

Conversion by Evaluation

Mathieu Boespflug*

École Polytechnique, INRIA
mboes@lix.polytechnique.fr

Abstract. We show how testing convertibility of two types in dependently typed systems can advantageously be implemented instead untyped normalization by evaluation, thereby reusing existing compilers and runtime environments for stock functional languages, without peeking under the hood, for a fast yet cheap system in terms of implementation effort.

Our focus is on performance of untyped normalization by evaluation. We demonstrate that with the aid of a standard optimization for higher order programs (namely uncurrying), the reuse of native datatypes and pattern matching facilities of the underlying evaluator, we may obtain a normalizer with little to no performance overhead compared to a regular evaluator.

1 Introduction

The objective here is to achieve efficient strong reduction (or full normalization) of terms in the λ -calculus. By *strong* reduction we mean the β -reduction of all redexes in a term, including inside functional values. By efficient we mean speedy execution on stock hardware.

Most implementations of the λ -calculus, such as those underpinning many functional languages, only implement *weak* reduction (also called *evaluation*). That is, reduction never occurs inside function bodies until these functions are applied to actual arguments. But for our purposes, weak reduction is not always enough.

Dependently typed theories underlie many proof assistants such as Agda, Coq, or Epigram. Such theories allow one to use a different type in lieu of another type so long as the two are convertible. Type checking a term therefore entails checking the convertibility of arbitrary terms (usually, this means deciding β -equivalence). This is typically captured by the following conversion rule:

$$\frac{\Gamma \vdash a : \tau \quad \tau \equiv \tau' : s}{\Gamma \vdash a : \tau'}$$

It is therefore the case that type checking (or equivalently proof checking) in such systems incurs the need to carry out arbitrary β -reductions. Efficient (full) normalization is particularly important when checking types entails a large amount of computation, as can often be the case, notably in proofs by reflection. Grégoire and Mahboubi [14] and Gonthier [12] provide ideal examples of such proofs.

* The research presented here was supported by a grant from Région Ile-de-France.

Other heavy users of normalization include partial evaluation, since specializing a function to statically known arguments amounts to fully normalizing this partially applied function.

Of late, functional languages have seen their influence considerably increase and their scope of application in the industry and in academia reach previously unforeseen niches. An enabling ingredient to this success has been the availability of efficient evaluation mechanisms for programs written in these languages, contending even with lower level imperative languages for the performance crown. A particularly elegant idea, normalization by evaluation (NbE), proposes to exploit off-the-shelf evaluators to implement normalization, rather than rolling out a custom built normalizer from scratch [2, 3, 4, 7, 10, 11]. All the better for speedy execution on stock hardware: some evaluators for functional languages have benefited from dozens of man years spent pouring over complex optimizations and tweaking the execution paths on a multitude of computer architectures.

Unfortunately, all flavors of NbE proposed so far have, to the best of our knowledge, achieved one or the other of the following two goals, but never both:

1. generalize to well typed terms in arbitrarily complex type systems.
2. Avoid making the cost of each reduction significantly higher than that of the underlying evaluator.

Starting from a normalizing interpreter for the λ -calculus with constants, we iteratively improve the performance of the evaluator through equational reasoning and the introduction of higher order abstract syntax (HOAS), ultimately deriving a form of normalization by evaluation. In contrast to usual approaches to NbE, where the normalization is type driven, and along the same lines as Aehlig et al. [1] and Filinski and Rohde [11], we shunt the first problem by deriving an *untyped* variant of NbE that finds the normal form of all λ -terms if there is one (Section 2). We then show how to improve on this naive implementation to the point where the time cost of β -reduction is typically within a few percentage points of that of the underlying evaluator. We demonstrate this using a few benchmarks whose results we discuss in Section 4.

Our main contribution is to show how to derive an efficient yet lightweight method for normalizing arbitrary λ -terms by enlisting the help of a few standard optimizations, further reaffirming that beyond the theoretical interest in NbE, it is also a realistic execution technique whose performance is on par with the best (albeit weak) reduction devices available.

2 Untyped NbE

2.1 The framework

Consider normalization of the pure λ -calculus with constants. By iteratively and exhaustively applying the β -rule one can of course find the normal form of some arbitrary term. This is a directed notion of normalization. But an alternative view of normalization is to consider normalization as a term equivalence relation.

Then, the normal form of a term is just a representative of the equational theory formed by the reflexive, transitive and symmetric closure of the β -reduction relation. A normalization function finds the normal form t' of a term t with t and t' equivalent. This is a reduction-free view of normalization [11].

