

Safe Zero-cost Coercions for Haskell

Joachim Breitner

Karlsruhe Institute of Technology
breitner@kit.edu

Richard A. Eisenberg

University of Pennsylvania
eir@cis.upenn.edu

Simon Peyton Jones

Microsoft Research
simonpj@microsoft.com

Stephanie Weirich

University of Pennsylvania
sweirich@cis.upenn.edu

Abstract

Generative type abstractions – present in Haskell, OCaml, and other languages – are useful concepts to help prevent programmer errors. They serve to create new types that are distinct at compile time but share a run-time representation with some base type. We present a new mechanism that allows for zero-cost conversions between generative type abstractions and their representations, even when such types are deeply nested. We prove type safety in the presence of these conversions and have implemented our work in GHC.

Categories and Subject Descriptors D.3.3 [Programming Languages]: Language Constructs and Features—abstract data types; F.3.3 [Logics and Meanings of Programs]: Studies of Program Constructs—Type structure

Keywords Haskell; Coercion; Type class; Newtype deriving

1. Introduction

Modular languages support *generative type abstraction*, the ability for programmers to define application-specific types, and rely on the type system to distinguish between these new types and their underlying representations. Type abstraction is a powerful tool for programmers, enabling both flexibility (implementors can change representations) and security (implementors can maintain invariants about representations). Typed languages provide these mechanisms with zero run-time cost – there should be no performance penalty for creating abstractions – using mechanisms such as ML’s module system [MTHM97] and Haskell’s **newtype** declaration [Mar10].

For example, a Haskell programmer might create an abstract type for HTML data, representing them as Strings (Figure 1). Although String values use the same patterns of bits in memory as HTML values, the two types are distinct. That is, a

```
module Html( HTML, text, unMk, ... ) where
  newtype HTML = Mk String
  unMk :: HTML → String
  unMk (Mk s) = s
  text :: String → HTML
  text s = Mk (escapeSpecialCharacters s)
```

Figure 1. An abstraction for HTML values

String will not be accepted by a function expecting an HTML. The constructor Mk converts a String to an HTML (see function text), while using Mk in a pattern converts in the other direction (see function unMk). By exporting the type HTML, but not its data constructor, module Html ensures that the type HTML is *abstract* – clients cannot make arbitrary strings into HTML – and thereby prevent cross-site scripting attacks.

Using **newtype** for abstraction in Haskell has always suffered from an embarrassing difficulty. Suppose in the module Html, the programmer wants to break HTML data into a list of lines:

```
linesH :: HTML → [HTML]
linesH h = map Mk (lines (unMk h))
```

To get the resulting [HTML] we are forced to map Mk over the list. Operationally, this map is the identity function – the run-time representation of [String] is identical to [HTML] – *but it will carry a run-time cost nevertheless*. The optimiser in the Glasgow Haskell Compiler (GHC) is powerless to fix the problem, because it works over a *typed* intermediate language; the Mk constructor changes the type of its operand, and hence cannot be optimised away. There is nothing that the programmer can do to prevent this run-time cost. What has become of the claim of zero-overhead abstraction?

In this paper we describe a robust, simple mechanism that programmers can use to solve this problem, making the following contributions:

- We describe the design of *safe coercions* (Section 2), which introduces the function

```
coerce :: Coercible a b ⇒ a → b
```

and a new type class Coercible. This function performs a zero-cost conversion between two types a and b that have the same representation. The crucial question becomes

[Copyright notice will appear here once ‘preprint’ option is removed.]

what instances of Coercible exist? We give a simple but non-obvious strategy (Sections 2.1–2.2), expressed largely in the familiar language of Haskell type classes.

- We formalise Coercible by translation into GHC’s intermediate language System FC, augmented with the concept of *roles* (Section 2.2), adapted from prior work [WVPZ11]. Our new contribution is a significant simplification of the roles idea in System FC; we formalise this simpler system and give the usual proofs of preservation and progress in Section 4.
- Adding safe coercions to the source language raises new issues for abstract types, and for the coherence of type elaboration. We articulate the issues, and introduce *role annotations* to solve them (Section 3).
- It would be too onerous to insist on programmer-supplied role annotations for every type, so we give a *role inference algorithm* in Section 5.
- To support our claim of practical utility, we have implemented the whole scheme in GHC (Section 6), and evaluated it against thousands of Haskell libraries (Section 9).

Our work finally resolves a notorious and long-standing bug in GHC (#1496), which concerns the interaction of newtype coercions with type families (Section 7). While earlier work [WVPZ11] was motivated by the same bug, it was too complicated to implement. Our new approach finds a sweet spot, offering a considerably simpler system in exchange for a minor loss of expressiveness (Sections 8 and 10).

As this work demonstrates, the interactions between type abstraction and advanced type system features, such as type families and GADTs, are subtle. The ability to create and enforce zero-cost type abstraction is not unique to Haskell – notably the ML module system also provides this capability, and more. As a result, OCaml developers are now grappling with similar difficulties. We discuss the connection between roles and OCaml’s variance annotations (Section 8), as well as other related work.

2. The design and interface of Coercible

We begin by focusing exclusively on the programmer’s-eye-view of safe coercions. We need no new syntax; rather, the programmer simply sees a new API, provided in just two declarations:

```
class Coercible a b
coerce :: Coercible a b => a -> b
```

The type class Coercible is abstract, i.e. its methods are not visible. It differs from other type classes in a few minor points: The user cannot create manual instances; instances are automatically generated by the compiler; and the visibility of instances is conditional. Generally, users can think of it as a normal type class, which is a nice property of the design.

The key principle is this: *If two types s and t are related by Coercible $s\ t$, then s and t have bit-for-bit identical run-time representations.* Moreover, as you can see from the type of coerce, if Coercible $s\ t$ holds then coerce can convert a value of type s to one of type t . And that’s it!

The crucial question, to which we devote the rest of this section and the next, becomes this: exactly when does Coercible $s\ t$ hold? To what your appetite consider these declarations:

```
newtype Age      = MkAge Int
newtype AgeRange = MkAR (Int,Int)
newtype BigAge   = MkBig Age
```

GHC generates the following instances of Coercible:

- (1) **instance** Coercible a a
- (2) For every **newtype** NT $x = \text{MkNT } (T\ x)$, the instances

```
instance Coercible (T x) b => Coercible (NT x) b
instance Coercible a (T x) => Coercible a (NT x)
```

which are visible if and only if the constructor MkNT is in scope.

- (3) For every type constructor TC $r\ p\ n$, where

- r stands for TC’s parameters at role representational,
- p for those at role phantom and
- n for those at role nominal,

the instance

```
instance Coercible r1 r2 =>
  Coercible (TC r1 p1 n) (TC r2 p2 n)
```

Figure 2. Coercible instances

Here are some coercions that hold, so that a single call to coerce suffices to convert between the two types:

- Coercible Int Age: we can coerce from Int to Age at zero cost; this is simply the MkAge constructor.
- Coercible Age Int: and the reverse; this is pattern matching on MkAge.
- Coercible [Age] [Int]: lifting the coercion over lists.
- Coercible (Either Int Age) (Either Int Int): lifting the coercion over Either.
- Coercible (Either Int Age) (Either Age Int): this is more complicated, because first argument of Either must be coerced in one direction, and the second in the other.
- Coercible (Int -> Age) (Age -> Int): all this works over function arrows too.
- Coercible (Age, Age) AgeRange: we have to unwrap the pair of Ages and then wrap with MkAR.
- Coercible [BigAge] [Int]: two levels of coercion.

In the rest of this section we will describe how Coercible constraints are solved or, equivalently, which instances of Coercible exist. (See Figure 2 for a concise summary.)

2.1 Coercing newtypes

Since Coercible relates a newtype with its base type, we need Coercible instance declarations for every such newtype. The naive **instance** Coercible Int Age does not work well, for reasons explained in the box on page 3, so instead we generate *two* instances for each newtype:

```
instance Coercible a Int => Coercible a Age  — (A1)
instance Coercible Int b => Coercible Age b  — (A2)
```

```
instance Coercible a Age => Coercible a BigAge — (B1)
instance Coercible Age b => Coercible BigAge b — (B2)
```

```
instance Coercible a AgeRange => Coercible a (Int,Int)
instance Coercible AgeRange b => Coercible (Int,Int) b
```

Notice that each instance unwraps just one layer of the newtype, so we call them the “unwrapping instances”.

If we now want to solve, say, a constraint `Coercible s Age`, for any type `s`, we can use (A1) to reduce it to the simpler goal `Coercible s Int`. A more complicated, two-layer coercion `Coercible BigAge Int` is readily reduced, in two such steps, to `Coercible Int Int`. All we need now is for GHC to have a built-in witness of reflexivity, expressing that any type has the same run-time representation as itself:

```
instance Coercible a a
```

This simple scheme allows coercions that involve arbitrary levels of wrapping or unwrapping, in either direction, with a single call to `coerce`. The solution path is not fully determined, but that does not matter. For example, here are two ways to solve `Coercible BigAge Age`:

```

Coercible BigAge Age
→ Coercible BigAge Int   — By (A1)
→ Coercible Age Int      — By (B2)
→ Coercible Int Int      — By (A2)
→ solved                  — By reflexivity

```

```

Coercible BigAge Age
→ Coercible Age Age      — By (B2)
→ solved                  — By reflexivity

```

Since `Coercible` constraints have no run-time behaviour (unlike normal type class constraints), we have no concerns about incoherence; any solution will do.

The newtype-unwrapping instances (i.e., (2) in Figure 2) are available *only if the corresponding newtype data constructor (Mk in our current example) is in scope*; this is required to preserve abstraction, as we explain in Section 3.1.

