



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 with Addendum

Main Authors:
Roberto M. Amadio, Nicolas Ayache, Yann Régis-Gianas

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 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.

Contents

1	Introduction	5
2	Labelling approach for the toy compiler	8
2.1	Labelled Imp	9
2.2	Labelled Vm	9
2.3	Labelled Mips	10
2.4	Labellings and instrumentations	10
2.5	Sound and precise labellings	11
3	Labelling approach for the C compiler	13
3.1	Labelled languages	13
3.2	Labelling of the source language	14
3.3	Verifications on the object code	15
3.4	Building the cost annotation	16
3.5	Testing	16
4	Conclusion and future work	16
A	A toy compiler	18
A.1	Imp: language and semantics	18
A.2	Vm: language and semantics	18
A.3	Compilation from Imp to Vm	19
A.4	Mips: language and semantics	20
A.5	Compilation from Vm to Mips	21
B	Proofs	23
B.1	Notation	23
B.2	Proof of proposition 1	23
B.3	Proof of proposition 3	23
B.4	Proof of proposition 4	23
B.5	Proof of proposition 6	24
B.6	Proof of proposition 7	24
B.7	Proof of proposition 10	24
B.8	Proof of proposition 12	24
B.9	Proof of proposition 13	25
B.10	Proof of proposition 15	25
B.11	Proof of proposition 19	25
C	A C compiler	27
C.1	Clight	27
C.2	Cminor	27
C.3	RTLabs	27
C.4	RTL	31
C.5	ERTL	31
C.6	LTL	33
C.7	LIN	34
C.8	Assembly	34
C.9	Benchmarks	36
D	A direct approach	37
D.1	Mips and Vm cost annotations	37
D.2	Imp cost annotation	38
D.3	Composition	40
D.4	Coq development	40
D.5	Limitations of the direct approach	40
E	Related approaches	41

F Assessment of the deliverable within the *CerCo* project, with hindsight

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 [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 mechanically verified 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] to infer more synthetic cost annotations. 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:

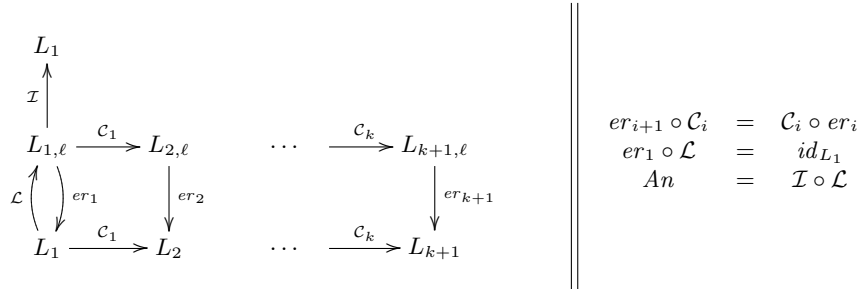
- enrich the program with a collection of *global cost variables* to measure resource consumption (time, stack size, heap size, . . .)
- inject suitable code at some critical points (procedures, loops, . . .) to keep track of the execution cost.

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 appendix D 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 $er_{L_1} \circ \mathcal{L}$ 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 summarized by the following diagram.

$$\text{Imp} \xrightarrow{c} \text{Vm} \xrightarrow{c'} \text{Mips}$$

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.

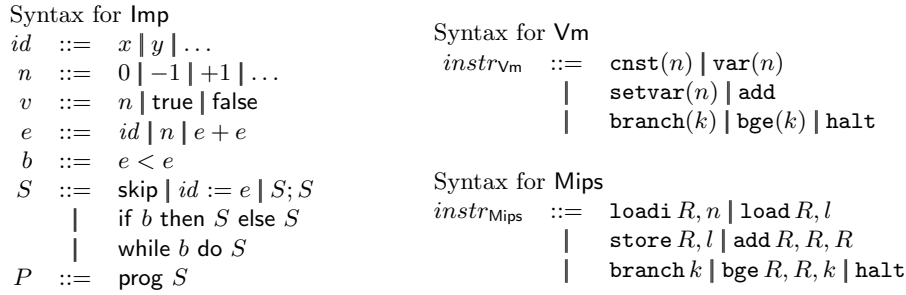
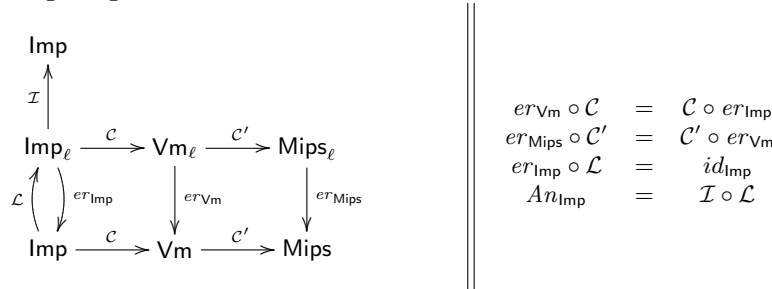