The normalization function does not have to be β -reduction based. Suppose we can construct a denotational model of the λ -calculus with the following two properties:

1. if $t_1 \leftrightarrow_{\beta\eta} t_2$ then $\llbracket t_1 \rrbracket = \llbracket t_2 \rrbracket$ (soundness);
2. if t_1 is in normal form then a term t_2 can be extracted from a denotation $\llbracket t_1 \rrbracket$, such that $t_1 \leftrightarrow_{\alpha} t_2$ (reproduction).

Then a normalization function taking as input a closed term t can be given as

$$\Downarrow t = \Downarrow (\llbracket t \rrbracket \emptyset),$$

where \Downarrow is the extraction function, which we will call *reification*, and \emptyset is the empty set. For any t_1 in normal form, by soundness of the model $\Downarrow (\llbracket t_1 \rrbracket \emptyset) = \Downarrow (\llbracket t_2 \rrbracket \emptyset)$ for all t_2 such that $t_1 \leftrightarrow_{\beta\eta} t_2$. Since by reproduction $\Downarrow (\llbracket t_1 \rrbracket \emptyset) \leftrightarrow_{\alpha} t_1$, we have $\Downarrow t_1 \leftrightarrow_{\beta\eta} t_2$ as expected.

2.2 Towards reduction-free normalization

Consider the following representation of the syntax¹ using de Bruijn levels. The grammar for the syntax is given by the `Term` production in Figure 1.

```
data Term = Var Int | App Term Term | Abs Term
```

A normal order normalization is usually implemented along the lines of²

```
norm1 :: Term -> Term
norm1 (App t1 t2) =
  case norm1 t1 of
    Abs x t'1 -> norm1 (subst x t2 t1)
    t'1 -> App t'1 (norm1 t2)
norm1 (Abs x t) = Abs x (norm1 t)
norm1 t = t
```

We can aim for a much simpler implementation by using higher order abstract syntax (HOAS), whereby binders of the term language are represented as functions in the metalanguage. This allows us to dispense with managing scopes, variables and capture avoiding substitutions ourselves. That work is offloaded to a contraption capable of doing it far more efficiently and correctly than we are:

¹ For notational clarity, we will underline in what follows the syntax of terms, writing applications explicitly as `@`, and denote the implementation language (or *metalanguage*) using the more convenient Haskell syntax.

² The definition of *subst* is elided for conciseness.

Var	\ni	x, y, z	
Term	\ni	t	$::= x \mid \lambda. t \mid t t$
Term \supset Term_N	\ni	t_e	$::= x \mid t_e t$
Term \supset Term_{NF}	\ni	t_n	$::= t_a \mid \lambda. t_n$
Term \supset Term_A	\ni	t_a	$::= x \mid t_a t_n$

Fig. 1. Grammar and subgrammars of terms. Variables are encoded using de Bruijn levels.

the metalanguage runtime. Moving to HOAS requires a few tweaks on the *Term* datatype:

```
data Term = Const String | Abs (Term → Term)
          | App Term Term
```

Syntax variables are represented by metalanguage variables. We can therefore dispense with the *Var* constructor and introduce in its place the *Const* constructor, which stands in lieu of uninterpreted constants — or equivalently, free variables. For example, the term using named variables $(\lambda x. (\lambda y. y x)) z$ parses to the expression

```
App (Abs (λx → Abs (λy → App y x))) (Const "0")
```

The datatype *Term* represents the universe of all λ -terms, normalization of which is achieved by the following code, taking meta-level terms to object-level terms:

```
norm2 n (App t1 t2) =
  case norm2 n t1 of
    Abs t'1 → norm2 n (t'1 t2)
    t'1 → t'1 @ (norm2 n t2)
norm2 n (Abs t) =
  λ. (norm2 (n + 1) (t (Const (show n))))
norm2 n (Const c) = c
```

One can see here how the problem with shifting bindings to the metalanguage is that we can no longer descend under abstractions; they have become black boxes. But descending under abstractions is needed to normalize, so let us deconstruct these abstractions, thus turning the variable bound by some abstraction free. Remember that we already have a way to represent free variables, using *Const*. So normalizing an abstraction simply requires applying the abstraction to a fresh³ (unbound) variable and normalizing the result.

³ In practice one can opt for one of a variety of strategies for freshness. For simplicity, in this paper we get away with a simple integer counter by using de Bruijn levels in the term syntax.

After deconstructing and normalizing under the abstraction comes the time to reconstruct this abstraction. Rather than reconstructing an opaque metalanguage term, we can simply reify the abstraction into a term of the syntax. Our normalization function is no longer an endomorphism on *Term*: its result is a syntactic term in normal form.

