



Certifying cost annotations in compilers

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

► **To cite this version:**

Roberto M. Amadio, Nicolas Ayache, Yann Régis-Gianas, Ronan Saillard. Certifying cost annotations in compilers. 2010. <hal-00524715>

HAL Id: hal-00524715

<https://hal.archives-ouvertes.fr/hal-00524715>

Submitted on 8 Oct 2010

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Certifying cost annotations in compilers*

Roberto M. Amadio⁽¹⁾ Nicolas Ayache⁽²⁾
Yann Régis-Gianas⁽²⁾ Ronan Saillard⁽²⁾

⁽¹⁾ Université Paris Diderot (UMR-CNRS 7126)

⁽²⁾ Université Paris Diderot (UMR-CNRS 7126) and INRIA (Team πr^2)

October 8, 2010

Abstract

We discuss 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 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 a so-called *labelling* approach to these three questions. As a first step, we examine its application to a toy compiler. This formal study suggests that the labelling approach 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.

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 optimizing compiler from a large subset of the C language to a typical assembly language of the kind used in embedded systems [9].

In the framework of the *Certified Complexity* (CerCo) project [3], 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 [4, 7] is that the *reaction time* of the program is estimated by means of abstract interpretation methods (such as those developed by AbsInt [1, 6]) 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

*This work was supported by the *Information and Communication Technologies (ICT) Programme* as Project FP7-ICT-2009-C-243881 CerCo.

determine the number of times a loop is iterated (see, *e.g.*, [15] 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 [5]. 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 a methodology addressing these three questions and to consider its concrete application to a simple toy compiler and to a moderately optimizing untrusted C compiler.

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.

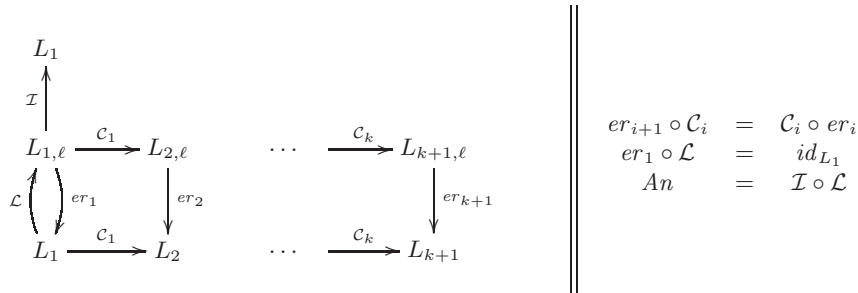
Thus producing a cost-annotation of a source program P amounts to build an *annotated program* $An(P)$ which behaves as P while self-monitoring its execution cost. In particular, if we do *not* observe the cost variables then we expect the annotated program $An(P)$ to be functionally equivalent to P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

Soundness and precision of cost annotations Suppose we have a functionally correct compiler \mathcal{C} that associates with a program P in the source language a program $\mathcal{C}(P)$ 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 $An(P)$ 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* with respect to the cost model if the *difference* between the predicted and real execution costs is bounded by a constant which depends on the program.

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. A previous experience described in [2] suggests that the process of pushing ‘hidden parameters’ in the definitions of cost annotations and of manipulating directly numerical cost is error prone and produces complex proofs. For this reason, we advocate next a ‘labelling approach’ where costs are handled at an abstract level and numerical values are produced at the very end of the construction.

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 \mathcal{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 $\text{er}_{L_1} \circ \mathcal{L}$ 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 global *cost* variable. Then an *annotation* An of the source program can be derived simply as the composition of the labelling and the instrumentation functions: $An = \mathcal{I} \circ \mathcal{L}$.

Suppose s is some adequate representation of the state of a program. Let P be a source program and suppose that its annotation satisfies the following property:

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

where c and δ are some non-negative numbers. 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

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

where $\mathcal{C} = \mathcal{C}_k \circ \dots \circ \mathcal{C}_1$ and λ is a sequence (or a multi-set) of labels whose ‘cost’ corresponds to the number δ 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 $er_{k+1}(\mathcal{C}(\mathcal{L}(P)))$ is actually equal to the compiled code $\mathcal{C}(P)$ we are interested in. Given this, the following question arises: under which conditions the sequence λ , *i.e.*, the increment δ , 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 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.

A toy compiler As a first case study for the labelling approach 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.

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 [9]. Notable differences are that some compilation steps are fused, 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 [13]. 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 [14]. The main optimisations it performs are common subexpression elimination, liveness analysis and register allocation, and graph compression. We ran some

benchmarks to ensure that our prototype implementation is realistic. The results are given in appendix B.9 and the compiler is available from the authors.