2.2 Coercing parameters of type constructors

As Figure 2 shows, as well as the unwrapping instances for a **newtype**, we also generate one instance for each type constructor, including data types, newtypes the function type, and built-in data types like tuples. We call this instance the “lifting instance” for the type, because it lifts coercions through the type. The shape of the instance depends on the so-called *roles* of the type constructor. Each type parameter of a type constructor has a role, determined by the way in which the parameter is used in the definition of the type constructor. In practice, the roles of a declared data type are determined by a role inference algorithm (Section 5) and can be modified by role annotations (Section 3.1). Once defined, the roles of a type constructor are the same in every scope, regardless of whether the concrete definition of that type is available in that scope.

Roles, a development of earlier work [WVPZ11] (Section 8), are a new concept for the programmer. In the following subsections, we discuss how the three possible roles, *representational*, *phantom* and *nominal*, ensure that lifting instances do not violate type safety by allowing coercions between types with different run-time representations.

2.2.1 Coercing representational type parameters

The most common role is *representational*. It is the role that is assigned to the type parameters of ordinary newtypes and data types like `Maybe`, the list type and `Either`. The `Coercible` instances for these type constructors are:

```

instance Coercible a b => Coercible (Maybe a) (Maybe b)
instance Coercible a b => Coercible [a] [b]
instance (Coercible a1 b1, Coercible a2 b2)
=> Coercible (Either a1 a2) (Either b1 b2)

```

Why a single instance is not enough

Why do we create two instances for every newtype, rather than just the single declaration

```
instance Coercible Int Age
```

to witness the fact that `Int` and `Age` have the same run-time representation?

That would indeed allow us to convert from `Int` to `Age`, using `coerce`, but what about the reverse direction? We then might need a second function

```
uncoerce :: Coercible a b => b -> a
```

although it would be tiresome for the programmer to remember which one to call. Alternatively, perhaps GHC should generate *two* instances:

```
instance Coercible Int Age
instance Coercible Age Int
```

But how would we get from `BigAge` to `Int`? We could try this:

```
down :: BigAge -> Int
down x = coerce (coerce x)
```

Our intent here is that each invocation of `coerce` unwraps one “layer” of newtype. But this is not good, because the type inference engine cannot figure out which type to use for the result of the inner `coerce`. To make the code typecheck we would have to add a type signature:

```
down :: BigAge -> Int
down x = coerce (coerce x :: Age)
```

Not very nice. Moreover we would prefer to do all this with a *single* call to `coerce`, implying that `Coercible BigAge Int` must hold. That might make us consider adding the instance declaration

```
instance (Coercible a b, Coercible b c) => Coercible a c
```

to express the transitivity of `Coercible`. But now the problem of the un-specified intermediate type `b` re-appears, and cannot be solved with a type signature.

All of these problems are nicely solved using the instances in Figure 2.

These instances are just as you would expect: for example, the type `Maybe t1` and `Maybe t2` have the same run-time representation if and only if `t1` and `t2` have the same representation.

Most primitive type constructors also have representational roles for their arguments. For example, the domain and co-domain of arrow types are representational, giving rise to the following `Coercible` instance:

```
instance (Coercible a1 b1, Coercible a2 b2)
=> Coercible (a1 -> a2) (b1 -> b2)
```

Likewise, the type `IORef` has a representational parameter, so expressions of type `IORef Int` can be converted to type `IORef Age` for zero cost (and outside of the `IO monad`).

Returning to the introduction, we can use these instances to write `linesH` very directly, thus:

```
linesH :: HTML → [HTML]
linesH = coerce lines
```

In this case, the call to `coerce` gives rise to a constraint `Coercible (String → [String]) (HTML → [HTML])`, which gets simplified to `Coercible String HTML` using the instances for arrow and list types. Then the instance for the newtype `HTML` reduces it to `Coercible String String`, which is solved by the reflexive instance.

2.2.2 Coercing phantom type parameters

A type parameter has a *phantom* role if it does not occur in the definition of the type, or if it does, then only as a phantom parameter of another type constructor. For example, these declarations

```
data Phantom b = Phantom
data NestedPhantom b = L [Phantom b] | SomethingElse
```

both have parameter `b` at a phantom role.

When do the types `Phantom t1` and `Phantom t2` have the same run-time representation? Always! Therefore, we have the instances

```
instance Coercible (Phantom a) (Phantom b)
instance Coercible (NestedPhantom a) (NestedPhantom b)
```

and `coerce` can be used to change the phantom parameter arbitrarily.

2.2.3 Coercing nominal type parameters

In contrast, the *nominal* role induces the strictest preconditions for `Coercible` instances. This role is assigned to a parameter that possibly affects the run-time representation of a type, commonly because it is passed to a type function. For example, consider the following code

```
type family EncData a where
  EncData String = (ByteString, Encoding)
  EncData HTML = ByteString
```

```
data Encoding = ...
data EncText a = MkET (EncData a)
```

Even though we have `Coercible HTML String`, it would be wrong to derive the instance `Coercible (EncText HTML) (EncText String)`, because these two types have quite different run-time representations! Therefore, there are no instances that change a nominal parameter of a type constructor.

All parameters of a type or data *family* have nominal role, because they could be inspected by the type family instances. For similar reasons, the non-uniform parameters to GADTs are also required to be nominal.

2.2.4 Coercing multiple type parameters

A type constructor can have multiple type parameters, each at a different role. In that case, an appropriate constraint for each type parameter is used:

```
data Params r p n = Con1 (Maybe r) | Con2 (EncData n)
```

yields the instance

```
instance Coercible r1 r2
  ⇒ Coercible (Params r1 p1 n) (Params r2 p2 n)
```

This instance expresses that the representational type parameters may change if there is a `Coercible` instance for them; the phantom type parameters may change arbitrarily; and the nominal type parameters must stay the same.

3. Abstraction and coherence

The purpose of the `HTML` type from the introduction is to prevent accidentally mixing up unescaped strings and `HTML` fragments. Rejecting programs that make this mistake is not a matter of type safety as traditionally construed, but rather of preserving a desired abstraction.

While the previous section described how the `Coercible` instances ensure that uses of `coerce` are type safe, this section discusses two other properties: *abstraction* and *class coherence*.

3.1 Preserving abstraction

When the constructors of a type are in scope then we can write code semantically equivalent to `coerce` by hand (although it might be less efficient). In this situation, the use of `coerce` should definitely be allowed. However, when the constructors are not in scope, it turns out that we sometimes want the lifting instance, and sometimes we do *not* want it.

The newtype unwrapping instance is directly controlled by the visibility of the constructor and can be used if and only if this is in scope. (See Section 2.1 for how this is accomplished.) For example, since the author of module `Html` did not export `Mk`, a client does not see the unwrapping instances for `HTML`, and the abstraction is preserved.

However, we permit the use of the coercion lifting instance for a type constructor even when the data constructors are not available. For example, built-in types like `IORef` or the function type `(→)` do not even have constructors that can be in scope. Nevertheless, coercing from `IORef HTML` to `IORef String` and from `HTML → HTML` to `String → String` should be allowed.

Therefore the rule for the lifting instance is that it can be used independent of the visibility of constructors. Instead, its form – what coercions it allows – is controlled by the roles of the type constructor’s parameters.

Library authors can control the roles assigned to type constructors using *role annotations*. In many cases, the role inferred by the type checker is sufficient, even for abstract types. Consider a library for non-empty lists:

```
module NonEmptyListLib( NE, singleton, ... ) where
  data NE a = MkNE [a]
  singleton :: a → NE a
  ... etc...
```

The type must be exported abstractly; otherwise, the non-empty property can be broken by its users. Nevertheless lifting a coercion through `NE`, i.e. coercing `NE HTML` to `NE String`, should be allowed. Therefore, the role of `NE`’s parameter should be representational. In this case, the library author does not have to actively set it: As it is the most permissive type-safe role, the role inference algorithm (Section 5.2) already chooses representational.

However, sometimes library authors must restrict the usage of the lifting coercion to ensure that the invariants of their abstract types can be preserved. For example, consider the data type `Map k v`, which implements an efficient finite map from keys of type `k` to values of type `v`, using an internal representation based on a balanced tree, something like this:

```
data Map k v = Leaf | Node k v (Map k v) (Map k v)
```

It would be disastrous if the user were allowed to coerce from $(\text{Map Age } v)$ to $(\text{Map Int } v)$, because a valid tree with regard to the ordering of Age might be completely bogus when using the ordering of Int.

To prevent that difficulty, the author specifies

type role Map nominal representational

As explained in Section 2.2, we now have the desirable and useful lifting instance

instance Coercible a b \Rightarrow Coercible (Map k a) (Map k b)

which allows the coercion from Map k HTML to Map k String.

Note that in the declaration of Map the parameters k and v are used in exactly the same way, so this distinction cannot be made by the compiler; it can only be specified by the programmer. However, the compiler ensures that programmer-specified role annotations cannot subvert the type system: if the annotation specifies an unsafe role, the compiler will reject the program.

3.2 Preserving class coherence

Another property of Haskell, independent of type-safety, is the coherence of type classes. There should only ever be one class instance for a particular class and type. We call this desirable property *coherence*. Without extra checks, Coercible could be used to create incoherence.

Consider this (non-Haskell98) data type, which reifies a Show instance as a value:

data HowToShow a **where**

MkHTS :: Show a \Rightarrow HowToShow a

showH :: HowToShow a \rightarrow a \rightarrow String

showH MkHTS x = show x

Here showH pattern-matches on a HowToShow value, and uses the instance stored inside it to obtain the show method. If we are not careful, the following code would break the coherence of the Show type class:

instance Show HTML **where**

show (Mk s) = "HTML:" ++ show s

stringShow :: HowToShow String

stringShow = MkHTS

htmlShow :: HowToShow HTML

htmlShow = MkHTS

badShow :: HowToShow HTML

badShow = coerce stringShow

$\lambda >$ showH stringShow "Hello"

"Hello"

$\lambda >$ showH htmlShow (Mk "Hello")

"HTML:Hello"

$\lambda >$ showH badShow (Mk "Hello")

"Hello"

In the final example we were applying show to a value of type HTML, but the Show instance for String (coerced to (Show HTML)) was used.