The next step is to split out of *norm₂* the code dealing with applications into an *app* function. By appeal to the semantics of the metalanguage, we can offload yet more work to the metalanguage runtime. Insofar as evaluation order of the normalizer and metalanguage correspond, all *App* nodes can be removed from terms and replaced with calls to the *app* function. The *App* constructor is still needed, but only to represent neutral terms⁴ (i.e. *Term_N* of Figure 1). The previous example then becomes

$$app (Abs (\lambda x \rightarrow Abs (\lambda y \rightarrow app y x))) (Const "0")$$

This leads to the final definition of our normalizer:

$$\begin{aligned} app (Abs t_1) t_2 &= t_1 t_2 \\ app t_1 t_2 &= App t_1 t_2 \\ norm n (App t_1 t_2) &= (norm n t_1) \underline{\@} (norm n t_2) \\ norm n (Abs t) &= \underline{\lambda}. (norm (n + 1) (t (Const (show n)))) \\ norm n (Const c) &= \underline{c} \end{aligned}$$

After this final step, notice that all forms in the syntax are now interpreted directly with their corresponding (tagged) forms in the metalanguage, as shown in Figure 2. *norm* matches the specification of a reification function. Indeed, parsing a term to the metalanguage, then unparsing the resulting construct with *norm*, is an untyped, reduction-free, normalization by evaluation function, in the sense of Section 2.

$$\begin{aligned} \llbracket x \rrbracket n &= \hat{x} && \text{if } x < n \\ \llbracket x \rrbracket n &= Const \underline{x} && \text{otherwise} \\ \llbracket \lambda. t \rrbracket n &= Abs (\lambda \hat{n} \rightarrow \llbracket t \rrbracket (n + 1)) \\ \llbracket t_1 t_2 \rrbracket n &= app (\llbracket t_1 \rrbracket n) (\llbracket t_2 \rrbracket n) \end{aligned}$$

Fig. 2. Translation of the syntax into the metalanguage. $\hat{\cdot}$ maps naturals to variable names.

3 Optimizations

In this section we will focus on offloading yet more work to the metalanguage runtime by exploiting intrinsic features of most higher order programming languages that go beyond the pure λ -calculus. One such feature is the uncurrying of function applications, the other is pattern matching on algebraic datatypes.

⁴ Neutral terms are variables or applications of a neutral term to a term. Substituting a neutral term anywhere in another term will not create additional redexes.

3.1 Minimizing closures

Functional values in functional programming languages are typically represented as *closures*, a pairing of code and an environment assigning values to all free variables appearing in the code. Consider a church encoding of lists and a right fold in a syntax where functions can be applied to multiple arguments in one go.

$$\begin{aligned} nil &\equiv \underline{\lambda f g. f} \\ cons &\equiv \underline{\lambda h t f g. g h (t f g)} \\ map &\equiv \underline{\lambda f l. l nil (\lambda h t. cons (f h) t)} \end{aligned}$$

Y is the usual call-by-name fixed-point combinator. The notation $\lambda x_1 \dots x_n. []$ is syntactic sugar for $(\lambda x_1. \dots (\lambda x_n. [])) \dots$. That is, the higher-order functions above take multiple arguments, but are encoded in terms of unary functions that return functions. This encoding is called *currying*.

Note however that currying has a cost. Applying a function to multiple arguments entails the creation of many short-lived intermediate closures, one for each function returned as a result of the application to one argument. In general, one will need to allocate (and then deallocate soon thereafter) $n - 1$ closures during the consecutive application of a function to n arguments. For instance,

$$\begin{aligned} \llbracket map\ id\ nil \rrbracket & \\ &= app\ (app\ map\ id)\ nil \\ &= app\ (app\ (Abs\ (\lambda f \rightarrow Abs\ (\lambda l \rightarrow \dots)))\ id)\ nil \\ &\rightarrow_{\beta} app\ (Abs\ (\lambda l \rightarrow \dots))\ nil \\ &\rightarrow_{\beta} nil \end{aligned}$$

Here, *map* is applied to two arguments, therefore one intermediate *Abs* structure is constructed. But an alternative encoding of n -ary functions could avoid this.

