

Proving Correctness of Compilers Using Structured Graphs

Patrick Bahr

Department of Computer Science, University of Copenhagen, Denmark
paba@di.ku.dk

Abstract. We present an approach to compiler implementation using Oliveira and Cook’s structured graphs that avoids the use of explicit jumps in the generated code. The advantage of our method is that it takes the implementation of a compiler using a tree type along with its correctness proof and turns it into a compiler implementation using a graph type along with a correctness proof. The implementation and correctness proof of a compiler using a tree type without explicit jumps is simple, but yields code duplication. Our method provides a convenient way of improving such a compiler without giving up the benefits of simple reasoning.

1 Introduction

Verification of compilers – like other software – is difficult [13]. In such an endeavour one typically has to balance the “cleverness” of the implementation with the simplicity of reasoning about it. A concrete example of this fact is given by Hutton and Wright [10] who present correctness proofs of compilers for a simple language with exceptions. The authors first present a naïve compiler implementation that produces a tree representing the possible control flow of the input program. The code that it produces is essentially the right code, but the compiler *loses information* since it duplicates code instead of sharing it. However, the simplicity of the implementation is matched with a clean and simple proof by equational reasoning. Hutton and Wright also present a more realistic compiler, which uses labels and explicit jumps, resulting in a target code in linear form and without code duplication. However, the cleverer implementation also requires a more complicated proof, in which one has to reason about the freshness and scope of labels.

In this paper we present an intermediate approach, which is still simple, both in its implementation and in its correctness proof, but which avoids the loss of information of the simple approach described by Hutton and Wright [10]. The remedy for the information loss of the simple approach is obvious: we use a graph instead of a tree structure to represent the target code. The linear representation with labels and jumps is essentially a graph as well – it is just a very inconvenient one for reasoning. Instead of using unique names to represent sharing, we use the *structured graphs* representation of Oliveira and Cook [18]. This representation

of graphs uses parametric higher-order abstract syntax [5] to represent binders, which in turn are used to represent sharing. This structure allows us to take the simple compiler implementation using trees, make a slight adjustment to it, and obtain a compiler implementation using graphs that preserves the sharing information.

In essence our approach teases apart two aspects that are typically combined in code generation: (1) the translation into the target language, and (2) generating fresh (label) names for representing jumps in the target language. By keeping the two aspects separate, we can implement further transformations, e.g. code optimisations, without having to deal with explicit jumps and names. Only in the final step, when the code is linearised, names have to be generated in order to produce explicit jump instructions. Consequently, the issues that ensue in this setting can be dealt with in isolation – separately from the actual translation and subsequent transformation steps.

Our main goal is to retain the simplicity of the correctness proof of the tree-based compiler. The key observation making this possible is that the semantics of the tree-based and the graph-based target language, i.e. their respective *virtual machines*, are equivalent after *unravelling* of the graph structure. More precisely, given the semantics of the tree-based and the graph-based target language as $exec_{\top}$ and $exec_{\mathcal{G}}$, respectively, we have the following equation:

$$exec_{\mathcal{G}} = exec_{\top} \circ unravel$$

We show that this correspondence is an inherent consequence of the recursion schemes that are used to define these semantics. In fact, this correspondence follows from the correctness of *short cut fusion* [8, 12]. That is, the above property is independent of the target language of the compiler. As a consequence, the correctness proof of the improved, graph-based compiler is reduced to a proof that its implementation is equivalent to the tree-based implementation modulo unravelling. More precisely, it then suffices to show that

$$comp_{\top} = unravel \circ comp_{\mathcal{G}}$$

which is achieved by a straightforward induction proof.

In sum, the technique that we propose here improves existing simple compiler implementations to more realistic ones using a graph representation for the target code. This improvement requires minimal effort – both in terms of the implementation and the correctness proof. The fact that we consider both the implementation and its correctness proof makes our technique the ideal companion to improve a compiler that has been obtained by calculation [16]. Such calculations derive a compiler from a specification, and produce not only an implementation of the compiler but also a proof of its correctness. The example compiler that we use in this paper has in fact been calculated in this way by Bahr and Hutton [3], and we have successfully applied our technique to other compilers derived by Bahr and Hutton [3], which includes compilers for languages with features such as (synchronous and asynchronous) exceptions, (global and local) state and non-determinism. Thus, despite its simplicity, our technique

is quite powerful, especially when combined with other techniques such as the abovementioned calculation techniques.

In short, the contributions of this paper are the following:

- From a compiler with code duplication we derive a compiler that avoids duplication using a graph representation.
- Using short cut fusion, we prove that folds over graphs are equal to corresponding folds over the unravelling of the input graphs.
- Using the above result, we derive the correctness of the graph-based compiler implementation from the correctness of the tree-based compiler.
- We further simplify the proof by using free monads to represent tree types together with a corresponding monadic graph type.

Throughout this paper we use Haskell [14] as the implementation language.

2 A Simple Compiler

The example language that we use throughout the paper is a simple expression language with integers, addition and exceptions:

```
data Expr = Val Int | Add Expr Expr
         | Throw | Catch Expr Expr
```

The semantics of this language is defined using an evaluation function that evaluates a given expression to an integer value or returns *Nothing* in case of an uncaught exception:

```
eval :: Expr → Maybe Int
eval (Val n)      = Just n
eval (Add x y)   = case eval x of
                    Nothing → Nothing
                    Just n  → case eval y of
                                Nothing → Nothing
                                Just m  → Just (n + m)
eval Throw       = Nothing
eval (Catch x h) = case eval x of
                    Nothing → eval h
                    Just n  → Just n
```

This is the same language and semantics used by Hutton and Wright [10]. Like Hutton and Wright, we chose a simple language in order to focus on the essence of the problem, which in our case is control flow in the target language and the use of duplication or sharing to represent it. Moreover, this choice allows us to compare our method to the original work of Hutton and Wright whose focus was on the simplicity of reasoning.

