

Evidence Normalization in System FC (Invited Talk)

Dimitrios Vytiniotis and Simon Peyton Jones

Microsoft Research, Cambridge
dimitris@microsoft.com, simonpj@microsoft.com

Abstract

System FC is an explicitly typed language that serves as the target language for Haskell source programs. System FC is based on System F with the addition of erasable but explicit type equality proof witnesses. Equality proof witnesses are generated from type inference performed on source Haskell programs. Such witnesses may be very large objects, which causes performance degradation in later stages of compilation, and makes it hard to debug the results of type inference and subsequent program transformations. In this paper we present an equality proof simplification algorithm, implemented in GHC, which greatly reduces the size of the target System FC programs.

1998 ACM Subject Classification D.3.3 Language Constructs and Features, F.3.3 Studies of Program Constructs

Keywords and phrases Haskell, type functions, system FC

Digital Object Identifier 10.4230/LIPIcs.RTA.2013.20

1 Introduction

A statically-typed intermediate language brings a lot of benefits to a compiler: it is free from the design trade-offs that come with source language features; types can inform optimisations; and type checking programs in the intermediate language provides a powerful consistency check on each stage of the compiler.

The Glasgow Haskell Compiler (GHC) has just such an intermediate language, which has evolved from System F to System FC [16, 20] to accommodate the source-language features of *GADTs* [6, 15, 13] and *type families* [9, 3]. The key feature that allows System FC to accommodate GADTs and type families is its use of explicit *coercions* that witness the equality of two syntactically-different types. Coercions are erased before runtime but, like types, serve as a static consistency proof that the program will not “go wrong”.

In GHC, coercions are produced by a fairly complex type inference (and proof inference) algorithm that elaborates source Haskell programs into FC programs [19]. Furthermore, coercions undergo major transformations during subsequent program optimization passes. As a consequence, they can become very large, making the compiler bog down. This paper describes how we fixed the problem:

- Our main contribution is a novel coercion simplification algorithm, expressed as a rewrite system, that allows the compiler to replace a coercion with an equivalent but much smaller one (Section 4).
- Coercion simplification is important in practice. We encountered programs whose unsimplified coercion terms grow to many times the size of the actual executable terms, to the point where GHC choked and ran out of heap. When the simplifier is enabled, coercions simplify to a small fraction of their size (Section 5).



© Dimitrios Vytiniotis and Simon Peyton Jones;
licensed under Creative Commons License CC-BY
24th International Conference on Rewriting Techniques and Applications (RTA'13).
Editor: Femke van Raamsdonk; pp. 20–38



Leibniz International Proceedings in Informatics
LIPIcs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



c	\in	Coercion variables	
x	\in	Term variables	
e, u	$::=$	$x \mid l \mid \lambda x:\sigma. e \mid e u$	
		$\mid \Lambda a:\eta. e \mid e \phi$	Type polymorphism
		$\mid \lambda c:\tau. e \mid e \gamma$	Coercion abstraction/application
		$\mid K \mid \mathbf{case} e \mathbf{of} \overline{p \rightarrow u}$	Constructors and case expressions
		$\mid \mathbf{let} x:\tau = e \mathbf{in} u$	Let binding
		$\mid e \triangleright \gamma$	Cast
p	$::=$	$K \overline{c:\tau} \overline{x:\tau}$	Patterns

■ **Figure 1** Syntax of System FC (Terms).

- To get these benefits, coercion simplification must take user-declared equality axioms into account, but the simplifier *must never loop* while optimizing a coercion – no matter which axioms are declared by users. Proof normalization theorems are notoriously hard, but we present such a theorem for our coercion simplification. (Section 6)

Equality proof normalization was first studied in the context of monoidal categories and we give pointers to early work in Section 7 – this work in addition addresses the simplification of open coercions containing variables and arbitrary user-declared axioms.

2 An overview of System FC

We begin by reviewing the role of an intermediate language. GHC desugars a rich, complex source language (Haskell) into a small, simple intermediate language. The source language, Haskell, is *implicitly typed*, and a type inference engine figures out the type of every binder and sub-expression. To make type inference feasible, Haskell embodies many somewhat ad-hoc design compromises; for example, λ -bound variables are assigned monomorphic types. By contrast, the intermediate language is simple, uniform, and *explicitly typed*. It can be typechecked by a simple, linear time algorithm. The type inference engine *elaborates* the implicitly-typed Haskell program into an explicitly-typed FC program.

To make this concrete, Figure 1 gives the syntax of System FC, the calculus implemented by GHC’s intermediate language. The term language is mostly conventional, consisting of System F, together with let bindings, data constructors and case expressions. The syntax of a term encodes its typing derivation: every binder carries its type, and type abstractions $\Lambda a:\eta. e$ and type applications $e \phi$ are explicit.

The types and kinds of the language are given in Figure 2. Types include variables (a) and constants H (such as `Int` and `Maybe`), type applications (such as `Maybe Int`), and polymorphic types $(\forall a:\eta. \phi)$. The syntax of types also includes *type functions* (or *type families* in the Haskell jargon), which are used to express type level computation. For instance the following declaration in source Haskell:

```
type family F (a :: *) :: a
type instance F [a] = a
```

introduces a type function F at the level of System FC. The accompanying `instance` line asserts that any expression of type `F [a]` can be viewed as having type `a`. We shall see in Section 2.2 how this fact is expressed in FC. Finally type constants include datatype constructors (T) but also arrow (\rightarrow) as well as a special type constructor $\sim\#$ whose role we

$$\boxed{\Gamma \Vdash^{\text{tm}} e : \tau}$$

$$\frac{(x:\tau) \in \Gamma}{\Gamma \Vdash^{\text{tm}} x : \tau} \text{EVAR} \qquad \frac{(K:\sigma) \in \Gamma}{\Gamma \Vdash^{\text{tm}} K : \sigma} \text{ECON}$$

$$\frac{\Gamma, (x:\sigma) \Vdash^{\text{tm}} e : \tau \quad \Gamma \Vdash^{\text{by}} \sigma : \star}{\Gamma \Vdash^{\text{tm}} \lambda x:\sigma. e : \sigma \rightarrow \tau} \text{EABS} \qquad \frac{\Gamma \Vdash^{\text{tm}} e : \sigma \rightarrow \tau \quad \Gamma \Vdash^{\text{tm}} u : \sigma}{\Gamma \Vdash^{\text{tm}} e u : \tau} \text{EAPP}$$

$$\frac{\Gamma, (c:\sigma) \Vdash^{\text{tm}} e : \tau \quad \Gamma \Vdash^{\text{by}} \sigma : \text{Constraint}_{\#}}{\Gamma \Vdash^{\text{tm}} \lambda c:\sigma. e : \sigma \rightarrow \tau} \text{ECABS} \qquad \frac{\Gamma \Vdash^{\text{tm}} e : (\sigma_1 \sim_{\#} \sigma_2) \rightarrow \tau \quad \Gamma \Vdash^{\text{co}} \gamma : \sigma_1 \sim_{\#} \sigma_2}{\Gamma \Vdash^{\text{tm}} e \gamma : \tau} \text{ECAPP}$$

$$\frac{\Gamma, (a:\eta) \Vdash^{\text{tm}} e : \tau}{\Gamma \Vdash^{\text{tm}} \Lambda a:\eta. e : \forall a:\eta. \tau} \text{ETABS} \qquad \frac{\Gamma \Vdash^{\text{tm}} e : \forall a:\eta. \tau \quad \Gamma \Vdash^{\text{by}} \phi : \eta}{\Gamma \Vdash^{\text{tm}} e \phi : \tau[\phi/a]} \text{ETAPP}$$

$$\frac{\Gamma, (x:\sigma) \Vdash^{\text{tm}} u : \sigma \quad \Gamma, (x:\sigma) \Vdash^{\text{tm}} e : \tau}{\Gamma \Vdash^{\text{tm}} \text{let } x:\sigma = u \text{ in } e : \tau} \text{ELET} \qquad \frac{\Gamma \Vdash^{\text{tm}} e : \tau \quad \Gamma \Vdash^{\text{co}} \gamma : \tau \sim_{\#} \phi}{\Gamma \Vdash^{\text{tm}} e \triangleright \gamma : \phi} \text{ECAST}$$

$$\begin{array}{l}
\Gamma \Vdash^{\text{tm}} e : T \bar{\kappa} \bar{\sigma} \\
\text{For each branch } K \bar{x}:\bar{\tau} \rightarrow u \\
(K:\forall(\bar{a}:\eta_{\bar{a}}). \bar{\sigma}_1 \sim_{\#} \bar{\sigma}_2 \rightarrow \bar{\tau} \rightarrow T \bar{a}) \in \Gamma \\
\phi_i = \tau_i[\bar{\sigma}/\bar{a}] \\
\phi_{1i} = \sigma_{1i}[\bar{\sigma}/\bar{a}] \\
\phi_{2i} = \sigma_{2i}[\bar{\sigma}/\bar{a}] \quad \Gamma, c:\phi_1 \sim_{\#} \phi_2 \quad x:\phi \Vdash^{\text{tm}} u : \sigma \\
\hline
\Gamma \Vdash^{\text{tm}} \text{case } e \text{ of } K (c:\sigma_1 \sim_{\#} \sigma_2) (\bar{x}:\bar{\tau}) \rightarrow u : \sigma \quad \text{ECASE}
\end{array}$$