The literature abounds with various encodings of n -ary functions (i.e. calling conventions) targeted by compilers to avoid costly closure allocation. Marlow and Peyton-Jones [19] proposes the Push/Enter and Eval/Apply dichotomy to describe them. We pick the Eval/Apply model here for its very cheap implementation cost and good performance in the common case [19]. That is, assuming a syntax where consecutive λ 's have been folded into multiple argument abstractions, we can forgo many *Abs* constructions by means of a family ap_n of application operators and the addition of a number of Abs_n constructors, as shown in Figure 3. Note that most functions appearing in terms of the syntax will typically have low arity, so that one could reap most of the benefit of this approach even if bounding the number of ap_n operators and Abs_n constructors to a small number such as 4 or 5. Though uncommon, applications of functions with higher arity is still possible, but at a slight performance cost due to extra closure construction.

Parsing the above terms to the metalanguage now gives:

$$\begin{aligned} nil &= Abs_2\ (\lambda f\ g \rightarrow f) \\ cons &= Abs_4\ (\lambda h\ t\ f\ g \rightarrow ap_2\ g\ h\ (ap_2\ t\ f\ g)) \end{aligned}$$

1. $ap_n (Abs_m f) t_1 \dots t_n = Abs_{m-n} (f t_1 \dots t_n)$
2. $ap_n (Abs_m f) t_1 \dots t_n = f t_1 \dots t_n$
3. $ap_n (Abs_m f) t_1 \dots t_n = ap_{n-m} (f t_1 \dots t_m) t_{m+1} \dots t_n$

where conditions on (1) are if $n < m$, on (2) if $n = m$, on (3) if $n > m$.

Fig. 3. A family of *ap* operators

$$map = Abs_2 (\lambda f l \rightarrow \\ ap_2 l nil (Abs_2 (\lambda h t \rightarrow ap_2 cons (ap_1 f h) t)))$$

For small n , n -ary functions in the syntax are encoded using n -ary functions in the metalanguage. Beyond economizing data structure allocations, this optimization permits us to reap the benefits of closure allocation strategies typically found in compilers to reduce the cost and frequency of extending closure environments. For example, many execution environments such as the OCaml interpreter can avoid any allocation of environments on the heap in the common case of n -ary functions applied to n arguments, instead pushing all arguments on the stack [16].

3.2 Specialized constructors

Representing all datatypes as functions via Church encodings induces needlessly many β -reductions and wastes opportunities for optimization. Haskell and many other statically typed functional programming languages feature algebraic datatypes and pattern matching facilities on these datatypes, enabling more natural and more efficient data manipulation. Compiling complex pattern matches to decision trees or to backtracking automata [15] can drastically reduce the amount of computation needed to access and manipulate algebraic structures.

With the current definition of *Term*, it is already possible to parse patterns in the syntax to case analysis constructs in the metalanguage, but currently a metalanguage representation of a pattern p_1 can become quite a bit larger than p_1 . Assume for instance constants *nil* and *cons*, constructors of the list type, and take the definition of *append* in the metalanguage:

$$append = Abs_2 (\lambda xs ys \rightarrow \mathbf{case} \ xs \ \mathbf{of} \\ \quad Const \ \mathbf{"nil"} \rightarrow ys \\ \quad App \ (App \ (Const \ \mathbf{"cons"}) \ x) \ xs' \rightarrow \\ \quad \quad ap_2 \ (Const \ \mathbf{"cons"}) \ x \ (ap_2 \ append \ xs' \ ys)$$

Replacing the constructor names with integers rather than strings to avoid string comparison cost does spare some computation, but it is better to avoid the *Const* constructor altogether. Rather than representing a datatype as an in-memory tree, with *App* constructors at branch nodes and *Const* constructors at the leaves, each in its own memory cell, it is much more memory efficient to add all data constructors found in the syntax as additional constructors to the

metalanguage interpretation, effectively flattening the representation in memory. That is, for constructors *nil* and *cons*, add

```
data Term = ... | Nil | Cons Term Term
```

As shall be detailed in Section 4, a flatter structure means less indirection when performing pattern matches, hence better performance.

The downside of mirroring syntax level constructors as constructors in *Term* is that doing so breaks modularity. Since the *Term* datatype is the universe of all syntax terms, breaking up definitions in the syntax into modules requires that all constructors in all modules need to be coalesced into the term *Term* datatype. Encoding modules in the syntax with modules in the metalanguage is useless, because introducing a new constructor means modifying *Term*, which in turn means recompiling all modules because they all depend on *Term*.