Figure 1: Syntax definitions.

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



2.1 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 is extended as described by the following rule.

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

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 $\lambda = \lambda_1 \dots \lambda_n$.

2.2 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 1 *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 2 *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.*

2.3 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 1. $m \Vdash \sigma, s$ means “the low-level Mips memory m realizes the Vm stack σ and state s ”.

Proposition 3 *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 $m \Vdash \sigma, s$ then $(\mathcal{C}'(C), m) \Downarrow (m', \lambda)$ and $m' \Vdash \sigma, s'$.

2.4 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 has follows.

Proposition 4 *Let S be an Imp_ℓ command. If $(\mathcal{I}(S), s[c/\text{cost}]) \Downarrow s'[c+\delta/\text{cost}]$ then $\exists \lambda \kappa(\lambda) = \delta$ and $(S, s[c/\text{cost}]) \Downarrow (s'[c/\text{cost}], \lambda)$.*

Definition 5 *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 6 *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 7 *Given a function κ for the labels and a labelling function \mathcal{L} , for all programs P of the source language if $(\mathcal{I}(\mathcal{L}(P)), s[c/\text{cost}]) \Downarrow s'[c + \delta/\text{cost}]$ and $m \Vdash -\epsilon, s[c/\text{cost}]$ then $(\mathcal{C}'(\mathcal{C}(\mathcal{L}(P))), m) \Downarrow (m', \lambda)$, $m' \Vdash -\epsilon, s'[c/\text{cost}]$ and $\kappa(\lambda) = \delta$.*

2.5 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 8 *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 9 *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 10 *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.

Definition 11 *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 12 *If M is precisely labelled and $(M, m) \Downarrow (m', \lambda)$ then the cost of the computation is $\kappa(\lambda)$.*

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

Table 1: Two labellings for the Imp language

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 1 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 13 *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.

Definition 14 *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 15 *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 summarize, 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 is enough to label commands while it is not necessary to label boolean conditions and expressions.

Besides soundness and precision, a third criterion to evaluate labellings is that they do not introduce too many unnecessary labels. We call this property *economy*. In practice, it seems that one can produce first a sound and possibly precise labelling and then apply heuristics to eliminate unnecessary labels.

We stress that our approach to labelling is based on the hypothesis that we can obtain precise informations on the execution time of each instruction of the generated binary code. This hypothesis is indeed realised in the processors of the 8051 family we are considering. On the other hand, the execution time of instructions running on more complex architectures including, *e.g.*, cache memory or pipelines are much less predictable. This lack of precision may somehow be compensated by analysing the worst-case execution time of (long) sequences of instructions. This point is further developed in appendix E.

3 Labelling approach for the C compiler

This section informally describes the labelled extensions of the languages in the compilation chain (see Appendix C 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 assembly 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 detailed in the following sections. It is summarized by the following steps.

1. Label the input Clight program.
2. Compile the labelled Clight program in the labelled world. This produces a labelled assembly code.
3. For each label of the labelled assembly code, compute the cost of the instructions under its scope and generate a *label-cost mapping*. An unlabelled assembly code — the result of the compilation — is obtained by removing the labels from the labelled assembly 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.

3.1 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 lmp 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* \rightarrow *node*. For LIN, Mips and 8051, it is *emit label*. The translation of these label instructions is immediate. In Mips and 8051, we also rely on a reserved label *begin_function* to pinpoint the beginning of a function code (cf. section 3.2).