The target for the compiler is a simple stack machine with the following instruction set:

data $Code = PUSH\ Int\ Code \mid ADD\ Code \mid HALT$
 $\mid UNMARK\ Code \mid MARK\ Code\ Code \mid THROW$

The intended semantics (which is made precise later) for the instructions is:

- $PUSH\ n$ pushes the integer value n on the stack,
- ADD expects two integers on the stack and replaces them with their sum,
- $MARK\ c$ pushes the exception handler code c on the stack,
- $UNMARK$ removes such a handler code from the stack,
- $THROW$ unwinds the stack until an exception handler code is found, which is then executed, and
- $HALT$ stops the execution.

For the implementation of the compiler we deviate slightly from the presentation of Hutton and Wright [10] and instead write the compiler in a style that uses an additional accumulation parameter c , which simplifies the proofs [9]:

$comp^A :: Expr \rightarrow Code \rightarrow Code$
 $comp^A (Val\ n) \quad c = PUSH\ n\ c$
 $comp^A (Add\ x\ y) \quad c = comp^A\ x\ (comp^A\ y\ (ADD\ c))$
 $comp^A\ Throw \quad c = THROW$
 $comp^A (Catch\ x\ h) \quad c = MARK\ (comp^A\ h\ c)\ (comp^A\ x\ (UNMARK\ c))$

Since the code generator is implemented in this code continuation passing style, function application corresponds to concatenation of code fragments. To stress this reading, we shall use the operator \triangleright , which is simply defined as function application and is declared to associate to the right with minimal precedence:

$(\triangleright) :: (a \rightarrow b) \rightarrow a \rightarrow b$
 $f \triangleright x = f\ x$

For instance, the equation for the *Add* case of the definition of $comp^A$ then reads:

$comp^A (Add\ x\ y)\ c = comp^A\ x \triangleright comp^A\ y \triangleright ADD \triangleright c$

To obtain the final code for an expression, we supply $HALT$ as the initial value of the accumulator of $comp^A$. The use of the \triangleright operator to supply the argument indicates the intuition that $HALT$ is placed at the end of the code produced by $comp^A$:

$comp :: Expr \rightarrow Code$
 $comp\ e = comp^A\ e \triangleright HALT$

The following examples illustrate the workings of the compiler $comp$:

$comp (Add (Val\ 2) (Val\ 3)) \rightsquigarrow PUSH\ 2 \triangleright PUSH\ 3 \triangleright ADD \triangleright HALT$
 $comp (Catch (Val\ 2) (Val\ 3)) \rightsquigarrow MARK (PUSH\ 3 \triangleright HALT)$
 $\quad \quad \quad \triangleright PUSH\ 2 \triangleright UNMARK \triangleright HALT$
 $comp (Catch\ Throw (Val\ 3)) \rightsquigarrow MARK (PUSH\ 3 \triangleright HALT) \triangleright THROW$

For the virtual machine that executes the code produced by the above compiler, we use the following type for the stack:

```
type Stack = [Item]
data Item = VAL Int | HAN (Stack → Stack)
```

This type deviates slightly from the one for the virtual machine defined by Hutton and Wright [10]. Instead of having the code of an exception handler on the stack (constructor *HAN*), we have the continuation of the virtual machine on the stack. This will simplify the proof as we shall see later on. However, this type and the accompanying definition of the virtual machine that is given below is exactly the result of the calculation given by Bahr and Hutton [3] just before the last calculation step (which then yields the virtual machine of Hutton and Wright [10]). The virtual machine that works on this stack is defined as follows:

```
exec :: Code → Stack → Stack
exec (PUSH n c) s = exec c (VAL n : s)
exec (ADD c) s = case s of
    (VAL m : VAL n : t) → exec c (VAL (n + m) : t)
exec THROW s = unwind s
exec (MARK h c) s = exec c (HAN (exec h) : s)
exec (UNMARK c) s = case s of (x : HAN _ : t) → exec c (x : t)
exec HALT s = s

unwind :: Stack → Stack
unwind [] = []
unwind (VAL _ : s) = unwind s
unwind (HAN h : s) = h s
```

The virtual machine does what is expected from the informal semantics that we have given above. The semantics of *MARK*, however, may seem counterintuitive at first: as mentioned above, *MARK* does not put the handler code on the stack but rather the continuation that is obtained by executing it. Consequently, when the unwinding of the stack reaches a handler *h* on the stack, this handler *h* is directly applied to the remainder of the stack. This slight deviation from the semantics of Hutton and Wright [10] makes sure that *exec* is in fact a fold.

We will not go into the details of the correctness proof for the compiler *comp*. One can show that it satisfies the following correctness property [3]:

Theorem 1 (compiler correctness).

$$\text{exec } (\text{comp } e) [] = \text{conv } (\text{eval } e) \quad \text{for all } e :: \text{Expr}$$

where $\text{conv } (\text{Just } n) = [\text{Val } n]$
 $\text{conv } \text{Nothing} = []$

That is, in particular, we have that

$$\text{exec } (\text{comp } e) [] = [\text{Val } n] \iff \text{eval } e = \text{Just } n$$

While the compiler has the nice property that it can be derived from the language semantics, the code that it produces is quite unrealistic. Note the duplication that occurs for generating the code for *Catch*: the continuation code c is inserted both after the handler code (in $comp^A h c$) and after the *UNMARK* instruction. This is necessary since the code c may have to be executed regardless whether an exception is thrown in the scope x of the *Catch* or not.

This duplication can be avoided by using explicit jumps in the code. Instead of duplicating code, jumps to a single copy of the code are inserted. However, this complicates both the implementation of the compiler and its correctness proof [10]. Also the derivation of such a compiler by calculation is equally cumbersome.

The approach that we suggest in this paper takes the above compiler and derives a slightly different variant that instead of a tree structure produces a graph structure. Along with the compiler we derive a virtual machine that also works on the graph structure. The two variants of the compiler and its companion virtual machine only differ in the sharing that the graph variant provides. This fact allows us to derive the correctness of the graph-based compiler very easily from the correctness of the original tree-based compiler.

3 From Trees to Graphs

Before we derive the graph-based compiler and the corresponding virtual machine, we restructure the definition of the original compiler and the corresponding virtual machine. This will smoothen the process and simplify the presentation.