A solution to recover modularity is to hardcode a set of constructors in the *Term* datatype, much as we hardcoded the set *Abs_n* of *n*-ary functions. This means that constructors with small arity in the source language can be represented using a single constructor in the metalanguage. Larger (less common) constructors in the source language can of course be represented as the composition of smaller constructors.

```
data Term = ... | Const0 Int | Const1 Int Term
           | ... | Constn Int Term ... Term
```

In languages that feature first class arrays, in particular allowing pattern matching on arrays (such as OCaml), one could also replace the definition of *Const* with

```
type term = ... | Const of name * term array
```

The effect of removing *Const* is to build in a closed world assumption on constructors of the syntax. Some languages allow the definition of extensible datatypes, which we can use to break the closed world assumption. Recent versions of OCaml feature polymorphic variants and Standard ML's *exn* exception datatype is extensible. Terms applied to a constant would simply be accumulated in the array. The array size is known in advance because all constructors have a fixed number of fields.

In summary, the appropriate option will be contingent on the runtime environment chosen to execute the normalizer. As always, the objective here is to make do with existing runtime environments without modification, whilst observing that the penalty of this constraint can be made close to negligible — an observation substantiated in the following section.

4 Benchmarks

Our use of untyped NbE is as a cheap contraption to efficiently perform the conversion test in dependent type theories. In this section we examine the effect

of various optimizations presented previously on a small set of benchmarks and compare them to earlier work on untyped NbE by Aehlig et al. [1]. In these benchmarks, the object language is Haskell. The interpretation stage of NbE then becomes a source-to-source transformation on programs, which we implement using Template Haskell. The transformed source is then compiled to native code by the GHC compiler.

We compare 6 flavors of NbE:

ahn This is untyped NbE as described in [1]. All functions are interpreted as unary functions. All function arguments are packed into lists that the function pattern matches over to extract individual arguments.

singularity This interpretation takes every function to a unary closure. Functions taking multiple arguments are curried and are represented using multiple embedded closures.

evalapply The optimization described in 3.1.

constructors Every constructor appearing in terms of the object language become additional constructors *Term*, as in 3.2.

ucea Combination of “evalapply” and “constructors”.

whnf The identify interpretation, where terms of the object language are interpreted as themselves.

We run the following benchmarks for each of the flavors:

append Concatenation of two large lists of integers of size 50,000.

even Test whether an input list is even or odd. Lists are represented using a Church encoding, so that no pattern matching occurs in this benchmark. It is meant to test performance of applications.

sort Sorting of large lists of integers encoded using constructors. This benchmark is meant to be rather more sensitive to pattern matching performance. The implementation is mergesort found in the base package of the Haskell libraries.

exp3-8 A tiny benchmark appearing in the nofib suite: taking 3 to the power of 8, in Peano arithmetic.

queens Enumerate the solutions to this classic constraint satisfaction problem: find a way to place 10 queens on a 10x10 chess board such that no two queens are on the same column or row.

The results are shown in Figure 4 and Table 1. Note immediately how the vast majority of the performance benefits comes from interpreting constructors as constructors; this greatly reduces the size of the patterns to match and help allocate fewer objects on the heap. An overview of the heap usage and garbage collection on each of the above benchmarks shows that using constructors typically halves total heap allocation during the lifetime of the program.

Currying functions, rather than grouping the arguments into lists that are frequently deconstructed and reconstructed, affords a gain in most benchmarks. The eval/apply optimization allows a further halving of execution time on benchmarks with functions with high arity, such as queens and its heavy use of *foldr*.

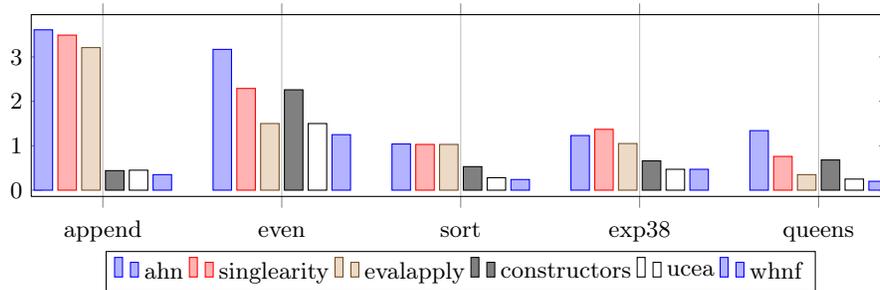


Fig. 4. Visual representation of the data in Table 1.

flavor	append	%	even	%	sort	%	exp3-8	%	queens	%
ahn	3.61	1031	3.17	253	1.04	433	1.23	261	1.34	670
evalapply	3.21	917	1.50	120	1.03	429	1.05	223	0.35	175
singularity	3.49	997	2.29	183	1.03	429	1.37	191	0.76	380
constructors	0.44	125	2.26	180	0.53	220	0.66	140	0.68	340
ucea	0.45	128	1.50	120	0.28	116	0.47	100	0.25	120
whnf	0.35	100	1.25	100	0.24	100	0.47	100	0.20	100

Table 1. Absolute execution times (seconds) and relative to execution time of whnf.

