



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Quantified Class Constraints

Citation for published version:

Bottu, G-J, Karachalias, G, Schrijvers, T, Oliveira, BCDS & Wadler, P 2017, Quantified Class Constraints. in *ICFP Workshops 2017 - 10th ACM SIGPLAN International Haskell Symposium (Haskell 2017)*. ACM SIGPLAN Notices, no. 10, vol. 52, ACM, pp. 148-161, 10th ACM SIGPLAN International Haskell Symposium, Oxford, United Kingdom, 3/09/17. <https://doi.org/10.1145/3122955.3122967>

Digital Object Identifier (DOI):

[10.1145/3122955.3122967](https://doi.org/10.1145/3122955.3122967)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

ICFP Workshops 2017 - 10th ACM SIGPLAN International Haskell Symposium (Haskell 2017).

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Quantified Class Constraints

Gert-Jan Bottu

KU Leuven
Belgium

gertjan.bottu@student.kuleuven.be

Georgios Karachalias

KU Leuven
Belgium

georgios.karachalias@cs.kuleuven.be

Tom Schrijvers

KU Leuven
Belgium

tom.schrijvers@cs.kuleuven.be

Bruno C. d. S. Oliveira

University of Hong Kong
China

bruno@cs.hku.hk

Philip Wadler

University of Edinburgh
UK

wadler@inf.ed.ac.uk

Abstract

Quantified class constraints have been proposed many years ago to raise the expressive power of type classes from Horn clauses to the universal fragment of Hereditary Harrop logic. Yet, while it has been much asked for over the years, the feature was never implemented or studied in depth. Instead, several workarounds have been proposed, all of which are ultimately stopgap measures.

This paper revisits the idea of quantified class constraints and elaborates it into a practical language design. We show the merit of quantified class constraints in terms of more *expressive modeling* and in terms of *terminating type class resolution*. In addition, we provide a declarative specification of the type system as well as a type inference algorithm that elaborates into System F. Moreover, we discuss termination conditions of our system and also provide a prototype implementation.

CCS Concepts • Theory of computation → Type structures; • Software and its engineering → Functional languages;

Keywords Haskell, type classes, type inference

ACM Reference Format:

Gert-Jan Bottu, Georgios Karachalias, Tom Schrijvers, Bruno C. d. S. Oliveira, and Philip Wadler. 2017. Quantified Class Constraints. In *Proceedings of 10th ACM SIGPLAN International Haskell Symposium, Oxford, UK, September 7-8, 2017 (Haskell'17)*, 14 pages.

<https://doi.org/10.1145/3122955.3122967>

1 Introduction

Since Wadler and Blott [38] originally proposed type classes as a means to make adhoc polymorphism less adhoc, the feature has become one of Haskell's cornerstone features. Over the years type classes have been the subject of many language extensions that increase their expressive power and enable new applications. Examples of such extensions include: multi-parameter type classes [19]; functional dependencies [18]; or associated types [3].

Several of these implemented extensions were inspired by the analogy between type classes and predicates in Horn clauses. Yet, Horn clauses have their limitations. As a small side-product of

their work on derivable type classes, Hinze and Peyton Jones [12] have proposed to raise the expressive power of type classes to essentially the universal fragment of Hereditary Harrop logic [9] with what they call *quantified class constraints*. Their motivation was to deal with higher-kinded types which seemed to require instance declarations that were impossible to express in the type-class system of Haskell at that time.

Unfortunately, Hinze and Peyton Jones never did elaborate on quantified class constraints. Later, Lämmel and Peyton Jones [21] found a workaround for the particular problem of the derivable type classes work that did not involve quantified class constraints. Nevertheless the idea of quantified class constraints has whet the appetite of many researchers and developers. GHC ticket #2893¹, requesting for quantified class constraints, was opened in 2008 and is still open today. Commenting on this ticket in 2009, Peyton Jones states that “*their lack is clearly a wart, and one that may become more pressing*”, yet clarifies in 2014 that “*(t)he trouble is that I don't know how to do type inference in the presence of polymorphic constraints*.” In 2010, 10 years after the original idea, Hinze [10] rues that the feature has not been implemented yet. As recently as 2016, Chauhan et al. [4] regret that “*Haskell does not allow the use of universally quantified constraints*” and now in 2017 Spivey [34] has to use pseudo-Haskell when modeling with quantified class constraints. While various workarounds have been proposed and are used in practice [20, 31, 36], none has stopped the clamor for proper quantified class constraints.

This paper finally elaborates the original idea of quantified class constraints into a fully fledged language design.

Specifically, the contributions of this paper are:

- We provide an overview of the two main advantages of quantified class constraints (Section 2):
 1. they provide a natural way to express more of a type class's specification, and
 2. they enable terminating type class resolution for a larger class of applications.
- We elaborate the type system sketch of Hinze and Peyton Jones [12] for quantified type class constraints into a full-fledged formalization (Section 3). Our formalization borrows the idea of focusing from COCHIS [32], a calculus for Scala-style implicits [26, 27], and adapts it to the Haskell setting. We account for two notable differences: a global set of non-overlapping instances and support for superclasses.
- We present a type inference algorithm that conservatively extends that of Haskell 98 (Section 4) and comes with a dictionary-passing elaboration into System F (Section 5).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Haskell'17, September 7-8, 2017, Oxford, UK

© 2017 Association for Computing Machinery.

ACM ISBN 978-1-4503-5182-9/17/09...\$15.00

<https://doi.org/10.1145/3122955.3122967>

¹<https://ghc.haskell.org/trac/ghc/ticket/2893>

- We discuss the termination conditions on a system with quantified class constraints (Section 6).
- We provide a prototype implementation, which incorporates higher-kinded datatypes and accepts all² examples in this paper, at <https://github.com/gkaracha/quantcs-impl>.

2 Motivation

This section illustrates the expressive power afforded by quantified class constraints to capture several requirements of type class instances more succinctly, and to provide terminating resolution for a larger group of applications.

2.1 Precise and Succinct Specifications

Monad Transformers Consider the MTL type class for monad transformers [15]:

```
class Trans t where
  lift :: Monad m => m a -> (t m) a
```

What is not formally expressed in the above type class declaration, but implicitly expected, is that for any type T that instantiates *Trans* there should also be a *Monad* instance of the form:

```
instance Monad m => Monad (T m) where ...
```

Because the type checker is not told about this requirement, it will not accept the following definition of monad transformer composition.

```
newtype (t1 * t2) m a = C { runC :: t1 (t2 m) a }
instance (Trans t1, Trans t2) => Trans (t1 * t2) where
  lift = C . lift . lift
```

The idea of this code is to *lift* from monad m to $(t_2 m)$ and then to *lift* from $(t_2 m)$ to $t_1 (t_2 m)$. However, the second *lift* is only valid if $(t_2 m)$ is a monad and the type checker has no way of establishing that this fact holds for all monad transformers t_2 . Workarounds for this problem do exist in current Haskell [13, 31, 36], but they clutter the code with heavy encodings.

Quantified class constraints allow us to state this requirement explicitly as part of the *Trans* class declaration:

```
class (forall m. Monad m => Monad (t m)) => Trans t where
  lift :: Monad m => m a -> (t m) a
```

The instance for transformer composition $t_1 * t_2$ now typechecks.

Second-Order Functors Another example can be found in the work of Hinze [11]. He represents parameterized datatypes, like polymorphic lists and trees, as the fixpoint *Mu* of a *second-order functor*:

```
data Mu h a = In { out :: h (Mu h) a }
data List2 f a = Nil | Cons a (f a)
type List = Mu List2
```

A second-order functor h is a type constructor that sends functors to functors. This can be concisely expressed with the quantified class constraint $\forall f. \text{Functor } f \Rightarrow \text{Functor } (h f)$, for example in the *Functor* instance of *Mu*:

```
instance (forall f. Functor f => Functor (h f)) => Functor (Mu h)
  where fmap f (In x) = In (fmap f x)
```

²except for the *HFunctor* example (Section 2.1), which needs higher-rank types [28].

Although this is Hinze's preferred formulation he remarks that:

Unfortunately, the extension has not been implemented yet. It can be simulated within Haskell 98 [36], but the resulting code is somewhat clumsy.

Johann and Ghani use essentially the same data-generic representation, the fixpoint of second-order functors, to represent so-called *nested datatypes* [2]. For instance, Hinze [10] represents perfect binary trees with the nested datatype

```
data Perfect a = Zero a | Succ (Perfect (a, a))
```

This can be expressed with the generic representation as *Mu HPerf*, the fixpoint of the second-order functor *HPerf*, defined as

```
data HPerf f a = HZero a | HSucc (f (a, a))
```

Johann and Ghani's notion of second-order functor differs slightly from Hinze's.³ Ideally, their notion would be captured by the following class declaration:

```
class (forall f. Functor f => Functor (h f)) => HFunctor h where
  hfmap :: (Functor f, Functor g)
    => (forall x. f x -> g x) -> (forall x. h f x -> h g x)
```

Like in Hinze's case, the quantified class constraint expresses that a second-order functor takes first-order functors to first-order functors. Additionally, second-order functors provide a second-order *fmap*, called *hfmap*, which replaces f by g , to take values of type $h f x$ to type $h g x$. Yet, in the absence of actual support for quantified class constraints, Johann and Ghani provide the following declaration instead:

```
class HFunctor h where
  ffmap :: Functor f => (a -> b) -> (h f a -> h f b)
  hfmap :: (Functor f, Functor g)
    => (forall x. f x -> g x) -> (forall x. h f x -> h g x)
```

In essence, they inline the *fmap* method provided by the quantified class constraint in the *HFunctor* class. This is unfortunate because it duplicates the *Functor* class's functionality.

2.2 Terminating Corecursive Resolution

Quantified class constraints were first proposed by Hinze and Peyton Jones [12] as a solution to a problem of diverging type class resolution. Consider their generalized rose tree datatype

```
data GRose f a = GBranch a (f (GRose f a))
```

and its *Show* instance

```
instance (Show a, Show (f (GRose f a))) => Show (GRose f a)
  where show (GBranch x xs) = unwords [show x, "-", show xs]
```

Notice the two constraints in the instance context which are due to the two *show* invocations in the method definition. Standard recursive type class resolution would diverge when faced with the constraint $(\text{Show } (\text{GRose } [] \text{ Bool}))$. Indeed, it would recursively resolve the instance context: *Show Bool* is easily dismissed, but *Show [GRose [] a]* requires resolving *Show (GRose [] Bool)* again. Clearly this process loops.