■ **Figure 4** Well-formed terms.

$$\boxed{\Gamma \Vdash^{\text{by}} \tau : \kappa}$$

$$\frac{(a:\eta) \in \Gamma}{\Gamma \Vdash^{\text{by}} a : \eta} \text{TVAR} \qquad \frac{(T:\kappa) \in \Gamma}{\Gamma \Vdash^{\text{by}} T : \kappa} \text{TDATA} \qquad \frac{(F:\kappa) \in \Gamma}{\Gamma \Vdash^{\text{by}} F : \kappa} \text{TFUN}$$

$$\frac{\kappa_1, \kappa_2 \in \{\text{Constraint}_{\#}, \star\}}{\Gamma \Vdash^{\text{by}} (\rightarrow) : \kappa_1 \rightarrow \kappa_2 \rightarrow \star} \text{TARR} \qquad \frac{}{\Gamma \Vdash^{\text{by}} (\sim_{\#}) : \kappa \rightarrow \kappa \rightarrow \text{Constraint}_{\#}} \text{TEQPRED}$$

$$\frac{\Gamma \Vdash^{\text{by}} \phi_1 : \kappa_1 \rightarrow \kappa_2 \quad \Gamma \Vdash^{\text{by}} \phi_2 : \kappa_1}{\Gamma \Vdash^{\text{by}} \phi_1 \phi_2 : \kappa_2} \text{TAPP} \qquad \frac{\Gamma, (a:\eta) \Vdash^{\text{by}} \tau : \star}{\Gamma \Vdash^{\text{by}} \forall a:\eta. \tau : \star} \text{TALL}$$

■ **Figure 5** Well-formed types.

$\Gamma \Vdash^\infty \gamma : \sigma_1 \sim_{\#} \sigma_2$		
$\frac{(c:\sigma_1 \sim_{\#} \sigma_2) \in \Gamma}{\Gamma \Vdash^\infty c : \sigma_1 \sim_{\#} \sigma_2}$ C _{VAR}	$\frac{(C \bar{a}:\bar{\eta} : \tau_1 \sim_{\#} \tau_2) \in \Gamma \quad \Gamma \Vdash^\infty \gamma_i : \sigma_i \sim_{\#} \phi_i}{\Gamma \Vdash^\infty C \bar{\gamma} : \tau_1[\bar{\sigma}/\bar{a}] \sim_{\#} \tau_2[\bar{\phi}/\bar{a}]}$ C _{AX}	$\frac{\Gamma \Vdash^\infty \phi : \kappa}{\Gamma \Vdash^\infty \langle \phi \rangle : \sigma \sim_{\#} \sigma}$ C _{REFL}
$\frac{\Gamma \Vdash^\infty \gamma_1 : \sigma_1 \sim_{\#} \sigma_2 \quad \Gamma \Vdash^\infty \gamma_2 : \sigma_2 \sim_{\#} \sigma_3}{\Gamma \Vdash^\infty \gamma_1; \gamma_2 : \sigma_1 \sim_{\#} \sigma_3}$ C _{TRANS}	$\frac{\Gamma \Vdash^\infty \gamma : \sigma_1 \sim_{\#} \sigma_2}{\Gamma \Vdash^\infty \text{sym } \gamma : \sigma_2 \sim_{\#} \sigma_1}$ C _{SYM}	$\frac{\Gamma \Vdash^\infty \gamma : H \bar{\sigma} \sim_{\#} H \bar{\tau}}{\Gamma \Vdash^\infty \text{nth } k \gamma : \sigma_k \sim_{\#} \tau_k}$ C _{TH}
$\frac{\Gamma, (a:\eta) \Vdash^\infty \gamma : \sigma_1 \sim_{\#} \sigma_2}{\Gamma \Vdash^\infty \forall a:\eta. \gamma : (\forall a:\eta. \sigma_1) \sim_{\#} (\forall a:\eta. \sigma_2)}$ C _{CALL}		
$\frac{\Gamma \Vdash^\infty \gamma_1 : \sigma_1 \sim_{\#} \sigma_2 \quad \Gamma \Vdash^\infty \gamma_2 : \phi_1 \sim_{\#} \phi_2 \quad \Gamma \Vdash^\infty \sigma_1 \phi_1 : \kappa}{\Gamma \Vdash^\infty \gamma_1 \gamma_2 : \sigma_1 \phi_1 \sim_{\#} \sigma_2 \phi_2}$ C _{APP}	$\frac{\Gamma \Vdash^\infty \gamma : (\forall a:\eta. \sigma_1) \sim_{\#} (\forall a:\eta. \sigma_2) \quad \Gamma \Vdash^\infty \phi : \eta}{\Gamma \Vdash^\infty \gamma @ \phi : \sigma_1[\phi/a] \sim_{\#} \sigma_2[\phi/a]}$ C _{INST}	

■ **Figure 6** Well-formed coercions.

```

data T a where
  T1 :: Int -> T Int
  T2 :: a -> T a

  f :: T a -> [a]
  f (T1 x) = [x+1]
  f (T2 v) = [v]

main = f (T1 4)

```

We regard the GADT data constructor T1 as having the type

$$T1 : \forall a. (a \sim_{\#} \text{Int}) \rightarrow \text{Int} \rightarrow T a$$

So in FC, T1 takes three arguments: a type argument to instantiate a , a coercion witnessing the equivalence of a and Int , and a value of type Int . Here is the FC elaboration of `main`:

```
main = f Int (T1 Int <Int> 4)
```

The coercion argument has kind $(\text{Int} \sim_{\#} \text{Int})$, for which the evidence is just $\langle \text{Int} \rangle$ (reflexivity). Similarly, pattern-matching on T1 binds two variables: a coercion variable, and a term variable. Here is the FC elaboration of function `f`:

```

f = /\(a:*) . \(x:T a) .
  case x of
    T1 (c:a ~# Int) (n:Int) -> (Cons (n+1) Nil) |> sym [c]
    T2 (v:a) -> Cons v Nil

```

The cast converts the type of the result from $[\text{Int}]$ to $[\mathbf{a}]$. The coercion $\text{sym } [c]$ is evidence for (or a proof of) the equality of these types, using coercion c , of type $(\mathbf{a} \sim_{\#} \text{Int})$.

2.2 Typing coercions

Figure 6 gives the typing rules for coercions. The rules include unsurprising cases for reflexivity (C_{REFL}), symmetry (C_{SYM}), and transitivity (C_{TRANS}). Rules C_{CALL} and C_{APP} allow us to construct coercions on more complex types from coercions on simpler types.

Rule CINST instantiates a coercion between two \forall -types, to get a coercion between two instantiated types. Rule CVAR allows us to use a coercion that has been introduced to the context by a coercion abstraction $(\lambda c:\tau \sim_{\#} \phi. e)$, or a pattern match against a GADT (as in the example above).

Rule CAX refers to instantiations of *axioms*. In GHC, axioms can arise as a result of *newtype* or *type family* declarations. Consider the following code:

```
newtype N a = MkN (a -> Int)

type family F (x :: *) :: *
type instance F [a] = a
type instance F Bool = Char
```

N is a *newtype* (part of the original Haskell 98 definition), and is desugared to the following FC coercion axiom:

$$C_N a : N a \sim_{\#} a \rightarrow \text{Int}$$

which provides evidence of the equality of types $(N a)$ and $(a \rightarrow \text{Int})$.

In the above Haskell code, F is a *type family* [4, 3], and the two `type instance` declarations above introduce two FC coercion axioms:

$$\begin{aligned} C_1 a & : F [a] \sim_{\#} a \\ C_2 & : F \text{ Bool} \sim_{\#} \text{Char} \end{aligned}$$

Rule CAX describes how these axioms may be used to create coercions. In this particular example, if we have $\gamma : \tau \sim_{\#} \sigma$, then we can prove that $C_1 \gamma : F [\tau] \sim_{\#} \sigma$. Using such coercions we can get, for example, that $(\exists \triangleright \text{sym } (C_1 \langle \text{Int} \rangle)) : F [\text{Int}]$.

Axioms always appear saturated in System FC, hence the syntax $C \bar{\gamma}$ in Figure 2.

3 The problem with large coercions

System FC terms arise as the result of elaboration of source language terms, through type inference. Type inference typically relies on a *constraint solver* [19] which produces System FC witnesses of equality (coercions), that in turn decorate the elaborated term. The constraint solver is not typically concerned with producing small or readable witnesses; indeed GHC's constraint solver can produce large and complex coercions. These complex coercions can make the elaborated term practically impossible to understand and debug.

Moreover, GHC's optimiser transforms well-typed FC terms. Insofar as these transformations involve coercions, the coercions *themselves* may need to be transformed. If you think of the coercions as little proofs that fragments of the program are well-typed, then the optimiser must maintain the proofs as it transforms the terms.