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 labelling approach to the toy compiler. Section 4 reports our experience in implementing and testing the labelling approach on the C compiler. Section 5 summarizes our contribution and outlines some perspectives for future work. Appendix A sketches the proofs that have not been mechanically checked in Coq and appendix B provides some details on the structure of the C compiler we have implemented.

2 A toy compiler

We formalise the toy compiler introduced in section 1.

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 $x \neq y$. The *big-step* operational semantics of Imp expressions and boolean conditions is defined as follows:

$$\frac{}{(v, s) \Downarrow v} \quad \frac{}{(x, s) \Downarrow s(x)} \quad \frac{(e, s) \Downarrow v \quad (e', s) \Downarrow v'}{(e + e', s) \Downarrow (v + \mathbf{z} v')} \quad \frac{(e, s) \Downarrow v \quad (e', s) \Downarrow v'}{(e < e', s) \Downarrow (v < \mathbf{z} v')}$$

A *continuation* K is a list of commands which terminates with a special symbol **halt**: $K ::= \text{halt} \mid S \cdot K$. Table 1 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: $(S, s) \Downarrow s'$ if $(S, \text{halt}, s) \rightarrow \dots \rightarrow (\text{skip}, \text{halt}, s')$.

2.2 Vm: language and semantics

Following [10], 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 σ .

Given a sequence C , we denote with $|C|$ its length and with $C[i]$ its i^{th} element (the leftmost element being the 0^{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

$(x := e, K, s)$	\rightarrow	$(\text{skip}, K, s[v/x])$	$\text{if } (e, s) \Downarrow v$
$(S; S', K, s)$	\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)	

Table 1: Small-step operational semantics of `Imp` commands

Rule	$C[i] =$
$C \vdash (i, \sigma, s) \rightarrow (i + 1, n \cdot \sigma, s)$	<code>cnst(n)</code>
$C \vdash (i, \sigma, s) \rightarrow (i + 1, s(x) \cdot \sigma, s)$	<code>var(x)</code>
$C \vdash (i, n \cdot \sigma, s) \rightarrow (i + 1, \sigma, s[n/x])$	<code>setvar(x)</code>
$C \vdash (i, n \cdot n' \cdot \sigma, s) \rightarrow (i + 1, (n +_{\mathbf{Z}} n') \cdot \sigma, s)$	<code>add</code>
$C \vdash (i, \sigma, s) \rightarrow (i + k + 1, \sigma, s)$	<code>branch(k)</code>
$C \vdash (i, n \cdot n' \cdot \sigma, s) \rightarrow (i + 1, \sigma, s)$	<code>bge(k)</code> and $n <_{\mathbf{Z}} n'$
$C \vdash (i, n \cdot n' \cdot \sigma, s) \rightarrow (i + k + 1, \sigma, s)$	<code>bge(k)</code> and $n \geq_{\mathbf{Z}} n'$

Table 2: Operational semantics `Vm` programs

2. 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 3 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

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

Table 3: Conditions for well-formed code

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

Table 4: Compilation from `Imp` to `Vm`

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 $h = h'$.

2.3 Compilation from `Imp` to `Vm`

In table 4, we define compilation functions \mathcal{C} from `Imp` to `Vm` which operate on expressions, boolean conditions, statements, and programs. We write $\mathit{sz}(e)$, $\mathit{sz}(b)$, $\mathit{sz}(S)$ for the number of instructions the compilation function associates with the expression e , the boolean condition b , and the statement S , respectively.

We follow [10] for the proof of soundness of the compilation function for expressions and boolean conditions (see also [11] 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 $(b, s) \Downarrow \mathbf{true}$ 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 $(b, s) \Downarrow \mathbf{false}$ then $C \cdot \mathcal{C}(b, k) \cdot C' \vdash (i, \sigma, s) \xrightarrow{*} (j, \sigma, s)$ where $i = |C|$ and $j = |C \cdot \mathcal{C}(b, k)|$.

Next we focus on the compilation of statements. We introduce a ternary relation $R(C, i, K)$ which relates a `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 . (A formal definition is available in appendix 4.) 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 $R(C, j, S' \cdot K')$.*

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

Proposition 5 *For any program P there is a unique h such that $\mathcal{C}(P) : h$.*

Rule	$M[i] =$
$M \vdash (i, m) \rightarrow (i + 1, m[n/R])$	loadi R, n
$M \vdash (i, m) \rightarrow (i + 1, m[l/R])$	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 5: Operational semantics Mips programs

2.4 Mips: language and semantics

We consider a Mips-like machine [8] 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 denote with M a list of instructions. The operational semantics is formalised in table 5 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 $M[j] = \text{halt}$.