To solve this problem, Hinze and Peyton Jones proposed to write the *GRose* instance with a quantified type class constraint as:

```
instance (Show a, forall x. Show x => Show (f x)) => Show (GRose f a)
  where show (GBranch x xs) = unwords [show x, "-", show xs]
```

³ It is more in line with the category theoretical notion of endofunctors over the category of endofunctors.

This would avoid the diverging loop in the type system extension they sketch, because the two recursive resolvents, $Show\ Bool$ and $\forall x. Show\ x \Rightarrow Show\ [x]$ are readily discharged with the available $Bool$ and $[a]$ instances.

When faced with the same looping issue in their *Scrap Your Boilerplate* work, Lämmel and Peyton Jones [22] implemented a different solution: *cycle-aware constraint resolution*. This approach detects that a recursive resolvent is identical to one of its ancestors and then ties the (co-)recursive knot at the level of the underlying type class dictionaries.

Unfortunately, cycle-aware resolution is not a panacea. It only deals with a particular class of diverging resolutions, those that cycle. The fixpoint of the second-order functor $HPerf$ presented above is beyond its capabilities.

```
instance (Show (h (Mu h) a))  $\Rightarrow$  Show (Mu h a) where
  show (In x) = show x
```

```
instance (Show a, Show (f (a, a)))  $\Rightarrow$  Show (HPerf f a) where
  show (HZero a) = "(Z " ++ show a ++ ")"
  show (HSucc xs) = "(S " ++ show xs ++ ")"
```

Resolving $Show\ (Mu\ HPerf\ Int)$ diverges without cycling back to the original constraint due to the nestedness of the perfect tree type:

```

  Show (Mu HPerf Int)
 $\rightarrow$  Show (HPerf (Mu HPerf) Int)
 $\rightarrow$  Show Int, Show (Mu HPerf (Int, Int))
 $\rightarrow$  Show (HPerf (Mu HPerf) (Int, Int))
 $\rightarrow$  Show (Int, Int), Show (Mu HPerf ((Int, Int), (Int, Int)))
 $\rightarrow$  ...
```

In contrast, with quantified type class constraints we can formulate the instances in a way that resolution does terminate.

```
instance (Show a,
   $\forall f\ x. (Show\ x, \forall y. Show\ y \Rightarrow Show\ (f\ y)) \Rightarrow Show\ (h\ f\ x)$ 
 $\Rightarrow Show\ (Mu\ h\ a)$  where show (In x) = show x
```

```
instance (Show a,  $\forall x. Show\ x \Rightarrow Show\ (f\ x)$ )  $\Rightarrow Show\ (HPerf\ f\ a)$ 
  where show (HZero a) = "(Z " ++ show a ++ ")"
  show (HSucc xs) = "(S " ++ show xs ++ ")"
```

2.3 Summary

In summary, quantified type class constraints enable 1) expressing more of a type class's specification in a natural and succinct manner, and 2) terminating type class resolution for a larger group of applications.

In the remainder of this paper we provide a declarative type system for a Haskell-like calculus with quantified class constraints (Section 3). Type inference is shown in Section 4 and Section 5 provides an elaboration into System F. Section 6 presents the conditions we require to ensure termination in the presence of quantified class constraints. Finally, Section 7 discusses related work and Section 8 concludes.

3 Declarative Type System

This section provides the declarative type system specification for our core Haskell calculus with quantified class constraints.

```
x, y, z, f ::= <term variable name>
a, b, c    ::= <type variable name>
TC         ::= <class name>
```

```
pgm ::= e | cls; pgm | inst; pgm           program
cls  ::= class A  $\Rightarrow$  TC a where { f ::  $\sigma$  }   class decl.
inst ::= instance A  $\Rightarrow$  TC  $\tau$  where { f = e }   instance decl.
```

```
e ::= x |  $\lambda x. e$  |  $e_1\ e_2$  | let x =  $e_1$  in  $e_2$            term
```

```
 $\tau$  ::= a |  $\tau_1 \rightarrow \tau_2$            monotype
 $\rho$  ::=  $\tau$  |  $C \Rightarrow \rho$            qualified type
 $\sigma$  ::=  $\rho$  |  $\forall a. \sigma$            type scheme
```

```
A ::= • | A, C           axiom set
C ::= Q |  $C_1 \Rightarrow C_2$  |  $\forall a. C$            constraint
Q ::= TC  $\tau$            class constraint
```

```
 $\Gamma$  ::= • |  $\Gamma, x : \sigma$  |  $\Gamma, a$            typing environment
P ::=  $\langle A_S, A_I, A_L \rangle$            program theory
```

Figure 1. Source Syntax

3.1 Syntax

Figure 1 presents the, mostly standard, syntax of our source language. A program pgm consists of class declarations cls , instance declarations $inst$ and a top-level expression e . For simplicity, each class has a single parameter and a single method.

Terms e comprise a λ -calculus extended with let-bindings. By convention, we use f to denote a method name and x, y, z to denote any kind of term variable name.

Types also appear in Figure 1. Like all extensions of the Damas-Milner system [5] with qualified types [14], we discriminate between monotypes τ , qualified types ρ and type schemes σ . Note that, to avoid clutter, our formalization does not feature higher-kinded types, but our prototype implementation does.

Our calculus differs from Haskell'98 in that it conservatively generalizes the language of constraints. In Haskell'98 the constraints that can appear in type signatures and in class and instance contexts are basic class constraints Q of the form $TC\ \tau$. As a consequence, the constraint schemes or axioms that are derived from instances (and for superclasses) are Horn clauses of the form:

$$\forall \bar{a}. Q_1 \wedge \dots \wedge Q_n \Rightarrow Q_0$$

These axioms are similar to rank-1 polymorphic types in the sense that the quantifiers (and the implication) only occur on the outside. We allow a more general form of constraints C where, in analogy with higher-rank types, quantifiers and implications occur in nested positions. This more expressive form of constraints can occur in signatures and class/instance contexts. Consequently, the syntactic sort C of constraints and axioms is one and the same.

Note that constraint schemes of the form $\forall \bar{a}. (Q_1 \wedge \dots \wedge Q_n) \Rightarrow Q_0$, used in earlier formalizations of type classes (e.g., [25]), are not valid syntax for our constraints C because we do not provide a notation for conjunction. Yet, we can easily see the scheme notation as syntactic sugar for a curried representation:

$$\forall \bar{a}. (Q_1 \wedge \dots \wedge Q_n) \Rightarrow Q_0 \quad \equiv \quad \forall \bar{a}. Q_1 \Rightarrow (\dots (Q_n \Rightarrow Q_0) \dots)$$

$P; \Gamma \varepsilon_m e : \sigma$	Term Typing
$\frac{(x : \sigma) \in \Gamma}{P; \Gamma \varepsilon_m x : \sigma}$	TMVAR
$\frac{x \notin \text{dom}(\Gamma) \quad P; \Gamma, x : \tau \varepsilon_m e_1 : \tau \quad P; \Gamma, x : \tau \varepsilon_m e_2 : \sigma}{P; \Gamma \varepsilon_m (\mathbf{let} \ x = e_1 \ \mathbf{in} \ e_2) : \sigma}$	TMLLET
$\frac{a \notin \Gamma \quad P; \Gamma, a \varepsilon_m e : \sigma}{P; \Gamma \varepsilon_m e : \forall a. \sigma}$	$(\forall I)$
$\frac{x \notin \text{dom}(\Gamma) \quad \Gamma \varepsilon_y \tau_1 \quad P; \Gamma, x : \tau_1 \varepsilon_m e : \tau_2}{P; \Gamma \varepsilon_m \lambda x. e : \tau_1 \rightarrow \tau_2}$	$(\rightarrow I)$
$\frac{P; \Gamma \varepsilon_m e_1 : \tau_1 \rightarrow \tau_2 \quad \Gamma \varepsilon_t C}{P; \Gamma \varepsilon_m e_2 : \tau_1}$	$(\rightarrow E)$
$\frac{P_{s,l} C; \Gamma \varepsilon_m e : \rho}{P; \Gamma \varepsilon_m e : C \Rightarrow \rho}$	$(\Rightarrow I)$
$\frac{P; \Gamma \varepsilon_m e : C \Rightarrow \rho \quad P; \Gamma \models C}{P; \Gamma \varepsilon_m e : \rho}$	$(\Rightarrow E)$
$\frac{P; \Gamma \varepsilon_m e : \forall a. \sigma \quad \Gamma \varepsilon_y \tau}{P; \Gamma \varepsilon_m e : [a \mapsto \tau] \sigma}$	$(\forall E)$
$\Gamma \varepsilon_{ls} \text{cls} : A_S; \Gamma_c$	Class Declaration Typing
$\frac{\Gamma, a \varepsilon_t C_i \quad \Gamma, a \varepsilon_y \sigma}{\Gamma \varepsilon_{ls} \mathbf{class} (C_1, \dots, C_n) \Rightarrow TC \ a \ \mathbf{where} \ \{ f :: \sigma \} : [\forall a. TC \ a \Rightarrow C_i]; [f : \forall a. TC \ a \Rightarrow \sigma]}$	CLASS
$P; \Gamma \varepsilon_{nst} \text{inst} : A_I$	Class Instance Typing
$\frac{\bar{b} = fv(\tau) \quad \Gamma, \bar{b} \varepsilon_{bx} A \quad \mathbf{class} (C_1, \dots, C_n) \Rightarrow TC \ a \ \mathbf{where} \ \{ f :: \sigma \} \quad P_{s,l} A; \Gamma, \bar{b} \models [\tau/a] C_i \quad P_{s,l} A_{s,l} TC \ \tau; \Gamma, \bar{b} \varepsilon_m e : [\tau/a] \sigma}{P; \Gamma \varepsilon_{nst} \mathbf{instance} \ A \Rightarrow TC \ \tau \ \mathbf{where} \ \{ f = e \} : [\forall \bar{b}. A \Rightarrow TC \ \tau]}$	INSTANCE

Figure 2. Declarative Type System (Selected Rules)

We denote a list of C -constraints as A , short for *axiom set* as we use them to represent, among others, axioms given through type class instances.