3.1 How big coercions arise

The trouble is that *term-level optimisation tends to make coercions bigger*. The full details of these transformations are given in the so called *push* rules in our previous work [20], but we illustrate them here with an example. Consider this term:

$$(\lambda x.e \triangleright \gamma) a$$

where

$$\begin{aligned}\gamma & : (\sigma_1 \rightarrow \tau_1) \sim_{\#} (\sigma_2 \rightarrow \tau_2) \\ a & : \sigma_2\end{aligned}$$

We would like to perform the beta reduction, but the cast is getting in the way. No matter! We can transform thus:

$$\begin{aligned}& (\lambda x. e \triangleright \gamma) a \\ = & ((\lambda x. e) (a \triangleright \mathit{sym} (\mathit{nth} 0 \gamma))) \triangleright \mathit{nth} 1 \gamma\end{aligned}$$

From the coercion γ we have derived two coercions whose syntactic form is larger, but whose types are smaller:

$$\begin{aligned}\gamma & : (\sigma_1 \rightarrow \tau_1) \sim_{\#} (\sigma_2 \rightarrow \tau_2) \\ \mathit{sym} (\mathit{nth} 0 \gamma) & : \sigma_2 \sim_{\#} \sigma_1 \\ \mathit{nth} 1 \gamma & : \tau_1 \sim_{\#} \tau_2\end{aligned}$$

Here we make use of the coercion combinators sym , which reverses the sense of the proof; and $\mathit{nth} i$, which from a proof of $T \bar{\sigma} \sim_{\#} T \bar{\tau}$ gives a proof of $\sigma_i \sim_{\#} \tau_i$. Finally, we use the derived coercions to cast the argument and result of the function separately. Now the lambda is applied directly to an argument (without a cast in the way), so β -reduction can proceed as desired. Since β -reduction is absolutely crucial to the optimiser, this ability to “push coercions out of the way” is fundamental. Without it, the optimiser is hopelessly compromised.

A similar situation arises with **case** expressions:

$$\mathit{case} (K e_1 \triangleright \gamma) \mathit{of} \{ \dots; K x \rightarrow e_2; \dots \}$$

where K is a data constructor. Here we want to simplify the **case** expression, by picking the correct alternative $K x \rightarrow e_2$, and substituting e_1 for x . Again the coercion gets in the way, but again it is possible to push the coercion out of way.

3.2 How coercions can be simplified

Our plan is to simplify complicated coercion terms into simpler ones, using rewriting. Here are some obvious rewrites we might think of immediately:

$$\begin{aligned}\mathit{sym} (\mathit{sym} \gamma) & \rightsquigarrow \gamma \\ \gamma; \mathit{sym} \gamma & \rightsquigarrow \langle \tau \rangle \quad \text{if } \gamma : \tau \sim_{\#} \phi\end{aligned}$$

But there are much more complicated rewrites to consider. Consider these coercions, where C_N is the axiom generated by the newtype coercion in Section 2.2:

$$\begin{aligned}\gamma_1 & : \tau_1 \sim_{\#} \tau_2 \\ \gamma_2 = \mathit{sym} (C_N \langle \tau_1 \rangle) & : (\tau_1 \rightarrow \mathbf{Int}) \sim_{\#} (N \tau_1) \\ \gamma_3 = N \langle \gamma_1 \rangle & : (N \tau_1) \sim_{\#} (N \tau_2) \\ \gamma_4 = C_N \langle \tau_2 \rangle & : (N \tau_2) \sim_{\#} (\tau_2 \rightarrow \mathbf{Int})\end{aligned}$$

$$\gamma_5 = \gamma_2; \gamma_3; \gamma_4 : (\tau_1 \rightarrow \mathbf{Int}) \sim_{\#} (\tau_2 \rightarrow \mathbf{Int})$$

Here γ_2 takes a function, and wraps it in the newtype; then γ_3 coerces that newtype from $N \tau_1$ to $N \tau_2$; and γ_4 unwraps the newtype. Composing the three gives a rather large, complicated coercion $\gamma_2; \gamma_3; \gamma_4$. *But its type is pretty simple*, and indeed the coercion $\gamma_1 \rightarrow \langle \mathbf{Int} \rangle$ is a much simpler witness of the same equality. The rewrite system we present shortly will rewrite the former to the latter.

Finally, here is an actual example taken from a real program compiled by GHC (don't look at the details!):

$$\begin{aligned} & \text{Mut } \langle v \rangle (\text{sym } (C_{\text{StateT}} \langle s \rangle)) \langle a \rangle \\ & ; \text{sym } (\text{nth } 0 ((\forall \text{wtb}. \text{Mut } \langle w \rangle (\text{sym } (C_{\text{StateT}} \langle t \rangle)) \langle b \rangle \rightarrow \langle \text{ST } t (w b) \rangle) @ v @ s @ a)) \\ \rightsquigarrow & \langle \text{Mut } v s a \rangle \end{aligned}$$

As you can see, the shrinkage in coercion size can be dramatic.

4 Coercion simplification

We now proceed to the details of our coercion simplification algorithm. We note that the design of the algorithm is guided by empirical evidence of its effectiveness on actual programs and that other choices might be possible. Nevertheless, we formally study the properties of this algorithm, namely we will show that it preserves validity of coercions and terminates – even when the rewrite system induced by the axioms is not strongly normalizing.

4.1 Simplification rules

Coercion simplification is given as a non-deterministic relation in Figure 7 and Figure 8. In these two figures we use some syntactic conventions: Namely, for sequences of coercions $\bar{\gamma}_1$ and $\bar{\gamma}_2$, we write $\bar{\gamma}_1; \bar{\gamma}_2$ for the sequence of pointwise transitive compositions and $\text{sym } \bar{\gamma}_1$ for pointwise application of symmetry. We write $\text{nontriv}(\gamma)$ iff γ contains some variable c or axiom application $C \bar{\gamma}$.

We define coercion evaluation contexts, \mathcal{G} , as coercion terms with holes inside them. The syntax of \mathcal{G} allows us to rewrite anywhere inside a coercion. The main coercion evaluation rule is COEVAL. If we are given a coercion γ , we first decompose it to some evaluation context \mathcal{G} with γ_1 in its hole. Rule COEVAL works up to associativity of transitive composition; for example, we will allow the term $(\gamma_1; \gamma_2); \gamma_3$ to be written as $\mathcal{G}[\gamma_2; \gamma_3]$ where $\mathcal{G} = \gamma_1; \square$. This treatment of transitivity is extremely convenient, but we must be careful to ensure that our argument for termination remains robust under associativity (Section 6). Once we have figured out a decomposition $\mathcal{G}[\gamma_1]$, COEVAL performs a single step of rewriting $\Delta \vdash \gamma_1 \rightsquigarrow \gamma_2$ and simply return $\mathcal{G}[\gamma_2]$. Since we are allowed to rewrite coercions under a type environment $(\forall a:\eta. \mathcal{G}$ is a valid coercion evaluation context), Δ (somewhat informally) enumerates the type variables bound by \mathcal{G} . For instance we should be allowed to rewrite $\forall a:\eta. \gamma_1$ to $\forall a:\eta. \gamma_2$. This can happen if $(a:\eta) | - \gamma_1 \rightsquigarrow \gamma_2$. The precondition $\Delta \Vdash \gamma_1 : \sigma \sim_{\#} \phi$ of rule COEVAL ensures that this context corresponds to the decomposition of γ into a context and γ_1 . Moreover, the Δ is passed on to the \rightsquigarrow relation, since some of the rules of the \rightsquigarrow relation that we will present later may have to consult the context Δ to establish preconditions for rewriting.

The soundness property for the \longrightarrow relation is given by the following theorem.

► **Theorem 1** (Coercion subject reduction). *If $\Vdash \gamma_1 : \sigma \sim_{\#} \phi$ and $\gamma_1 \longrightarrow \gamma_2$ then $\Vdash \gamma_2 : \sigma \sim_{\#} \phi$.*

The rewriting judgement $\Delta \vdash \gamma_1 \rightsquigarrow \gamma_2$ satisfies a similar property.

► **Lemma 2.** *If $\Delta \Vdash \gamma_1 : \sigma \sim_{\#} \phi$ and $\Delta \vdash \gamma_1 \rightsquigarrow \gamma_2$ then $\Delta \Vdash \gamma_2 : \sigma \sim_{\#} \phi$.*

To explain coercion simplification, we now present the reaction rules for the \rightsquigarrow relation, organized in several groups.