To avoid this confusion, the parameters of a type class are all assigned a *nominal* role by default. Accordingly, the parameter of HowToShow is also assigned a nominal role by default, preventing the coercion between (HowToShow HTML) and (HowToShow String).

Metavariables:

x	term	α, β	type	c	coercion
C	axiom	D	data type	N	newtype
F	type family	K	data constructor		
e	$::= \lambda c:\phi.e \mid e \gamma \mid e \triangleright \gamma \mid \dots$				terms
τ, σ	$::= \alpha \mid \tau_1 \tau_2 \mid \forall \alpha:\kappa.\tau \mid H \mid F(\bar{\tau})$				types
κ	$::= * \mid \kappa_1 \rightarrow \kappa_2$				kinds
H	$::= (\rightarrow) \mid (\Rightarrow) \mid (\sim_{\rho}^{\kappa}) \mid T$				type constants
T	$::= D \mid N$				algebraic data types
ϕ	$::= \tau \sim_{\rho}^{\kappa} \sigma$				proposition
γ, η	$::=$				coercions
	$\mid \langle \tau \rangle \mid \langle \tau, \sigma \rangle_P \mid \mathbf{sym} \gamma \mid \gamma_1 \ddagger \gamma_2$ $\mid H(\bar{\gamma}) \mid F(\bar{\gamma}) \mid \gamma_1 \gamma_2 \mid \forall \alpha:\kappa.\gamma$ $\mid c \mid C(\bar{\tau})$				equivalence congruence assumptions
	$\mid \mathbf{nth}^i \gamma \mid \mathbf{left} \gamma \mid \mathbf{right} \gamma \mid \gamma @ \tau$ $\mid \mathbf{sub} \gamma$				decomposition sub-rolling
ρ	$::= N \mid R \mid P$				roles
Γ	$::= \emptyset \mid \Gamma, \alpha:\kappa \mid \Gamma, c:\phi \mid \Gamma, x:\tau$				typing contexts
Ω	$::= \emptyset \mid \Omega, \alpha:\rho$				role contexts

Figure 3. An excerpt of the grammar of System FC

4. Ensuring type safety: System FC with roles

Haskell is a large and complicated language. How do we know that the ideas sketched above in source-language terms are actually sound? What, precisely, do roles mean, and when precisely are two types equal? In this section we answer these questions for GHC's small, statically-typed intermediate language, GHC Core. Every Haskell program is translated into Core, and we can typecheck Core to reassure ourselves that the (large, complicated) front end accepts only good programs.

Core is an implementation of a calculus called System FC, itself an extension of the classical Girard/Reynolds System F. The version of FC that we develop in this paper derives from much prior work.¹ However, for clarity we give a self-contained description of the system and do not assume familiarity with previous versions.

Figure 3 gives the syntax of System FC. The starting point is an entirely conventional lambda calculus in the style of System F. We therefore elide most of the syntax of terms e , giving the typing judgement for terms in the extended version of this paper [BEPW14]. Types τ are also conventional, except that we add (saturated) type-family applications $F(\bar{\tau})$, to reflect their addition to source Haskell [CKP05, CKPM05]. Types are classified by kinds κ in the usual way; the kinding judgement $\Gamma \vdash \tau : \kappa$ on types is conventional and appears in the extended version of this paper. To avoid clutter we use only monomorphic kinds, but it is easy to add kind polymorphism along the lines of [YWC⁺12], and our implementation does so.

¹ Several versions of System FC are described in published work. Some of these variants have had decorations to the FC name, such as FC₂ or FC_C[†]. We do not make these distinctions in the present work, referring instead to all of these systems – in fact, one evolving system – as “FC”.

4.1 Roles and casts

FC's distinctive feature is a type-safe cast ($e \triangleright \gamma$) (Figure 3), which uses a *coercion* γ to cast a term from one type to another. A coercion γ is a witness or proof of the equality of two types. Coercions are classified by the judgement

$$\Gamma \vdash \gamma : \tau \sim_{\rho}^{\kappa} \sigma$$

given in Figure 4, and pronounced “in type environment Γ the coercion γ witnesses that the types τ and σ both have kind κ , and are equal at role ρ ”. The notion of being “equal at role ρ ” is the important feature of this paper; it is a development of earlier work, as Section 8 describes. There are precisely three roles (see Figure 3), written N, R, and P, with the following meaning:

Nominal equality, written \sim_N , is the equality that the type checker reasons about. When a Haskell programmer says that two Haskell types are the “same”, we mean that the types are nominally equal. Thus, we can say that $\text{Int} \sim_N \text{Int}$. Type families introduce new nominal equalities. So, if we have **type instance** $F \text{ Int} = \text{Bool}$, then $F \text{ Int} \sim_N \text{Bool}$.

Representational equality, written \sim_R , holds between two types that share the same run-time representation. Because all types that are nominally equal also share the same representation, nominal equality is a subset of representational equality. Continuing the example from the introduction, $\text{HTML} \sim_R \text{String}$.

Phantom equality, written \sim_P , holds between any two types, whatsoever. It may seem odd that we produce and consume proofs of this “equality”, but doing so keeps the system uniform and easier to reason about. The idea of phantom equality is new in this work, and it allows for zero-cost conversions among types with phantom parameters.

We can now give the typing judgement for type-safe cast:

$$\frac{\Gamma \vdash e : \tau_1 \quad \Gamma \vdash \gamma : \tau_1 \sim_R \tau_2}{\Gamma \vdash e \triangleright \gamma : \tau_2} \text{ TM_CAST}$$

The coercion γ must be a proof of *representational* equality, as witnessed by the R subscript to the result of the coercion typing premise. This makes good sense: we can treat an expression of one type τ_1 as an expression of some other type τ_2 if and only if those types share a representation.

4.2 Coercions

Coercions (Figure 3) and their typing rules (Figure 4) are the heart of System FC. The basic typing judgement for coercions is $\Gamma \vdash \gamma : \tau \sim_{\rho}^{\kappa} \sigma$. When this judgement holds, it is easy to prove that τ and σ must have the same kind κ . However, kinds are not very relevant to the focus of this work, and so we often omit the kind annotation in our presentation. It can always be recovered by using the (syntax-directed) kinding judgement on types.

We can understand the typing rules in Figure 4, by thinking about the equalities that they define.

4.2.1 Nominal implies representational

If we have a proof that two types are nominally equal, then they are certainly representationally equal. This intuition is expressed by the **sub** operator, and the rule CO_SUB.

$$\boxed{\Gamma \vdash \gamma : \phi}$$

$$\frac{\Gamma \vdash \tau : \kappa}{\Gamma \vdash \langle \tau \rangle : \tau \sim_N \tau} \text{ CO_REFL}$$

$$\frac{\Gamma \vdash \gamma : \sigma \sim_{\rho} \tau}{\Gamma \vdash \mathbf{sym} \gamma : \tau \sim_{\rho} \sigma} \text{ CO_SYM}$$

$$\frac{\Gamma \vdash \gamma_1 : \tau_1 \sim_{\rho} \tau_2 \quad \Gamma \vdash \gamma_2 : \tau_2 \sim_{\rho} \tau_3}{\Gamma \vdash \gamma_1 \circ \gamma_2 : \tau_1 \sim_{\rho} \tau_3} \text{ CO_TRANS}$$

$$\frac{\overline{\Gamma \vdash \gamma : \tau \sim_{\rho} \sigma} \quad \bar{\rho} \text{ is a prefix of roles}(H) \quad \Gamma \vdash H\bar{\tau} : \kappa}{\Gamma \vdash H(\bar{\gamma}) : H\bar{\tau} \sim_R H\bar{\sigma}} \text{ CO_TYCONAPP}$$

$$\frac{\overline{\Gamma \vdash \gamma : \tau \sim_N \sigma} \quad \Gamma \vdash F(\bar{\sigma}) : \kappa}{\Gamma \vdash F(\bar{\gamma}) : F(\bar{\tau}) \sim_N F(\bar{\sigma})} \text{ CO_TYFAM}$$

$$\frac{\Gamma \vdash \gamma_1 : \tau_1 \sim_{\rho} \sigma_1 \quad \Gamma \vdash \gamma_2 : \tau_2 \sim_N \sigma_2 \quad \Gamma \vdash \tau_1 \tau_2 : \kappa \quad \Gamma \vdash \sigma_1 \sigma_2 : \kappa}{\Gamma \vdash \gamma_1 \gamma_2 : \tau_1 \tau_2 \sim_{\rho} \sigma_1 \sigma_2} \text{ CO_APP}$$

$$\frac{\Gamma, \alpha : \kappa \vdash \gamma : \tau \sim_{\rho} \sigma}{\Gamma \vdash \forall \alpha : \kappa. \gamma : \forall \alpha : \kappa. \tau \sim_{\rho} \forall \alpha : \kappa. \sigma} \text{ CO_FORALL}$$

$$\frac{\Gamma \vdash \tau : \kappa \quad \Gamma \vdash \sigma : \kappa}{\Gamma \vdash \langle \tau, \sigma \rangle_P : \tau \sim_P \sigma} \text{ CO_PHANTOM}$$

$$\frac{c : \tau \sim_{\rho} \sigma \in \Gamma}{\Gamma \vdash c : \tau \sim_{\rho} \sigma} \text{ CO_VAR}$$