3.1 Preparations

Instead of defining the type *Code* directly, we represent it as the initial algebra of a functor. To distinguish this representation from the graph representation we introduce later, we use the name *Tree* for the initial algebra construction.

data $Tree\ f = In\ (f\ (Tree\ f))$

The functor that induces the initial algebra that we shall use for representing the target language is easily obtained from the original *Code* data type:

data $Code_F\ a = PUSH_F\ Int\ a \mid ADD_F\ a \quad \mid HALT_F$
 $\mid MARK_F\ a\ a \mid UNMARK_F\ a \mid THROW_F$

The type representing the target code is thus $Tree\ Code_F$, which is isomorphic to *Code* modulo non-strictness. We proceed by reformulating the definition of *comp* to work on the type $Tree\ Code_F$ instead of *Code*:

$comp_{\top}^A :: Expr \rightarrow Tree\ Code_F \rightarrow Tree\ Code_F$
 $comp_{\top}^A\ (Val\ n) \quad c = PUSH_{\top}\ n \triangleright c$
 $comp_{\top}^A\ (Add\ x\ y) \quad c = comp_{\top}^A\ x \triangleright comp_{\top}^A\ y \triangleright ADD_{\top} \triangleright c$

$$\begin{aligned}
\text{comp}_{\top}^{\Delta} \text{ Throw } & c = \text{THROW}_{\top} \\
\text{comp}_{\top}^{\Delta} (\text{Catch } x \ h) & c = \text{MARK}_{\top} (\text{comp}_{\top}^{\Delta} h \triangleright c) \triangleright \text{comp}_{\top}^{\Delta} x \triangleright \text{UNMARK}_{\top} \triangleright c \\
\text{comp}_{\top} & :: \text{Expr} \rightarrow \text{Tree Code}_{\mathbb{F}} \\
\text{comp}_{\top} e & = \text{comp}_{\top}^{\Delta} e \triangleright \text{HALT}_{\top}
\end{aligned}$$

Note that we do not use the constructors of $\text{Code}_{\mathbb{F}}$ directly, but instead we use *smart constructors* that also apply the constructor In of the type constructor Tree . These smart constructors serve as drop-in replacements for the constructors of the original Code data type. For example, PUSH_{\top} is defined as follows:

$$\begin{aligned}
\text{PUSH}_{\top} & :: \text{Int} \rightarrow \text{Tree Code}_{\mathbb{F}} \rightarrow \text{Tree Code}_{\mathbb{F}} \\
\text{PUSH}_{\top} i \ c & = \text{In} (\text{PUSH}_{\mathbb{F}} i \ c)
\end{aligned}$$

Lastly, we also reformulate the semantics of the target language, i.e. we define the function exec on the type $\text{Tree Code}_{\mathbb{F}}$. To do this, we use the following definition of a fold on an initial algebra:

$$\begin{aligned}
\text{fold} & :: \text{Functor } f \Rightarrow (f \ r \rightarrow r) \rightarrow \text{Tree } f \rightarrow r \\
\text{fold } \text{alg} (\text{In } t) & = \text{alg} (\text{fmap } (\text{fold } \text{alg}) t)
\end{aligned}$$

The definition of the semantics is a straightforward transcription of the definition of exec into an algebra:

$$\begin{aligned}
\text{execAlg} & :: \text{Code}_{\mathbb{F}} (\text{Stack} \rightarrow \text{Stack}) \rightarrow \text{Stack} \rightarrow \text{Stack} \\
\text{execAlg} (\text{PUSH}_{\mathbb{F}} n \ c) & \quad s = c (\text{VAL } n : s) \\
\text{execAlg} (\text{ADD}_{\mathbb{F}} c) & \quad s = \mathbf{case} \ s \ \mathbf{of} \\
& \quad \quad (\text{VAL } m : \text{VAL } n : t) \rightarrow c (\text{VAL } (n + m) : t) \\
\text{execAlg } \text{THROW}_{\mathbb{F}} & \quad s = \text{unwind } s \\
\text{execAlg} (\text{MARK}_{\mathbb{F}} h \ c) & \quad s = c (\text{HAN } h : s) \\
\text{execAlg} (\text{UNMARK}_{\mathbb{F}} c) & \quad s = \mathbf{case} \ s \ \mathbf{of} \ (x : \text{HAN } _ : t) \rightarrow c (x : t) \\
\text{execAlg } \text{HALT}_{\mathbb{F}} & \quad s = s \\
\text{exec}_{\top} & :: \text{Tree Code}_{\mathbb{F}} \rightarrow \text{Stack} \rightarrow \text{Stack} \\
\text{exec}_{\top} & = \text{fold } \text{execAlg}
\end{aligned}$$

From the correctness of the original compiler from Section 2, as expressed in Theorem 1, we obtain the correctness of our reformulation of the implementation:

Corollary 1 (correctness of comp_{\top}).

$$\text{exec}_{\top} (\text{comp}_{\top} e) [] = \text{conv} (\text{eval } e) \quad \text{for all } e :: \text{Expr}$$

Proof. Let $\phi :: \text{Code} \rightarrow \text{Tree Code}_{\mathbb{F}}$ be the function that recursively maps each constructor of Code to the corresponding smart constructor of $\text{Tree Code}_{\mathbb{F}}$. We can easily check that comp_{\top} and exec_{\top} are equivalent to the original functions comp respectively exec via ϕ , i.e.

$$\text{comp}_{\top} = \phi \circ \text{comp} \quad \text{and} \quad \text{exec}_{\top} \circ \phi = \text{exec}$$

Consequently, we have that $\text{exec}_{\top} \circ \text{comp}_{\top} = \text{exec} \circ \text{comp}$, and thus the corollary follows from Theorem 1. \square

3.2 Deriving a Graph-Based Compiler

Finally, we turn to the graph-based implementation of the compiler. Essentially, this implementation is obtained from $comp_{\top}$ by replacing the type $Tree\ Code_{\mathbb{F}}$ with a type $Graph\ Code_{\mathbb{F}}$, which instead of a tree structure has a graph structure, and using explicit sharing instead of duplication.