Coercion evaluation contexts		$\mathcal{G} ::= \square \mid \mathcal{G} \ \gamma \mid \gamma \ \mathcal{G} \mid C \ \bar{\gamma}_1 \mathcal{G} \bar{\gamma}_2 \mid \text{sym} \ \mathcal{G} \mid \forall a:\eta. \mathcal{G} \mid \mathcal{G} @ \tau \mid \mathcal{G}; \gamma \mid \gamma; \mathcal{G}$
$\gamma \cong \mathcal{G}[\gamma_1]$ modulo associativity of ($;$)		$\Delta \vdash^\circ \gamma_1 : \sigma \sim_{\#} \phi \quad \Delta \vdash \gamma_1 \rightsquigarrow \gamma_2$
		COEVAL
$\gamma \longrightarrow \mathcal{G}[\gamma_2]$		
$\Delta \vdash \gamma_1 \rightsquigarrow \gamma_2$		
Reflexivity rules		
REFLAPP	$\Delta \vdash \langle \phi_1 \rangle \langle \phi_2 \rangle \rightsquigarrow \langle \phi_1 \ \phi_2 \rangle$	
REFLALL	$\Delta \vdash \forall a:\eta. \langle \phi \rangle \rightsquigarrow \langle \forall a:\eta. \phi \rangle$	
REFLELIML	$\Delta \vdash \langle \phi \rangle; \gamma \rightsquigarrow \gamma$	
REFLELIMR	$\Delta \vdash \gamma; \langle \phi \rangle \rightsquigarrow \gamma$	
Eta rules		
ETAALLL	$\Delta \vdash (\forall a:\eta. \gamma_1); \gamma_2 @ \phi \rightsquigarrow \gamma_1[\phi/a]; (\gamma_2 @ \phi)$	
ETAALLR	$\Delta \vdash (\gamma_1; (\forall a:\eta. \gamma_2)) @ \phi \rightsquigarrow \gamma_1 @ \phi; \gamma_2[\phi/a]$	
ETANTHL	$\Delta \vdash \text{nth } k \ (\langle H \ \bar{\tau}^{1..l} \rangle \bar{\gamma}; \gamma) \rightsquigarrow \begin{cases} \text{nth } k \ \gamma & \text{if } k \leq l \\ \gamma_{k-l}; \text{nth } k \ \gamma & \text{otherwise} \end{cases}$	
ETANTHR	$\Delta \vdash \text{nth } k \ (\gamma; \langle H \ \bar{\tau}^{1..l} \rangle \bar{\gamma}) \rightsquigarrow \begin{cases} \text{nth } k \ \gamma & \text{if } k \leq l \\ \text{nth } k \ \gamma; \gamma_{k-l} & \text{otherwise} \end{cases}$	
Symmetry rules		
SYMREFL	$\Delta \vdash \text{sym} \langle \phi \rangle \rightsquigarrow \langle \phi \rangle$	
SYMALL	$\Delta \vdash \text{sym} (\forall a:\eta. \gamma) \rightsquigarrow \forall a:\eta. \text{sym} \ \gamma$	
SYMAPP	$\Delta \vdash \text{sym} (\gamma_1 \ \gamma_2) \rightsquigarrow (\text{sym} \ \gamma_1) (\text{sym} \ \gamma_2)$	
SYMTRANS	$\Delta \vdash \text{sym} (\gamma_1; \gamma_2) \rightsquigarrow (\text{sym} \ \gamma_2); (\text{sym} \ \gamma_1)$	
SYMSYM	$\Delta \vdash \text{sym} (\text{sym} \ \gamma) \rightsquigarrow \gamma$	
Reduction rules		
REDNTH	$\Delta \vdash \text{nth } k \ (\langle H \ \bar{\tau}^{1..l} \rangle \bar{\gamma}) \rightsquigarrow \begin{cases} \langle \tau_k \rangle & \text{if } k \leq l \\ \gamma_{k-l} & \text{otherwise} \end{cases}$	
REDINSTCO	$\Delta \vdash (\forall a:\eta. \gamma) @ \phi \rightsquigarrow \gamma[\phi/a]$	
REDINSTTY	$\Delta \vdash (\forall a:\eta. \tau) @ \phi \rightsquigarrow \langle \tau[\phi/a] \rangle$	
Push transitivity rules		
PUSHAPP	$\Delta \vdash (\gamma_1 \ \gamma_2); (\gamma_3 \ \gamma_4) \rightsquigarrow (\gamma_1; \gamma_3) (\gamma_2; \gamma_4)$	
PUSHALL	$\Delta \vdash (\forall a:\eta. \gamma_1); (\forall a:\eta. \gamma_2) \rightsquigarrow \forall a:\eta. \gamma_1; \gamma_2$	
PUSHINST	$\Delta \vdash (\gamma_1 @ \tau); (\gamma_2 @ \tau) \rightsquigarrow (\gamma_1; \gamma_2) @ \tau$	when $\Delta \vdash^\circ \gamma_1; \gamma_2 : \sigma_1 \sim_{\#} \sigma_2$
PUSHNTH	$\Delta \vdash (\text{nth } k \ \gamma_1); (\text{nth } k \ \gamma_2) \rightsquigarrow \text{nth } k \ (\gamma_1; \gamma_2)$	when $\Delta \vdash^\circ \gamma_1; \gamma_2 : \sigma_1 \sim_{\#} \sigma_2$

■ **Figure 7** Coercion simplification (I).

4.1.1 Pulling reflexivity up

Rules REFLAPP, REFLALL, REFLELIML, and REFLELIMR, deal with uses of reflexivity. Rules REFLAPP and REFLALL “swallow” constructors from the coercion language (coercion application, and quantification respectively) into the type language (type application, and quantification respectively). Hence they pull reflexivity as high as possible in the tree structure of a coercion term. Rules REFLELIML and REFLELIMR simply eliminate reflexivity uses that are composed with other coercions.

4.1.2 Pushing symmetry down

Uses of symmetry, contrary to reflexivity, are pushed as close to the leaves as possible or eliminated, (rules SYMREFL, SYMALL, SYMAPP, SYMTRANS, and SYMSYM) only getting stuck at terms of the form $\text{sym } x$ and $\text{sym } (C \bar{\gamma})$. The idea is that by pushing uses of symmetry towards the leaves, the rest of the rules may completely ignore symmetry, except where symmetry-pushing gets stuck (variables or axiom applications).

4.1.3 Reducing coercions

Rules REDNTH, REDINSTCO, and REDINSTTY comprise the first interesting group of rules. They eliminate uses of injectivity and instantiation. Rule REDNTH is concerned with the case where we wish to decompose a coercion of type $H \bar{\phi} \sim_{\#} H \bar{\sigma}$, where the coercion term contains H in its head. Notice that H is a type and may already be applied to some type arguments $\bar{\tau}^{1..l}$, and hence the rule has to account for selection from the first l arguments, or a later argument. Rule REDINSTCO deals with instantiation of a polymorphic coercion with a type. Notice that in rule REDINSTCO the quantified variable may only appear “protected” under some $\langle \sigma \rangle$ inside γ , and hence simply substituting $\gamma[\phi/a]$ is guaranteed to produce a syntactically well-formed coercion. Rule REDINSTTY deals with the instantiation of a polymorphic coercion that is *just* a type.

4.1.4 Eta expanding and subsequent reducing

Redexes of REDNTH and REDINSTCO or REDINSTTY may not be directly visible. Consider $\text{nth } k \ (\langle H \bar{\tau}^{1..l} \rangle \bar{\gamma}; \gamma)$. The use of transitivity stands in our way for the firing of rule REDNTH. To the rescue, rules ETAALLL, ETAALLR, ETANTHL, and ETANTHR, push decomposition or instantiation through transitivity and eliminate such redexes. We call these rules “eta” because in effect we are η -expanding and immediately reducing one of the components of the transitive composition. Here is a decomposition of ETAALLL in smaller steps that involve an η -expansion (of γ_2 in the second line):

$$\begin{aligned} & ((\forall a:\eta.\gamma_1); \gamma_2)@ \phi \\ \rightsquigarrow & ((\forall a:\eta.\gamma_1); (\forall a:\eta.\gamma_2@a))@ \phi \\ \rightsquigarrow & (\forall a:\eta.\gamma_1; \gamma_2@a)@ \phi \rightsquigarrow \gamma_1[\phi/a]; \gamma_2@ \phi \end{aligned}$$

We have merged these steps in a single rule to facilitate the proof of termination. In doing this, we do not lose any reactions, since all of the intermediate terms can also reduce to the final coercion.

There are many design possibilities for rules that look like our η -rules. For instance one may wonder why we are not always expanding terms of the form $\gamma_1; (\forall a:\eta.\gamma_2)$ to $\forall a:\eta.\gamma_1@a; \gamma_2$, whenever γ_1 is of type $\forall a:\eta.\tau \sim_{\#} \forall a:\eta.\phi$. We experimented with several

$$\begin{array}{c}
\frac{\Delta \Vdash^{\circ} c : \tau \sim_{\#} v}{\Delta \vdash c; \text{sym } c \rightsquigarrow \langle \tau \rangle} \text{VARSYM} \qquad \frac{\Delta \Vdash^{\circ} c : \tau \sim_{\#} v}{\Delta \vdash \text{sym } c; c \rightsquigarrow \langle v \rangle} \text{SYMVAR} \\
\\
\frac{(C(\overline{a:\eta}) : \tau \sim_{\#} v) \in \Gamma \quad \overline{a} \subseteq \text{ftv}(v)}{\Delta \vdash C \overline{\gamma}_1; \text{sym}(C \overline{\gamma}_2) \rightsquigarrow [\overline{a} \mapsto \overline{\gamma}_1; \text{sym } \overline{\gamma}_2] \uparrow(\tau)} \text{AXSYM} \qquad \frac{(C(\overline{a:\eta}) : \tau \sim_{\#} v) \in \Gamma \quad \overline{a} \subseteq \text{ftv}(\tau)}{\Delta \vdash \text{sym}(C \overline{\gamma}_1); C \overline{\gamma}_2 \rightsquigarrow [\overline{a} \mapsto \text{sym } \overline{\gamma}_1; \overline{\gamma}_2] \uparrow(v)} \text{SYMAX} \\
\\
\frac{(C(\overline{a:\eta}) : \tau \sim_{\#} v) \in \Gamma \quad \overline{a} \subseteq \text{ftv}(v) \quad \text{nontriv}(\delta) \quad \delta = [\overline{a} \mapsto \overline{\gamma}_2] \uparrow(v)}{\Delta \vdash (C \overline{\gamma}_1); \delta \rightsquigarrow C \overline{\gamma}_1; \overline{\gamma}_2} \text{AXSUCKR} \qquad \frac{(C(\overline{a:\eta}) : \tau \sim_{\#} v) \in \Gamma \quad \overline{a} \subseteq \text{ftv}(\tau) \quad \text{nontriv}(\delta) \quad \delta = [\overline{a} \mapsto \overline{\gamma}_1] \uparrow(\tau)}{\Delta \vdash \delta; (C \overline{\gamma}_2) \rightsquigarrow C \overline{\gamma}_1; \overline{\gamma}_2} \text{AXSUCKL} \\
\\
\frac{(C(\overline{a:\eta}) : \tau \sim_{\#} v) \in \Gamma \quad \overline{a} \subseteq \text{ftv}(\tau) \quad \text{nontriv}(\delta) \quad \delta = [\overline{a} \mapsto \overline{\gamma}_2] \uparrow(\tau)}{\Delta \vdash \text{sym}(C \overline{\gamma}_1); \delta \rightsquigarrow \text{sym}(C \overline{\gamma}_2; \overline{\gamma}_1)} \text{SYMAXSUCKR} \\
\\
\frac{(C(\overline{a:\eta}) : \tau \sim_{\#} v) \in \Gamma \quad \overline{a} \subseteq \text{ftv}(v) \quad \text{nontriv}(\delta) \quad \delta = [\overline{a} \mapsto \overline{\gamma}_1] \uparrow(v)}{\Delta \vdash \delta; \text{sym}(C \overline{\gamma}_2) \rightsquigarrow \text{sym}(C \overline{\gamma}_2; \text{sym } \overline{\gamma}_1)} \text{SYMAXSUCKL}
\end{array}$$