Finally, Figure 1 presents typing environments Γ , which are entirely standard, and the program theory P . The latter is a triple of three axiom sets: the superclass axioms A_S , the instance axioms A_I and local axioms A_L . We use the notation $P_{s,l} C$ to denote that we extend the local component of the triple, and similar notation for the other components. In earlier type class formalizations these separate kinds of axioms are typically conflated into a single axiom set. However, in this paper it is convenient to distinguish them for accurately stating the different restrictions imposed on them. Moreover, it is instructive for contrasting with regular Haskell. In Haskell, the local constraints are basic type class constraints Q only, while the instance and superclass axioms have the more expressive Horn clause form. In contrast, in our setting all three components support the same (and more general) form of *Harrop formulae*.

3.2 The Type System

Figure 2 presents the main judgments of our declarative type system for the language of Figure 1, namely term typing and typing of class and instance declarations.

Type & Constraint Well-Scopedness The judgments for well-scopedness of types, constraints and axiom sets are denoted $\Gamma \varepsilon_y \sigma$, $\Gamma \varepsilon_t C$ and $\Gamma \varepsilon_{bx} A$ respectively. Their definitions are straightforward and can be found in Appendix A.

Term Typing Term typing takes the form $P; \Gamma \varepsilon_m e : \sigma$ and can be read as “under program theory P and typing environment Γ , expression e has type σ ”. The rules are almost literally those of Chakravarty et al. [3]. There is only one difference, which is a simplification for the sake of convenience: following Vytiniotis et al. [37] we have opted for recursive let-bindings that are not generalized. Nevertheless, we generalize the type of top-level bindings (see Appendix A).

Apart from that, there are no noticeable differences with conventional Haskell in the typing rules. All the interesting differences are concentrated in the definition of the constraint entailment judgment $P; \Gamma \models C$, which is used in the constraint elimination Rule

$(\Rightarrow E)$. The definition of this auxiliary judgment is discussed in detail in Section 3.3.

Class Declaration Typing Typing for class declarations takes the form $\Gamma \varepsilon_{ls} \text{cls} : A_S; \Gamma_c$ and is given by Rule CLASS, presented in Figure 2.

In addition to checking the well-formedness of the method type, we ensure that the class context (C_1, \dots, C_n) is also well-formed, extending the environment with the local variable a . In turn, this implies that $fv(C_i) \subseteq \{a\}$, in line with the Haskell standard.

As usual, typing a class declaration extends the typing environment with the method typing, and the program’s theory with the superclass axioms. For instance, the extended monad transformer class yields the superclass axiom:

$$\forall t. \text{Trans } t \Rightarrow (\forall m. \text{Monad } m \Rightarrow \text{Monad } (t \ m))$$

Class Instance Typing Instance typing takes the form $P; \Gamma \varepsilon_{nst} \text{inst} : A_I$ and is given by Rule INSTANCE, also presented in Figure 2.

We check the well-formedness of the instance context A under the extended typing environment, and that each superclass constraint C_i is entailed by the instance context.

Finally, we check that the method implementation e has the type indicated by the class declaration, appropriately instantiated for the instance in question.

Program Typing The judgment for program typing ties everything together and takes the form $P; \Gamma \varepsilon_{pgm} \text{pgm} : \sigma$. Its definition is straightforward and can be found in Appendix A.

3.3 Constraint Entailment

Following the approach of Schrijvers et al. [32] for their COCHIS calculus, we present constraint entailment in two steps. First, we provide an easy-to-understand and expressive, yet also highly ambiguous, specification. Then we present a syntax-directed, semi-algorithmic variant that takes the ambiguity away, but has a more complicated formulation inspired by the *focusing* technique used in proof search [1, 23, 24].

Declarative Specification Constraint entailment takes the form $P; \Gamma \models C$, and its high-level declarative specification is given by the following rules:

$$\begin{array}{c} \frac{C \in P}{P; \Gamma \models C} \text{ (SPEC)} \quad \frac{a \notin \Gamma \quad P; \Gamma, a \models C}{P; \Gamma \models \forall a. C} \text{ (VIC)} \quad \frac{P; \Gamma \models \forall a. C \quad \Gamma \vDash_{\text{ty}} \tau}{P; \Gamma \models [\tau/a]C} \text{ (VEC)} \\ \\ \frac{P_{\text{.L}} C_1; \Gamma \models C_2}{P; \Gamma \models C_1 \Rightarrow C_2} \text{ (\Rightarrow IC)} \quad \frac{P; \Gamma \models C_1 \Rightarrow C_2}{P; \Gamma \models C_1} \text{ (\Rightarrow EC)} \end{array}$$

If we interpret constraints C as logical formulas, the above rules are nothing more than the rules of the universal fragment of Hereditary Harrop logic [9]. Rule (SPEC) is the standard axiom rule. Rules (\Rightarrow IC) and (\Rightarrow EC) correspond to implication introduction and elimination, respectively. Similarly, Rules (VIC) and (VEC) correspond to introduction and elimination of universal quantification, respectively. These are also essentially the rules Hinze and Peyton Jones [12] propose.

While compact and elegant, there is a serious downside to these rules: They are highly ambiguous and give rise to many trivially different proofs for the same constraint. For instance, assuming $\Gamma = \bullet, a$ and $P = \langle \bullet, \bullet, Eq a \rangle$, here are only two of the infinitely many proofs of $P; \Gamma \models Eq a$:

$$\frac{Eq a \in P}{P; \Gamma \models Eq a} \text{ (SPEC)}$$

versus

$$\frac{\frac{Eq a \in P'}{P'; \Gamma \models Eq a} \text{ (SPEC)}}{P; \Gamma \models Eq a \Rightarrow Eq a} \text{ (\Rightarrow IC)} \quad \frac{Eq a \in P}{P; \Gamma \models Eq a} \text{ (SPEC)} \quad \frac{Eq a \in P \Rightarrow Eq a}{P; \Gamma \models Eq a} \text{ (\Rightarrow EC)}$$

where $P' = P_{\text{.L}} Eq a$. Observe that the latter proof makes an unnecessary appeal to implication introduction.

Type-Directed Specification To avoid the trivial forms of ambiguity like in the example, we adopt a solution from proof search known as *focusing* [1]. This solution was already adopted by the COCHRIS calculus, for the same reason. The key idea of focusing is to provide a syntax-directed definition of constraint entailment where only one inference rule applies at any given time.

Figure 3 presents our definition of constraint entailment with focusing. The main judgment $P; \Gamma \models C$ is defined in terms of two auxiliary judgments, $P; \Gamma \models [C]$ and $\Gamma; [C] \models Q \rightsquigarrow A$, each of which is defined by structural induction on the constraint enclosed in square brackets.

The main entailment judgment is equivalent to the first auxiliary judgment $P; \Gamma \models [C]$. This auxiliary judgment focuses on the constraint C whose entailment is checked – we call this constraint the “goal”. There are three rules, for the three possible syntactic forms of C . Rules (\Rightarrow R) and (\forall R) decompose the goal by applying implication and quantifier introductions respectively. Once the goal is stripped down to a simple class constraint Q , Rule (QR) selects an axiom C from the theory P to discharge it. The selected axiom must *match* the goal, a notion that is captured by the second auxiliary judgment. Matching gives rise to a sequence A of new (and hopefully simpler) goals whose entailment is checked recursively.

$P; \Gamma \models C$ Constraint Entailment

$$\frac{P; \Gamma \models [C]}{P; \Gamma \models C}$$

$P; \Gamma \models [C]$ Constraint Resolution

$$\frac{P_{\text{.L}} C_1; \Gamma \models [C_2]}{P; \Gamma \models [C_1 \Rightarrow C_2]} \text{ (\Rightarrow R)} \quad \frac{P; \Gamma, b \models [C]}{P; \Gamma \models [\forall b. C]} \text{ (\forall R)}$$

$$\frac{C \in P : \Gamma; [C] \models Q \rightsquigarrow A \quad \forall C_i \in A : P; \Gamma \models [C_i]}{P; \Gamma \models [Q]} \text{ (QR)}$$

$\Gamma; [C] \models Q \rightsquigarrow A$ Constraint Matching

$$\frac{\Gamma; [C_2] \models Q \rightsquigarrow A}{\Gamma; [C_1 \Rightarrow C_2] \models Q \rightsquigarrow A, C_1} \text{ (\Rightarrow L)}$$

$$\frac{\Gamma; [[\tau/b]C] \models Q \rightsquigarrow A \quad \Gamma \vDash_{\text{ty}} \tau}{\Gamma; [\forall b. C] \models Q \rightsquigarrow A} \text{ (\forall L)} \quad \frac{}{\Gamma; [Q] \models Q \rightsquigarrow \bullet} \text{ (QL)}$$

Figure 3. Tractable Constraint Entailment

The second auxiliary judgment $\Gamma; [C] \models Q \rightsquigarrow A$ focuses on the axiom C and checks whether it matches the simple goal Q . Again, there are three rules for the three possible forms the axiom can take. Rule (QL) expresses the base case where the axiom is identical to the goal and there are no new goals. Rule (\Rightarrow L) handles an implication axiom $C_1 \Rightarrow C_2$ by recursively checking whether C_2 matches the goal. At the same time it yields a new goal C_1 which needs to be entailed in order for the axiom to apply. Finally, Rule (\forall L) handles universal quantification by instantiating the quantified variable in a way that recursively yields a match.

It is not difficult to see that this type-directed formulation of entailment greatly reduces the number of proofs for given goal.⁴ For instance, for the example above there is only one proof:

$$\frac{Eq a \in P \quad \Gamma; [Eq a] \models Eq a \rightsquigarrow \bullet \text{ (QL)}}{P; \Gamma \models [Eq a]} \text{ (QR)} \quad \frac{}{P; \Gamma \models Eq a}$$

3.4 Remaining Nondeterminism

While focusing makes the definition of constraint entailment type-directed, there are still two sources of nondeterminism. As a consequence, the specification is still ambiguous and not an algorithm.