$$\frac{C : [\bar{\alpha} : \bar{\kappa}]. \sigma_1 \sim_{\rho} \sigma_2 \quad \overline{\Gamma \vdash \tau : \kappa}}{\Gamma \vdash C(\bar{\tau}) : \sigma_1[\bar{\tau}/\bar{\alpha}] \sim_{\rho} \sigma_2[\bar{\tau}/\bar{\alpha}]} \text{ CO_AXIOM}$$

$$\frac{\Gamma \vdash \gamma : H\bar{\tau} \sim_R H\bar{\sigma} \quad \bar{\rho} \text{ is a prefix of roles}(H) \quad H \text{ is not a newtype}}{\Gamma \vdash \mathbf{nth}^i \gamma : \tau_i \sim_{\rho_i} \sigma_i} \text{ CO_NTH}$$

$$\frac{\Gamma \vdash \gamma : \tau_1 \tau_2 \sim_N \sigma_1 \sigma_2 \quad \Gamma \vdash \tau_1 : \kappa \quad \Gamma \vdash \sigma_1 : \kappa}{\Gamma \vdash \mathbf{left} \gamma : \tau_1 \sim_N \sigma_1} \text{ CO_LEFT}$$

$$\frac{\Gamma \vdash \gamma : \tau_1 \tau_2 \sim_N \sigma_1 \sigma_2 \quad \Gamma \vdash \tau_2 : \kappa \quad \Gamma \vdash \sigma_2 : \kappa}{\Gamma \vdash \mathbf{right} \gamma : \tau_2 \sim_N \sigma_2} \text{ CO_RIGHT}$$

$$\frac{\Gamma \vdash \gamma : \forall \alpha : \kappa. \tau_1 \sim_{\rho} \forall \alpha : \kappa. \sigma_1 \quad \Gamma \vdash \tau : \kappa}{\Gamma \vdash \gamma @ \tau : \tau_1[\tau/\alpha] \sim_{\rho} \sigma_1[\tau/\alpha]} \text{ CO_INST}$$

$$\frac{\Gamma \vdash \gamma : \tau \sim_N \sigma}{\Gamma \vdash \mathbf{sub} \gamma : \tau \sim_R \sigma} \text{ CO_SUB}$$

Figure 4. Formation rules for coercions

4.2.2 Equality is an equivalence relation

Equality is an equivalence relation at all three roles. Symmetry (rule CO_SYM) and transitivity (CO_TRANS) work for any role ρ . Reflexivity is more interesting: CO_REFL is a proof of nominal equality only. From this we can easily get representational reflexivity using **sub**. But what does “phantom” reflexivity mean? It is a proof term that any two types τ and σ are equal at role P , and we need a new coercion form to express that, written as $\langle \tau, \sigma \rangle_P$ (rule CO_PHANTOM).

4.2.3 Axioms for equality

Each newtype declaration, and each type-family instance, gives rise to an FC *axiom*; newtypes give rise to representational axioms, and type-family instances give rise to nominal axioms.² For example, the declarations

```
newtype HTML = Mk String
type family F [a] = Maybe a
```

produce the axioms

$$C_1 : \text{HTML} \sim_R \text{String}$$

$$C_2 : [\alpha : *]. F([\alpha]) \sim_N \text{Maybe } \alpha$$

Axiom C_1 states that HTML is *representationally* equal to String (since they are distinct types, but share a common representation), while C_2 states that $F([\sigma])$ is *nominally* equal to $\text{Maybe } \sigma$ (meaning that the two are considered to be the same type by the type checker). In C_2 , the notation “[$\alpha : *$].” binds α in the types being equated. Uses of these axioms are governed by the rule CO_AXIOM. Axioms must always appear fully applied, and we assume that they live in a global context, separate from the local context Γ .

4.2.4 Equality can be abstracted

Just as one can abstract over types and values in System F, one can also abstract over equality proofs in FC. To this end, FC terms (Figure 3) include coercion abstraction $\lambda c : \phi . e$ and application $e \gamma$. These are the introduction and elimination forms for the coercion-abstraction arrow (\Rightarrow), just as ordinary value abstraction and application are the introduction and elimination forms for ordinary arrow (\rightarrow) (see the extended version of this paper).

A coercion abstraction binds a coercion variable $c : \phi$. These variables can occur only in coercions; see the entirely conventional rule CO_VAR. Coercion variables can also be bound in the patterns of a **case** expression, which supports the implementation of generalised algebraic data types (GADTs).

4.2.5 Equality is congruent

Several rules witness that, ignoring roles, equality is *congruent* – for example, if $\sigma \sim_\rho \tau$ then $\text{Maybe } \sigma \sim_\rho \text{Maybe } \tau$. However, the roles in these rules deserve some study, as they are the key to understanding the whole system.

Congruence of type application Before diving into the rules themselves, it is helpful to consider some examples of how we want congruence and roles to interact. Let’s consider the definitions in Figure 5. With these definitions in hand, what equalities should be derivable? (Recall the intuitive meanings of the different roles in Section 4.1.)

1. Should $\text{Maybe HTML} \sim_R \text{Maybe String}$ hold? Yes, it should. The type parameter to **Maybe** has a representational role, so it makes sense that two **Maybe**s built

²For simplicity, we are restricting ourselves to *open* type families. Closed type families [EVPW14] are readily accommodated.

```
newtype HTML = Mk String
```

```
type family F a
type instance F String = Int
type instance F HTML = Bool
```

```
data T a = MkT (F a)
```

Figure 5. Congruence and roles example code

out of representationally equal types should be representationally equal.

2. Should $\text{Maybe HTML} \sim_N \text{Maybe String}$ hold? Certainly not. These two types are entirely distinct to Haskell programmers and its type checker.
3. Should $T \text{ HTML} \sim_R T \text{ String}$ hold? Certainly not. We can see, by unfolding the definition for T , that the representations of the two types are different.
4. Should $\alpha \text{ HTML} \sim_R \alpha \text{ String}$ hold, for a type variable α ? It depends on the instantiation of α ! If α becomes **Maybe**, then “yes”; if α becomes **T**, then “no”. Since we may be abstracting over α , we do not know which of the two will happen, so we take the conservative stance and say that $\alpha \text{ HTML} \sim_R \alpha \text{ String}$ does *not* hold.

This last point is critical. The alternative is to express α ’s argument roles in its kind, but that leads to a much more complicated system; see related work in Section 8. A distinguishing feature of this paper is the substantial simplification we obtain by attributing roles only to the arguments to type constants (H , in the grammar), and not to abstracted type variables. We thereby lose a little expressiveness, but we have not found that to be a big problem in practice. See Section 8.1 for an example of an easily fixed problem case.

To support both (1) and (4) requires two coercion forms and corresponding typing rules:

- The coercion form $H(\bar{\gamma})$ has an explicit type constant at its head. This form always proves a representational equality, and it requires input coercions of the roles designated by the roles of H ’s parameters (rule CO_TYCONAPP). The *roles* function gives the list of roles assigned to H ’s parameters, as explained in Section 2.2. We allow $\bar{\rho}$ to be a prefix of *roles*(H) to accommodate partially-applied type constants.
- The coercion form $\gamma_1 \gamma_2$ does not have an explicit type constant, so we must use the conservative treatment of roles discussed above. Rule CO_APP therefore requires γ_2 to be a nominal coercion, though the role of γ_1 carries through to $\gamma_1 \gamma_2$.

What if we wish to prove a nominal equality such as $\text{Maybe (F String)} \sim_N \text{Maybe Int}$? We can’t use the $H(\bar{\gamma})$ form, which proves only representational equality, but we can use the $\gamma_1 \gamma_2$ form. The leftmost coercion would just be $\langle \text{Maybe} \rangle$.

Congruence of type family application Rule CO_TYFAM proves the equality of two type-family applications. It requires nominal coercions among all the arguments. Why? Because type families can inspect their (type) arguments and branch on them. We would not want to be able to prove any equality between $F \text{ String}$ and $F \text{ HTML}$.

Congruence of polymorphic types The rule CO_FORALL works for any role ρ ; polymorphism and roles do not interact.

4.2.6 Equality can be decomposed

If we have a proof of $\text{Maybe } \sigma \sim_{\rho} \text{Maybe } \tau$, should we be able to get a proof of $\sigma \sim_{\rho} \tau$, by decomposing the equality? Yes, in this case, but we must be careful here as well.

Rule `CO_NTH` is almost an inverse to `CO_TYCONAPP`. The difference is that `CO_NTH` prohibits decomposing equalities among newtypes. Why? Because `nth` witnesses injectivity and newtypes are not injective! For example, consider these definitions:

```
data Phant a = MkPhant
newtype App a b = MkApp (a b)
```

Here, $\text{roles}(\text{App}) = R, N$. (The roles are inferred during compilation; see Section 5.) Yet, we can see the following chain of equalities:

$\text{App Phant Int} \sim_R \text{Phant Int} \sim_R \text{Phant Bool} \sim_R \text{App Phant Bool}$
By transitivity, we can derive a coercion γ witnessing

$$\text{App Phant Int} \sim_R \text{App Phant Bool}$$

If we could use `nth`² on γ , we would get $\text{Int} \sim_N \text{Bool}$: disaster! We eliminate this possibility by preventing `nth` on newtypes.