■ **Figure 8** Coercion simplification (II).

variations like this, but we found that such expansions either complicated the termination argument, or did not result in smaller coercion terms. Our rules in effect perform η -expansion *only* when there is a firing reduction directly after the expansion.

4.1.5 Pushing transitivity down

Rules PUSHAPP, PUSHALL, PUSHNTH, and PUSHINST push uses of transitivity *down* the structure of a coercion term, towards the leaves. These rules aim to reveal more redexes at the leaves, that will be reduced by the next (and final) set of rules. Notice that rules PUSHINST and PUSHNTH impose side conditions on the transitive composition $\gamma_1; \gamma_2$. Without these conditions, the resulting coercion may not be well-formed. Take $\gamma_1 = \forall a:\eta. \langle T a a \rangle$ and $\gamma_2 = \forall a:\eta. \langle T a \text{Int} \rangle$. It is certainly the case that $(\gamma_1 @ \text{Int}); (\gamma_2 @ \text{Int})$ is well formed. However, $\Vdash^{\circ} \gamma_1 : \forall a:\eta. T a a \sim_{\#} \forall a:\eta. T a a$ and $\Vdash^{\circ} \gamma_2 : \forall a:\eta. T a \text{Int} \sim_{\#} \forall a:\eta. T a \text{Int}$, and hence $(\gamma_1; \gamma_2) @ \text{Int}$ is not well-formed. A similar argument applies to rule PUSHNTH.

4.1.6 Leaf reactions

When transitivity and symmetry have been pushed as low as possible, new redexes may appear, for which we introduce rules VARSYM, SYMVAR, AXSYM, SYMAX, AXSUCKR, AXSUCKL, SYMAXSUCKR, SYMAXSUCKL. (Figure 8)

- Rules VARSYM and SYMVAR are entirely straightforward: a coercion variable (or its symmetric coercion) meets its symmetric coercion (or the variable) and the result is the identity.
- Rules AXSYM and SYMAX are more involved. Assume that the axiom $(C(\overline{a:\eta}) : \tau \sim_{\#} v) \in \Gamma$, and a well-formed coercion of the form: $C \overline{\gamma}_1; \text{sym}(C \overline{\gamma}_2)$. Moreover $\Delta \Vdash^{\circ}$

$[a \mapsto \gamma] \uparrow(\tau) = \gamma'$	
$[a \mapsto \gamma] \uparrow(a)$	$= \gamma$
$[a \mapsto \gamma] \uparrow(b)$	$= \langle b \rangle$
$[a \mapsto \gamma] \uparrow(H)$	$= \langle H \rangle$
$[a \mapsto \gamma] \uparrow(F)$	$= \langle F \rangle$
$[a \mapsto \gamma] \uparrow(\tau_1 \tau_2)$	$= \begin{cases} \langle \phi_1 \phi_2 \rangle & \text{when } [a \mapsto \gamma] \uparrow(\tau_i) = \langle \phi_i \rangle \\ ([a \mapsto \gamma] \uparrow(\tau_1)) ([a \mapsto \gamma] \uparrow(\tau_2)) & \text{otherwise} \end{cases}$
$[a \mapsto \gamma] \uparrow(\forall b:\eta. \tau)$	$= \begin{cases} \langle \forall a:\eta. \phi \rangle & \text{when } [a \mapsto \gamma] \uparrow(\tau) = \langle \phi \rangle \\ \forall b:\eta. ([a \mapsto \gamma] \uparrow(\tau)) & \text{otherwise } (b \notin \text{ftv}(\gamma), b \neq a) \end{cases}$

■ **Figure 9** Lifting.

$\bar{\gamma}_1 : \bar{\sigma}_1 \sim_{\#} \bar{\phi}_1$ and $\Delta \Vdash \bar{\gamma}_2 : \bar{\sigma}_2 \sim_{\#} \bar{\phi}_2$. Then we know that $\Delta \Vdash C \bar{\gamma}_1; \text{sym}(C \bar{\gamma}_2) : \tau[\bar{\sigma}_1/\bar{a}] \sim_{\#} \tau[\bar{\sigma}_2/\bar{a}]$. Since the composition is well-formed, it must be the case that $v[\bar{\phi}_1/\bar{a}] = v[\bar{\phi}_2/\bar{a}]$. If $\bar{a} \subseteq \text{ftv}(v)$ then it must be $\bar{\phi}_1 = \bar{\phi}_2$. Hence, the pointwise composition $\bar{\gamma}_1; \text{sym} \bar{\gamma}_2$ is well-formed and of type $\bar{\sigma}_1 \sim_{\#} \bar{\sigma}_2$. Consequently, we may replace the original coercion with the *lifting* of τ over a substitution that maps \bar{a} to $\bar{\gamma}_1; \text{sym} \bar{\gamma}_2 : [\bar{a} \mapsto \bar{\gamma}_1; \text{sym} \bar{\gamma}_2] \uparrow(\tau)$.

What is this lifting operation, of a substitution from type variables to coercions, over a type? Its result is a new coercion, and the definition of the operation is given in Figure 9.

The easiest way to understand it is by its effect on a type:

► **Lemma 3** (Lifting). *If $\Delta, (a:\eta) \Vdash \tau : \eta$ and $\Delta \Vdash \gamma : \sigma \sim \phi$ such that $\Delta \Vdash \sigma : \eta$ and $\Delta \Vdash \phi : \eta$, then $\Delta \Vdash [a \mapsto \gamma] \uparrow(\tau) : \tau[\sigma/a] \sim_{\#} \tau[\phi/a]$*

Notice that we have made sure that lifting pulls reflexivity as high as possible in the syntax tree – the only significance of this on-the-fly normalization was that it appeared to simplify the argument we have given for termination of coercion normalization.

Returning to rules AXSYM and SYMAX, we stress that the side condition is essential for the rule to be sound. Consider the following example:

$$C(a:\star) : F [a] \sim_{\#} \text{Int} \in \Gamma$$

Then $(C \langle \text{Int} \rangle); \text{sym}(C \langle \text{Bool} \rangle)$ is well-formed and of type $F [\text{Int}] \sim_{\#} F [\text{Bool}]$, but $\langle F \rangle (\langle \text{Int} \rangle; \text{sym} \langle \text{Bool} \rangle)$ is not well-formed! Rule SYMAX is symmetric and has a similar soundness side condition on the free variables of τ this time.

- The rest of the rules deal with the case when an axiom meets a lifted type – the reaction swallows the lifted type inside the axiom application. For instance, here is rule AXSUCKR:

$$\frac{(C(\bar{a}:\eta):\tau \sim_{\#} v) \in \Gamma \quad \bar{a} \subseteq \text{ftv}(v) \quad \text{nontriv}(\delta) \quad \delta = [\bar{a} \mapsto \bar{\gamma}_2] \uparrow(v)}{\Delta \vdash (C \bar{\gamma}_1); \delta \rightsquigarrow C \bar{\gamma}_1; \bar{\gamma}_2} \text{AXSUCKR}$$

This time let us assume that $\Delta \Vdash \bar{\gamma}_1 : \bar{\sigma}_1 \sim_{\#} \bar{\phi}_1$. Consequently $\Delta \Vdash C \bar{\gamma}_1 : \tau[\bar{\sigma}_1/\bar{a}] \sim_{\#} v[\bar{\phi}_1/\bar{a}]$. Since $\bar{a} \subseteq \text{ftv}(v)$ it must be that $\Delta \Vdash \bar{\gamma}_2 : \bar{\phi}_1 \sim_{\#} \bar{\phi}_3$ for some $\bar{\phi}_3$ and we can pointwise compose $\bar{\gamma}_1; \bar{\gamma}_2$ to get coercions between $\bar{\sigma}_1 \sim_{\#} \bar{\phi}_3$. The resulting coercion $C \bar{\gamma}_1; \bar{\gamma}_2$ is well-formed and of type $\tau[\bar{\sigma}_1/\bar{a}] \sim_{\#} v[\bar{\phi}_3/\bar{a}]$. Rules AXSUCKL, SYMAXSUCKL, and SYMAXSUCKR involve a similar reasoning.