2.5 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), and (3) we store the first b elements of the stack σ 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 σ and the store s , and write $m \Vdash \sigma, s$, if the following conditions are satisfied: (1) $s(x) = m(l_x)$, and (2) if $0 \leq i < |\sigma|$ then $\sigma[i] = m(R_i)$ if $i < b$, and $\sigma[i] = m(l_i)$ if $i \geq b$.

The compilation function \mathcal{C}' from Vm to Mips is described in table 6. 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 $\mathcal{C}'(0, C) \dots \mathcal{C}'(|C| - 1, C)$. 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 as¹ $p(i, C) = \sum_{0 \leq j < i} d(j, C)$, where the function $d(i, C)$ is defined as $d(i, C) = |\mathcal{C}'(i, C)|$. Hence $d(i, C)$ is the number of Mips instructions associated with the i^{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 $m \Vdash \sigma, s$ then $\mathcal{C}'(C) \vdash (p(i, C), m) \xrightarrow{*} (p(j, C), m')$ and $m' \Vdash \sigma', s'$.*

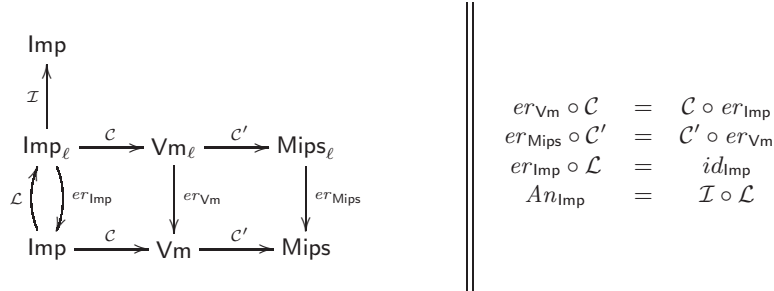
¹There is an obvious circularity in this definition that can be easily eliminated by defining first the function d following the case analysis in table 6, then the function p , and finally the function \mathcal{C}' as in table 6.

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

Table 6: Compilation from Vm to Mips

3 Labelling approach for the toy compiler

We apply the labelling approach introduced in section 1 to the toy compiler which results in the following diagram.



3.1 Labelled Imp

We extend the syntax so that statements can be labelled: $S ::= \dots \mid \ell : S$. For instance, $\ell : (\text{while } (n < x) \text{ do } \ell : S)$ is a labelled command. The small step semantics of statements defined in table 1 is extended as follows.

$$(\ell : S, K, s) \xrightarrow{\ell} (S, K, s)$$

We denote with λ, λ', \dots finite sequences of labels. In particular, we denote with ϵ the empty sequence and identify an unlabelled transition with a transition labelled with ϵ . Then the small step reduction relation we have defined on statements becomes a *labelled transition system*. There is an obvious *erasure* function er_{Imp} from the labelled language to the unlabelled one which is the identity on expressions and boolean conditions, and traverses commands removing all labels. We derive a *labelled* big-step semantics as follows: $(S, s) \Downarrow (s', \lambda)$ if $(S, \text{halt}, s) \xrightarrow{\lambda_1} \dots \xrightarrow{\lambda_n} (\text{skip}, \text{halt}, s')$ and $\lambda = \lambda_1 \cdots \lambda_n$.

3.2 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_{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 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 Imp_ℓ 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 7 *For all commands S in Imp_ℓ 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 8 *In the current formulation, a sequence of transitions λ 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 λ is a multi-set of labels.*

3.3 Labelled Mips

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

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

The evaluation predicate for labelled Mips is defined as $(M, m) \Downarrow (m', \lambda)$ if $M \vdash (0, m) \xrightarrow{\lambda_1} \dots \xrightarrow{\lambda_n} (j, m')$, $\lambda = \lambda_1 \dots \lambda_n$ and $M[j] = \text{halt}$. The following proposition relates Vm_ℓ code and its compilation and it is similar to proposition 7.

Proposition 9 *Let C be a Vm_ℓ 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 $m' \Vdash -\epsilon, s'$.*

3.4 Labellings and instrumentations

Assuming a function κ 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 κ and $cost$) as follows:

$$\mathcal{I}(\ell : S) = inc(\ell); \mathcal{I}(S) .$$

The function \mathcal{I} just distributes over the other operators of the language. We extend the function κ 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 as follows.

Proposition 10 *Let S be an Imp_ℓ 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 11 *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 12 *For any labelling function \mathcal{L} , and `Imp` program P , the following holds:*

$$er_{\text{Mips}}(\mathcal{C}'(\mathcal{C}(\mathcal{L}(P)))) = \mathcal{C}'(\mathcal{C}(P)) . \quad (3)$$

Proposition 13 *Given a function κ 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 $m \Vdash -\epsilon, s[c/cost]$ then $(\mathcal{C}'(\mathcal{C}(\mathcal{L}(P))), m) \Downarrow (m', \lambda)$, $m' \Vdash -\epsilon, s'[c/cost]$ and $\kappa(\lambda) = \delta$.*

3.5 Sound and precise labellings

With any Mips_ℓ 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 `nop` ℓ .