Overlapping Axioms The first source of non-determinism is that in Rule (QR) there may be multiple matching axioms that make the entailment go through. For applications of logic where proofs are irrelevant this is not a problem, but in Haskell where the proofs have computational content (namely the method implementations) this is a cause for concern. Haskell'98 also faces this problem. Consider two instances for the same type:

```
class Default a where { default :: a }
instance Default Bool where { default = True }
instance Default Bool where { default = False }
```

⁴Without loss of expressive power. See for example [30].

The two instances give rise to two different proofs for *Default Bool*, with distinct computational content (*True* vs. *False*). We steer away from this problem in the same way as Haskell'98, by requiring that instance declarations do not overlap. This does not rule out the possibility of distinct proofs for the same goal, but at least distinct proofs have the same computational content. Consider a class hierarchy where C is the superclass of both D and E .

```

class  $C$   $a$  where { ... }
class  $C$   $a$   $\Rightarrow$   $D$   $a$  where { ... }
class  $C$   $a$   $\Rightarrow$   $E$   $a$  where { ... }

```

This gives rise to the superclass axioms $\forall a. D a \Rightarrow C a$ and $\forall a. E a \Rightarrow C a$. Given additionally two local constraints $D \tau$ and $E \tau$, we have two ways to establish $C \tau$. The proofs are distinct, yet ultimately the computational content is the same. This is easy to see as only instances supply the computational content and there can be at most one instance for any given type τ .

In summary, non-overlap of instances is sufficient to ensure coherence.

Guessing Polymorphic Instantiation A second source of ambiguity is that Rule (VL) requires guessing an appropriate type τ for substituting the type variable b . Guessing is problematic because there are an infinite number of types to choose from and more than one of those choices can make the entailment work out. Choosing an appropriate type is a problem for the type inference algorithm in the next section. Different choices leading to different proofs is a more fundamental problem that also manifests itself in Haskell'98. Consider the following instances.

```

instance  $C$   $Char$  where { ... }
instance  $C$   $Bool$  where { ... }
instance  $C$   $a$   $\Rightarrow$   $D$   $Int$  where { ... }

```

The third instance gives rise to the axiom $\forall a. C a \Rightarrow D Int$. When resolving $D Int$ with this axiom we can choose a to be either $Char$ or $Bool$ and thus select a different C instance.

Haskell'98 avoids this problem by requiring that all quantified type variables, like a in the example, appear in the head of the axiom. Because our axioms have a more general, recursively nested form, we generalize this requirement in a recursively nested fashion. The predicate $unamb(C)$ in Figure 4 formalizes the requirement in terms of the auxiliary judgment $\bar{a} \vdash_{unamb} C$, where \bar{a} are type variables that need to be determined by the head of C . Rule (QU) constitutes the base case where Q is the head and contains the determinable type variables \bar{a} . Rule (VU) processes a quantifier by adding the new type variable to the list of determinable type variables \bar{a} . Finally, Rule (\Rightarrow U) checks whether the head C_2 of the implication determines the type variables \bar{a} . It also recursively checks whether C_1 is unambiguous on its own. The latter check is necessary because left-hand sides of implications are themselves added as axioms to the theory in Rule (\Rightarrow R); hence they must be well-behaved on their own.

The predicate $unamb(C)$ must be imposed on all constraints that are added to the theory. This happens in four places: the instance axioms added in Rule INSTANCE, the superclass axioms added in Rule CLASS, the local axioms added when checking against a given signature in Rule (\Rightarrow I) and the local axioms added during constraint entailment checking in Rule (\Rightarrow R). These four places can be traced back to three places in the syntax: class and instance heads, and (method) signatures.

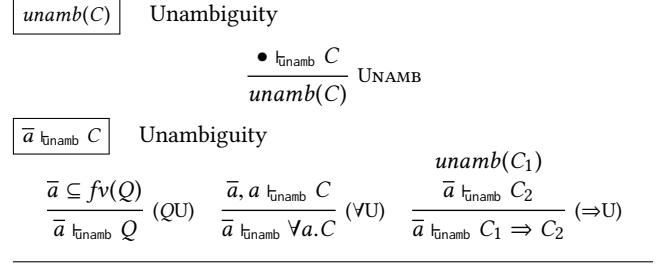


Figure 4. Unambiguity

4 Type Inference

We provide a type inference algorithm with elaboration into System F [7]. To simplify the presentation, this section focuses solely on type inference. The parts of the rules highlighted in gray concern elaboration and are discussed in Section 5.

To make the connection to the relations of the declarative specification (Section 3.2) more clear, corresponding rules share the same name.

4.1 Preliminaries

Before diving into the details of the algorithm, we first introduce some additional notation and constructs.

Variable-Annotated Constraints & Type Equalities Since our goal is to perform type inference and elaboration to System F simultaneously, we annotate all constraints with their corresponding System F evidence term (dictionary variable d). We keep the notational burden minimal by reusing the same letters as in Figure 1, yet with a calligraphic font:

$\mathcal{P} ::= \langle \mathcal{A}_S, \mathcal{A}_I, \mathcal{A}_L \rangle$ variable-annotated theory
 $\mathcal{A} ::= \bullet \mid \mathcal{A}, C$ variable-annotated axiom set
 $C ::= d : C$ variable-annotated constraint
 $Q ::= d : Q$ variable-annotated class constraint

Additionally, like every HM(X)-based system, our type-inference algorithm proceeds by first generating type constraints from the program text (constraint generation) and then solving these constraints independently of the program text (constraint solving).

During constraint generation, our algorithm gives rise to both (variable-annotated) constraints \mathcal{A} , as well as type equalities E :

$E ::= \bullet \mid E, \tau_1 \sim \tau_2$ type equalities

Type & Evidence Substitutions Furthermore, we introduce two kinds of substitutions: type substitutions θ and dictionary substitutions η :

$\theta ::= \bullet \mid \theta \cdot [\tau/a]$ type substitution
 $\eta ::= \bullet \mid \eta \cdot [t/d]$ evidence substitution

A type substitution θ maps type variables to monotypes, while an evidence substitution η maps dictionary variables d to System F terms t (see Section 5.1 for the formal syntax of System F terms).

4.2 Constraint Generation For Terms

Figure 5 presents constraint generation for terms. The relation takes the form $\Gamma \vdash_{\text{tm}} e : \tau \rightsquigarrow t \mid \mathcal{A}; E$. Given a typing environment Γ and a term e we infer (1) a monotype τ , (2) a set of wanted constraints \mathcal{A} , and (3) a set of wanted equalities E . Its definition is standard.

$$\boxed{\Gamma \varepsilon_{\text{tm}} e : \tau \rightsquigarrow t \mid \mathcal{A}; E} \quad \text{Term Typing}$$

$$\frac{\bar{b}, \bar{d} \text{ fresh} \quad (x : \forall \bar{a}. \bar{C} \Rightarrow \tau) \in \Gamma}{\Gamma \varepsilon_{\text{tm}} x : [\bar{b}/\bar{a}] \tau \rightsquigarrow x \bar{b} \bar{d} \mid (d : [\bar{b}/\bar{a}] C); \bullet} \text{TMVAR} \quad \frac{a \text{ fresh} \quad \Gamma, x : a \varepsilon_{\text{tm}} e_1 : \tau_1 \rightsquigarrow t_1 \mid \mathcal{A}_1; E_1 \quad \Gamma, x : \tau_1 \varepsilon_{\text{tm}} e_2 : \tau_2 \rightsquigarrow t_2 \mid \mathcal{A}_2; E_2 \quad \varepsilon_{\text{ty}} \tau_1 \rightsquigarrow v_1}{\Gamma \varepsilon_{\text{tm}} \text{let } x = e_1 \text{ in } e_2 : \tau_2 \rightsquigarrow \text{let } x : v_1 = t_1 \text{ in } t_2 \mid (\mathcal{A}_1, \mathcal{A}_2); (E_1, E_2, a \sim \tau_1)} \text{TMLET}$$

$$\frac{a \text{ fresh} \quad \Gamma, x : a \varepsilon_{\text{tm}} e : \tau \rightsquigarrow t \mid \mathcal{A}; E}{\Gamma \varepsilon_{\text{tm}} \lambda x. e : a \rightarrow \tau \rightsquigarrow \lambda(x : a). t \mid \mathcal{A}; E} \text{TMABS} \quad \frac{a \text{ fresh} \quad \Gamma \varepsilon_{\text{tm}} e_1 : \tau_1 \rightsquigarrow t_1 \mid \mathcal{A}_1; E_1 \quad \Gamma \varepsilon_{\text{tm}} e_2 : \tau_2 \rightsquigarrow t_2 \mid \mathcal{A}_2; E_2}{\Gamma \varepsilon_{\text{tm}} e_1 e_2 : a \rightsquigarrow t_1 t_2 \mid (\mathcal{A}_1, \mathcal{A}_2); (E_1, E_2, \tau_1 \sim \tau_2 \rightarrow a)} \text{TMAPP}$$

Figure 5. Constraint Generation for Terms with Elaboration

Rule **TMVAR** handles variables. We instantiate the polymorphic type $\forall \bar{a}. \bar{C} \Rightarrow \tau$ of a term variable x with fresh unification variables \bar{b} , introducing \bar{C} as wanted constraints, instantiated likewise. Rule **TMABS** assigns a fresh unification variable to the abstracted term variable x , and adds it to the context for checking the body of the abstraction. Rule **TMAPP** handles applications ($e_1 e_2$). We collect wanted class and equality constraints from each subterm, we generate a fresh type variable a for the result and record that the type of e_1 is a function type ($\tau_1 \sim \tau_2 \rightarrow a$). Rule **TMLET** handles (possibly recursive) let bindings.

4.3 Constraint Solving

The type class and equality constraints derived from terms are solved with the following two algorithms.

Solving Equality Constraints We solve a set of equality constraints E by means of unification. The function $\text{unify}(\bar{a}; E) = \theta_{\perp}$ takes the set of equalities and a set of “untouchable” type variables, and returns either the most general unifier θ of the equalities or fails if none exists. The untouchable type variables \bar{a} originate from type signatures; all other type variables are unification variables. The unifier is of course only allowed to substitute unification variables.

