commands of the shape *while b do S* and *while lb do S* introduce loops in the compiled code. In the second case, the compilation ensures that a label is placed in the loop. In the first case, we notice that both labelling introduce a label in the loop (though at different places). Thus all loops go through a label and the compiled code is always sound.

To show the precision of the second labelling  $\mathcal{L}_p$ , we note the following property.

**Lemma 29** *A soundly labelled graph is precise if each label occurs at most once in the graph and if the immediate successors of the bge nodes are either halt (no successor) or labelled nodes.*

Indeed, in a such a graph starting from a labelled node we can follow a unique path up to a leaf, another labelled node, or a bge node. In the last case, the hypotheses in the lemma 29 guarantee that the two simple paths one can follow from the bge node have the same length/cost.  $\square$

### B.13 Proof of proposition 25

By applying consecutively proposition 17 and propositions 20 or 22.  $\square$