Definition 14 *A simple path in a Mips_ℓ 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 15 *A Mips_ℓ 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 `Mips` code M , we can associate with every label ℓ a number $\kappa(\ell)$ which is the maximum (estimated) cost of executing a simple path whose first node is labelled with ℓ . We stress that in the following we assume that the cost of a simple path is proportional to the number of `Mips` instructions that are crossed in the path.

Proposition 16 *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 `Mips` code the sequence of labels associated with a computation is a significant information on the execution cost.

$\mathcal{L}_s(\text{prog } S)$	$= \text{prog } \ell : \mathcal{L}_s(S)$
$\mathcal{L}_s(\text{skip})$	$= \text{skip}$
$\mathcal{L}_s(x := e)$	$= x := e$
$\mathcal{L}_s(S; S')$	$= \mathcal{L}_s(S); \mathcal{L}_s(S')$
$\mathcal{L}_s(\text{if } b \text{ then } S_1 \text{ else } S_2)$	$= \text{if } b \text{ then } \mathcal{L}_s(S_1) \text{ else } \mathcal{L}_s(S_2)$
$\mathcal{L}_s(\text{while } b \text{ do } S)$	$= \text{while } b \text{ do } \ell : \mathcal{L}_s(S)$
$\mathcal{L}_p(\text{prog } S)$	$= \text{prog } \mathcal{L}_p(S)$
$\mathcal{L}_p(S)$	$= \text{let } \ell = \text{new}, (S', d) = \mathcal{L}'_p(S) \text{ in } \ell : S'$
$\mathcal{L}'_p(S)$	$= (S, 0) \text{ if } S = \text{skip} \text{ or } S = (x := e)$
$\mathcal{L}'_p(\text{if } b \text{ then } S_1 \text{ else } S_2)$	$= (\text{if } b \text{ then } \mathcal{L}_p(S_1) \text{ else } \mathcal{L}_p(S_2), 1)$
$\mathcal{L}'_p(\text{while } b \text{ do } S)$	$= (\text{while } b \text{ do } \mathcal{L}_p(S), 1)$
$\mathcal{L}'_p(S_1; S_2)$	$= \text{let } (S'_1, d_1) = \mathcal{L}'_p(S_1), (S'_2, d_2) = \mathcal{L}'_p(S_2) \text{ in}$
	$\text{case } d_1$
	$0 : (S'_1; S'_2, d_2)$
	$1 : \text{let } \ell = \text{new} \text{ in } (S'_1; \ell : S'_2, d_2)$

Table 7: Two labellings for the Imp language

Definition 17 We say that a soundly labelled code is precise if for every label ℓ in the code, the simple paths starting from a node labelled with ℓ have the same cost.

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

Proposition 18 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 Mips code. To discuss this point, we introduce in table 7 two labelling functions \mathcal{L}_s and \mathcal{L}_p for the Imp language. The first labelling relies on just one label while the second one relies on a function “new” which is meant to return fresh labels and on an auxiliary function \mathcal{L}'_p which returns a labelled command and a binary directive $d \in \{0, 1\}$. If $d = 1$ then the command that follows (if any) must be labelled.

Proposition 19 For all Imp programs P :

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

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 $\ell : (\text{if } 0 < x \text{ then } x := x + 1 \text{ else skip})$. In the compiled code, we find two simple paths associated with the label ℓ 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.

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 8: Building the annotation of a Clight program in the labelling approach

Definition 20 *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:*

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

Proposition 21 *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 $m \Vdash -\epsilon, s[c/\text{cost}]$ entails that $(\mathcal{C}'(\mathcal{C}(P)), m) \Downarrow m', m' \Vdash -\epsilon, s'[c/\text{cost}]$ and the cost of the execution is bound (is exactly) δ .*

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, it enough to label commands while it is not necessary to label boolean conditions and expressions.

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 modern 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.

4 Labelling approach for the C compiler

This section informally describes the labelled extensions of the languages in the compilation chain (see appendix B for details), 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 8 and is detailed in the following sections.

4.1 Labelled languages

Both the Clight and Cminor languages are extended in the same way by labelling both statements and expressions (by comparison, in the toy language `Imp` we just labelled statements). 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 those of Cminor.

The labelled versions of RTLabs and the languages in the back-end simply consist in adding a new instruction whose semantics is to emit a label without modifying the state. For the CFG based languages (RTLabs 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 4.2).