The definition of this unification function is folklore, following Damas and Milner [5] and accounting for signatures; it can be found in Appendix A.

Solving Type Class Constraints Figure 6 defines the judgment for solving type class constraints; it takes the form $\bar{a}; \mathcal{P} \models \mathcal{A}_1 \rightsquigarrow \mathcal{A}_2; \eta$. Given a set of untouchable type variables \bar{a} and a theory \mathcal{P} , it (exhaustively) replaces a set of constraints \mathcal{A}_1 with a set of simpler, residual constraints \mathcal{A}_2 , via the auxiliary judgment $\bar{a}; \mathcal{P} \models [C] \rightsquigarrow \mathcal{A}; \eta$, explained below.

This form differs from the specification in Figure 3: we allow constraints to be partially entailed, which in turn allows us to perform *simplification* [17] of top-level signatures. This is standard practice in Haskell when inferring types. For instance, when inferring the signature for

$$f x = [x] == [x]$$

Haskell simplifies the derived constraint $\text{Eq } [a]$ to $\text{Eq } a$, yielding the signature $\forall a. \text{Eq } a \Rightarrow a \rightarrow \text{Bool}$.

Simplification Auxiliary judgment $\bar{a}; \mathcal{P} \models [C] \rightsquigarrow \mathcal{A}; \eta$ uses the theory \mathcal{P} to simplify a single constraint C to a set of simpler constraints without instantiating any of the untouchable type variables \bar{a} . Following the focusing approach, the judgment is defined by three rules, one for each of the syntactic forms of the goal C .

Rules $(\Rightarrow R)$ and $(\forall R)$ recursively simplify the head of the goal. Observe that we add the bound variable b to the untouchables \bar{a}

$$\boxed{\bar{a}; \mathcal{P} \models \mathcal{A}_1 \rightsquigarrow \mathcal{A}_2; \eta} \quad \text{Constraint Solving Algorithm}$$

$$\frac{\nexists C \in \mathcal{A}_1 : \bar{a}; \mathcal{P} \models [C] \rightsquigarrow \mathcal{A}_2; \eta}{\bar{a}; \mathcal{P} \models \mathcal{A}_1 \rightsquigarrow \mathcal{A}_1; \bullet} \text{STOP}$$

$$\frac{\bar{a}; \mathcal{P} \models [C] \rightsquigarrow \mathcal{A}_2; \eta_1 \quad \bar{a}; \mathcal{P} \models \mathcal{A}_1, \mathcal{A}_2 \rightsquigarrow \mathcal{A}_3; \eta_2}{\bar{a}; \mathcal{P} \models \mathcal{A}_1, C \rightsquigarrow \mathcal{A}_3; (\eta_2 \cdot \eta_1)} \text{STEP}$$

$$\boxed{\bar{a}; \mathcal{P} \models [C] \rightsquigarrow \mathcal{A}; \eta} \quad \text{Constraint Simplification}$$

$$\frac{\varepsilon_{\text{ct}} C_1 \rightsquigarrow v_1 \quad \bar{a}; \mathcal{P}, \perp (d_1 : C_1) \models [d_2 : C_2] \rightsquigarrow (d : C); \eta \quad \bar{d}', d_1, d_2 \text{ fresh} \quad \eta' = [\lambda(d_1 : v_1). [d' d_1/d](\eta(d_2))/d_0]}{\bar{a}; \mathcal{P} \models [d_0 : C_1 \Rightarrow C_2] \rightsquigarrow (d' : C_1 \Rightarrow C); \eta'} (\Rightarrow R)$$

$$\frac{\bar{d}', d_C \text{ fresh} \quad \bar{a}, b; \mathcal{P} \models [d_C : C_0] \rightsquigarrow (d : C); \eta \quad \eta' = [\Delta b. [d' b/d](\eta(d_C))/d_0]}{\bar{a}; \mathcal{P} \models [d_0 : \forall b. C_0] \rightsquigarrow (d' : \forall b. C); \eta'} (\forall R)$$

$$\frac{C \in \mathcal{P} : \bar{a}; [C] \models Q \rightsquigarrow \mathcal{A}; \theta; \eta}{\bar{a}; \mathcal{P} \models [Q] \rightsquigarrow \mathcal{A}; \eta} (\text{QR})$$

$$\boxed{\bar{a}; [C] \models Q \rightsquigarrow \mathcal{A}; \theta; \eta} \quad \text{Constraint Matching}$$

$$\frac{d_1, d_2 \text{ fresh} \quad \bar{a}; [d_2 : C_2] \models Q \rightsquigarrow \mathcal{A}; \theta; \eta}{\bar{a}; [d : C_1 \Rightarrow C_2] \models Q \rightsquigarrow \mathcal{A}, d_1 : \theta(C_1); \theta; [d d_1/d_2] \cdot \eta} (\Rightarrow L)$$

$$\frac{d' \text{ fresh} \quad \bar{a}; [d' : C] \models Q \rightsquigarrow \mathcal{A}; \theta; \eta}{\bar{a}; [d : \forall b. C] \models Q \rightsquigarrow \mathcal{A}; \theta; [d(\theta(b))/d'] \cdot \eta} (\forall L)$$

$$\frac{\theta = \text{unify}(\bar{a}; \tau_1 \sim \tau_2)}{\bar{a}; [d' : TC \tau_1] \models d : TC \tau_2 \rightsquigarrow \bullet; \theta; [d'/d]} (\text{QL})$$

Figure 6. Constraint Entailment with Dictionary Construction

when going under a binder in Rule $(\forall R)$. Once the goal is stripped down to a simple class constraint Q , Rule (QR) selects an axiom C whose head matches the goal, and uses it to replace the goal with a set of simpler constraints \mathcal{A} (a process known as *context reduction* [16]). Goal matching is performed by judgment $\bar{a}; [C] \models Q \rightsquigarrow \mathcal{A}; \theta; \eta$, discussed below.

Matching Auxiliary judgment $\bar{a}; [C] \models Q \rightsquigarrow \mathcal{A}; \theta; \eta$ focuses on the axiom C and checks whether it matches the simple goal Q . The main difference between this algorithmic relation and its declarative specification in Figure 3 lies in the type substitution θ . Instead of guessing a type for instantiating a polymorphic axiom

$\Gamma \varepsilon_{\text{cls}} \text{ cls} : \mathcal{A}_S; \Gamma_c \rightsquigarrow \text{fdata}; \overline{\text{fval}}$

Class Declaration Typing

$$\frac{\Gamma, a \varepsilon_y \sigma \quad \varepsilon_y \sigma \rightsquigarrow v \quad \Gamma, a \varepsilon_t C_i \quad \varepsilon_t C_i \rightsquigarrow v_i \quad \overline{d, \overline{d}^n} \text{ fresh} \quad \text{fdata} = \mathbf{data} \ T_{TC} a = K_{TC} \overline{v}^n v \quad \text{fval}_1^i = \mathbf{let} f : (\forall a. T_{TC} a \rightarrow v) = \Lambda a. \lambda (d : T_{TC} a). \text{proj}_{TC}^{i+1}(d) \quad \text{fval}_2^i = \mathbf{let} d_i : (\forall a. T_{TC} a \rightarrow v_i) = \Lambda a. \lambda (d : T_{TC} a). \text{proj}_{TC}^i(d)}{\Gamma \varepsilon_{\text{cls}} (\mathbf{class} (C_1, \dots, C_n) \Rightarrow TC a \ \mathbf{where} \ { f :: \sigma }) : [\overline{d_i} : \forall a. TC a \Rightarrow C_i]; [f : \forall a. TC a \Rightarrow \sigma] \rightsquigarrow \text{fdata}; \text{fval}_1, \overline{\text{fval}_2}^n} \text{CLASS}$$

$\mathcal{P}; \Gamma \varepsilon_{\text{inst}} \text{ inst} : \mathcal{A}_I \rightsquigarrow \text{fval}$

Class Instance Typing

$$\frac{\mathbf{class} (C'_1, \dots, C'_m) \Rightarrow TC a \ \mathbf{where} \ { f :: \sigma} \quad \overline{b} = fv(\tau) \quad \overline{d}, \overline{d}', d_I \text{ fresh} \quad \mathcal{P}_I = \mathcal{P}_{\cdot, \text{L}} \ \overline{d} : C \quad \Gamma_I = \Gamma, \overline{b} \quad \Gamma_I \varepsilon_t C_i \quad \overline{b}; \mathcal{P}_{I, \text{L}} (\overline{d}_I : \sqrt{\overline{b}}. \overline{C}^n \Rightarrow TC \tau); \Gamma_I \varepsilon_m e : [\tau/a] \sigma \rightsquigarrow t \quad \varepsilon_t C_i \rightsquigarrow v_i \quad \overline{b}; \mathcal{P}_I \models \overline{d}' : [\tau/a] C' \rightsquigarrow \bullet; \eta}{\mathcal{P}; \Gamma \varepsilon_{\text{inst}} (\mathbf{instance} (C_1, \dots, C_n) \Rightarrow TC \tau \ \mathbf{where} \ { f = e }) : [\overline{d}_I : \sqrt{\overline{b}}. \overline{C}^n \Rightarrow TC \tau] \rightsquigarrow \mathbf{let} d_I : (\sqrt{\overline{b}}. \overline{v} \rightarrow T_{TC} \tau) = \Lambda \overline{b}. \lambda (d : \overline{v}). K_{TC} \tau \ \eta(\overline{d}') t} \text{INSTANCE}$$

Figure 7. Declaration Elaboration

in Rule (VL) (top-down), we defer the choice until the head of the axiom is met, in Rule (QL) (bottom-up). Observe that Rule (VL) does not record b as untouchable, effectively turning it into a unification variable. Thus, by unifying the head of the axiom with the goal we can determine without guessing an instantiation for all top-level quantifiers, captured by the type substitution θ .