In order to define graphs over a functor, we use the representation of Oliveira and Cook [18] called *structured graphs*. Put simply, a structured graph is a tree with added sharing facilitated by let bindings. In turn, let bindings are represented using parametric higher-order abstract syntax [5].

$$\begin{aligned} \mathbf{data}\ Graph'\ f\ v = & \mathit{GIn}\ (f\ (Graph'\ f\ v)) \\ & | \mathit{Let}\ (Graph'\ f\ v)\ (v \rightarrow Graph'\ f\ v) \\ & | \mathit{Var}\ v \end{aligned}$$

The first constructor has the same structure as the constructor of the $Tree$ type constructor. The other two constructors will allow us to express let bindings: $\mathit{Let}\ g\ (\lambda x \rightarrow h)$ binds g to the metavariable x in h . Metavariables bound in a let binding have type v ; the only way to use them is with the constructor Var . To enforce this invariant, the type variable v is made polymorphic:

$$\mathbf{newtype}\ Graph\ f = \mathit{MkGraph}\ (\forall v.\ Graph'\ f\ v)$$

We shall use the type constructor $Graph$ (and $Graph'$) as a replacement for $Tree$. For the purposes of our compiler we only need acyclic graphs. That is why we only consider non-recursive let bindings as opposed to the more general structured graphs of Oliveira and Cook [18]. This restriction to non-recursive let bindings is crucial for the reasoning principle that we use to prove correctness.

We can use the graph type almost as a drop-in replacement for the tree type. The only thing that we need to do is to use smart constructors that use the constructor GIn instead of In , e.g.

$$\begin{aligned} \mathit{PUSH}_{\mathbb{G}} &:: \mathit{Int} \rightarrow Graph'\ Code_{\mathbb{F}}\ v \rightarrow Graph'\ Code_{\mathbb{F}}\ v \\ \mathit{PUSH}_{\mathbb{G}}\ i\ c &= \mathit{GIn}\ (\mathit{PUSH}_{\mathbb{F}}\ i\ c) \end{aligned}$$

From the type of the smart constructors we can observe that graphs are constructed using the type constructor $Graph'$, not $Graph$. Only after the construction of the graph is completed, the constructor $\mathit{MkGraph}$ is applied in order to obtain a graph of type $Graph\ Code_{\mathbb{F}}$.

The definition of $comp_{\top}^{\mathbb{A}}$ can be transcribed into graph style by simply using the abovementioned smart constructors instead:

$$\begin{aligned} comp_{\mathbb{G}}^{\mathbb{A}} &:: \mathit{Expr} \rightarrow Graph'\ Code_{\mathbb{F}}\ a \rightarrow Graph'\ Code_{\mathbb{F}}\ a \\ comp_{\mathbb{G}}^{\mathbb{A}}\ (\mathit{Val}\ n) &\quad c = \mathit{PUSH}_{\mathbb{G}}\ n \triangleright c \\ comp_{\mathbb{G}}^{\mathbb{A}}\ (\mathit{Add}\ x\ y) &\quad c = comp_{\mathbb{G}}^{\mathbb{A}}\ x \triangleright comp_{\mathbb{G}}^{\mathbb{A}}\ y \triangleright \mathit{ADD}_{\mathbb{G}} \triangleright c \\ comp_{\mathbb{G}}^{\mathbb{A}}\ (\mathit{Throw}) &\quad c = \mathit{THROW}_{\mathbb{G}} \\ comp_{\mathbb{G}}^{\mathbb{A}}\ (\mathit{Catch}\ x\ h) &\quad c = \mathit{MARK}_{\mathbb{G}}\ (comp_{\mathbb{G}}^{\mathbb{A}}\ h \triangleright c) \triangleright comp_{\mathbb{G}}^{\mathbb{A}}\ x \triangleright \mathit{UNMARK}_{\mathbb{G}} \triangleright c \end{aligned}$$

The above is a one-to-one transcription of $comp_{\top}^A$. But this is not what we want. We want to make use of the fact that the target language allows sharing. In particular, we want to get rid of the duplication in the code generated for *Catch*.

We can avoid this duplication by simply using a let binding to replace the two occurrences of c with a metavariable c' that is then bound to c . The last equation for $comp_{\mathbb{G}}^A$ is thus rewritten as follows:

$$comp_{\mathbb{G}}^A (Catch\ x\ h)\ c = Let\ c\ (\lambda c' \rightarrow MARK_{\mathbb{G}} (comp_{\mathbb{G}}^A\ h \triangleright Var\ c') \\ \triangleright comp_{\mathbb{G}}^A\ x \triangleright UNMARK_{\mathbb{G}} \triangleright Var\ c')$$

The right-hand side for the case *Catch* $x\ h$ has now only one occurrence of c .

The final code generator function $comp_{\mathbb{G}}^A$ is then obtained by supplying $HALT_{\mathbb{G}}$ as the initial value of the code continuation and wrapping the result with the *MkGraph* constructor so as to return a result of type *Graph Code_F*:

$$comp_{\mathbb{G}} :: Expr \rightarrow Graph\ Code_F \\ comp_{\mathbb{G}}\ e = MkGraph (comp_{\mathbb{G}}^A\ e \triangleright HALT_{\mathbb{G}})$$

To illustrate the difference between $comp_{\mathbb{G}}$ and $comp_{\top}$, we apply both of them to an example expression $e = Add (Catch (Val\ 1) (Val\ 2)) (Val\ 3)$:

$$comp_{\top}\ e \rightsquigarrow MARK_{\top} (PUSH_{\top}\ 2 \triangleright PUSH_{\top}\ 3 \triangleright ADD_{\top} \triangleright HALT_{\top}) \\ \triangleright PUSH_{\top}\ 1 \triangleright UNMARK_{\top} \triangleright PUSH_{\top}\ 3 \triangleright ADD_{\top} \triangleright HALT_{\top} \\ comp_{\mathbb{G}}\ e \rightsquigarrow MkGraph (Let (PUSH_{\mathbb{G}}\ 3 \triangleright ADD_{\mathbb{G}} \triangleright HALT_{\mathbb{G}}) (\lambda v \rightarrow \\ MARK_{\mathbb{G}} (PUSH_{\mathbb{G}}\ 2 \triangleright Var\ v) \triangleright PUSH_{\mathbb{G}}\ 1 \triangleright UNMARK_{\mathbb{G}} \triangleright Var\ v))$$

Note that $comp_{\top}$ duplicates the code fragment $PUSH_{\top}\ 3 \triangleright ADD_{\top} \triangleright HALT_{\top}$, which is supposed to be executed after the catch expression, whereas $comp_{\mathbb{G}}$ binds this code fragment to a metavariable v , which is then used as a substitute.