3.2 Labelling of the source language

As for the toy compiler (cf. end of section 2), 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 assembly code.

Sequential instructions A sequence of `Clight` instructions that compile to sequential assembly code, such as a sequence of assignments, can be handled by a single label which covers the unique execution path.

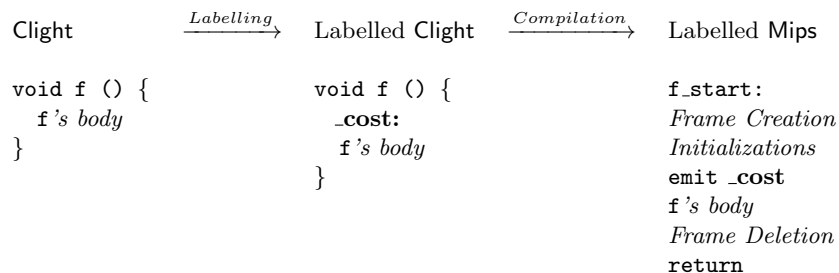
Ternary expressions and conditionals Most `Clight` expressions compile to sequential assembly code. *Ternary expressions*, that introduce a branching in the control flow, are one exception. In this case, we achieve precision by associating a label with each branch. This is similar to the treatment of the conditional we have already discussed in section 2. 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 2 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.

<code>Clight</code>	$\xrightarrow{\text{Labelling}}$	Labelled <code>Clight</code>	$\xrightarrow{\text{Compilation}}$	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 assembly code are performed by indirect jumps, the address of the callee being in the memory. In the general case, this address cannot be inferred statically. Even though the destination point of a function call is unknown, when the considered assembly 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:



We notice that some instructions in assembly code 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 behavior 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.

3.3 Verifications on the object code

The labelling previously described has been designed so that the compiled assembly 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 3.2. 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.

In the particular case of the 8051, which is a very basic micro-controller, the number of cycles required by each instruction is fully specified. Therefore, if the labelling is precise, the cost associated to a label is the exact number of cycles necessary to execute every simple execution path starting from that label.

3.4 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 initialized 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.

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

3.5 Testing

It is desirable to test the coherence of the labelling from Clight to assembly code. 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 optimizations that may re-order computations, the weaker condition mentioned in remark 2 could be considered.

Furthermore, in the case of the 8051 back-end, we can run our compiled code into a dedicated simulator and we check that the statically computed number of cycles for each simple path is the same as the real one.

4 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 plan to formalize and validate in the *Calculus of Inductive Constructions* the prototype implementation of the labelling approach for the C compiler described in section 3. This requires a major implementation effort which will be carried on in collaboration with our partners of the CerCo project [2].

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 *Embedded Real Time Software (ERTS)*, 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 hw/sw interface*, by Hennessy and Patterson, 2005.
- [8] MCS 51 Microcontroller Family User’s Manual. Publication number 121517, by Intel Corporation, 1994.
- [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] 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.
- [14] Microcontroller market and Technology. Analysis report 2008. EMITT solutions.

$$\begin{array}{ll}
(x := e, K, s) & \rightarrow (\text{skip}, K, s[v/x]) \quad \text{if } (e, s) \Downarrow v \\
(S; S', K, s) & \rightarrow (S, S' \cdot K, s) \\
(\text{if } b \text{ then } S \text{ else } S', K, s) & \rightarrow \begin{cases} (S, K, s) & \text{if } (b, s) \Downarrow \text{true} \\ (S', K, s) & \text{if } (b, s) \Downarrow \text{false} \end{cases} \\
(\text{while } b \text{ do } S, K, s) & \rightarrow \begin{cases} (S, (\text{while } b \text{ do } S) \cdot K, s) & \text{if } (b, s) \Downarrow \text{true} \\ (\text{skip}, K, s) & \text{if } (b, s) \Downarrow \text{false} \end{cases} \\
(\text{skip}, S \cdot K, s) & \rightarrow (S, K, s)
\end{array}$$
Table 2: Small-step operational semantics of `Imp` commands

A A toy compiler

We formalize the toy compiler introduced in section 1 ².