As an example, consider the derivation of one-step simplification of $\forall b. Eq b \Rightarrow Eq [a]$, when $(\forall a. Eq a \Rightarrow Eq [a]) \in \mathcal{P}$:⁵

$$\frac{\frac{\text{unify}(b; a \sim b) = \theta = [b/a]}{b; [Eq [a]] \models Eq [b] \rightsquigarrow \bullet; \theta} \text{(QL)}}{b; [Eq a \Rightarrow Eq [a]] \models Eq [b] \rightsquigarrow Eq b; \theta} \text{(}\Rightarrow\text{L)}}{\frac{b; [\forall a. Eq a \Rightarrow Eq [a]] \models Eq [b] \rightsquigarrow Eq b; \theta}{b; \mathcal{P}, Eq b \models [Eq [b]] \rightsquigarrow Eq b} \text{(QR)}}{b; \mathcal{P} \models [Eq b \Rightarrow Eq [b]] \rightsquigarrow (Eq b \Rightarrow Eq b)} \text{(}\Rightarrow\text{R)}}{\bullet; \mathcal{P} \models [\forall b. Eq b \Rightarrow Eq [b]] \rightsquigarrow (\forall b. Eq b \Rightarrow Eq b)} \text{(}\forall\text{R)}$$

Search As Section 3.4 has remarked, there may be multiple matching axioms, e.g., due to overlapping superclass axioms. The straightforward algorithmic approach to the involved nondeterminism is search, possibly implemented by backtracking. The GHC Haskell implementation can employ a heuristic to keep this search shallow. It does so by using the superclass constraints very selectively: whenever a new local constraint is added to the theory, it pro-actively derives all its superclasses and adds them as additional local axioms. When looking for a match, it does not consider the superclass axioms and prefers the local axioms over the instance axioms. If a matching local axiom exists, it immediately discharges the entire goal without further recursive resolution. This is the case because in regular Haskell local axioms are always simple class constraints Q .

In our setting, we can also implement a (modified version) of GHC's heuristic, but this does not obviate the need for deep search. The reason is that our local axioms are not necessarily simple axioms, and matching against them may leave residual goals that require further recursive resolution. When that recursive resolution gets stuck, we have to backtrack over the choice of axiom. Consider the following example.

$$\begin{array}{l} \mathbf{class} (E a \Rightarrow C a) \Rightarrow D a \\ \mathbf{class} (G a \Rightarrow C a) \Rightarrow F a \end{array}$$

⁵We omit the evidence substitutions for brevity.

Given local axioms $D a$, $F a$ and $G a$, consider what happens when we resolve the goal $C a$. The superclasses $E a \Rightarrow C a$ and $G a \Rightarrow C a$ of respectively $D a$ and $F a$ both match this goal. If we pick the first one, we get stuck when recursively resolving $E a$. However, if we backtrack and consider the second one instead, we can recursively resolve $G a$ against the given local constraint.

In summary, because we do not see a general way to avoid search, our prototype implementation uses backtracking for choosing between the different axioms.⁶

Implementation Our prototype implementation is available at <https://github.com/gkaracha/quantcs-impl>. It incorporates higher-kinded datatypes and performs type inference, elaboration into System F (as explained in the next section), and type checking of the generated code.

The examples we have tested with the prototype provide confidence that our system is sound and that the elaboration is type preserving. The formal proof of the metatheory is future work.

4.4 Checking Declarations

Figure 7 defines type checking of class and instance declarations.

Class Declaration Typing Typing for class declarations is given by Rule CLASS. For the purposes of type inference, Rule CLASS is identical to the corresponding rule of Figure 2, so we defer its analysis to Section 5.5 which discusses elaboration.

Instance Declaration Typing Typing for instance declarations takes the form $\mathcal{P}; \Gamma \varepsilon_{\text{inst}} \text{ inst} : \mathcal{A}_I \rightsquigarrow \text{fval}$ and is given by Rule INSTANCE. For the most part it is identical to the corresponding rule of Figure 2.

The most notable difference is the handling of the method implementation e : method implementations have their type imposed by the method signature in the class declaration. Hence, we need to *check* rather than *infer* their type.

This operation is expressed succinctly by relation $\overline{a}; \mathcal{P}; \Gamma \varepsilon_m e : \sigma \rightsquigarrow t$, presented in Figure 8. Essentially, it ensures that the inferred type for e subsumes the expected type σ . A type σ_1 is said to subsume type σ_2 if any expression that can be assigned type σ_1 can also be assigned type σ_2 .

⁶It is worth mentioning that the rules of Figure 6 conservatively extend standard Haskell resolution, both in terms of expressivity and performance.

$$\boxed{\bar{a}; \mathcal{P}; \Gamma \vdash_{\text{tm}} e : \sigma \rightsquigarrow t} \quad \text{Explicitly Annotated Term Typing}$$

$$\frac{\Gamma \vdash_{\text{tm}} e : \tau_1 \rightsquigarrow t \mid \mathcal{A}_e; E_e \quad \bar{d} \text{ fresh} \quad \theta = \text{unify}(\bar{a}, \bar{b}; E_e, \tau_1 \sim \tau_2) \quad \vdash_{\text{ct}} C_i \rightsquigarrow v_i \quad \bar{a}, \bar{b}; \mathcal{P}, \perp \vdash d : C \models \theta(\mathcal{A}_e) \rightsquigarrow \bullet; \eta}{\bar{a}; \mathcal{P}; \Gamma \vdash_{\text{tm}} e : (\forall \bar{b}. \bar{C} \Rightarrow \tau_2) \rightsquigarrow \Lambda \bar{b}. \lambda(d : v). \eta(\theta(t))} (\leq)$$

Figure 8. Subsumption Rule

$$\begin{array}{ll}
fpgm ::= t \mid fval; fpgm \mid fdata; fpgm & \text{program} \\
fval ::= \text{let } x : v = t & \text{value binding} \\
fdata ::= \text{data } T a = K \bar{v} & \text{datatype} \\
\\
t ::= x \mid K \mid \lambda(x : v). t \mid t_1 t_2 \mid \Lambda a. t \mid t v & \text{term} \\
\mid \text{let } x : v = t_1 \text{ in } t_2 \mid \text{case } t_1 \text{ of } K \bar{x} \rightarrow t_2 & \\
\\
v ::= a \mid v_1 \rightarrow v_2 \mid \forall a. v \mid T v & \text{type}
\end{array}$$

Figure 9. System F Syntax

Rule (\leq) performs type inference and type subsumption checking simultaneously: First, it infers a monotype τ_1 for expression e , as well as wanted constraints \mathcal{A}_e and type equalities E_e . Type equalities E_e should have a unifier and the inferred type τ_1 should also be unifiable with the expected type τ_2 . Finally, the given constraints \bar{C} should completely entail the wanted constraints \mathcal{A}_e .

4.5 Program Typing

Type inference and elaboration for programs is straightforward and can be found in Appendix A.

5 Translation to System F

This section discusses the elaboration aspect of the algorithm presented in Section 4.

5.1 Target Language: System F

Syntax The syntax of System F [7] – extended with data types and recursive let-bindings – is presented in Figure 9 and is entirely standard. Like in the source language, we elide all mention of kinds. Without loss of generality, we simplify matters by allowing only data types with a single type parameter and a single data constructor and case expressions with a single branch; this is sufficient for our dictionary-passing translation of type classes.

Semantics & Typing Since the operational semantics and typing for System F with data types are entirely standard and do not contribute to the novelty of this paper, we omit them from our main presentation. They can be found in Appendix B.

5.2 Elaboration of Types & Constraints

Our system follows the traditional approach of translating source type class constraints into explicitly-passed System F terms, the

so-called *dictionaries* [8, 38]. This transition is reflected in the translation of types, performed by judgment $\vdash_{\text{ty}} \sigma \rightsquigarrow v$:

$$\frac{}{\vdash_{\text{ty}} a \rightsquigarrow a} \text{TYVAR} \quad \frac{\vdash_{\text{ty}} \tau_1 \rightsquigarrow v_1 \quad \vdash_{\text{ty}} \tau_2 \rightsquigarrow v_2}{\vdash_{\text{ty}} \tau_1 \rightarrow \tau_2 \rightsquigarrow v_1 \rightarrow v_2} \text{TYARR} \\
\frac{\vdash_{\text{ct}} C \rightsquigarrow v_1 \quad \vdash_{\text{ty}} \rho \rightsquigarrow v_2}{\vdash_{\text{ty}} C \Rightarrow \rho \rightsquigarrow v_1 \rightarrow v_2} \text{TYQUAL} \quad \frac{\vdash_{\text{ty}} \sigma \rightsquigarrow v}{\vdash_{\text{ty}} \forall a. \sigma \rightsquigarrow \forall a. v} \text{TYALL}$$

Rules TYVAR, TYARR and TYALL are straightforward. Rule TYQUAL elaborates a qualified type into a System F arrow type: the constraint C is translated into the dictionary type v_1 , via relation $\vdash_{\text{ct}} C \rightsquigarrow v_1$ which performs elaboration of constraints:

$$\frac{\vdash_{\text{ty}} \tau \rightsquigarrow v}{\vdash_{\text{ct}} TC \tau \rightsquigarrow T_{TC} v} (CQ) \quad \frac{\vdash_{\text{ct}} C \rightsquigarrow v}{\vdash_{\text{ct}} \forall a. C \rightsquigarrow \forall a. v} (CV) \\
\frac{\vdash_{\text{ct}} C_1 \rightsquigarrow v_1 \quad \vdash_{\text{ct}} C_2 \rightsquigarrow v_2}{\vdash_{\text{ct}} C_1 \Rightarrow C_2 \rightsquigarrow v_1 \rightarrow v_2} (C\Rightarrow)$$

Rule (CQ) elaborates a class constraint ($TC \tau$) into a type constructor application ($T_{TC} v$), which corresponds to the type of dictionaries that witness ($TC \tau$). Rule (CV) is straightforward. Rule (C \Rightarrow) elaborates implication constraints of the form ($C_1 \Rightarrow C_2$) into System F arrow types ($v_1 \rightarrow v_2$), that is, types of *dictionary transformers*. As a concrete example, the constraint corresponding to the *Show* instance for type *HPerf* (Section 2.2):

$$\forall f a. \text{Show } a \Rightarrow (\forall x. \text{Show } x \Rightarrow \text{Show } (f x)) \Rightarrow \text{Show } (\text{HPerf } f a)$$

is elaborated into the type

$$\forall f a. T_{\text{Show}} a \rightarrow (\forall x. T_{\text{Show}} x \rightarrow T_{\text{Show}} (f x)) \rightarrow T_{\text{Show}} (\text{HPerf } f a)$$

5.3 Elaboration of Terms

Term elaboration is straightforward. Rule TMVAR handles term variables. The instantiation of the type scheme $\forall a. \bar{C} \Rightarrow \tau$ to $[\bar{b}/\bar{a}]\tau$ becomes explicit in the System F representation, by the application of x to type variables \bar{b} , as well as the fresh dictionary variables \bar{d} , corresponding one-to-one to the implicit constraints \bar{C} . Rule TMABS elaborates λ -abstractions. Since in System F all bindings are explicitly typed, in the elaborated term we annotate the binding of x with its type a . Similarly, Rule TMLET elaborates let bindings, again explicitly annotating x with its type v_1 in the elaborated term. Rule TMAPP is straightforward.