The recursion schemes on structured graphs make use of the parametricity in the metavariable type as well. The general fold over graphs as given by Oliveira and Cook [18] is defined as follows:¹

$$gfold :: Functor\ f \Rightarrow (v \rightarrow r) \rightarrow (r \rightarrow (v \rightarrow r) \rightarrow r) \rightarrow (f\ r \rightarrow r) \rightarrow \\ Graph\ f \rightarrow r \\ gfold\ v\ l\ i\ (MkGraph\ g) = trans\ g \\ \mathbf{where}\ trans\ (Var\ x) = v\ x \\ trans\ (Let\ e\ f) = l\ (trans\ e)\ (trans\ \circ\ f) \\ trans\ (GIn\ t) = i\ (fmap\ trans\ t)$$

The combinator takes three functions, which are used to interpret the three constructors of *Graph'*. This general form is needed for example if we want to transform the graph representation into a linearised form [2], but for our purposes we only need a simple special case of it:

¹ Oliveira and Cook [18] considered the more general case of cyclic graphs, the definition of *gfold* given here is specialised to the case of acyclic graphs.

$$\begin{aligned}
ufold &:: \text{Functor } f \Rightarrow (f \ r \rightarrow r) \rightarrow \text{Graph } f \rightarrow r \\
ufold &= \text{gfold } id \ (\lambda e \ f \rightarrow f \ e)
\end{aligned}$$

Note that the type signature is identical to the one for *fold* except for the use of *Graph* instead of *Tree*. Thus, we can reuse the algebra *execAlg* from Section 3.1, which defines the semantics of *Tree Code_F*, in order to define the semantics of *Graph Code_F*:

$$\begin{aligned}
exec_G &:: \text{Graph Code}_F \rightarrow \text{Stack} \rightarrow \text{Stack} \\
exec_G &= ufold \ execAlg
\end{aligned}$$

4 Correctness Proof

In this section we shall prove that the graph-based compiler that we defined in Section 3.2 is indeed correct. This turns out to be rather simple: we derive the correctness property for *comp_G* from the correctness property for *comp_T*. The simplicity of the argument is rooted in the fact that *comp_T* is the same as *comp_G* followed by unravelling. In other words, *comp_G* only differs from *comp_T* in that it adds sharing – as expected.

4.1 Compiler Correctness by Unravelling

Before we prove this relation between *comp_T* and *comp_G*, we need to specify what unravelling means:

$$\begin{aligned}
unravel &:: \text{Functor } f \Rightarrow \text{Graph } f \rightarrow \text{Tree } f \\
unravel &= ufold \ In
\end{aligned}$$

While this definition is nice and compact, we gain more insight into what it actually does by unfolding it:

$$\begin{aligned}
unravel &:: \text{Functor } f \Rightarrow \text{Graph } f \rightarrow \text{Tree } f \\
unravel \ (MkGraph \ g) &= unravel' \ g \\
unravel' &:: \text{Functor } f \Rightarrow \text{Graph}' \ f \ (\text{Tree } f) \rightarrow \text{Tree } f \\
unravel' \ (Var \ x) &= x \\
unravel' \ (Let \ e \ f) &= unravel' \ (f \ (unravel' \ e)) \\
unravel' \ (GIn \ t) &= In \ (fmap \ unravel' \ t)
\end{aligned}$$

We can see that *unravel* simply replaces *GIn* with *In*, and applies the function argument *f* of a let binding to the bound value *e*. For example, we have that

$$\begin{aligned}
&MkGraph \ (Let \ (PUSH_G \ 2 \triangleright \ HALT_G) \ (\lambda v \rightarrow MARK_G \ (Var \ v) \triangleright \ Var \ v)) \\
&\overset{unravel}{\rightsquigarrow} MARK_T \ (PUSH_T \ 2 \triangleright \ HALT_T) \triangleright \ PUSH_T \ 2 \triangleright \ HALT_T
\end{aligned}$$

We can now formulate the relation between *comp_T* and *comp_G*:

Lemma 1. $comp_{\top} = unravel \circ comp_{\mathbb{G}}$

This lemma, which we shall prove at the end of this section, is one half of the argument for deriving the correctness property for $comp_{\mathbb{G}}$. The other half is the property that $exec_{\top}$ and $exec_{\mathbb{G}}$ have the converse relationship, viz.

$$exec_{\mathbb{G}} = exec_{\top} \circ unravel$$

Proving this property is much simpler, though, because it follows from a more general property of *fold*.

Theorem 2. *Given a strictly positive functor f , a type c , and $alg :: f\ c \rightarrow c$, we have the following:*

$$unfold\ alg = fold\ alg \circ unravel$$

The equality $exec_{\mathbb{G}} = exec_{\top} \circ unravel$ is an instance of Theorem 2 where $alg = execAlg$. We defer discussion of the proof of this theorem until Section 4.2.

We derive the correctness of $comp_{\mathbb{G}}$ by combining Lemma 1 and Theorem 2:

Theorem 3 (correctness of $comp_{\mathbb{G}}$).

$$exec_{\mathbb{G}} (comp_{\mathbb{G}}\ e)\ [] = conv (eval\ e) \quad \text{for all } e :: Expr$$

Proof. $exec_{\mathbb{G}} (comp_{\mathbb{G}}\ e)\ [] = exec_{\top} (unravel (comp_{\mathbb{G}}\ e)\ [])$ (Theorem 2)
 $= exec_{\top} (comp_{\top}\ e)\ []$ (Lemma 1)
 $= conv (eval\ e)$ (Corollary 1) \square

We conclude this section by giving the proof of Lemma 1.