4.2 Labelling of the source language

As for the toy compiler (cf. end of section 3), the goals of a labelling are soundness, precision, and possibly economy. We explain our labelling by considering the constructions of Clight and their compilation to Mips.

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.

Ternary expressions and conditionals Most Clight expressions compile to sequential Mips code. *Ternary expressions*, that introduce a branching in the control flow, are one exception. In this case, we achieve precision by associating a label with each branch. This is similar to the treatment of the conditional we have already discussed in section 3. As for the Clight operations `&&` and `||` which have a lazy semantics, they are transformed to ternary expressions *before* computing the labelling.

Loops Loops in Clight are guarded by a condition. Following the arguments for 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 similar to the treatment of `while` loops discussed in section 3 and it is enough to guarantee that the loop in the compiled code goes through a label.

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
<pre> lbl: i++; ... goto lbl; </pre>		<pre> lbl: _cost: i++; ... goto lbl; </pre>		<pre> lbl: emit _cost li \$v0, 1 add \$a0, \$a0, \$v0 ... j lbl </pre>

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 global invariants of the compilation: (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. Invariant (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
<pre> void f () { f's body } </pre>		<pre> void f () { _cost: f's body } </pre>		<pre> f_start: Frame Creation Initializations emit _cost f's body Frame Deletion return </pre>

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. Invariant (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.

4.3 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 4.2. 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 bound 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.

4.4 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.

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

4.5 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, the weaker condition mentioned in remark 8 could be considered.

5 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 *labelling* approach and discussed its formal application to a toy compiler. Based on this experience, we have argued that the approach has good scalability properties, and 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 are currently testing the current compiler on the kind of C code produced for embedded applications by a `Lustre` compiler. Starting from the annotated C code, we are relying on the `Frama - C` tool to produce automatically meaningful information on, say, the reaction time of a given synchronous program. Second, we are porting the current compiler to other assembly languages. In particular, we are interested in targeting 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 B. This requires a major implementation effort which will be carried on in collaboration with our partners of the `CerCo` project [3].

References

- [1] AbsInt Angewandte Informatik. <http://www.absint.com/>.
- [2] R.M. Amadio, N. Ayache, K. Memarian, R. Saillard, Y. Régis-Gianas. Compiler Design and Intermediate Languages. Deliverable 2.1 of [3].
- [3] Certified Complexity (Project description). ICT-2007.8.0 FET Open, Grant 243881. <http://cerco.cs.unibo.it>.
- [4] Esterel Technologies. <http://www.esterel-technologies.com>.
- [5] Frama - C software analysers. <http://frama-c.com/>.
- [6] 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 *Embedded Real Time Software (ERTS)*, 2008.
- [7] X. Fornari. Understanding how SCADE suite KCG generates safe C code. White paper, Esterel Technologies, 2010.
- [8] J. Larus. Assemblers, linkers, and the SPIM simulator. Appendix of *Computer Organization and Design: the hw/sw interface*, by Hennessy and Patterson, 2005.
- [9] X. Leroy. Formal verification of a realistic compiler. *Commun. ACM*, 52(7):107-115, 2009.
- [10] X. Leroy. Mechanized semantics, with applications to program proof and compiler verification. *Marktoberdorf summer school*, 2009.
- [11] J. McCarthy and J. Painter. Correctness of a compiler for arithmetic expressions. In *Math. aspects of Comp. Sci. 1*, vol. 19 of Symp. in Appl. Math., AMS, 1967.
- [12] K. Memarian. Complexité Certifiée. *Travail d'étude et de recherche*, Master Informatique, Université Paris Diderot, 2010. <http://www.pps.jussieu.fr/~yrg/miniCerCo/>
- [13] 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.
- [14] F. Pottier. Compilation (INF 564), École Polytechnique, 2009-2010. <http://www.enseignement.polytechnique.fr/informatique/INF564/>.
- [15] 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 Proofs

We omit the proofs that have been checked by K. Memarian with the Coq proof assistant [12].

A.1 Notation

Let \xrightarrow{t} be a family of reduction relations where t ranges over the set of labels and ϵ . 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.

A.2 Proof of proposition 4

Given a Vm code C , we define an ‘accessibility relation’ \xrightarrow{C} as the least binary relation on $\{0, \dots, |C| - 1\}$ such that:

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

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

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

The following properties are useful.

Lemma 22 (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 $C \overset{i}{\cdot} C'$ 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:

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

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:

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

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 $C \vdash (i, \sigma, s) \overset{*}{\rightarrow} (i', \sigma, s) \overset{*}{\rightarrow} (a, \sigma, s')$, and it follows from (3'), the fact that $(e_1 < e_2) \Downarrow \text{true}$, and proposition 3(2).
- Property (7), follows from lemma 22(2), fact (5), and the following proof tree:

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

□

A.3 Proof of proposition 5

We actually prove that 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 $\mathcal{C}(P) : h$.