5.4 Dictionary Construction

The entailment algorithm of Figure 6 constructs explicit witness proofs (in the form of dictionary substitutions) while entailing a constraint.

Simplification The evidence substitution η in the simplification relation shows how to construct a witness for the wanted constraint C from the simpler constraints \mathcal{A}' and program theory \mathcal{P} .

The goal of Rule (\Rightarrow R) is to build an evidence substitution η' , which constructs a proof for $(d_0 : C_1 \Rightarrow C_2)$ from the proofs \bar{d}' for the simpler constraints $\bar{C}_1 \Rightarrow \bar{C}_2$. It is instructive to consider the generated evidence substitution in parts, also taking the types into account:

1. η illustrates how to generate a proof for $(d_2 : C_2)$, from the local assumption $(d_1 : C_1)$ and local residual constraints $(\bar{d} : \bar{C})$.

2. $\overline{[d' \ d_1/d]}$ generates proofs for the (local) residual constraints $(\overline{d} : \overline{C})$, by applying the residual constraints $(d' : C_1 \Rightarrow \overline{C})$ to the local assumption $(d_1 : C_1)$.
3. $(\overline{[d' \ d_1/d]} \cdot \eta)(d_2)$ is a proof for C_2 , under assumptions $(d_1 : C_1)$ and $(\overline{d'} : C_1 \Rightarrow \overline{C})$.
4. Finally, we construct the proof for $(d_0 : C_1 \Rightarrow C_2)$ by explicitly abstracting over $d_1 : \lambda(d_1 : v_1). \overline{[d' \ d_1/d]}(\eta(d_2))$

Rule $(\forall R)$ proceeds similarly. Finally, Rule (QR) generates the evidence substitution via constraint matching, which we discuss next.

Matching Similarly, the evidence substitution η in the matching relation shows how to construct a witness for the wanted constraint Q from the simpler constraints \mathcal{A} and program theory \mathcal{P} .

Rule $(\Rightarrow L)$ generates two fresh dictionary variables, d_1 for the residual constraint $\theta(C_1)$, and d_2 for the local assumption C_2 . Finally, dictionary d_2 is replaced by the application of the dictionary transformer d to the residual dictionary d_1 . Rule $(\forall L)$ behaves similarly. The instantiation of the axiom d becomes explicit, by applying it to the chosen type $\theta(b)$. Finally, Rule (QL) is straightforward: since the wanted and the given constraints are identical (given that they unify), the wanted dictionary d is replaced by the given d' .

5.5 Declaration Elaboration

Figure 7 presents the elaboration of both class and instance declarations into System F.

Elaboration of Class Declarations A declaration for a class TC is encoded in System F as a dictionary type T_{TC} , with a single data constructor K_{TC} and $n + 1$ arguments: n arguments for the superclass dictionaries (of type \overline{v}^n) and one more for the method implementation (of type v). For example, the *Trans* declaration of Section 2.1 gives rise to the following dictionary type:

$$\text{data } T_{Trans} \ t = K_{Trans} \ (\forall m. T_{Monad} \ m \rightarrow T_{Monad} \ (t \ m)) \\ (\forall m \ a. T_{Monad} \ m \rightarrow m \ a \rightarrow (t \ m) \ a)$$

Accordingly, we generate $n + 1$ projection functions that extract each of the arguments (d_i extracts the i -th superclass dictionary and f the method implementation). We use $proj_{TC}^i(d)$ to denote pattern matching against d and extracting the i -th argument:

$$proj_{TC}^i(d) \equiv \text{case } d \text{ of } K_{TC} \ \overline{x}^k \rightarrow x_i \quad , \overline{x}^k \text{ fresh}$$

where k denotes the arity of data constructor K_{TC} . E.g., the superclass projection function for class *Trans* takes the form:

$$d_{sc} : \forall t. T_{Trans} \ t \rightarrow (\forall m. T_{Monad} \ m \rightarrow T_{Monad} \ (t \ m)) \\ d_{sc} = \Lambda t. \lambda(d : T_{Trans} \ t). \text{case } d \text{ of } \{ K_{Trans} \ d' _ \rightarrow d' \}$$

Elaboration of Class Instances A class instance is elaborated into a System F dictionary transformer d_I :

$$\text{let } d_I : (\forall \overline{b}. \overline{v} \rightarrow T_{TC} \ \tau) = \Lambda \overline{b}. \lambda(\overline{d} : \overline{v}). K_{TC} \ \tau \ \overline{\eta(d')} \ t$$

Given dictionaries \overline{d} – corresponding to the given context constraints – we need to provide all arguments of the data constructor K_{TC} : (a) the instantiation of the class type parameter, (b) the superclass dictionaries, and (c) the method implementation. The first argument is trivial. We obtain the superclass dictionaries by applying the evidence substitution η on the dictionary variables \overline{d} that abstract over the required superclass constraints. The method implementation t is elaborated via premise

$$\overline{b}; \mathcal{P}_I; \Gamma_I \vdash_{\text{tm}} e : [\tau/a]\sigma \rightsquigarrow t$$

which elaborates type subsumption in a similar manner.

6 Termination of Resolution

Termination of resolution is the cornerstone of the overall termination of type inference. This section discusses how to enforce termination by means of syntactic conditions on the axioms. These conditions are adapted from those of COCHIS [32] and generalize the earlier conditions for Haskell by Sulzmann et al. [35].

Overall Strategy We show termination by characterising the resolution process as a (resolution) tree with goals in the nodes and axioms on the (multi-)edges. The initial goal sits at the root of the tree. A multi-edge from a parent node to its children presents an axiom that matches the parent node's goal and its children are the residual goals. Resolution terminates iff the tree is finite. Hence, if it does not terminate, there is an infinite path from the root in the tree, that denotes an infinite sequence of axiom applications.

To show that there cannot be such an infinite path, we use a norm $\| \cdot \|$ that maps the head⁷ of every goal C to a natural number, its size. The size of a class constraint $TC \ \tau$ is the size of its type parameter τ , which is given by the following equations:

$$\|a\| = 1 \\ \|\tau_1 \rightarrow \tau_2\| = 1 + \|\tau_1\| + \|\tau_2\|$$

If we can show that this size strictly decreases from any parent goal to its children, then we know that, because the order on the natural numbers is well-founded, on any path from the root there is eventually a goal that has no children.

Termination Condition It is trivial to show that the size strictly decreases, if we require that every axiom makes it so. This requirement is formalised as the termination condition of axioms $term(C)$:

$$\frac{}{term(Q)} (QT) \quad \frac{term(C)}{term(\forall a. C)} (\forall T) \\ \frac{term(C_1) \quad term(C_2) \quad \|Q_1\| < \|Q_2\| \quad \forall a \in fv(C_1) \cup fv(C_2) : occ_a(Q_1) \leq occ_a(Q_2)}{term(C_1 \Rightarrow C_2)} (\Rightarrow T)$$

Rule $(\Rightarrow T)$ for $C_1 \Rightarrow C_2$ enforces the main condition, that the size of the residual constraint's head Q_1 is strictly smaller than the head Q_2 of C_2 . In addition, the rule ensures that this property is stable under type substitution. Consider for instance the axiom $\forall a. C(a \rightarrow a) \Rightarrow C(a \rightarrow Int \rightarrow Int)$. The head's size 5 is strictly greater than the context constraint's size 3. Yet, if we instantiate a to $(Int \rightarrow Int \rightarrow Int)$, then the head's size becomes 10 while the context constraint's size becomes 11. Declaratively, we can formulate stability as:

$$\forall \theta. dom(\theta) \subseteq fv(C_1) \cup fv(C_2) \Rightarrow \|\theta(Q_1)\| < \|\theta(Q_2)\|$$

The rule uses instead an equivalent algorithmic formulation which states that the number of occurrences of any free type variable a may not be larger in Q_1 than in Q_2 . Here the number of occurrences of a type variable a in a class constraint $TC \ \tau$ (denoted as

⁷ The head of a constraint is defined as: $head(Q) = Q$; $head(\forall a. C) = head(C)$; and $head(C_1 \Rightarrow C_2) = head(C_2)$.

$occ_a(TC \tau)$ is the same as the number of free occurrences of a in the parameter τ , where function $occ_a(\tau)$ is defined as:

$$occ_a(b) = \begin{cases} 1 & , \text{ if } a = b \\ 0 & , \text{ if } a \neq b \end{cases}$$

$$occ_a(\tau_1 \rightarrow \tau_2) = occ_a(\tau_1) + occ_a(\tau_2)$$

Finally, as the constraints have a recursive structure whereby their components are themselves used as axioms, the rules also enforce the termination condition recursively on the components.

Superclass Condition If we could impose the termination condition above on all axioms in the theory P , we would be set. Unfortunately, this condition is too strong for the superclass axioms. Consider the superclass axiom $\forall a. Ord\ a \Rightarrow Eq\ a$ of the standard Haskell'98 *Ord* type class. Here both *Ord* a and *Eq* a have size 1; in other words, the size does not strictly decrease and so the axiom does not satisfy the termination condition.

To accommodate this and other examples, we impose an alternative condition for superclass axioms. This superclass condition relaxes the strict size decrease to a non-strict size decrease and makes up for it by requiring that the superclass relation forms a *directed acyclic graph* (DAG). The superclass relation is defined as follows on type classes.