A.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 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: $(S, s) \Downarrow s'$ if $(S, \text{halt}, s) \rightarrow \dots \rightarrow (\text{skip}, \text{halt}, s')$.

A.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 .

²A mechanized proof in the Coq proof assistant can be found in K. Memarian’s *Travail d’étude et de recherche* entitled *Complexité Certifiée* which can be downloaded at <http://www.pps.jussieu.fr/~yrg/miniCerCo/>.

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 3: Operational semantics Vm programs

$C[i] =$	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 4: Conditions for well-formed code

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 formalized 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 16 *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 17 (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 $h = h'$.**

A.3 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 $sz(e)$, $sz(b)$, $sz(S)$ for the number of

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

Table 5: Compilation from Imp to Vm

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.

Proposition 18 *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 $j = |C \cdot \mathcal{C}(b, k)|$.*

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

Proposition 19 *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 20 *For any program P there is a unique h such that $\mathcal{C}(P) : h$.*

A.4 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 denote with M a list of instructions. The operational semantics is formalized 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

$C[i] =$	$C'(i, C) =$
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)	(branch k') 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 7: Compilation from Vm to Mips

A.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 $m \parallel \sigma, s$, if the following conditions are satisfied: (1) $s(x) = m(l_x)$, and (2) if $0 \leq i < |\sigma|$ then $\sigma[i] = m(R_i)$ if $i < b$, and $\sigma[i] = m(l_i)$ if $i \geq b$.

The compilation function 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 20(3), there is a unique h such that $C : h$. We denote with $C'(C)$ the concatenation $C'(0, C) \dots C'(|C| - 1, C)$. Given a well formed Vm code C with $i < |C|$ we denote with $p(i, C)$ the position of the first instruction in $C'(C)$ which corresponds to the compilation of the instruction with position i in C . This is defined as³ $p(i, C) = \sum_{0 \leq j < i} d(j, C)$, where the function $d(i, C)$ is defined as $d(i, C) = |C'(i, C)|$. 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

³There 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 , and finally the function C' as in table 7.

can then be stated as follows.

Proposition 21 *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 $C'(C) \vdash (p(i, C), m) \xrightarrow{*} (p(j, C), m')$ and $m' \Vdash_{-\sigma', s'}$.*

B Proofs

We omit the proofs that have been mechanically checked by K. Memarian and R. Saillard with the Coq proof assistant ⁴.

B.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.

B.2 Proof of proposition 1

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

Proposition 22 *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 19 and it is proven in the same way with an additional case for labelled commands. \square

B.3 Proof of proposition 3

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

Proposition 23 *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 $C'(C) \vdash (p(i, C), m) \xrightarrow{t} (p(j, C), m')$ and $m' \Vdash \sigma', s'$.*

B.4 Proof of proposition 4

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

- 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:

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

⁴ These proofs can be downloaded at <http://www.pps.jussieu.fr/~yrg/miniCerCo/> and at <http://www.pps.jussieu.fr/~saillard>.

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

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

□

B.5 Proof of proposition 6

By diagram chasing using propositions 1(1), 3(1), and the definition 5 of labelling.

□

B.6 Proof of proposition 7

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 4, for some λ :

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

Finally, by propositions 1(2) and 3(2) :

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

□

B.7 Proof of proposition 10

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 $\kappa(\lambda) = \kappa(\ell_1) + \cdots + \kappa(\ell_n)$.

□

B.8 Proof of proposition 12

Same proof as proposition 10, by replacing the word *bounds* by *is exactly* and the words *bounded by* by *exactly*.

□

B.9 Proof of proposition 13

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 24 *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 24 guarantee that the two simple paths one can follow from the bge node have the same length/cost. \square

B.10 Proof of proposition 15

By applying consecutively proposition 7 and propositions 10 or 12. \square

B.11 Proof of proposition 19

Given a Vm code C , we define an ‘accessibility relation’ $\overset{C}{\rightsquigarrow}$ 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 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.

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

The following properties are useful.