The side condition $\text{nontriv}(\delta)$ is not restrictive in any way – it merely requires that

δ contains some variable c or axiom application. If not, then δ can be converted to reflexivity:

► **Lemma 4.** *If $\vdash^\infty \delta : \sigma \sim_{\#} \phi$ and $\neg \text{nontriv}(\delta)$, then $\delta \longrightarrow^* \langle \phi \rangle$.*

Reflexivity, when transitively composed with any other coercion, is eliminable via REFLELIML/R or and consequently the side condition is not preventing any reactions from firing. It will, however, be useful in the simplification termination proof in Section 6.

The purpose of rules AXSUCKL/R and SYMAXSUCKL/R is to eliminate intermediate coercions in a big transitive composition chain, to give the opportunity to an axiom to meet its symmetric version and react with rules AXSYM and SYMAX. In fact this rule is *precisely* what we need for the impressive simplifications from Section 3. Consider that example again:

$$\begin{aligned}
\gamma_5 &= \gamma_2; \gamma_3; \gamma_4 \\
&= \text{sym}(C_N \langle \tau_1 \rangle); ((N) \gamma_1); (C_N \langle \tau_2 \rangle) && (\text{AXSUCKL with } \delta := ((N) \gamma_1)) \\
\longrightarrow & \text{sym}(C_N \langle \tau_1 \rangle); (C_N (\gamma_1; \langle \tau_2 \rangle)) && (\text{REFLELIMR with } \gamma := \gamma_1, \phi := \tau_2) \\
\longrightarrow & \text{sym}(C_N \langle \tau_1 \rangle); (C_N \gamma_1) && (\text{SYMAX}) \\
\longrightarrow & \langle \rightarrow \rangle (\langle \tau_1 \rangle; \gamma_1) \langle \text{Int} \rangle && (\text{REFLELIML with } \phi := \tau_1, \gamma := \gamma_1) \\
\longrightarrow & \langle \rightarrow \rangle \gamma_1 \langle \text{Int} \rangle
\end{aligned}$$

Notably, rules AXSUCKL/R and SYMAXSUCKL/R generate axiom applications of the form $C \bar{\gamma}$ (with a coercion as argument). In our previous papers, the syntax of axiom applications was $C \bar{\tau}$, with *types* as arguments. But we need the additional generality to allow coercions rewriting to proceed without getting stuck.

5 Coercion simplification in GHC

To assess the usefulness of coercion simplification we added it to GHC. For Haskell programs that make no use of GADTs or type families, the effect will be precisely zero, so we took measurements on two bodies of code. First, our regression suite of 151 tests for GADTs and type families; these are all very small programs. Second, the `Data.Accelerate` library that we know makes use of type families [5]. This library consists of 18 modules, containing 8144 lines of code.

We compiled each of these programs with and without coercion simplification, and measured the percentage reduction in size of the coercion terms with simplification enabled. This table shows the minimum, maximum, and aggregate reduction, taken over the 151 tests and 18 modules respectively. The “aggregate reduction” is obtained by combining all the programs in the group (testsuite or `Accelerate`) into one giant “program”, and computing the reduction in coercion size.

	Testsuite	Accelerate
Minimum	−97%	−81%
Maximum	+14%	0%
Aggregate	−58%	−69%

There is a substantial aggregate decrease of 58% in the testsuite and 69% in `Accelerate`, with a massive 97% decrease in special cases. These special cases should not be taken lightly: in one program the types and coercions taken together were five times bigger than the term they decorated; after simplification they were “only” twice as big. The coercion simplifier makes the compiler less vulnerable to falling off a cliff.

Only one program showed an increase in coercion size, of 14%, which turned out to be the effect of this rewrite:

$$\text{sym}(C; D) \longrightarrow (\text{sym } D); (\text{sym } C)$$

Smaller coercion terms make the compiler faster, but the normalization algorithm itself consumes some time. However, the effect on compile time is barely measurable (less than 1%), and we do not present detailed figures.

Of course none of this would matter if coercions were always tiny, so that they took very little space in the first place. And indeed that is often the case. But for programs that make heavy use of type functions, un-optimised coercions can dominate compile time. For example, the `Accelerate` library makes heavy use of type functions. The time and memory consumption of compiling all 21 modules of the library are as follows:

	Compile time	Memory allocated	Max residency
With coercion optimisation	68s	31 Gbyte	153 Mbyte
Without coercion optimisation	291s	51 Gbyte	2,000 Mbyte

As you can see, the practical effects can be extreme; the cliff is very real.

6 Termination and confluence

We have demonstrated the effectiveness of the algorithm in practice, but we must also establish termination. This is important, since it would not be acceptable for a compiler to loop while simplifying a coercion, no matter what axioms are declared by users. Since the rules fire non-deterministically, and some of the rules (such as `REDINSTCO` or `AXSYM`) create potentially larger coercion trees, termination is not obvious.

6.1 Termination

To formalize a termination argument, we introduce several definitions in Figure 10. The *axiom polynomial* of a coercion over a distinguished variable z , $p(\cdot)$, returns a polynomial with natural number coefficients that can be compared to any other polynomial over z . The *coercion weight* of a coercion is defined as the function $w(\cdot)$ and the *symmetry weight* of a coercion is defined with the function $sw(\cdot)$ in Figure 10. Unlike the polynomial and coercion weights of a coercion, $sw(\cdot)$ does take symmetry into account. Finally, we will also use the *number of coercion applications and coercion \forall -introductions*, denoted with $intros(\cdot)$ in what follows.

Our termination argument comprises of the lexicographic left-to-right ordering of:

$$\mu(\cdot) = \langle p(\cdot), w(\cdot), intros(\cdot), sw(\cdot) \rangle$$

We will show that each of the \rightsquigarrow reductions reduces this tuple. For this to be a valid termination argument for (\longrightarrow) we need two more facts about *each* component measure, namely that (i) $(=)$ and $(<)$ are preserved under arbitrary contexts, and (ii) each component is invariant with respect to the associativity of $(;)$.

► **Lemma 5.** *If $\Delta \Vdash \gamma_1 : \tau \sim_{\#} \sigma$ and $\gamma_1 \cong \gamma_2$ modulo associativity of $(;)$, then $p(\gamma_1) = p(\gamma_2)$, $w(\gamma_1) = w(\gamma_2)$, $intros(\gamma_1) = intros(\gamma_2)$, and $sw(\gamma_1) = sw(\gamma_2)$.*

Proof. This is a simple inductive argument, the only interesting case is the case for $p(\cdot)$ where the reader can calculate that $p(\gamma_1; (\gamma_2; \gamma_3)) = p((\gamma_1; \gamma_2); \gamma_3)$ and by induction we are done. ◀

Axiom polynomial	Coercion weight
$p(\text{sym } \gamma) = p(\gamma)$	$w(\text{sym } \gamma) = w(\gamma)$
$p(C \bar{\gamma}) = z \cdot \Sigma p(\gamma_i) + z + 1$	$w(C \bar{\gamma}) = \Sigma w(\gamma_i) + 1$
$p(c) = 1$	$w(c) = 1$
$p(\gamma_1; \gamma_2) = p(\gamma_1) + p(\gamma_2) + p(\gamma_1) \cdot p(\gamma_2)$	$w(\gamma_1; \gamma_2) = 1 + w(\gamma_1) + w(\gamma_2)$
$p(\langle \phi \rangle) = 0$	$w(\langle \phi \rangle) = 1$
$p(\text{nth } k \ \gamma) = p(\gamma)$	$w(\text{nth } k \ \gamma) = 1 + w(\gamma)$
$p(\gamma @ \phi) = p(\gamma)$	$w(\gamma @ \phi) = 1 + w(\gamma)$
$p(\gamma_1 \ \gamma_2) = p(\gamma_1) + p(\gamma_2)$	$w(\gamma_1 \ \gamma_2) = 1 + w(\gamma_1) + w(\gamma_2)$
$p(\forall a:\eta. \gamma) = p(\gamma)$	$w(\forall a:\eta. \gamma) = 1 + w(\gamma)$
Symmetry weight	
$sw(\text{sym } \gamma) = w(\gamma) + sw(\gamma)$	
$sw(C \bar{\gamma}) = \Sigma sw(\gamma_i)$	
$sw(c) = 0$	
$sw(\gamma_1; \gamma_2) = sw(\gamma_1) + sw(\gamma_2)$	
$sw(\langle \phi \rangle) = 0$	
$sw(\text{nth } k \ \gamma) = sw(\gamma)$	
$sw(\gamma @ \phi) = sw(\gamma)$	
$sw(\gamma_1 \ \gamma_2) = sw(\gamma_1) + sw(\gamma_2)$	
$sw(\forall a:\eta. \gamma) = sw(\gamma)$	

■ **Figure 10** Metrics on coercion terms.