The rules `CO_LEFT` and `CO_RIGHT` are almost inverses to `CO_APP`. The difference is that both `CO_LEFT` and `CO_RIGHT` require and produce only nominal coercions. We need a new newtype to see why this must be so:

```
newtype EitherInt a = MkEI (Either a Int)
```

This definition yields an axiom showing that, for all a , $\text{EitherInt } a \sim_R (\text{Either } a \text{ Int})$. Suppose we could apply `left` and `right` to coercions formed from this axiom. Using `left` would get us a proof of $\text{EitherInt} \sim_R (\text{Either } a)$, which could then be used to show, say, $(\text{Either Char}) \sim_R (\text{Either Bool})$ and then (using `nth`) $\text{Char} \sim_N \text{Bool}$. Using `right` would get us a proof of $a \sim_R \text{Int}$, for *any* a . These are both clearly disastrous. So, we forbid using these coercion formers on representational coercions.³

Thankfully, polymorphism and roles play well together, and the `CO_INST` rule (inverse to `CO_FORALL`) shows quite straightforwardly that, if two polytypes are equal, then so are the instantiated types.

There is no decomposition form for type family applications: knowing that $F(\bar{\tau})$ is equal to $F(\bar{\sigma})$ tells us nothing whatsoever about the relationship between $\bar{\tau}$ and $\bar{\sigma}$.

4.3 Role attribution for type constants

In System FC we assume an unwritten global environment of top-level constants: data types, type families, axioms, and so on. For a data type H , for example, this environment will give the kind of H , the types of H 's data constructors, and the roles of H 's parameters. Clearly this global environment must be internally consistent. For example, a data constructor K must return a value of type $D \bar{\tau}$ where D is a data type; K 's type must be well-kinded, and that kind must be consistent with D 's kind.

³We note in passing that the forms `left` and `right` are present merely to increase expressivity. They are not needed anywhere in the metatheory to prove type soundness. Though originally part of FC, they were omitted in previous versions [WVPZ11] and even in the implementation. Haskell users then found that some desirable program were no longer type-checking. Thus, these forms were re-introduced.

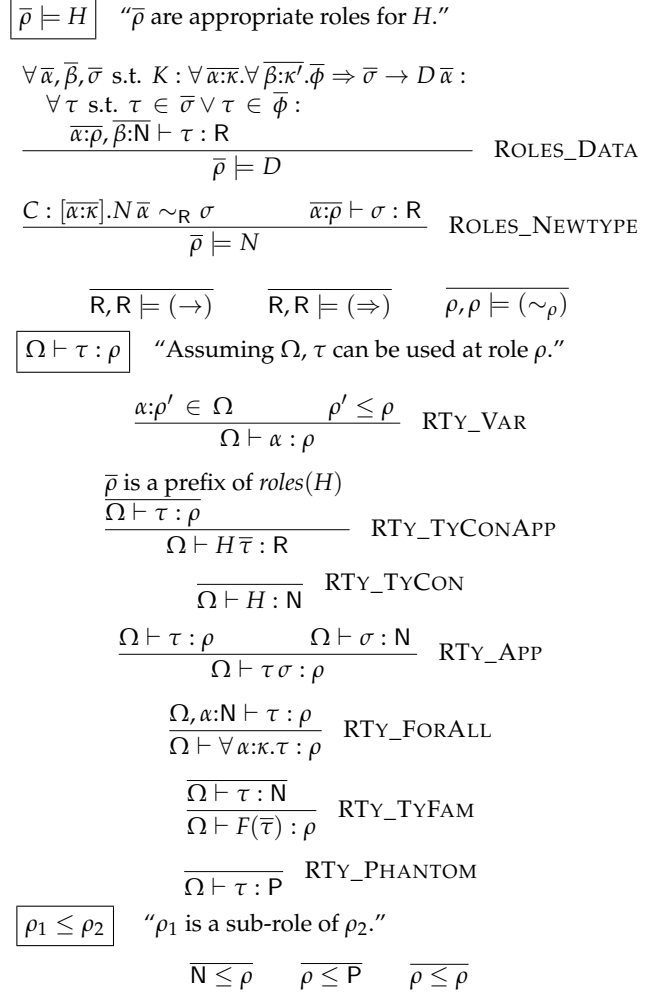


Figure 6. Rules asserting a correct assignment of roles to data types

All of this is standard except for roles. It is essential that the roles of D 's parameters, $\text{roles}(D)$, are consistent with D 's definition. For example, it would be utterly wrong for the global environment to claim that $\text{roles}(\text{Maybe}) = P$, because then we could prove that $\text{Maybe Int} \sim_R \text{Maybe Bool}$ using `CO_TYCONAPP`.

We use the judgement $\bar{\rho} \models H$, to mean " $\bar{\rho}$ are suitable roles for the parameters of H ", and in our proof of type safety, we assume that $\text{roles}(H) \models H$ for all H . The rules for this judgement and two auxiliary judgements appear in Figure 6. Note that this judgement defines a *relation* between roles and data types. Our role inference algorithm (Section 5) determines the most permissible roles for this relation, but often other, less permissive roles, such as those specified by role annotations, are also included by this relation.

Start with `ROLES_NEWTYPER`. Recall that a newtype declaration for N gives rise to an axiom $C : [\bar{\alpha}:\bar{\kappa}]. N \bar{\alpha} \sim_R \sigma$. The rule says that roles $\bar{\rho}$ are acceptable for N if each parameter α_i is used in σ in a way consistent with ρ_i , expressed using the auxiliary judgement $\bar{\alpha}:\bar{\rho} \vdash \sigma : R$.

The key auxiliary judgement $\Omega \vdash \tau : \rho$ checks that the type variables in τ are used in a way consistent with their roles specified in Ω , when considered at role ρ . More pre-

cisely, if $\alpha:\rho' \in \Omega$ and if $\sigma_1 \sim_{\rho'} \sigma_2$ then $\tau[\sigma_1/\alpha] \sim_{\rho} \tau[\sigma_2/\alpha]$. Unlike in many typing judgements, the role ρ (as well as Ω) is an *input* to this judgement, not an output. With this in mind, the rules for the auxiliary judgement are straightforward. For example, `RTY_TYFAM` says that the argument types of a type family application are used at nominal role. The variable rule, `RTY_VAR`, allows a variable to be assigned a more restrictive role (via the sub-role judgement) than required, which is needed both for multiple occurrences of the same variable, and to account for role signatures. Note that rules `RTY_TYCONAPP` and `RTY_APP` overlap – this judgement is not syntax-directed.

Returning to our original judgement $\bar{\rho} \models H$, `ROLES_DATA` deals with algebraic data types D , by checking roles in each of its data constructors K . The type of a constructor is parameterised by universal type variables $\bar{\alpha}$, existential type variables $\bar{\beta}$, coercions (with types $\bar{\phi}$), and term-level arguments (with types $\bar{\sigma}$). For each constructor, we must examine each proposition ϕ and each term-level argument type σ , checking to make sure that each is used at a representational role. Why check for a representational role specifically? Because *roles* is used in `CO_TYCONAPP`, which produces a representational coercion. In other words, we must make sure that each term-level argument appears at a representational role within the type of each constructor K for `CO_TYCONAPP` to be sound.

Finally (\rightarrow) and (\Rightarrow) have representational roles: functions care about representational equality but never branch on the nominal identity of a type. (For example, functions always treat `HTML` and `String` identically.) We also see that the roles of the arguments to an equality proposition match the role of the proposition. This fact comes from the congruence of the respective equality relations.

These definitions lead to a powerful theorem:

Theorem (Roles assignments are flexible). *If $\bar{\rho} \models H$, where H is a data type or newtype, and $\bar{\rho}'$ is such that $\rho'_i \leq \rho_i$ (for $\rho_i \in \bar{\rho}$ and $\rho'_i \in \bar{\rho}'$), then $\bar{\rho}' \models H$.*

Proof. Straightforward induction on $\Omega \vdash \tau : \rho$. \square

This theorem states that, given a sound role assignment for H , any more restrictive role assignment is also sound. This property of our system here is one of its distinguishing characteristics from our prior work on roles – see Section 10 for discussion.

4.4 Metatheory

The preceding discussion gave several non-obvious examples where admitting *too many* coercions would lead to unsoundness. However, we must have *enough* coercions to allow us to make progress when evaluating a program. (We do not have space to elaborate, but a key example is the use of `nth` in rule `S_KPUSH`, presented in the extended version of this paper.) Happily, we can be confident that we have enough coercions, but not too many, because we prove the usual progress and preservation theorems for System FC. The structure of the proofs follows broadly that in previous work, such as [WVPZ11] or [YWC⁺12].

A key step in the proof of progress is to prove *consistency*; that is, that no coercion can exist between, say, `Int` and `Bool`. This is done by defining a non-deterministic, role-directed rewrite relation on types and showing that the rewrite system is confluent and preserves type constants (other than newtypes) appearing in the heads of types. We then prove that, if a coercion exists between two types τ_1 and τ_2 , these two types both rewrite to a type σ . We conclude then that τ_1

and τ_2 , if headed by a non-newtype type constant, must be headed by the same such constant.

Alas, the rewrite relation is *not* confluent! The non-linear patterns allowed in type families (that is, with a repeated variable on the left-hand side), combined with non-termination, break the confluence property (previous work gives full details [EVPW14]). However, losing confluence does not necessarily threaten consistency – it just threatens the particular proof technique we use. However, a more powerful proof appears to be an open problem in the term rewriting community.⁴ For the purposes of our proof we dodge this difficulty by restricting type families to have only linear patterns, thus leading to confluence; consistency of the full system remains an open problem.

The full proof of type safety appears in the extended version of this paper; it exhibits no new proof techniques.

5. Roles on type constructors

In System FC we assume that, for every type constant H , the global environment specifies *roles(H)*, the roles of H 's parameters. However, there is some flexibility about this role assignment; the only requirement for type soundness is that *roles(H) $\models H$* .

In GHC, the roles of a type constructor are determined first by any role annotations provided by the programmer. If these are missing, the type checker calculates the default roles using the inference algorithm described below.