A.4 Proof of proposition 7

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

Proposition 23 *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. □

A.5 Proof of proposition 9

- (1) The compilation of the `Vm` instruction `nop(ℓ)` is the `Mips` instruction `(nop ℓ)`.
- (2) By iterating the following proposition.

Proposition 24 *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. □

A.6 Proof of proposition 10

We extend the instrumentation to the continuations by defining:

$$\mathcal{I}(S \cdot K) = \mathcal{I}(S) \cdot \mathcal{I}(K) \quad \mathcal{I}(\text{halt}) = \text{halt} .$$

Then we examine the possible reductions of a configuration $(\mathcal{I}(S), \mathcal{I}(K), s[c/cost])$.

- If S is an unlabelled statement such as `while b do S'` then $\mathcal{I}(S) = \text{while } b \text{ do } \mathcal{I}(S')$ and assuming $(b, s) \Downarrow \text{true}$ the reduction step is:

$$(\mathcal{I}(S), \mathcal{I}(K), s[c/cost]) \rightarrow (\mathcal{I}(S'), \mathcal{I}(S) \cdot \mathcal{I}(K), s[c/cost]) .$$

Noticing that $\mathcal{I}(S) \cdot \mathcal{I}(K) = \mathcal{I}(S \cdot K)$, this step is matched in the labelled language as follows:

$$(S, K, s[c/cost]) \rightarrow (S', S \cdot K, s[c/cost]) .$$

- On the other hand, if $S = \ell : S'$ is a labelled statement then $\mathcal{I}(S) = \text{inc}(\ell); \mathcal{I}(S')$ and, by a sequence of reductions steps, we have:

$$(\mathcal{I}(S), \mathcal{I}(K), s[c/cost]) \xrightarrow{*} (\mathcal{I}(S'), \mathcal{I}(K), s[c + \kappa(\ell)/cost]) .$$

This step is matched by the labelled reduction:

$$(S, K, s[c/cost]) \xrightarrow{\ell} (S', K, s[c/cost]) .$$

□

A.7 Proof of proposition 12

By diagram chasing using propositions 7(1), 9(1), and the definition 11 of labelling. □

A.8 Proof of proposition 13

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 10, for some λ :

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

Finally, by propositions 7(2) and 9(2) :

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

□

A.9 Proof of proposition 16

If $\lambda = \ell_1 \cdots \ell_n$ then the computation is the concatenation of simple paths labelled with ℓ_1, \dots, ℓ_n . Since $\kappa(\ell_i)$ bounds the cost of a simple path labelled with ℓ_i , the cost of the overall computation is bounded by $\kappa(\lambda) = \kappa(\ell_1) + \cdots + \kappa(\ell_n)$. □

A.10 Proof of proposition 18

Same proof as proposition 16, by replacing the word *bounds* by *is exactly* and the words *bounded by* by *exactly*. \square

A.11 Proof of proposition 19

In both labellings under consideration the root node is labelled. An obvious observation is that only commands of the shape *while b do S* introduce loops in the compiled code. 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 25 *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 25 guarantee that the two simple paths one can follow from the bge node have the same length/cost. \square

A.12 Proof of proposition 21

By applying consecutively proposition 13 and propositions 16 or 18. \square

B 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.

B.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 [13] 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 [9] for a formal definition of the Clight language. Here we just recall in figure B.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.

B.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 B.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.