► **Lemma 6.** *If $\Gamma, \Delta \Vdash \gamma_i : \tau \sim_{\#} \sigma$ (for $i = 1, 2$) and $p(\gamma_1) < p(\gamma_2)$ then $p(\mathcal{G}[\gamma_1]) < p(\mathcal{G}[\gamma_2])$ for any \mathcal{G} with $\Gamma \Vdash \mathcal{G}[\gamma_i] : \phi \sim_{\#} \phi'$. Similarly if we replace $(<)$ with $(=)$.*

Proof. By induction on the shape of \mathcal{G} . The only interesting case is the transitivity case again. Let $\mathcal{G} = \gamma; \mathcal{G}'$. Then $p(\gamma; \mathcal{G}'[\gamma_1]) = p(\gamma) + p(\mathcal{G}'[\gamma_1]) + p(\gamma) \cdot p(\mathcal{G}'[\gamma_1])$ whereas $p(\gamma; \mathcal{G}'[\gamma_2]) = p(\gamma) + p(\mathcal{G}'[\gamma_2]) + p(\gamma) \cdot p(\mathcal{G}'[\gamma_2])$. Now, either $p(\gamma) = 0$, in which case we are done by induction hypothesis for $\mathcal{G}'[\gamma_1]$ and $\mathcal{G}'[\gamma_2]$, or $p(\gamma) \neq 0$ in which case again induction hypothesis gives us the result since we are multiplying $p(\mathcal{G}'[\gamma_1])$ and $p(\mathcal{G}'[\gamma_2])$ by the same polynomial. The interesting “trick” is that the polynomial for transitivity contains both the product of the components *and* their sum (since product alone is not preserved by contexts!). ◀

► **Lemma 7.** *If $\Gamma, \Delta \Vdash \gamma_i : \tau \sim_{\#} \sigma$ and $w(\gamma_1) < w(\gamma_2)$ then $w(\mathcal{G}[\gamma_1]) < w(\mathcal{G}[\gamma_2])$ for any \mathcal{G} with $\Gamma \Vdash \mathcal{G}[\gamma_i] : \phi \sim_{\#} \phi'$. Similarly if we replace $(<)$ with $(=)$.*

► **Lemma 8.** *If $\Gamma, \Delta \Vdash \gamma_i : \tau \sim_{\#} \sigma$ and $\text{intros}(\gamma_1) < \text{intros}(\gamma_2)$ then $\text{intros}(\mathcal{G}[\gamma_1]) < \text{intros}(\mathcal{G}[\gamma_2])$ for any \mathcal{G} with $\Gamma \Vdash \mathcal{G}[\gamma_i] : \phi \sim_{\#} \phi'$. Similarly if we replace $(<)$ with $(=)$.*

► **Lemma 9.** *If $\Gamma, \Delta \Vdash \gamma_i : \tau \sim_{\#} \sigma$, $w(\gamma_1) \leq w(\gamma_2)$, and $sw(\gamma_1) < sw(\gamma_2)$ then $sw(\mathcal{G}[\gamma_1]) < sw(\mathcal{G}[\gamma_2])$ for any \mathcal{G} with $\Gamma \Vdash \mathcal{G}[\gamma_i] : \phi \sim_{\#} \phi'$.*

Proof. The only interesting case is when $\mathcal{G} = \text{sym } \mathcal{G}'$ and hence we have that $sw(\mathcal{G}[\gamma_1]) = sw(\text{sym } \mathcal{G}'[\gamma_1]) = w(\mathcal{G}'[\gamma_1]) + sw(\mathcal{G}'[\gamma_1])$. Similarly $sw(\mathcal{G}[\gamma_2]) = w(\mathcal{G}'[\gamma_2]) + sw(\mathcal{G}'[\gamma_2])$. By the precondition for the weights and induction hypothesis we are done. The precondition on the weights is not restrictive, since $w(\cdot)$ has higher precedence than $sw(\cdot)$ inside $\mu(\cdot)$. ◀

The conclusion is the following theorem.

► **Theorem 10.** *If $\gamma \cong \mathcal{G}[\gamma_1]$ modulo associativity of $(;)$ and $\Delta \Vdash \gamma_1 : \sigma \sim_{\#} \phi$, and $\Delta \vdash \gamma_1 \rightsquigarrow \gamma_2$ such that $\mu(\gamma_2) < \mu(\gamma_1)$, it is the case that $\mu(\mathcal{G}[\gamma_2]) < \mu(\gamma)$.*

► **Corollary 11.** (\longrightarrow) *terminates on well-formed coercions if each of the \rightsquigarrow transitions reduces $\mu(\cdot)$.*

Note that often the term rewrite literature requires similar conditions (preservation under contexts and associativity), but also *stability under substitution* (e.g. see [1], Chapter 5). In our setting, variables are essentially treated as constants and this is the reason that we do not rely on stability under substitutions. For instance the rule REFLELIMR $\Delta | - \gamma; \langle \phi \rangle \rightsquigarrow \gamma$ is *not* expressed as $\Delta | - c; \langle \phi \rangle \rightsquigarrow c$, as would be customary in a more traditional term-rewrite system presentation.

We finally show that indeed each of the \rightsquigarrow steps reduces $\mu(\cdot)$.

► **Theorem 12 (Termination).** *If $\Delta \Vdash^{\text{co}} \gamma_1 : \sigma \sim_{\#} \phi$ and $\Delta \vdash \gamma_1 \rightsquigarrow \gamma_2$ then $\mu(\gamma_2) < \mu(\gamma_1)$.*

Proof. It is easy to see that the reflexivity rules, the symmetry rules, the reduction rules, and the η -rules preserve or reduce the polynomial component $p(\cdot)$. The same is true for the push rules but the proof is slightly more interesting. Let us consider PUSHAPP, and let us write p_i for $p(\gamma_i)$. We have that $p((\gamma_1 \ \gamma_2); (\gamma_3 \ \gamma_4)) = p_1 + p_2 + p_3 + p_4 + p_1 p_3 + p_2 p_3 + p_1 p_4 + p_2 p_4$. On the other hand $p((\gamma_1; \gamma_3) (\gamma_2; \gamma_4)) = p_1 + p_3 + p_1 p_3 + p_2 + p_4 + p_2 p_4$ which is a smaller or equal polynomial than the left-hand side polynomial. Rule PUSHALL is easier. Rules PUSHINST and PUSHNTH have exactly the same polynomials on the left-hand and the right-hand side so they are ok. Rules VARSYM and SYMVAR reduce $p(\cdot)$. The interesting bit is with rules AXSYM, SYMAX, and AXSUCKR/L and SYMAXSUCKR/L. We will only show the cases for AXSYM and AXSUCKR as the rest of the rules involve very similar calculations:

- Case SYMAX. We will use the notational convention \bar{p}_1 for $p(\bar{\gamma}_1)$ (a vector of polynomials) and similarly \bar{p}_2 for $p(\bar{\gamma}_2)$. Then the left-hand side polynomial is:

$$\begin{aligned} & (z\Sigma\bar{p}_1 + z + 1) + (z\Sigma\bar{p}_2 + z + 1) + \\ & \quad (z\Sigma\bar{p}_1 + z + 1) \cdot (z\Sigma\bar{p}_2 + z + 1) = \\ & (z^2 + 2z)\Sigma\bar{p}_1 + (z^2 + 2z)\Sigma\bar{p}_2 + z^2\Sigma\bar{p}_1\Sigma\bar{p}_2 + (z^2 + 4z + 3) \end{aligned}$$

For the right-hand side polynomial we know that each $\gamma_{1i}; \text{sym } \gamma_{2i}$ will have polynomial $p_{1i} + p_{2i} + p_{1i}p_{2i}$ and it cannot be repeated inside the lifted type more than a finite number of times (bounded by the maximum number of occurrences of a type variable from \bar{a} in type τ), call it k . Hence the right-hand side polynomial is smaller or equal to:

$$k\Sigma\bar{p}_1 + k\Sigma\bar{p}_2 + k\Sigma(p_{1i}p_{2i}) \leq k\Sigma\bar{p}_1 + k\Sigma\bar{p}_2 + k\Sigma\bar{p}_1\Sigma\bar{p}_2$$

But that polynomial is strictly smaller than the left-hand side polynomial, hence we are done.

- Case AXSUCKR. In this case the left-hand side polynomial is going to be greater or equal to (because of reflexivity inside δ and because some of the \bar{a} variables may appear more than once inside v it is not exactly equal to) the following:

$$\begin{aligned} & (z\Sigma\bar{p}_1 + z + 1) + \Sigma\bar{p}_2 + (z\Sigma\bar{p}_1 + z + 1)\Sigma\bar{p}_2 = \\ & \quad z\Sigma\bar{p}_1\Sigma\bar{p}_2 + z\Sigma\bar{p}_1 + z\Sigma\bar{p}_2 + 2\Sigma\bar{p}_2 + z + 1 \end{aligned}$$

On the other hand, the right-hand side polynomial is:

$$z\Sigma(p_{1i} + p_{2i} + p_{1i}p_{2i}) + z + 1 \leq z\Sigma\bar{p}_1 + z\Sigma\bar{p}_2 + z\Sigma\bar{p}_1\Sigma\bar{p}_2 + z + 1$$

We observe that there is a difference of $2\Sigma\bar{p}_2$, but we know that δ satisfies $\text{nontriv}(\delta)$, and consequently there must exist some variable or axiom application inside one of the $\bar{\gamma}_2$. Therefore, $\Sigma\bar{p}_2$ is *non-zero* and the case is finished.