5.1 Role inference

A type constructor's roles are assigned depending on its nature:

- Primitive type constructors like (\rightarrow) and (\sim_{ρ}^x) have pre-defined roles (Figure 6).
- Type families (Section 2.2.3) and type classes (Section 3.2) have nominal roles for all parameters.
- For a **data** type or **newtype** T GHC *infers* the roles for T 's type parameters, possibly modified by role annotations (Section 3.1).

The role inference algorithm is quite straightforward. At a high level, it simply starts with the role information of the built-in constants (\rightarrow) , (\Rightarrow) , and (\sim_{ρ}) , and propagates the roles until it finds a fixpoint. In the description of the algorithm, we assume a mutable environment; *roles(H)* pulls a list of roles from this environment. Only after the algorithm is complete will *roles(H) $\models H$* hold.

1. Populate *roles(T)* (for all T) with user-supplied annotations; omitted role annotations default to phantom. (See Section 5.2 for discussion about this choice of default.)
2. For every data type D , every constructor for that data type K , and every coercion type and term-level argument type σ to that constructor: run `walk(D, σ)`.
3. For every newtype N with representation type σ , run `walk(N, σ)`.
4. If the role of any parameter to any type constant changed in the previous steps, go to step 2.

⁴Specifically, we believe that a positive answer to open problem #79 of the Rewriting Techniques and Applications (RTA) conference would lead to a proof of consistency; see <http://www.win.tue.nl/rtaloop/problems/79.html>.

- For every T , check $roles(T)$ against a user-supplied annotation, if any. If these disagree, reject the program. Otherwise, $roles(T) \models T$ holds.

The procedure $walk(T, \sigma)$ is defined as follows, matching from top to bottom:

```
walk(T, α)      := mark the α parameter to T as R.
walk(T, H  $\bar{\tau}$ ) := let  $\bar{\rho} = roles(H)$ ;
                  for every  $i, 0 < i \leq \text{length}(\bar{\tau})$ :
                    if  $\rho_i = N$ , then
                      mark all variables free in  $\tau_i$  as N;
                    else if  $\rho_i = R$ , then  $walk(T, \tau_i)$ .
walk(T,  $\tau_1 \tau_2$ ) := walk(T,  $\tau_1$ );
                  mark all variables free in  $\tau_2$  as N.
walk(T,  $F(\bar{\tau})$ ) := mark all variables free in the  $\bar{\tau}$  as N.
walk(T,  $\forall \beta: \kappa. \tau$ ) := walk(T,  $\tau$ ).
```

When marking, we must follow these two rules:

- If a variable to be marked does not appear as a type-level argument to the data type T in question, ignore it.
- Never allow a variable previously marked N to be marked R. If such a mark is requested, ignore it.

The first rule above deals with existential and local (\forall -bound) type variables, and the second one deals with the case where a variable is used both in a nominal and in a representational context. In this case, we wish the variable to be marked N, not R.

Theorem. *The role inference algorithm always terminates.*

Theorem (Role inference is sound). *After running the role inference algorithm, $roles(H) \models H$ will hold for all H .*

Theorem (Role inference is optimal). *After running the role inference algorithm, any loosening of roles (a change from ρ to ρ' , where $\rho \leq \rho'$ and $\rho \neq \rho'$) would violate $roles(H) \models H$.*

Proofs of these theorems appear in the extended version of this paper.

5.2 The role of role inference

According to the specification of sound role assignments in Figure 6, a type constructor H can potentially have several different sound role assignments. For example, assigning Maybe’s parameter to have a representational role is type-safe, but assigning a nominal role would be, too. Note that nominal roles are always sound for data types, according to the definition in Figure 6. However, as we saw in the description of the role inference algorithm, we choose default roles for data types to be as permissive as possible – in other words, the default role for a data type constructor parameter starts at phantom and only change when constrained by the algorithm. Here, we discuss this design decision and its consequences.

What if we had no role inference whatsoever and required programmers to annotate every data type? In this case, the burden on programmers seems drastic and migration to this system overwhelming, requiring all existing data type declarations to be annotated with roles.

Alternatively, we could specify that all unannotated roles default to nominal (thus removing the need for role inference). This choice would lead to greater abstraction safety by default – we would not have to worry that the implementor of Map is unaware of roles and forgets a critical role annotation.

However, we choose to use the most permissive roles by default for several reasons. First, for convenience: this choice increases the availability of coerce (as only those types with annotations would be Coercible otherwise), and it supports backward compatibility with the Generalized Newtype Deriving (GND) feature (see Section 7).

Furthermore, our choice of using phantom as the default also means that the majority of programmers do not need to learn about roles. They will not need role annotations in their code. Users of coerce will need to consider roles, as will library implementors who use class-based invariants (see Section 3.1). Other users are unaffected by roles and will not be burdened by them.

Our choices in the design of the role system, and the default of phantom in particular, has generated vigorous debate.⁵ This discussion is healthy for the Haskell community. The difficulty with abstraction is not new: with GND, it has always been possible to lift coercions through data types, potentially violating their class-based invariants. The features described in this paper make this subversion both more convenient (through the use of coerce) and, more importantly, now preventable (through the use of role annotations).

6. Implementing Coercible

We have described the source-language view of Coercible (Sections 2, 3), and System FC, the intermediate language into which the source language is elaborated (Section 4). In this section we link the two by describing how the source-language use of Coercible is translated into Core.

6.1 Coercible and coerce

When the compiler transforms Haskell to Core, type classes become ordinary types and type class constraints turn into ordinary value arguments [WB89]. In particular, type classes typically become simple product types with one field per method.

The same holds for the type class Coercible $a \ b$, which has one method, namely the witness of representational equality $a \sim_R b$. As that type cannot be expressed in Haskell, the actual definition of Coercible is built in:

```
data Coercible a b = MkCoercible (a  $\sim_R$  b)
```

The definition of coerce, which is also only possible in Core, pattern-matches on MkCoercible to get hold of the equality witness, and then uses Core’s primitive cast operation:

```
coerce :: forall α β. Coercible α β → α → β
coerce =  $\Lambda$  α β.  $\lambda$  (c :: Coercible α β) (x :: α). case c of
  MkCoercible eq → x  $\triangleright$  eq
```

Since type applications are explicit in Core, coerce now takes four arguments: the types to cast from and to, the coercion witness, and finally the value to cast.

The data type Coercible also serves to box the primitive, unboxed type \sim_R , just as Int serves to box the primitive, unboxed type Int#:

```
data Int = I# Int#
```

All boxed types are represented uniformly by a heap pointer. In GHC all constraints (such as Eq a or Coercible a b) are boxed, so that they can be treated uniformly, and even polymorphically [YWC⁺12]. In contrast, an unboxed type is rep-

⁵To read some of this debate, see the thread beginning with this post: <http://www.haskell.org/pipermail/libraries/2014-March/022321.html>

resented by a non-pointer bit field, such as a 32 or 64-bit int in the case of `Int#` [PL91].

A witness of (unboxed) type \sim_R carries no information: we never actually inspect an equality proof at run-time. So the type \sim_R can be represented by a *zero-width* bit-field – that is, by nothing at all. This implementation trick, of boxing a zero-bit witness, is exactly analogous to the wrapping of boxed nominal equalities used to implement deferred type errors [VPMa12].

Since `Coercible` is a regular data type, you might worry about bogus programs like this, which uses recursion to construct an unsound witness `co` whose value is bottom:

```
looksUnsound :: forall α β. α → β
looksUnsound = \α β x →
  let co :: Coercible α β = co in
  coerce α β co x
```

However, since `coerce` evaluates the `Coercible` argument (see the definition of `coerce` above), `looksUnsound` will simply diverge. Again, this follows the behaviour of deferred type errors [VPMa12].

In uses of `coerce`, the `Coercible` argument will be constructed from the instances which, as described below (Section 6.4), are guaranteed to be acyclic. The usual simplification machinery of GHC then ensures that these are inlined, causing the `case` to cancel with the `MkCoercible` constructor, leaving only the cast `x > eq`, which is operationally free.

6.2 On-demand instance generation

The language of Section 2 suggests that we generate Haskell instance declarations for `Coercible`, based on type declarations. Although this is a useful way to explain the design to a programmer (who is already familiar with type classes and instance declarations), GHC’s implementation is much simpler and more direct.

Rather than generate and compile instance declarations, the constraint solver treats `Coercible` constraints specially: to solve a `Coercible` constraint, the solver uses the rules of Section 2 directly to decompose the constraint into simpler sub-goals. This approach makes it easy to implement the non-standard visibility rules of `Coercible` instances (see Section 3.1), by simply not applying the `newtype-unwrapping` rule if the constructor is not in scope.

6.3 The higher rank instance

Consider this declaration, whose constructor uses a higher-rank type:

```
newtype Sel = MkSel (forall a. [a] → a)
```

We would expect its `newtype-unwrapping` instance to take the form

```
instance Coercible (forall a. [a] → a) b => Coercible Sel b
instance Coercible a (forall a. [a] → a) => Coercible a Sel
```

These declarations are illegal in source Haskell, even with all GHC extensions enabled. Nevertheless, we can generate internally and work with them in the solver just fine. This leads to constraints of the form

```
Coercible (forall a. s) (forall b. t)
```

which need special support in the solver. It already supports solving (nominal) type equalities of the form $(\text{forall } a. s) \sim (\text{forall } b. t)$, by generating a fresh type variable `c` and solving $s[c/a] \sim t[c/b]$. We generalised this functionality to handle representational type equalities as well.