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

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"> id n $\text{addrsymbol}(id)$ $\text{addrstack}(\delta)$ $op_1 a$ $op_2 a a$ $\kappa[a]$ $a?a : a$ 	<ul style="list-style-type: none"> local variable integer constant address of global symbol address within stack data unary arithmetic operation binary arithmetic operation memory read conditional expression
Statements:	$s ::=$ <ul style="list-style-type: none"> skip $id = a$ $\kappa[a] = a$ $id^? = a(\vec{a}) : sig$ $\text{tailcall } a(\vec{a}) : sig$ $\text{return}(a^?)$ $s; s$ $\text{if } a \text{ then } s \text{ else } s$ $\text{loop } s$ $\text{block } s$ $\text{exit } n$ $\text{switch } a \text{ tbl}$ $\text{lbl} : s$ $\text{goto } \text{lbl}$ 	<ul style="list-style-type: none"> empty statement assignment memory write function call function tail call function return sequence conditional infinite loop block delimiting exit constructs terminate the $(n + 1)^{th}$ enclosing block multi-way test and exit labelled statement jump to a label
Switch tables:	$tbl ::=$ <ul style="list-style-type: none"> default:exit(n) $\text{case } i: \text{exit}(n);tbl$ 	
Functions:	$Fd ::=$ <ul style="list-style-type: none"> internal $sig \vec{id} \vec{id} n s$ $\text{external } id \text{ sig}$ 	<ul style="list-style-type: none"> internal function: signature, parameters, local variables, stack size and body external function
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 $\rightarrow node$	(no instruction)
$psd_reg := op(psd_reg^*) \rightarrow node$	(operation)
$psd_reg := \&var_name \rightarrow node$	(address of a global)
$psd_reg := \&locals[n] \rightarrow node$	(address of a local)
$psd_reg := fun_name \rightarrow node$	(address of a function)
$psd_reg := memq(psd_reg[psd_reg]) \rightarrow node$	(memory load)
$memq(psd_reg[psd_reg]) := psd_reg \rightarrow node$	(memory store)
$psd_reg := fun_ref(psd_reg^*) : signature \rightarrow node$	(function call)
$fun_ref(psd_reg^*) : signature$	(function tail call)
test $op(psd_reg^*) \rightarrow node, node$	(branch)
return $psd_reg?$	(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 9: 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 9. Local variables are now represented by *pseudo registers* that are available in unbounded number. The grammar rule *op* that is not detailed in table 9 defines usual arithmetic and boolean operations (+, xor, ≤, 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.

$size ::= Byte \mid HalfWord \mid Word$	$fun_ref ::= fun_name \mid psd_reg$
$instruction ::=$	
$skip \rightarrow node$	(no instruction)
$psd_reg := n \rightarrow node$	(constant)
$psd_reg := unop(psd_reg) \rightarrow node$	(unary operation)
$psd_reg := binop(psd_reg, psd_reg) \rightarrow node$	(binary operation)
$psd_reg := \&globals[n] \rightarrow node$	(address of a global)
$psd_reg := \&locals[n] \rightarrow node$	(address of a local)
$psd_reg := fun_name \rightarrow node$	(address of a function)
$psd_reg := size(psd_reg[n]) \rightarrow node$	(memory load)
$size(psd_reg[n]) := psd_reg \rightarrow node$	(memory store)
$psd_reg := fun_ref(psd_reg^*) \rightarrow node$	(function call)
$fun_ref(psd_reg^*)$	(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)
$return\ psd_reg?$	(return)
$fun_def ::= fun_name(psd_reg^*)$	$program ::= globals : n$
$result : psd_reg?$	fun_def^*
$locals : psd_reg^*$	
$stack : n$	
$entry : node$	
$exit : node$	
$(node : instruction)^*$	

Table 10: Syntax of the RTL language

B.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 explicits 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 to an offset. The syntax of the language can be found in table 10. 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.

B.5 ERTL

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

$size ::= \text{Byte} \mid \text{HalfWord} \mid \text{Word}$	$fun_ref ::= fun_name \mid psd_reg$
$instruction ::=$	
skip $\rightarrow node$	(no instruction)
NewFrame $\rightarrow node$	(frame creation)
DelFrame $\rightarrow node$	(frame deletion)
$psd_reg := \text{stack}[slot, n] \rightarrow node$	(stack load)
$\text{stack}[slot, n] := psd_reg \rightarrow node$	(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)
$fun_def ::=$	$program ::=$
$fun_name(n)$	globals : n
locals : psd_reg^*	fun_def^*
stack : n	
entry : $node$	
($node : instruction$)*	

Table 11: Syntax of the ERTL language

Syntax. The syntax of the language is given in table 11. 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