Lemma 25 (1) *The relation $\overset{C}{\rightsquigarrow}$ is transitive.*

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

The first property can be proven by induction on the definition of $\overset{C}{\rightsquigarrow}$ 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:

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

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 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:

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

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 $C \vdash (i, \sigma, s) \xrightarrow{*} (i', \sigma, s) \xrightarrow{*} (a, \sigma, s')$, and it follows from (3'), the fact that $(e_1 < e_2) \Downarrow \text{true}$, and proposition 18(2).
- Property (7), follows from lemma 25(2), fact (5), and the following proof tree:

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

□

C 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.

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

There is a notable difficulty to compile the C language into 8051 assembly code due to the fact that 8051's machine word are 8bits long and C has 32bits primitive datatypes. To deal with the dissimilarity in that case, we first translate the Clight input program into a Clight input program that only uses 8bits primitive datatypes.

C.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 C.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.

C.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

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

Figure 2: Syntax of the Clight language

Signatures:	$sig ::= sig \vec{int} \ (int void)$	arguments and result
Expressions:	$a ::= id$ $ n$ $ \text{addrsymbol}(id)$ $ \text{addrstack}(\delta)$ $ op_1 a$ $ op_2 a a$ $ \kappa[a]$ $ a?a : a$	local variable integer constant address of global symbol address within stack data unary arithmetic operation binary arithmetic operation memory read conditional expression
Statements:	$s ::= skip$ $ 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}$ $ lbl : s$ $ \text{goto } lbl$	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:	$tbl ::= \text{default} : \text{exit}(n)$ $ \text{case } i : \text{exit}(n);tbl$	
Functions:	$Fd ::= \text{internal } sig \vec{id} \vec{id} \ n \ s$ $ \text{external } id \ sig$	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 3: 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 8: Syntax of the RTLabs language

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.

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 8. Local variables are now represented by *pseudo registers* that are available in unbounded number. The grammar rule *op* that is not detailed in table 8 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>fun_name</i> (<i>psd_reg</i> *)	<i>program</i> ::= globals : <i>n</i>
	result : <i>psd_reg</i> ?
	locals : <i>psd_reg</i> *
	stack : <i>n</i>
	entry : <i>node</i>
	exit : <i>node</i>
	(node : <i>instruction</i>)*

Table 9: Syntax of the RTL language

C.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 and 8051 assembly languages.

Syntax. RTL is very close to RTLabs. It is based on CFGs and explicits the Mips or the 8051 instructions corresponding to the RTLabs instructions. Type information disappears: everything is represented using machine words. Moreover, each global of the program is associated to an offset. The syntax of the language can be found in table 9. The grammar rules *unop*, *binop*, *uncon*, and *bincon*, respectively, represent the sets of unary operations, binary operations, unary conditions and binary conditions of the target assembly language.

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

C.5 ERTL

As in RTL, the structure of ERTL programs is based on CFGs. ERTL explicits the calling conventions of the Mips assembly language. In the back-end for 8051, we defined our own

<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)
NewFrame → <i>node</i>	(frame creation)
DelFrame → <i>node</i>	(frame deletion)
<i>psd_reg</i> := stack[slot, <i>n</i>] → <i>node</i>	(stack load)
stack[slot, <i>n</i>] := <i>psd_reg</i> → <i>node</i>	(stack store)
<i>hdw_reg</i> := <i>psd_reg</i> → <i>node</i>	(pseudo to hardware)
<i>psd_reg</i> := <i>hdw_reg</i> → <i>node</i>	(hardware to pseudo)
<i>psd_reg</i> := <i>n</i> → <i>node</i>	(constant)
<i>psd_reg</i> := unop(<i>psd_reg</i>) → <i>node</i>	(unary operation)
<i>psd_reg</i> := binop(<i>psd_reg</i> , <i>psd_reg</i>) → <i>node</i>	(binary operation)
<i>psd_reg</i> := <i>fun_name</i> → <i>node</i>	(address of a function)
<i>psd_reg</i> := size(<i>psd_reg</i> [<i>n</i>]) → <i>node</i>	(memory load)
size(<i>psd_reg</i> [<i>n</i>]) := <i>psd_reg</i> → <i>node</i>	(memory store)
<i>fun_ref</i> (<i>n</i>) → <i>node</i>	(function call)
<i>fun_ref</i> (<i>n</i>)	(function tail call)
test uncon(<i>psd_reg</i>) → <i>node</i> , <i>node</i>	(branch unary condition)
test bincon(<i>psd_reg</i> , <i>psd_reg</i>) → <i>node</i> , <i>node</i>	(branch binary condition)
return <i>b</i>	(return)
<i>fun_def</i> ::= <i>fun_name</i> (<i>n</i>)	<i>program</i> ::= globals : <i>n</i>
locals : <i>psd_reg</i> *	<i>fun_def</i> *
stack : <i>n</i>	
entry : <i>node</i>	
(<i>node</i> : <i>instruction</i>)*	