6.4 Preventing circular reasoning and diverging instances

For most type classes, like `Show`, it is perfectly fine (and useful) to use a not-yet solved type class constraint to solve another, even though this can lead to cycles [LP05]. Consider the following code and execution:

```
newtype Fix a = MkFix (a (Fix a))
deriving instance Show (a (Fix a)) => Show (Fix a)
```

```
λ> show (MkFix (Just (MkFix (Just (MkFix Nothing))))))
"MkFix (Just (MkFix (Just (MkFix Nothing))))"
```

There are two `Show` instances at work: one for `Show (Maybe a)`, which uses the instance of `Show a`; and one for `Show (Fix a)`, which uses the the instance `Show (a (Fix a))`. Plugging them together to solve `Show (Fix Maybe)`, we see that this instance calls, by way of `Show (Maybe (Fix Maybe))`, itself. Nevertheless, the result is perfectly well-behaved and indeed terminates.

But with `Coercible`, such circular reasoning would be problematic; we could then seemingly write the bogus function `looksUnsoundH`:

```
newtype Id a = MkId a
c1 :: a → Fix Id
c1 = coerce
c2 :: Fix Id → b
c2 = coerce
looksUnsoundH :: a → b
looksUnsoundH = c2 ∘ c1
```

With the usual constraint solving, this code would type check: to solve the constraint `Coercible a (Fix Id)`, we need to solve `Coercible a (Id (Fix Id))`, which requires `Coercible a (Fix Id)`. This is a constraint we already looked at, so the constraint solver would normally consider all required constraints solved and accept the program.

Fortunately, there is no soundness problem here. Circular constraint-solving leads to a recursive definition of the `Coercible` constraints, exactly like the (Core) `looksUnsound` in Section 6.1, and `looksUnsoundH` will diverge just like `looksUnsound`. Nevertheless, unlike normal type classes, a recursive definition of `Coercible` is *never* useful, so it is more helpful to reject it statically. GHC therefore uses the existing depth-counter of the solver to spot and reject recursion of `Coercible` constraints.

6.5 Coercible and rewrite rules

What if a client of module `Html` writes this?

```
....(map unMk hs)...
```

She cannot use `coerce` because `HTML` is an abstract type, so the type system would (rightly) reject an attempt to use `coerce` (Section 3.1). However, since `HTML` is a `newtype`, one might hope that GHC’s optimiser would transform `(map unMk)` to `coerce`. The optimiser must respect type soundness, but (by design) it does not respect abstraction boundaries: dissolving abstractions is one key to high performance.

The correctness of transforming `(map unMk)` to `coerce` depends on a theorem about `map`, which a compiler can hardly be expected to identify and prove all by itself. Fortunately GHC already comes with a mechanism that allows a library author to specify *rewrite rules* for their code [PTH01]. The author takes the proof obligation that the rewrite is semantics-preserving, while GHC simply applies the rewrite whenever possible. In this case the programmer could write

```
{-# RULES "map/co" map coerce = coerce #-}
```

In our example, the programmer wrote `(map unMk)`. The definition `unMk` in module `Html` does not mention `coerce`, but both produce the same System FC code (a cast). So via cross-module inlining (more dissolution of abstraction boundaries) `unMk` will be inlined, transforming the call to the equivalent of `(map coerce)`, and that in turn fires the rewrite rule. Indeed even a nested call like `map (map unMk)` will also be turned into a single call of `coerce` by this same process applied twice.

The bottom line is this: the author of a map-like function `someMap` can accompany `someMap` with a RULE, and thereby optimise calls of `someMap` that do nothing into a simple call to `coerce`.

Could we dispense with a user-visible `coerce` function altogether, instead using map-like functions and RULEs as above? No: doing so would replace the zero-cost guarantee with best-effort optimisation; it would burden the author of every map-like function with the obligation to write a suitable RULE; it would be much less convenient to use in deeply-nested cases; and there might simply *be* no suitable map-like function available.

7. Generalized Newtype Deriving done right

As mentioned before, `newtype` is a great tool to make programs more likely to be correct, by having the type checker enforce certain invariants or abstractions. But newtypes can also lead to tedious boilerplate. Assume the programmer needs an instance of the type class `Monoid` for her type `HTML`. The underlying type `String` already comes with a suitable instance for `Monoid`. Nevertheless, she has to write quite a bit of code to convert that instance into one for `HTML`:

```
instance Monoid HTML where
  mempty = Mk mempty
  mappend (Mk a) (Mk b) = Mk (mappend a b)
  mconcat xs = Mk (mconcat (map unMk xs))
```

Note that this definition is not only verbose, but also non-trivial, as invocations of `Mk` and `unMk` have to be put in the right places, possibly via some higher order functions like `map` – all just to say “just use the underlying instance”!

This task is greatly simplified with `Coercible`: Instead of wrapping and unwrapping arguments and results, she can directly coerce the method of the base type’s instance itself:

```
instance Monoid HTML where
  mempty = coerce (mempty :: String)
  mappend = coerce (mappend :: String → String → String)
  mconcat = coerce (mconcat :: [String] → String)
```

The code is pure boilerplate: apply `coerce` to the method, instantiated at the base type by a type signature. And because it is boilerplate, the compiler can do it for her; all she has to do is to declare which instances of the base type should be lifted to the new type by listing them in the `deriving` clause:

```
newtype HTML = Mk String deriving Monoid
```

This is not a new feature: GHC has provided this *Generalized Newtype Deriving* (GND) for many years. But, the implementation was “magic” – GND would produce code that a user could not write herself. Now, the feature can be explained easily and fully via `coerce`.

Furthermore, GND was previously unsound [WVPZ11]. When combined with other extensions of GHC, such as type families [CKP05, CKPM05] or GADTs [CH03], GND could be exploited to completely break the type system: Figure 7

```
newtype Id1 a = MkId1 a
newtype Id2 a = MkId2 (Id1 a) deriving (UnsafeCast b)
```

```
type family Discern a b
type instance Discern (Id1 a) b = a
type instance Discern (Id2 a) b = b
```

```
class UnsafeCast to from where
  unsafe :: from → Discern from to
```

```
instance UnsafeCast b (Id1 a) where
  unsafe (MkId1 x) = x
```

```
unsafeCoerce :: a → b
unsafeCoerce x = unsafe (MkId2 (MkId1 x))
```

Figure 7. The above implementation of `unsafeCoerce` compiles (with appropriate flags) in GHC 7.6.3 but does not in GHC 7.8.1.

shows how this notorious bug can allow any type to be coerced to any other. The clause “`deriving (UnsafeCast b)`” is the bogus use of GND, and now will generate the instance

```
instance UnsafeCast b c ⇒ UnsafeCast b (Id2 c) where
  unsafe = coerce (unsafe :: c → Discern c b)
```

which will rightly be rejected because `Discern`’s first parameter has a nominal role. Indeed, preventing abuse of GND was the entire subject of the previous work [WVPZ11] the current paper is based on.

Similarly, it was possible to use GND to break invariants of abstract data types. The addition of `coerce` makes it yet easier to break such abstractions. As discussed in Section 3.1, these abuses can now be prevented via role annotations.

8. Related work

Prior work discusses the relationship between roles in FC and languages with generativity and abstraction, type-indexed constructs, and universes in dependent type theory. We do not repeat that discussion here. Instead we use this section to clarify the relationship between this paper and [WVPZ11], as well as make connections to other systems.

8.1 Prior version of roles

The idea of *roles* was initially developed in [WVPZ11] as a solution to the Generalized Newtype Deriving problem. That work introduces the equality relations \sim_R and \sim_N (called “type equality” and “code equality” resp. in [WVPZ11]). However, the system presented in [WVPZ11] was quite invasive: it required annotating every sub-tree of every kind with a role. Kinds in GHC are already quite complicated because of kind polymorphism, and a new form of role-annotated kinds would be more complex still.

In this paper, we present a substantially simplified version of the roles system of [WVPZ11], requiring role information only on the parameters to data types. Our new design keeps roles and kinds modularly separate, so that roles can be handled almost entirely separately (both intellectually and in the implementation) from kinds. The key simplification is to “assume the worst” about higher-kinded parameters, by assuming that their arguments are all nominal. In exchange we give up some expressiveness; specifically, we give up the ability to abstract over type constructors with non-nominal argument roles (see Section 10).

Furthermore, the observation that it is sound to “assume the worst” and use parameterised types with less permissive roles opens the door to role annotations. In this work, programmers are allowed to deliberately specify less permissive roles, giving them the ability to preserve type abstractions.

Surprisingly, this flexibility means that our version of roles actually *increases* expressiveness compared to [WVPZ11] in some places. In [WVPZ11] a role is part of a type’s kind, so a type expecting a higher-kinded argument (such as `Monad`) would also have to specify the roles expected by its argument. Therefore if `Monad` is applicable to `Maybe`, it would not also be applicable to a type `T` whose parameter has a nominal role. In the current work, however, there is no problem because `Maybe` and `T` have the same kind.

Besides the simplification discussed above, this paper makes two other changes to the specification of roles presented in [WVPZ11].

- The treatment of the phantom role is entirely novel; the rule `CO_PHANTOM` has no analogue in prior work.
- The coercion formation rules (Figure 4) are refactored so that the role on the coercion is an *output* of the (syntax-directed) judgement instead of an input. This is motivated by the implementation (which does not know the role at which coercions should be checked) and requires the addition of the `CO_SUB` rule.

There are, of course, other minor differences between this system and [WVPZ11] in keeping with the evolution of System FC. The main significant change, unrelated to roles, is the re-introduction of **left** and **right** coercions; see Section 4.2.6.

One important non-difference relates to the linear-pattern requirement. Section 4.4 describes that our language is restricted to have only *linear* patterns in its type families. (GHC, on the other hand, allows non-linear patterns as well.) This restriction exists in the language in [WVPZ11] as well. Section 4.2.2 of [WVPZ11] defines so-called Good contexts as having certain properties. Condition 1 in this definition subtly implies that all type families have linear patterns – if a type family had a non-linear pattern, it would be impossible, in general, to establish this condition. The fact that the definition of Good implies linear patterns came as a surprise, further explored in [EVPW14]. The language described in the present paper clarifies this restriction, but it is not a new restriction.