$size ::= \text{Byte} \mid \text{HalfWord} \mid \text{Word}$	$fun_ref ::= fun_name \mid hdw_reg$																												
$instruction ::=$ <table style="width: 100%; border-collapse: collapse;"> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$skip \rightarrow node$</td><td>(no instruction)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$NewFrame \rightarrow node$</td><td>(frame creation)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$DelFrame \rightarrow node$</td><td>(frame deletion)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$hdw_reg := n \rightarrow node$</td><td>(constant)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$hdw_reg := unop(hdw_reg) \rightarrow node$</td><td>(unary operation)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$hdw_reg := binop(hdw_reg, hdw_reg) \rightarrow node$</td><td>(binary operation)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$hdw_reg := fun_name \rightarrow node$</td><td>(address of a function)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$hdw_reg := size(hdw_reg[n]) \rightarrow node$</td><td>(memory load)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$size(hdw_reg[n]) := hdw_reg \rightarrow node$</td><td>(memory store)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$fun_ref() \rightarrow node$</td><td>(function call)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$fun_ref()$</td><td>(function tail call)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$test uncon(hdw_reg) \rightarrow node, node$</td><td>(branch unary condition)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$test bincon(hdw_reg, hdw_reg) \rightarrow node, node$</td><td>(branch binary condition)</td></tr> <tr><td style="border-right: 1px solid black; padding-right: 10px;">$return$</td><td>(return)</td></tr> </table>	$skip \rightarrow node$	(no instruction)	$NewFrame \rightarrow node$	(frame creation)	$DelFrame \rightarrow node$	(frame deletion)	$hdw_reg := n \rightarrow node$	(constant)	$hdw_reg := unop(hdw_reg) \rightarrow node$	(unary operation)	$hdw_reg := binop(hdw_reg, hdw_reg) \rightarrow node$	(binary operation)	$hdw_reg := fun_name \rightarrow node$	(address of a function)	$hdw_reg := size(hdw_reg[n]) \rightarrow node$	(memory load)	$size(hdw_reg[n]) := hdw_reg \rightarrow node$	(memory store)	$fun_ref() \rightarrow node$	(function call)	$fun_ref()$	(function tail call)	$test uncon(hdw_reg) \rightarrow node, node$	(branch unary condition)	$test bincon(hdw_reg, hdw_reg) \rightarrow node, node$	(branch binary condition)	$return$	(return)	
$skip \rightarrow node$	(no instruction)																												
$NewFrame \rightarrow node$	(frame creation)																												
$DelFrame \rightarrow node$	(frame deletion)																												
$hdw_reg := n \rightarrow node$	(constant)																												
$hdw_reg := unop(hdw_reg) \rightarrow node$	(unary operation)																												
$hdw_reg := binop(hdw_reg, hdw_reg) \rightarrow node$	(binary operation)																												
$hdw_reg := fun_name \rightarrow node$	(address of a function)																												
$hdw_reg := size(hdw_reg[n]) \rightarrow node$	(memory load)																												
$size(hdw_reg[n]) := hdw_reg \rightarrow node$	(memory store)																												
$fun_ref() \rightarrow node$	(function call)																												
$fun_ref()$	(function tail call)																												
$test uncon(hdw_reg) \rightarrow node, node$	(branch unary condition)																												
$test bincon(hdw_reg, hdw_reg) \rightarrow node, node$	(branch binary condition)																												
$return$	(return)																												
$fun_def ::= fun_name(n)$ $locals : n$ $stack : n$ $entry : node$ $(node : instruction)^*$	$program ::=$ <table style="width: 100%; border-collapse: collapse;"> <tr><td style="padding-right: 20px;">$globals : n$</td><td></td></tr> <tr><td>fun_def^*</td><td></td></tr> </table>	$globals : n$		fun_def^*																									
$globals : n$																													
fun_def^*																													

Table 12: Syntax of the LTL language

a *register 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.

B.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 explicit in ERTL; thus, function calls and returns do not need parameters in LTL. The syntax is defined in table 12.

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.

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

Table 13: Syntax of the LIN language

B.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 13.

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.

B.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 14.

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

B.9 Benchmarks

To ensure that our prototype compiler is realistic, we performed some preliminary benchmarks on a 183MHz MIPS 4KEc processor, running a linux based distribution. We compared the

$load ::= lb \mid lhw \mid lw$ $store ::= sb \mid shw \mid sw$ $fun_ref ::= fun_name \mid hdw_reg$

$instruction ::=$

nop	(empty instruction)
li hdw_reg, n	(constant)
unop hdw_reg, hdw_reg	(unary operation)
binop $hdw_reg, hdw_reg, hdw_reg$	(binary operation)
la hdw_reg, fun_name	(address of a function)
load $hdw_reg, n(hdw_reg)$	(memory load)
store $hdw_reg, n(hdw_reg)$	(memory store)
call fun_ref	(function call)
uncon $hdw_reg, node$	(branch unary condition)
bincon $hdw_reg, hdw_reg, node$	(branch binary condition)
mips_label :	(Mips label)
j mips_label	(goto)
return	(return)

$program ::=$

globals : n
entry : mips_label*
instruction*

Table 14: Syntax of the Mips language

	gcc -00	acc	gcc -01
badsort	55.93	34.51	12.96
fib	76.24	34.28	45.68
mat_det	163.42	156.20	54.76
min	12.21	16.25	3.95
quicksort	27.46	17.95	9.41
search	463.19	623.79	155.38

Figure 3: Benchmarks results (execution time is given in seconds).

wall clock execution time of several simple C programs compiled with our compiler against the ones produced by GCC set up with optimization levels 0 and 1. As shown by Figure 3, our prototype compiler produces executable programs that are on average faster than GCC's without optimizations.