Proof (of Lemma 1). Instead of proving the equation directly, we prove the following equation for all $e :: Expr$ and $c :: Graph'\ Code_{\mathbb{F}} (Tree\ Code_{\mathbb{F}})$:

$$comp_{\top}^{\mathbb{A}}\ e \triangleright unravel'\ c = unravel'\ (comp_{\mathbb{G}}^{\mathbb{A}}\ e \triangleright c) \quad (1)$$

In particular, the above equation holds for all $c :: \forall v . Graph'\ Code_{\mathbb{F}}\ v$. Thus, the lemma follows from the above equation as follows:

$$\begin{aligned} & comp_{\top}\ e \\ = & \{ \text{definition of } comp_{\top} \} \\ & comp_{\top}^{\mathbb{A}}\ e \triangleright HALT_{\top} \\ = & \{ \text{definition of } unravel' \} \\ & comp_{\top}^{\mathbb{A}}\ e \triangleright unravel'\ HALT_{\mathbb{G}} \\ = & \{ \text{Equation (1)} \} \\ & unravel'\ (comp_{\mathbb{G}}^{\mathbb{A}}\ e \triangleright HALT_{\mathbb{G}}) \\ = & \{ \text{definition of } unravel \} \\ & unravel (MkGraph (comp_{\mathbb{G}}^{\mathbb{A}}\ e \triangleright HALT_{\mathbb{G}})) \\ = & \{ \text{definition of } comp_{\mathbb{G}} \} \\ & unravel (comp_{\mathbb{G}}\ e) \end{aligned}$$

We prove (1) by induction on e :

– Case $e = \text{Val } n$:

$$\begin{aligned}
& \text{unravel}' (comp_G^A (\text{Val } n) \triangleright c) \\
= & \{ \text{definition of } comp_G^A \} \\
& \text{unravel}' (\text{PUSH}_G n \triangleright c) \\
= & \{ \text{definition of } unravel' \} \\
& \text{PUSH}_T n \triangleright unravel' c \\
= & \{ \text{definition of } comp_T^A \} \\
& comp_T^A (\text{Val } n) \triangleright unravel' c
\end{aligned}$$

– Case $e = \text{Throw}$:

$$\begin{aligned}
& \text{unravel}' (comp_G^A \text{Throw} \triangleright c) \\
= & \{ \text{definition of } comp_G^A \} \\
& \text{unravel}' \text{THROW}_G \\
= & \{ \text{definition of } unravel' \} \\
& \text{THROW}_T \\
= & \{ \text{definition of } comp_T^A \} \\
& comp_T^A \text{Throw} \triangleright unravel' c
\end{aligned}$$

– Case $e = \text{Add } x y$:

$$\begin{aligned}
& \text{unravel}' (comp_G^A (\text{Add } x y) \triangleright c) \\
= & \{ \text{definition of } comp_G^A \} \\
& \text{unravel}' (comp_G^A x \triangleright comp_G^A y \triangleright \text{ADD}_G \triangleright c) \\
= & \{ \text{induction hypothesis} \} \\
& comp_T^A x \triangleright unravel' (comp_G^A y \triangleright \text{ADD}_G \triangleright c) \\
= & \{ \text{induction hypothesis} \} \\
& comp_T^A x \triangleright comp_T^A y \triangleright unravel' (\text{ADD}_G \triangleright c) \\
= & \{ \text{definition of } unravel' \} \\
& comp_T^A x \triangleright comp_T^A y \triangleright \text{ADD}_T \triangleright unravel' c \\
= & \{ \text{definition of } comp_T^A \} \\
& comp_T^A (\text{Add } x y) \triangleright unravel' c
\end{aligned}$$

– Case $e = \text{Catch } x h$:

$$\begin{aligned}
& \text{unravel}' (comp_G^A (\text{Catch } x h) \triangleright c) \\
= & \{ \text{definition of } comp_G^A \} \\
& \text{unravel}' (\text{Let } c (\lambda c' \rightarrow \text{MARK}_G (comp_G^A h \triangleright \text{Var } c') \\
& \quad \triangleright comp_G^A x \triangleright \text{UNMARK}_G \triangleright \text{Var } c')) \\
= & \{ \text{definition of } unravel' \text{ and } \beta\text{-reduction} \} \\
& \text{unravel}' (\text{MARK}_G (comp_G^A h \triangleright \text{Var } (unravel' c)) \\
& \quad \triangleright comp_G^A x \triangleright \text{UNMARK}_G \triangleright \text{Var } (unravel' c)) \\
= & \{ \text{definition of } unravel' \} \\
& \text{MARK}_T (\text{unravel}' (comp_G^A h \triangleright \text{Var } (unravel' c))) \\
& \quad \triangleright unravel' (comp_G^A x \triangleright \text{UNMARK}_G \triangleright \text{Var } (unravel' c)) \\
= & \{ \text{induction hypothesis} \} \\
& \text{MARK}_T (comp_T^A h \triangleright unravel' (\text{Var } (unravel' c))) \\
& \quad \triangleright comp_T^A x \triangleright unravel' (\text{UNMARK}_G \triangleright \text{Var } (unravel' c)) \\
= & \{ \text{definition of } unravel' \} \\
& \text{MARK}_T (comp_T^A h \triangleright unravel' c) \triangleright comp_T^A x \triangleright \text{UNMARK}_T \triangleright unravel' c \\
= & \{ \text{definition of } comp_T^A \} \\
& comp_T^A (\text{Catch } x h) \triangleright unravel' c
\end{aligned}$$

□

4.2 Proof of Theorem 2

Theorem 2 states that folding a structured graph $g :: \text{Graph } f$ over a strictly positive functor f with an algebra alg yields the same result as first unravelling g and then folding the resulting tree with alg , i.e.

$$unfold\ alg = fold\ alg \circ unravel$$

Since $unravel$ is defined as $unfold\ In$, the above equality follows from a more general law of folds over algebraic data types, known as *short cut fusion* [8]:

$$b\ alg = fold\ alg\ (b\ In) \quad \text{for all } b :: \forall c. (f\ c \rightarrow c) \rightarrow c$$

This law holds for all strictly positive functors f as proved by Johann [12]. Essential for its correctness is the polymorphic type of b .