The main observation, however, is that untyped normalization by evaluation with the addition of the eval/apply optimization and the use of metalanguage constructors is hardly any slower on these benchmarks than the execution of these benchmarks by evaluation alone. In pathological cases where none of the execution time is spent in pattern matching, such as the “even” benchmark, we observe a penalty of about 20%. However, pattern matching or garbage collection and heap allocation dominates the runtime of many functional programs. In such cases the extra cost of tagging closures is often negligible.

4.1 Proofs by reflection

A popular style of proof consists in reusing the proof language provided by the theorem prover as a programming language. For some predicate P ranging over terms of type T , rather than proving directly the property

```
forall t : T, P t
```

one instead introduces a decision procedure f , along with a proof $f_{correct}$ that f is correct. In COQ, this would go something along the lines of

```
Variable P : T -> Prop.
Variable f : T -> bool.
Variable x : T.
Variable f_correct : forall x:T, f x = true -> P x.
```

Now we have that the term

variables	1	%	2	%	3	%	4	%	5	%
no conv	0.63	94	0.68	94	1.40	93	2.25	77	3.92	3.11
nbe	0.64	95	0.70	97	1.42	94	2.30	79	27.27	20.02
Coq VM	0.67	100	0.72	100	1.50	100	2.92	100	136.2	100

Table 2. Solving formulae of n variables with Cooper’s quantifier elimination.

```
fun t:T => f_correct t (refl_equal true):P t
```

is a proof of `forall t : T, P t`. For some $x : T$, the conversion test here consists in verifying that the function f applied to $x : T$ reduces to `true`. The SSREFLECT proof language encourages this style of proof in the small as well as in the large, so that typical properties such as the symmetry of the proposition disjunction operator might be proved more efficiently and concisely using reflection. In effect, reflection rephrases the problem so as to shift much of the burden of proof to mere calculation, avoiding tedious deductive reasoning.

As more proofs adopt this style of reasoning, computation starts dominating the time needed for proof checking. Using a prototype implementation inside the kernel of COQ of the normalization scheme of Section 2, we briefly report on the impact of using normalization by evaluation for the conversion test on a tactic that generates proofs in the reflexive style: Cooper’s quantifier elimination for Presburger arithmetic (unpublished work by Salil Joshi and Assia Mahboubi).

Figure 2 shows the computational blowup as the formulae to solve increase in the number of variables. Starting from 6 variables, the problem is so large that the runtime exhausts all available memory after over 30 minutes. For each formula, we record two markers for our performance measurements. The reference time is the time needed by the kernel to verify the proof generated by the tactic when compiling the proof to COQ’s existing virtual machine. The best we can hope to do is the time required to check the proof when the conversion test is unplugged, i.e. the time spent in other proof checking tasks save conversion. As evidenced by the last column, speedup compared to the already existing virtual machine based reduction scheme is a fivefold increase in the purely computational part of the proof (which dominates the entire proof checking time on even short formulae), as can be expected from moving from a bytecode based environment to execution of native code. We expect similar gains for other (large and small) proofs by reflection, such as [12]. However, for very small proofs the overhead associated with compiling everything in the environment might not pay its worth. The default conversion routine of COQ should fare better in these cases.

5 A note on correctness

A detailed treatment of the correctness of the normalization algorithm presented here is beyond the scope of this paper. We note, however, that the conversion tested is implemented in the trusted base of most any theorem prover. High

assurance of correctness is hence a very desirable property. Previous work on other variants of untyped normalization by evaluation has already established partial correctness properties [1] and soundness (the output term, if any, is β -equivalent to the input term), standardization (β -equivalent terms are mapped to the same result) and completeness (normal forms are found for all terms that have a normal form) [11]. For instance, a meaning preserving embedding of terms as represented in Section 2 into terms of the form found in [1] is straightforward (arguments to functions are boxed into lists), by which means we may port the results found therein.

Correctness may alternatively be derived via meaning preserving transformations from preexisting normalizers, in the style of [6].