Table 10: Syntax of the ERTL language

calling convention since there were none.

Syntax. The syntax of the language is given in table 10. 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.

<i>size</i> ::= Byte HalfWord Word	<i>fun_ref</i> ::= <i>fun_name</i> <i>hdlw_reg</i>
<i>instruction</i> ::=	
skip → <i>node</i>	(no instruction)
NewFrame → <i>node</i>	(frame creation)
DelFrame → <i>node</i>	(frame deletion)
<i>hdlw_reg</i> := <i>n</i> → <i>node</i>	(constant)
<i>hdlw_reg</i> := <i>unop</i> (<i>hdlw_reg</i>) → <i>node</i>	(unary operation)
<i>hdlw_reg</i> := <i>binop</i> (<i>hdlw_reg</i> , <i>hdlw_reg</i>) → <i>node</i>	(binary operation)
<i>hdlw_reg</i> := <i>fun_name</i> → <i>node</i>	(address of a function)
<i>hdlw_reg</i> := <i>size</i> (<i>hdlw_reg</i> [<i>n</i>]) → <i>node</i>	(memory load)
<i>size</i> (<i>hdlw_reg</i> [<i>n</i>]) := <i>hdlw_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>hdlw_reg</i>) → <i>node</i> , <i>node</i>	(branch unary condition)
test <i>bincon</i> (<i>hdlw_reg</i> , <i>hdlw_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> ::= globals : <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 11: Syntax of the LTL language

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 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.

C.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 11.

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>asm_label</i> :	(assembly 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 12: Syntax of the LIN language

C.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 12.

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.

C.8 Assembly

Mips is a rather simple assembly language whereas 8051's is rather complex [8]. We only describe Mips. 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 13.

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

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

$instruction ::=$

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

$program ::=$

<code>globals : n</code>
<code>entry : mips_label*</code>
<code>instruction*</code>

Table 13: 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 4: Benchmarks results (execution time is given in seconds).

C.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 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 4, our prototype compiler produces executable programs that are on average faster than GCC's without optimizations.

D A direct approach

Our first attempt at building cost annotations followed a ‘direct’ approach which is summarised by the following diagram.

$$\begin{array}{ccccc}
 L_1 & \xrightarrow{C_1} & L_2 & \cdots & \xrightarrow{C_k} & L_{k+1} \\
 \downarrow An_1 & & \downarrow An_1 & & & \downarrow An_{k+1} \\
 L_1 & & L_2 & & & L_{k+1}
 \end{array}$$

With respect to the introductive discussion in section 1, 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.

Applying the direct approach to the toy compiler, results in the following diagram.

$$\begin{array}{ccccc}
 \text{Imp} & \xrightarrow{C} & \text{Vm} & \xrightarrow{C'} & \text{Mips} \\
 \downarrow An_{\text{Imp}} & & \downarrow An_{\text{Vm}} & & \downarrow An_{\text{Mips}} \\
 \text{Imp} & & \text{Vm} & & \text{Mips}
 \end{array}$$

D.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_{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_{cost} :

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

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

Table 14: Upper bounds on the cost of Vm instructions

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 .
2. Before the instruction i^{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 14):

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

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 26 *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 21 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 27 *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*

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

D.2 Imp cost annotation

The definition of the cost annotation $An_{\text{Imp}}(P)$ for an Imp program P is defined in table 15 and it relies on an auxiliary function κ 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*:

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

Table 15: Annotation for Imp programs

- 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 D.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 18.