Finally, because this system has been implemented in GHC, this paper discusses more details related to compilation from source Haskell. In particular, the role inference algorithm of Section 5 is a new contribution of this work.

8.2 OCaml and variance annotations

The interactions between sub-typing, type abstraction, and various type system extensions such as GADTs and parameter constraints also appear in the OCaml language. In that context, *variance annotations* act like roles; they ensure that subtype coercions between compatible types are safe. For example, the type α `list` of immutable lists is covariant in the parameter α : if $\sigma \leq \tau$ then σ `list` $\leq \tau$ `list`. Variances form a lattice, with *invariant*, the most restrictive, at the bottom; *covariant* and *contravariant* incomparable; and *bivariant* at the top, allowing sub-typing in both directions. It is tempting to identify invariant with nominal and bivariant with phantom, but the exact connection is unclear. Scherer and Rémy [SR13] show that GADT parameters are not always invariant.

Exploration of the interactions between type abstraction, GADTs, and other features have recently revealed a soundness issue in OCaml⁶ that has been confirmed to date back several years. Garrigue discusses these issues [Gar13]. His proposed solution is to “assume that nothing is known about abstract types when they are used in parameter constraints and GADT return types” – akin to assigning nominal roles. However, this solution is too conservative, and in practice the OCaml 4.01 compiler relies on no fewer than *six* flags to describe the variance of type parameters. However, lacking anything equivalent to Core and its tractable metatheory, the OCaml developers cannot demonstrate the soundness of their solution in the way that we have done here.

What is clear, however, is that generative type abstraction interacts in interesting and non-trivial ways with type equality and sub-typing. Roles and type-safe coercion solve an immediate practical problem in Haskell, but we believe that the ideas have broader applicability in advanced type systems.

9. Roles in Practice

We have described a mechanism to allow safe coercions among distinct types, and we have reimplemented GHC’s previously unsafe `GeneralizedNewtypeDeriving` extension in terms of these safe coercions. Naturally, this change causes some code that was previously accepted to be rejected. Given that Haskell has a large user base and a good deal of production code, how does this change affect the community?

Advance testing During the development of this feature, we tested it against several popular Haskell packages available through Hackage, an online Haskell open-source distribution site. These tests were all encouraging and did not find any instances of hard-to-repair code in the wild.

Compiling all of Hackage As of 30 September 2013, 3,234 packages on Hackage compiled with GHC 7.6.3, the last released version without roles. The development version of GHC at that time included roles. A total of only four packages failed to compile directly due to GND failure.⁷ Of these, three of the failures were legitimate – the use of GND was indeed unsafe. For example, one case involved coercing a type variable passed into a type family; the author implicitly assumed that a newtype and its representation type were always considered equivalent with respect to the type family. Only one package failed to compile because of the gap in expressiveness between the roles in [WVPZ11] and those here. No other Hackage package depends on this one, indicating it is not a key part of the Haskell open-source fabric. See Section 10 for discussion of the failure.

These data were gathered almost two months after the implementation of roles was pushed into the development version of GHC, so active maintainers may have made changes to their packages before the study took place. Indeed, we are aware of a few packages that needed manual updates. In these cases, instances previously derived using GND had to be written by hand, but quite straightforwardly.

⁶ <http://caml.inria.fr/mantis/view.php?id=5985>

⁷ These data come from Bryan O’Sullivan’s work, described here: <http://www.haskell.org/pipermail/ghc-devs/2013-September/002693.html> That posting includes 3 additional GND failures; these were due to an implementation bug, since fixed.

10. Future directions

As of the date of writing (May 2014), roles seem not to have caused an undue burden to the community. The first release candidate for GHC 7.8 was released on 3 February 2014, followed by the full release on 9 April, and package authors have been updating their work to be compatible for some time. The authors of this paper are unaware of any major problems that Haskellers have had in updating existing code, despite hundreds of packages being available for GHC 7.8.⁸

However, we are aware that some users wish to use roles in higher-order scenarios that are currently impossible. We focus on one such scenario, as it is representative of all examples we have seen, including the package that did not compile when testing all of Hackage (Section 9).

Imagine adding the join method to the Monad class, as follows:

```
class Monad m where
  ...
  join :: forall a. m (m a) -> m a
```

With this definition, GND would still work in many cases. For example, if we define

```
newtype M a = Mk (Maybe a)
  deriving Monad
```

GND will work without a problem. We would need to show Coercible (Maybe (Maybe a) -> Maybe a) (M (M a) -> M a), which is straightforward.

More complicated constructions run into trouble, though. Take this definition, written to restrict a monad’s interface:

```
newtype Restr m a = Mk (m a)
  deriving Monad
```

To perform GND in this scenario, we must prove Coercible (m (m a) -> m a) (Restr m (Restr m a) -> Restr m a). In solving for this constraint, we eventually simplify to Coercible (m (m a)) (m (Restr m a)). At this point, we are stuck, because we do not have any information about the role of m’s parameter, so we must assume it is nominal. The GND feature is thus not available here. Similar problems arise when trying to use GND on monad transformers, a relatively common idiom.

How would this scenario play out under the system proposed in [WVPZ11]? This particular problem wouldn’t exist – m’s kind could have the right roles – but a different problem would. A type’s kind also stores its roles in [WVPZ11]. This means that Monad instances could be defined only for types that expect a representational parameter. Yet, it is sometimes convenient to define a Monad instance for a data type whose parameter is properly assigned a nominal role. The fact that the system described in this paper can accept Monad instances both for types with representational parameters and nominal parameters is a direct consequence of the *Role assignments are flexible* theorem (Section 4.3), which does not hold of the system in [WVPZ11].

Looking forward, there is a proposal to indeed add join to Monad, and so we want to be able to allow the use of GND on this enhanced Monad class. We have started to formulate solutions to this problem and have hope that we can overcome this barrier without modifications to the core language.

⁸Package authors have the option of specifying which compilers their package is known to work with. Of the 555 packages listed as working with one of the GHC 7.6 versions, 183 also are listed as compatible with GHC 7.8. These packages include 43 that use the GND extension.

11. Conclusion

Our focus has been on Haskell, for the sake of concreteness, but we believe that this work is important beyond the Haskell community. Any language that offers *both* generative type abstraction *and* type-level computation must deal with their interaction, and those interactions are extremely subtle. We have described one sound and tractable way to combine the two, including the source language changes, type inference, core calculus, and metatheory. In doing so we have given a concrete foundation for others to build upon.

Acknowledgments

Thanks to Antal Spector-Zabusky for contributing to this version of FC; and to Edward Kmett and Dimitrios Vytiniotis for discussion and feedback. This material is based upon work supported by the National Science Foundation under grant nos. CCF-1116620 and CCF-1319880. The first author was supported by the Deutsche Telekom Stiftung.

References

- [BEPW14] Joachim Breitner, Richard A. Eisenberg, Simon Peyton Jones, and Stephanie Weirich, *Safe zero-cost coercions for Haskell (extended version)*, Tech. Report MS-CIS-14-07, University of Pennsylvania, 2014.
- [CH03] James Cheney and Ralf Hinze, *First-class phantom types*, Tech. report, Cornell University, 2003.
- [CKP05] Manuel M. T. Chakravarty, Gabriele Keller, and Simon Peyton Jones, *Associated type synonyms*, ICFP, ACM, 2005, pp. 241–253.
- [CKPM05] Manuel M. T. Chakravarty, Gabriele Keller, Simon Peyton Jones, and Simon Marlow, *Associated types with class*, POPL, ACM, 2005, pp. 1–13.
- [EVPW14] Richard A. Eisenberg, Dimitrios Vytiniotis, Simon Peyton Jones, and Stephanie Weirich, *Closed type families with overlapping equations*, POPL, ACM, 2014, pp. 671–683.
- [Gar13] Jacques Garrigue, *On variance, injectivity, and abstraction*, OCaml Meeting, Boston., September 2013.
- [LP05] Ralf Lämmel and Simon Peyton Jones, *Scrap your boilerplate with class: Extensible generic functions*, ICFP, 2005.
- [Mar10] Simon Marlow (editor), *Haskell 2010 language report*, 2010.
- [MTHM97] Robin Milner, Mads Tofte, Robert Harper, and David MacQueen, *The definition of Standard ML (revised)*, 1997.
- [PL91] Simon Peyton Jones and J Launchbury, *Unboxed values as first class citizens*, FPCA, LNCS, vol. 523, 1991, pp. 636–666.
- [PTH01] Simon Peyton Jones, Andrew Tolmach, and Tony Hoare, *Playing by the rules: rewriting as a practical optimisation technique in GHC*, Haskell Workshop, 2001, pp. 203–233.
- [SR13] Gabriel Scherer and Didier Rémy, *GADTs meet subtyping*, ESOP, 2013, pp. 554–573.
- [VPMa12] Dimitrios Vytiniotis, Simon Peyton Jones, and José Pedro Magalhães, *Equality proofs and deferred type errors: A compiler pearl*, ICFP, ACM, 2012, pp. 341–352.
- [WB89] Philip Wadler and Stephen Blott, *How to make ad-hoc polymorphism less ad-hoc*, POPL, ACM, 1989, pp. 60–76.
- [WVPZ11] Stephanie Weirich, Dimitrios Vytiniotis, Simon Peyton Jones, and Steve Zdancewic, *Generative type abstraction and type-level computation*, POPL, ACM, 2011, pp. 227–240.
- [YWC⁺12] Brent A. Yorgey, Stephanie Weirich, Julien Cretin, Simon Peyton Jones, Dimitrios Vytiniotis, and José Pedro Magalhães, *Giving Haskell a promotion*, TLDI, ACM, 2012, pp. 53–66.