6 Related Work

Our work is a continuation of many other contributions regarding normalization by evaluation and its applications. Whilst many treatments of NbE do discuss computational efficiency, few quantify empirically performance on select benchmarks. [1] is one work on which we build upon, being closely related both in its attention to the performance side of the coin and in the essence of their scheme. They too map terms of the object language to tagged equivalents in the metalanguage by embedding functions, free variables and constants into a datatype. Our approach differs from theirs in that we treat functions of arbitrary arity uniformly by currying. In their approach functions of the object language are mapped to single arity functions within the metalanguage, encapsulating all arguments of the functions inside lists. The body of the functions then pattern match on the input list to extract arguments. Whilst appealing in its simplicity, their approach suffers performance-wise from allocating many lists during function application time that are then immediately deconstructed. In addition, encapsulating arguments inside lists breaks the optimization described in Section 3.1. For simplicity, constructors in the object language are not translated to constructors in the metalanguage but rather represented with a special constructor for constants. Lindley [18] also considers untyped normalization by evaluation in a performance sensitive context, giving a quantitative analysis of the performance of a number of algorithms and variants compared to reduction based approaches. Optimizations for higher order programs and data constructors are not considered, however.

Filinski and Rohde [11] propose a similar algorithm for untyped normalization by evaluation. Whilst Aehlig et al. prove only partial correctness, namely that if their algorithm returns a term then that term is in normal form and convertible to the input (*soundness* and *standardization* properties), Filinski and Rohde further prove *completeness*. However, the focus there is on a precise semantic study, rather than an evaluation of performance.

Of particular note in the work of Aehlig et al. [1] is their generalization of NbE to the symbolic normalization of terms with regards to arbitrary user-provided rewrite rules. For conciseness, we do not discuss this matter further in

this paper, but their translation of rewrite rules as pattern matching functions in the metalanguage can readily be adapted to the normalization scheme presented here. This generalization is not required for the conversion test in the Calculus of Inductive Constructions used by COQ for instance, but it is useful for reduction in Isabelle/HOL and for the conversion test in formalisms such as λII -modulo [8]. Blanqui et al. [5] independently propose a similar translation of rewrite rules into OCaml though in the context of finding canonical forms for non-free algebraic datatypes rather than applied to normalization.

A variety of virtual machines have been proposed for normalization. Notably, Crégut [9] proves correct a normalizer for the λ -calculus. The code can be executed by expansion to Motorola 68000 assembly code, resulting in an efficient but more heavyweight (in the sense of implementation effort) and less portable execution model compared to NbE based approaches. The machine of Grégoire and Leroy [13] that COQ sometimes uses for the conversion test should also be mentioned here. Theirs is a modified and formalized version of a bytecode interpreter for OCaml (the ZAM), to do normalization via reduction to weak head normal form along with a *readback* phase to restart weak reduction under binders. Whilst offering striking similarities to NbE, including in its reuse of existing evaluators, one important difference lies in the fact that the implementation of the underlying evaluator needs to be modified, whereas the objective of NbE, here and elsewhere, is to get away without looking under the hood. As a side effect, NbE affords more freedom of choice regarding which evaluator to choose, allowing for instance to trade off minimizing the trusted base for better performance.

The principal extension made to the ZAM to normalize COQ terms is the introduction of *accumulators*, which represent applications of free variables to a number of terms. Embedding this construct within the virtual machine avoids having to do case analysis at every application to discriminate between function applications and applications of neutral terms. We show that with the simple optimization of Section 3.1, the overhead of this case analysis is very small in practise.

These approaches can be seen as complementary to the one exposed here in that these normalizers are abstract machines whose correctness is more readily established, hence avoiding extending the trusted base of a theorem prover with code as large as that of a full scale compiler and the associated runtime environment for the chosen metalanguage. They may also reduce the cost of compilation, which for small terms can far exceed the time needed to normalize them.

7 Conclusion

Just as moving from a naive interpreter to an optimizing compiler can mean moving from the intractable to the feasible for the evaluation of programs, so too does compiling the costly components of the type checking problem in dependent type theories may reap enormous benefits. Others have shown how it is possible

to bring to bear the power of existing compiler technology in proof assistants with little implementation effort. We have shown that to get excellent performance rivaling that of stock runtime systems for popular programming languages, the implementation effort is nearly trivial: parse the object language and pretty print it to tagged terms in the form of a functional program. We can have our cake and eat it too.

A limitation of normalization by evaluation is that terms are always evaluated to weak head normal forms before normalizing under binders. When strongly normalizing a term, this may not be the best strategy: in fact [17] has shown that this could lead to redundant copying of exponentially many λ -terms, which an optimal strategy might avoid. But seeking the optimal strategy may introduce far too much overhead to be viable in practice. As in [13], the approach presented here seeks to minimize the cost of each reduction, at some expense on the total number of reductions performed. It would be interesting however, to allow for short-circuiting of normalization when reduction so far has yielded enough information to decide the convertibility of two terms, whilst retaining the conceptual and implementation simplicity of normalization by evaluation.