Proposition 28 *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 $(d' - d) \leq (c' - c) + \kappa(S)$.

D.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 29 *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 $m'[c'/l_{\text{cost}}] \Vdash_{-\epsilon, s'[c'/\text{cost}]}$.

D.4 Coq development

We have formalised and mechanically checked in COQ the application of the direct approach to the toy compiler (but for propositions 27 and 29 for which we produced ‘paper proofs’). 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 [10] 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.

D.5 Limitations of the direct approach

As mentioned in section D.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 15 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 κ 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 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, we manipulate costs at an abstract level as labels and produce numerical values only at the very end of the construction.

E Related approaches

Our fine-grained approach to labelling is based on the hypothesis that we can obtain precise information on the execution time of each instruction of the generated binary code. This hypothesis is indeed realised in the processors of the 8051 family we are considering. On the other hand, the execution time of instructions running on more complex architectures including, *e.g.*, cache memory or pipelines are much less predictable. This lack of precision may somehow be compensated by analysing the worst-case execution time of sequences of instructions. As mentioned in the introduction, tools such as those developed by AbsInt require explicit annotations of the binaries to determine the number of times a loop is iterated and apply abstract interpretation methods in order to provide an estimation of the execution time of the code. From the point of view of a *formal* certification this approach rises at least two questions.

First, one needs a formal certification of the fact that the abstract interpretation method does indeed produce correct results for a given processor. This in turn supposes a formal modelling of the processor and a proof that the abstract interpretation method does indeed abstract the processor's behaviour. Taking a pragmatic (and informal) approach, one could argue that the results of the tool could be trusted, in the same way we trust the execution times given by the manufacturers.

Second, one needs a formal certification of the fact that the annotations on the loops are indeed correct. For instance, suppose we can identify in the source code *for-loops* that iterate a *constant* number of times some sequence of statements. In our approach, we must insert a label in the body of the loop so as to avoid that the control flow graph associated with the binary code contains a loop that does not cross any label. Having determined the cost of one iteration of the body of the loop, we can then reason at the level of the source code to obtain the cost of several iterations.

In an alternative approach, we could try to propagate to the control flow graph the information that the body of the loop is iterated a constant number of times and then invoke tools such as those mentioned above to get an estimation of the cost of running the whole iteration. As mentioned above, the rationale for this approach is to get a more precise estimation of the cost of executing the iteration on a complex architecture. While the approach seems easily implementable, it is not clear how one should proceed in order to formally certify its correctness. Specifically, one needs to propagate enough information at the level of the binary code to be able to verify automatically that the predicted number of iterations of the body of the loop is indeed correct. In turn, this means to be able to track down the initial value of a certain number of memory locations (registers and/or stack locations) and to guarantee that their contents is modified at each iteration leading eventually to the exit of the loop. While this seems within the reach of a *proof-carrying code* approach, the challenge is to devise a more lightweight approach which targets the specific properties of interest.

F Assessment of the deliverable within the *CerCo* project, with hindsight

The port of the compiler to the 8051 processor has required a major and rather unrewarding implementation effort which is described in deliverable D2.2.

One problem that has emerged is that the compilation of general purpose C programs with recursion and/or 16-32 integers bits produces binaries whose size approaches quickly the limits of the architecture (program memory is limited to 64Kbytes).

This is a problem for instance for the case study which is planned in deliverable D5.3 (Case study: analysis of synchronous code) whose purpose is to produce automatically cost annotations for the C programs generated by a compiler of a synchronous language such as LUSTRE. The generated C programs are rather large (though their control flow is simple) and assume 32 bits integers.

At the time being, two options are being considered as far as deliverable D5.3 is concerned.

1. Improve the efficiency of the compilation of 16/32 bits integers for the 8051 compiler.
2. Conduct the case study relying on a compiler to a 32 bits processor.

If the second option is chosen, one possibility is to target the MIPS processor for which a prototype compiler has already been developed. More generally, one advantage of the second option is that it allows to conduct realistic case studies on a larger class of programming constructs. For instance, we are interested in studying the extension of the labelling approach to a compilation chain from a functional language of the ML family to C.