For any given graph $g :: \text{Graph } f$, we can instantiate b with the function $\lambda a \rightarrow unfold\ a\ g$, which yields that

$$(\lambda a \rightarrow unfold\ a\ g)\ alg = fold\ alg\ ((\lambda a \rightarrow unfold\ a\ g)\ In)$$

Note that $\lambda a \rightarrow unfold\ a\ g$ has indeed the required polymorphic type. After applying beta-reduction, we obtain the equation

$$unfold\ alg\ g = fold\ alg\ (unfold\ In\ g)$$

Since g was chosen arbitrarily, and $unravel$ is defined as $unfold\ In$, we thus obtain the equation as stated in Theorem 2:

$$unfold\ alg = fold\ alg \circ unravel$$

5 Other Approaches

5.1 Other Graph Representations

The technique presented here is not necessarily dependent on the particular representation of graphs that we chose. However, while other representations are conceivable, structured graphs have two properties that make them a suitable choice for this application: (1) they have a simple representation in Haskell and (2) they provide a convenient interface for introducing sharing, viz. variable binding in the host language.

Nevertheless, in other circumstances a different representation may be advantageous. For example the use of higher-order abstract syntax may have a negative impact on performance in practical applications. Moreover, the necessity of reasoning over parametricity may be inconvenient for a formalisation of the proofs in a proof assistant.

Therefore, we also studied an alternative representation of graphs that uses de Bruijn indices for encoding binders instead of parametric higher-order abstract syntax (PHOAS). To this end, we have used the technique proposed by

Bernardy and Pouillard [4] to provide a PHOAS interface to this graph representation. This allows us to use essentially the same simple definition of the graph-based compiler as presented in Section 3.2. Using this representation of graphs – PHOAS interface on the outside, de Bruijn indices under the hood – we formalised the proofs presented here in the Coq theorem prover².

5.2 A Monadic Approach

We briefly describe a variant of our technique that is based on free monads and a corresponding monadic graph structure. The general approach of this variant is similar to what we have seen thus far; however, the monadic structure simplifies some of the proofs. The details can be found in the companion report [2].

The underlying idea, originally developed by Matsuda et al. [15], is to replace a function f with accumulation parameters by a function f' that produces a *context* with the property that

$$f\ x\ a_1 \dots a_n = (f'\ x)(a_1, \dots, a_n)$$

That is, we obtain the result of the original function f by plugging in the accumulation arguments a_1, \dots, a_n in to the context that f' produces.

In order to represent contexts, we use a free monad type $Tree_M$ instead of a tree type $Tree$, where $Tree_M$ is obtained from $Tree$ by adding a constructor of type $a \rightarrow Tree_M\ f\ a$. A context with n holes is represented by a type $Tree_M\ f\ (Fin\ n)$ – where $Fin\ n$ is a type with exactly n distinct inhabitants – and context application is represented by the monadic bind operator \gg . The compiler is then reformulated as follows – using the shorthand $hole = return\ ()$:

$$\begin{aligned} comp_M^C &:: Expr \rightarrow Tree_M\ Code_F\ () \\ comp_M^C\ (Val\ n) &= PUSH_M\ n\ hole \\ comp_M^C\ (Add\ x\ y) &= comp_M^C\ x \gg comp_M^C\ y \gg ADD_M\ hole \\ comp_M^C\ (Throw) &= THROW_M \\ comp_M^C\ (Catch\ x\ h) &= MARK_M\ (comp_M^C\ h)\ (comp_M^C\ x \gg UNMARK_M\ hole) \end{aligned}$$

As we only have a single accumulator for the compiler, we use the type $() \simeq Fin\ 1$ to express that there is exactly one type of hole.

Also graphs can be given monadic structure by adding a constructor of type $a \rightarrow Graph'_M\ f\ v\ a$ to the data type $Graph'$. And the compiler $comp_G^A$ can be reformulated in terms of this type accordingly.

We can define fold combinators for the monadic structures as well. The virtual machines are thus easily adapted to this monadic style by simply reusing the same algebra $execAlg$. Again, one half of the correctness proof follows from a generic theorem about folds corresponding to Theorem 2. The other half of the proof can be simplified. In the corresponding proof of Lemma 1 it suffices to show the following simpler equation, in which $unravel'$ only appears once:

$$comp_T^A = unravel' \circ comp_G^A$$

² Available from the author's web site.

This simplifies the induction proof. While this proof requires an additional lemma, viz. that unravelling distributes over \gg , this lemma can be proved (once and for all) for any strictly positive functor f :

$$\text{unravel}' (g_1 \gg g_2) = \text{unravel}' g_1 \gg \text{unravel}' g_2$$

Unfortunately, we cannot exploit short cut fusion to prove this lemma because it involves a genuine graph transformation, viz. \gg on graphs. However, with the representation mentioned in Section 5.1, we can prove it by induction.

Note that the full monadic structure of Tree_M and Graph_M is not needed for our example compiler since we only use the simple bind operator \gg , not $\gg\equiv$. However, a different compiler implementation may use more than one accumulation parameter (for example an additional code continuation that contains the current exception handler), for which we need the more general bind operator.

6 Concluding Remarks

6.1 Related Work

Compiler verification is still a hard problem and in this paper we only cover one – but arguably the central – part of a compiler, viz. the translation of a high-level language to a low-level language. The literature on the topic of compiler verification is vast (e.g. see the survey of Dave [7]). More recent work has shown impressive results in verification of a realistic compiler for the C language [13]. But there are also efforts in verifying compilers for higher-level languages (e.g. by Chlipala [6]).

This paper, however, focuses on identifying simple but powerful techniques for reasoning about compilers rather than engineering massive proofs for full-scale compilers. Our contributions thus follow the work on calculating compilers [21, 16, 1] as well as Hutton and Wright’s work on equational reasoning about compilers [10, 11].