It is the arbitrary copying of coercions $\bar{\gamma}_1$ and $\bar{\gamma}_2$ in rules AXSYM and SYMAX that prevents simpler measures that only involve summation of coercions for axioms or transitivity. Other reasonable measures such as the height of transitivity uses from the leaves would not be preserved from contexts, due to AXSYM again.

So far we've shown that all rules but the axiom rules preserve the polynomials, and the axiom rules reduce them. We next show that in the remaining rules, some other component reduces, lexicographically. Reflexivity rules reduce $w(\cdot)$. Symmetry rules preserve $w(\cdot)$ and $intros(\cdot)$ but reduce $sw(\cdot)$. Reduction rules and η -rules reduce $w(\cdot)$. Rules PUSHAPP and PUSHALL preserve or reduce $w(\cdot)$ but certainly reduce $intros(\cdot)$. Rules PUSHINST and PUSHNTH reduce $w(\cdot)$. ◀

We conclude that (\longrightarrow) terminates.

6.2 Confluence

Due to the arbitrary types of axioms and coercion variables in the context, we do not expect confluence to be true. Here is a short example that demonstrates the lack of confluence; assume we have the following in our context:

$$\begin{aligned} C_1 (a:\star \rightarrow \star) &: F \ a \sim_{\#} a \\ C_2 (a:\star \rightarrow \star) &: G \ a \sim_{\#} a \end{aligned}$$

Consider the coercion:

$$(C_1 \langle \sigma \rangle); sym (C_2 \langle \sigma \rangle)$$

of type $F \sigma \sim_{\#} G \sigma$. In one reduction possibility, using rule AXSUCKR, we may get

$$C_1 (sym (C_2 \langle \sigma \rangle))$$

In another possibility, using SYMAXSUCKL, we may get

$$sym (C_2 (sym (C_1 \langle \sigma \rangle)))$$

Although the two normal forms are different, it is unclear if one of them is “better” than the other.

Despite this drawback, confluence or syntactic characterization of normal forms is, for our purposes, of secondary importance (if possible at all for open coercions in such an under-constrained problem!), since we never reduce coercions for the purpose of comparing their normal forms. That said, we acknowledge that experimental results may vary with respect to the actual evaluation strategy, but we do not expect wild variations.

7 Related and future work

Traditionally, work on proof theory is concerned with proof normalization theorems, namely cut-elimination. Category and proof theory has studied the commutativity of diagrams in *monoidal categories* [12], establishing coherence theorems. In our setting Lemma 4 expresses such a result: any coercion that does not include axioms or free coercion variables is equivalent to reflexivity. More work on proof theory is concerned with cut-elimination theorems – in our setting eliminating transitivity completely is plainly impossible due to the presence of axioms. Recent work on *2-dimensional type theory* [10] provides an equivalence relation on equality proofs (and terms), which suffices to establish that types enjoy canonical forms.

Although that work does not provide an algorithm for checking equivalence (this is harder to do because of actual computation embedded with isomorphisms), that definition shares many rules with our normalization algorithm. Finally there is a large literature in associative commutative rewrite systems [7, 2].

To our knowledge, most programming languages literature on coercions is not concerned with coercion simplification but rather with inferring the placement of coercions in source-level programs. Some recent examples are [11] and [17]. A comprehensive study of coercions *and their normalization* in programming languages is that of [8], motivated by coercion placement in a language with *type dynamic*. Henglein’s coercion language differs to ours in that (i) coercions there are not symmetric, (ii) do not involve polymorphic axiom schemes and (iii) may have computational significance. Unlike us, Henglein is concerned with characterizations of minimal coercions and confluence, fixes an equational theory of coercions, and presents a normalization algorithm for that equational theory. In our case, in the absence of a denotational semantics for System FC and its coercions, such an axiomatization would be no more ad-hoc than the algorithm and hence not particularly useful: for instance we could consider adding type-directed equations like $\Delta \vdash \gamma \rightsquigarrow \langle \tau \rangle$ when $\Delta \Vdash \gamma : \tau \sim_{\#} \tau$, or other equations that only hold in consistent or confluent axiom sets. It is certainly an interesting direction for future work to determine whether there even exists a maximal syntactic axiomatization of equalities between coercions with respect to some denotational semantics of System FC.

In the space of typed intermediate languages, xMLF[14] is a calculus with coercions that capture *instantiation* instead of equality, and which serves as target for the MLF language. Although the authors are not directly concerned with normalization as part of an intermediate language simplifier, their translation of the graph-based instantiation witnesses does produce xMLF normal proofs.

Finally, another future work direction would be to determine whether we can encode coercions as λ -terms, and derive coercion simplification by normalization in some suitable λ -calculus.

Acknowledgments

Thanks to Tom Schrijvers for early discussions and for contributing a first implementation. We would particularly like to thank Thomas Ströder for his insightful and detailed feedback in the run-up to submitting the final paper.

References

- 1 Franz Baader and Tobias Nipkow. *Term rewriting and all that*. Cambridge University Press, New York, NY, USA, 1998.
- 2 Leo Bachmair and David A. Plaisted. Termination orderings for associative-commutative rewriting systems. *J. Symb. Comput.*, 1(4):329–349, December 1985.
- 3 Manuel M. T. Chakravarty, Gabriele Keller, and Simon Peyton Jones. Associated type synonyms. In *ICFP’05: Proceedings of the Tenth ACM SIGPLAN International Conference on Functional Programming*, pages 241–253, New York, NY, USA, 2005. ACM.
- 4 Manuel M. T. Chakravarty, Gabriele Keller, Simon Peyton Jones, and Simon Marlow. Associated types with class. *SIGPLAN Not.*, 40(1):1–13, 2005.
- 5 Manuel M.T. Chakravarty, Gabriele Keller, Sean Lee, Trevor L. McDonell, and Vinod Grover. Accelerating Haskell array codes with multicore GPUs. In *Proceedings of the sixth workshop on Declarative Aspects of Multicore Programming, DAMP’11*, pages 3–14, New York, NY, USA, 2011. ACM.

- 6 James Cheney and Ralf Hinze. First-class phantom types. CUCIS TR2003-1901, Cornell University, 2003.
- 7 Nachum Dershowitz, Jien Hsiang, N. Alan Josephson, and David A. Plaisted. Associative-commutative rewriting. In *Proceedings of the Eighth international joint conference on Artificial intelligence – Volume 2*, IJCAI'83, pages 940–944, San Francisco, CA, USA, 1983. Morgan Kaufmann Publishers Inc.
- 8 Fritz Henglein. Dynamic typing: syntax and proof theory. *Sci. Comput. Program.*, 22:197–230, June 1994.
- 9 Oleg Kiselyov, Simon Peyton Jones, and Chung-chieh Shan. Fun with type functions. In Cliff Jones and Bill Roscoe, editors, *Reflections on the work of CAR Hoare*. Springer, 2010.
- 10 Daniel R. Licata and Robert Harper. Canonicity for 2-dimensional type theory. In *Proceedings of the 39th annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL'12, pages 337–348, New York, NY, USA, 2012. ACM.
- 11 Zhaohui Luo. Coercions in a polymorphic type system. *Mathematical Structures in Computer Science*, 18(4):729–751, 2008.
- 12 Saunders Mac Lane. *Categories for the Working Mathematician*. Number 5 in Graduate Texts in Mathematics. Springer-Verlag, 1971.
- 13 Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Geoffrey Washburn. Simple unification-based type inference for GADTs. In *ICFP'06: Proceedings of the Eleventh ACM SIGPLAN International Conference on Functional Programming*, pages 50–61, New York, NY, USA, 2006. ACM Press.
- 14 Didier Rémy and Boris Yakobowski. A Church-style intermediate language for MLF. In Matthias Blume, Naoki Kobayashi, and German Vidal, editors, *Functional and Logic Programming*, volume 6009 of *Lecture Notes in Computer Science*, pages 24–39. Springer Berlin / Heidelberg, 2010.
- 15 Tim Sheard and Emir Pasalic. Meta-programming with built-in type equality. In *Proc 4th International Workshop on Logical Frameworks and Meta-languages (LFM'04)*, Cork, pages 106–124, July 2004.
- 16 Martin Sulzmann, Manuel M. T. Chakravarty, Simon Peyton Jones, and Kevin Donnelly. System F with type equality coercions. In *TLDI'07: Proceedings of the 2007 ACM SIGPLAN International Workshop on Types in Languages Design and Implementation*, pages 53–66, New York, NY, USA, 2007. ACM.
- 17 Nikhil Swamy, Michael Hicks, and Gavin M. Bierman. A theory of typed coercions and its applications. In *Proceedings of the 14th ACM SIGPLAN International Conference on Functional Programming*, ICFP'09, pages 329–340, New York, NY, USA, 2009. ACM.
- 18 Dimitrios Vytiniotis, Simon Peyton Jones, and Pedro Magalhaes. Equality proofs and deferred type errors. In *Proceedings of ACM SIGPLAN International Conference on Functional Programming (ICFP'12)*, pages 341–352, 2012.
- 19 Dimitrios Vytiniotis, Simon Peyton Jones, Tom Schrijvers, and Martin Sulzmann. Outsidein(x): modular type inference with local assumptions. *Journal of Functional Programming*, 21, 2011.
- 20 Stephanie Weirich, Dimitrios Vytiniotis, Simon Peyton Jones, and Steve Zdancewic. Generative type abstraction and type-level computation. In *Proceedings of the 38th annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL'11, pages 227–240, New York, NY, USA, 2011. ACM.