The normalization algorithm presented here is at the heart of a new proof checker for the λII -calculus modulo called DEDUKI⁵, but transferring this technology to full-fledged proof assistants would be of benefit. We have also implemented this scheme inside the kernel of COQ that works in the common case of comparing non-functional closed values, but a full treatment of terms of the Calculus of Inductive Constructions requires careful attention to the reductions rules of that calculus when in the presence of free variables.

8 Acknowledgements

Many thanks to Klaus Aehlig, Olivier Danvy and Benjamin Grégoire for fruitful discussions on normalization by evaluation, to Bruno Barras for enlightening discussions of the implementation of COQ and to Assia Mahboubi for her kind encouragements. Chantal Keller and Assia Mahboubi were of great help in elaborating benchmarks. The author is greatly indebted to Arnaud Spiwack for making his time and expertise available to understand the COQ system's kernel.

References

- [1] Aehlig, K., Haftmann, F., Nipkow, T.: A Compiled Implementation of Normalization by Evaluation. *Theorem Proving in Higher Order Logics (TPHOLs 2008)*, Lecture Notes in Computer Science. Springer-Verlag (2008)
- [2] Altenkirch, T., Dybjer, P., Hofmann, M., Scott, P.: Normalization by evaluation for typed lambda calculus with coproducts. In: *Proceedings of the Sixteenth Annual IEEE Symposium on Logic in Computer Science*. pp. 203–210 (2001)
- [3] Berger, U., Eberl, M., Schwichtenberg, H.: Normalization by evaluation. *Prospects for Hardware Foundations* pp. 117–137 (1998)

⁵ <http://www.lix.polytechnique.fr/dedukti>

- [4] Berger, U., Eberl, M., Schwichtenberg, H.: Term rewriting for normalization by evaluation. *Information and Computation* 183(1), 19–42 (2003)
- [5] Blanqui, F., Hardin, T., Weis, P.: On the Implementation of Construction Functions for Non-free Concrete Data Types. *Lecture Notes In Computer Science* 4421, 95 (2007)
- [6] Boespflug, M.: From self-interpreters to normalization by evaluation. In: *Informal proceedings of the 2009 Workshop on Normalization by Evaluation*. pp. 35–38 (August 2009)
- [7] Coquand, T., Dybjer, P.: Intuitionistic model constructions and normalization proofs. *Mathematical Structures in Computer Science* 7(01), 75–94 (1997)
- [8] Cousineau, D., Dowek, G.: Embedding pure type systems in the lambda-pi-calculus modulo. *Lecture Notes in Computer Science* 4583, 102–117 (2007)
- [9] Crégut, P.: Strongly reducing variants of the Krivine abstract machine. *Higher-Order and Symbolic Computation* 20(3), 209–230 (2007)
- [10] Danvy, O.: Type-directed partial evaluation. *POPL’96* pp. 242–257 (1996)
- [11] Filinski, A., Rohde, H.: A denotational account of untyped normalization by evaluation (2004)
- [12] Gonthier, G.: The four colour theorem: Engineering of a formal proof. In: Kapur, D. (ed.) *ASCM. Lecture Notes in Computer Science*, vol. 5081, p. 333. Springer (2007)
- [13] Grégoire, B., Leroy, X.: A compiled implementation of strong reduction. *Proceedings of the seventh ACM SIGPLAN international conference on Functional programming* pp. 235–246 (2002)
- [14] Grégoire, B., Mahboubi, A.: Proving equalities in a commutative ring done right in coq. *Lecture notes in computer science* 3603, 98 (2005)
- [15] Le Fessant, F., Maranget, L.: Optimizing pattern matching. *Proceedings of the sixth ACM SIGPLAN international conference on Functional programming* pp. 26–37 (2001)
- [16] Leroy, X.: The ZINC experiment: an economical implementation of the ML language. *Tech. rep., INRIA* (1990)
- [17] Lévy, J.: Réductions correctes et optimales dans le Lambda-Calcul. *Université Paris 7* (1978)
- [18] Lindley, S.: Normalisation by evaluation in the compilation of typed functional programming languages (2005)
- [19] Marlow, S., Peyton-Jones, S.: Making a fast curry: push/enter vs. eval/apply for higher-order languages. *Journal of Functional Programming* 16(4-5), 415–449 (2006)