Structured graphs have been used in the setting of programming language implementation before: Oliveira and Löh [17] used structured graphs to represent embedded domain-specific languages. That is, graphs are used for the representation of the source language. Graph structures used for representing intermediate languages in a compiler typically employ pointers (e.g. Ramsey and Dias [20]) or labels (e.g. Ramsey et al. [19]). We are not aware of any work that makes use of higher-order abstract syntax or de Bruijn indices in the representation of graph structures in this setting.

6.2 Discussion and Future Work

The underlying goal of our method is to separate the transformation to the target language from the need to generate fresh names for representing jumps. For a full compiler, we still have to deal with explicit jumps eventually, but we can do so in isolation. That is, (1) we have to define a function

$linearise :: Graph\ Code_F \rightarrow Code_L$

that transforms the graph-based representation into a linear representation of the target language, and (2) we have to prove that it preserves the semantics. The proof can focus solely on the aspect of fresh names and explicit jumps. Since *linearise* is trivial for all cases except for the let bindings of the graph representation, we expect that the proof can be made independently of the actual language under consideration.

While our method reduces the proof obligations for the graph-based compiler considerably, there is still room for improvement. Indeed, we only require a simple induction proof showing the equality $comp_T = unravel \circ comp_G$. But since the two compiler variants differ only in the sharing they produce, one would hope the proof obligation could be further reduced to the only interesting case, i.e. the case for *Catch* in our example. In a proof assistant such as Coq, we can indeed take care of all the other cases with a single tactic and focus on the interesting case. However, it would be desirable to have a more systematic approach that captures this intuitive understanding.

A shortcoming of our method is its limitation to acyclic graphs. Nevertheless, the implementation part of our method easily generalises to cyclic structures, which permits compilation of cyclic control structures like loops. Corresponding correctness proofs, however, need a different reasoning principle.

Acknowledgements

The author is indebted to Janis Voigtländer for his many helpful suggestions to improve both the substance and the presentation of this paper. In particular, the idea to use short cut fusion to prove Theorem 2 was his. The author would like to thank Nicolas Pouillard and Daniel Gustafsson for their assistance in the accompanying Coq development.

This work was supported by the Danish Council for Independent Research, Grant 12-132365, “Efficient Programming Language Development and Evolution through Modularity”.

References

- [1] Ager, M.S., Biernacki, D., Danvy, O., Midtgaard, J.: From interpreter to compiler and virtual machine: A functional derivation. Tech. Rep. RS-03-14, BRICS, Department of Computer Science, University of Aarhus (2003)
- [2] Bahr, P.: Proving correctness of compilers using structured graphs (extended version). Tech. rep., University of Copenhagen (2014)
- [3] Bahr, P., Hutton, G.: Calculating correct compilers (2014), unpublished manuscript
- [4] Bernardy, J.P., Pouillard, N.: Names for free: polymorphic views of names and binders. In: Proceedings of the 2013 ACM SIGPLAN symposium on Haskell. pp. 13–24. ACM (2013)

- [5] Chlipala, A.: Parametric higher-order abstract syntax for mechanized semantics. In: *Proceeding of the 13th ACM SIGPLAN International Conference on Functional Programming*. pp. 143–156. ACM (2008)
- [6] Chlipala, A.: A verified compiler for an impure functional language. In: *Proceedings of the 37th annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. pp. 93–106. ACM (2010)
- [7] Dave, M.A.: *Compiler verification: a bibliography*. SIGSOFT Softw. Eng. Notes 28(6), 2–2 (Nov 2003)
- [8] Gill, A., Launchbury, J., Peyton Jones, S.L.: A short cut to deforestation. In: *Proceedings of the Conference on Functional Programming Languages and Computer Architecture*. pp. 223–232. ACM (1993)
- [9] Hutton, G.: *Programming in Haskell*, vol. 2. Cambridge University Press Cambridge (2007)
- [10] Hutton, G., Wright, J.: Compiling exceptions correctly. In: Kozen, D. (ed.) *Mathematics of Program Construction. Lecture Notes in Computer Science*, vol. 3125, pp. 211–227. Springer Berlin / Heidelberg (2004)
- [11] Hutton, G., Wright, J.: What is the meaning of these constant interruptions? *J. Funct. Program.* 17(06), 777–792 (2007)
- [12] Johann, P.: A generalization of short-cut fusion and its correctness proof. *Higher Order Symbol. Comput.* 15(4), 273–300 (2002)
- [13] Leroy, X.: Formal certification of a compiler back-end or: programming a compiler with a proof assistant. In: *Conference record of the 33rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. pp. 42–54. ACM (2006)
- [14] Marlow, S.: *Haskell 2010 language report* (2010)
- [15] Matsuda, K., Inaba, K., Nakano, K.: Polynomial-time inverse computation for accumulative functions with multiple data traversals. In: *Proceedings of the ACM SIGPLAN 2012 workshop on Partial evaluation and program manipulation*. pp. 5–14. ACM (2012)
- [16] Meijer, E.: *Calculating Compilers*. Ph.D. thesis, Katholieke Universiteit Nijmegen (1992)
- [17] Oliveira, B.C.d.S., Löb, A.: Abstract syntax graphs for domain specific languages. In: *Proceedings of the ACM SIGPLAN 2013 Workshop on Partial Evaluation and Program Manipulation*. pp. 87–96. ACM (2013)
- [18] Oliveira, B.C., Cook, W.R.: Functional programming with structured graphs. In: *Proceedings of the 17th ACM SIGPLAN International Conference on Functional Programming*. pp. 77–88. ACM (2012)
- [19] Ramsey, N., Dias, J.a., Peyton Jones, S.: Hoopl: a modular, reusable library for dataflow analysis and transformation. In: *Proceedings of the third ACM Haskell symposium on Haskell*. pp. 121–134. ACM (2010)
- [20] Ramsey, N., Dias, J.: An applicative control-flow graph based on huet’s zipper. In: *Proceedings of the ACM-SIGPLAN Workshop on ML*. pp. 105 – 126 (2006)
- [21] Wand, M.: Deriving target code as a representation of continuation semantics. *ACM Trans. Program. Lang. Syst.* 4(3), 496–517 (Jul 1982)