Definition 6.1 (Superclass Relation). Given a class declaration

$$\text{class } (C_1, \dots, C_n) \Rightarrow TC\ a \text{ where } \{ f :: \sigma \}$$

each type class TC_i is a superclass of TC , where $head(C_i) = TC_i \tau_i$.

Observe that the DAG induces a well-founded partial order on type classes. Hence, on any path in the resolution tree, any uninterrupted sequence of superclass axiom applications has to be finite. For the length of such a sequence, the size of the goal does not increase (but might not decrease either). Yet, after a finite number of steps the sequence has to come to an end. If the path still goes on at that point, it must be due to the application of an instance or local axiom, which strictly decreases the goal size. Hence, overall we have preserved the variant that the goal size decreases after a bounded number⁸ of steps.

Termination & Soundness Finally, although we have not proven it formally yet, we are confident that soundness of type inference and preservation of typing under elaboration hold independently of termination (and thus are not affected by whether the termination conditions are met). Such a property is crucial for integrating our algorithm within GHC in the future, where flags such as `UndecidableInstances` are heavily used.

7 Related Work

This section discusses related work, focusing mostly on comparing our approach with existing encodings/workarounds in Haskell. The history of quantified class constraints and their demand in previous research was already discussed in Section 1.

The Coq Proof Assistant Coq provides very flexible support for type classes [33] and allows for arbitrary formulas in class and instance contexts – actually the contexts are just parameters. For instance, we can model the *Trans* class as:

```
Class Trans (T : (Type -> Type) -> Type -> Type)
  {forall M, {Monad M} -> Monad (T M)} :=
  { lift : forall A M, {Monad M} -> M A -> (T M) A }.
```

The downside of Coq's flexibility is that resolution can be ambiguous and non-terminating. The accepted workaround is for the programmer to perform resolution manually when necessary. This is acceptable in the context of Coq's interactive approach to proving, but would mean a great departure from Haskell's non-interactive type inference.

Trifonov's Workaround and Monatron Trifonov [36] gives an encoding of quantified class constraints in terms of regular class constraints. The encoding introduces a new type class that encapsulates the quantified constraint, e.g. *Monad_t t* for $\forall m. Monad\ m \Rightarrow Monad\ (t\ m)$, and that provides the implied methods under a new name. This expresses the *Trans* problem as follows:

```
class Monad_t t where
  treturn :: Monad m => a -> t m a
  tbind   :: Monad m => t m a -> (a -> t m b) -> t m b

class Monad_t t => Trans t where
  lift :: Monad m => m a -> t m a
```

While this approach captures the intention of the quantified constraint, it does not enable the type checker to see that *Monad (t m)* holds for any transformer t and monad m . While the monad methods are available for $t\ m$, they do not have the usual name.

For this reason, Trifonov presents a further (non-Haskell'98) refinement of the encoding, which was adopted by the Monatron [13] library⁹ among others. A non-essential difference is that Monatron merges the above *Monad_t* and *Trans* into a single class:

```
class MonadT t where
  lift   :: Monad m => m a -> t m a
  treturn :: Monad m => a -> t m a
  tbind  :: Monad m => t m a -> (a -> t m b) -> t m b
```

The key novelty is that it also makes the methods *treturn* and *tbind* available under their usual name with a single *Monad* instance for all monad transformers.

```
instance (Monad m, MonadT t) => Monad (t m) where
  return = treturn
  (>>=) = tbind
```

With these definitions the monad transformer composition does type check. Unfortunately, the head of the *Monad (t m)* instance is highly generic and easily overlaps with other instances.

The MonadZipper Because they found Monatron's overlapping instances untenable, Schrijvers and Oliveira [31] presented a different workaround for this problem in the context of their monad zipper datatype, which is an extended form of transformer composition. Their solution adds a method *mw* to the *Trans* type class:

```
class Trans t where
  lift :: Monad m => m a -> t m a
  mw  :: Monad m => MonadWitness t m
```

⁸bounded by the height of the superclass DAG

⁹For the implementation see <https://hackage.haskell.org/package/Monatron>

For any monad m this method returns a GADT [29] witness for the fact that $t m$ is a monad. This is possible because with GADTs, type class instances can be stored in the data constructors.

```
data MonadWitness (t :: (* -> *) -> (* -> *)) m where
  MW :: Monad (t m) => MonadWitness t m
```

By pattern matching on the witness of the appropriate type the programmer can bring the required *Monad* ($t_2 m$) constraint into scope to satisfy the type checker.

```
instance (Trans t1, Trans t2) => Trans (t1 * t2) where
  lift :: ∀m a. Monad m => m a -> (t1 * t2) m a
  lift = case (mw :: MonadWitness t2 m) of
    MW -> C · lift · lift

  mw = . . .
```

The downside of this approach is that it offloads part of the type checker's work on the programmer. As a consequence, the code becomes cluttered with witness manipulation.

The constraint Library Kmett's constraint library [20] provides generic infrastructure for reifying quantified constraints in terms of GADTs, generalizing the *MonadZipper* solution above. Additionally, it complements the encoding with ample utilities for the manipulation of such constraints. Unfortunately, it suffers from the same problem: passing, construction and deconstruction of dictionaries needs to be manually performed by the programmer.

Corecursive Resolution Fu et al. [6] address the divergence problem that arises for generic nested datatypes. They turn the diverging resolution with user-supplied instances into a terminating resolution in terms of automatically derived instances. These auxiliary instances are derived specifically to deal with the query at hand; they shift the pattern of divergence to the term-level in the form of co-recursively defined dictionaries. The authors do point out that the class of divergent cases they support is limited and that deriving quantified instances would be beneficial.

COCHIS The calculus of coherent implicits, COCHIS [32], and its focusing-based resolution in particular, have been a major inspiration of this work. Just like this work, COCHIS supports recursive resolution of quantified constraints. Yet, there are a number of significant differences. Firstly, COCHIS does not feature a separate syntactic sort for type classes, but implicitly resolves regular terms in the Scala tradition. As a consequence, it does not distinguish between instance and superclass axioms, e.g., for the sake of enforcing termination and coherence. Perhaps more significantly, COCHIS features local "instances" as opposed to our globally scoped instances. Local instances may overlap with one another and coherence is obtained by prioritizing those instances that are introduced in the innermost scope. This way COCHIS's resolution is entirely deterministic, while ours is non-deterministic (yet coherent) due to overlapping local and superclass axioms.

8 Conclusion

This paper has presented a fully fledged design of quantified class constraints. We have shown that this feature significantly increases the modelling power of type classes, while at the same enables a terminating type class resolution for a larger class of applications. Interesting future work we aim to pursue includes (a) establishing the metatheory, (b) extending the system with quantification over

predicates¹⁰, raising the power of type classes to (a fragment of) second-order logic, and (c) studying the interaction of quantified class constraints with commonly used type-level features like *functional dependencies* [18] or *associated type families* [3], allowing us to integrate the new feature in Haskell's ecosystem.

A Additional Judgments

A.1 Well-formedness of Types & Constraints

Well-formedness of types takes the form $\Gamma \Vdash_{\text{ty}} \sigma$ and is given by the following rules:

$$\frac{a \in \Gamma}{\Gamma \Vdash_{\text{ty}} a} \text{TYVAR} \quad \frac{\Gamma \Vdash_{\text{ty}} \tau_1 \quad \Gamma \Vdash_{\text{ty}} \tau_2}{\Gamma \Vdash_{\text{ty}} \tau_1 \rightarrow \tau_2} \text{TYARR}$$

$$\frac{\Gamma \Vdash_{\text{ct}} C \quad \Gamma \Vdash_{\text{ty}} \rho}{\Gamma \Vdash_{\text{ty}} C \Rightarrow \rho} \text{TYQUAL} \quad \frac{a \notin \Gamma \quad \Gamma, a \Vdash_{\text{ty}} \sigma}{\Gamma \Vdash_{\text{ty}} \forall a. \sigma} \text{TYALL}$$

It is entirely straightforward and ensures that type terms are well-scoped. Rule *TYQUAL* requires checking the well-formedness of our new form of constraints C , via relation $\Gamma \Vdash_{\text{ct}} C$, given by the following rules:

$$\frac{\Gamma \Vdash_{\text{ty}} \tau}{\Gamma \Vdash_{\text{ct}} TC \tau} \text{(CQ)} \quad \frac{\Gamma \Vdash_{\text{ct}} C_1 \quad \Gamma \Vdash_{\text{ct}} C_2}{\Gamma \Vdash_{\text{ct}} C_1 \Rightarrow C_2} \text{(C}\Rightarrow\text{)} \quad \frac{a \notin \Gamma \quad \Gamma, a \Vdash_{\text{ct}} C}{\Gamma \Vdash_{\text{ct}} \forall a. C} \text{(CV)}$$

Finally, an axiom set A is well-formed if all constraints it contains are well-formed:

$$\frac{}{\Gamma \Vdash_{\text{ax}} \bullet} \text{AxNIL} \quad \frac{\Gamma \Vdash_{\text{ax}} A \quad \Gamma \Vdash_{\text{ct}} C}{\Gamma \Vdash_{\text{ax}} A, C} \text{AxCONS}$$

A.2 Program Typing

The judgment for program typing takes the form $P; \Gamma \Vdash_{\text{pgm}} \text{pgm} : \sigma$ and is given by the following rules:

$$\frac{\Gamma \Vdash_{\text{cls}} \text{cls} : AS; \Gamma_C \quad P_{\text{ss}} AS; \Gamma, \Gamma_C \Vdash_{\text{pgm}} \text{pgm} : \sigma}{P; \Gamma \Vdash_{\text{pgm}} (\text{cls}; \text{pgm}) : \sigma} \text{PGMCLS}$$

$$\frac{P; \Gamma \Vdash_{\text{inst}} \text{inst} : AI \quad P_{\text{sl}} AI; \Gamma \Vdash_{\text{pgm}} \text{pgm} : \sigma}{P; \Gamma \Vdash_{\text{pgm}} (\text{inst}; \text{pgm}) : \sigma} \text{PGMINST}$$

$$\frac{P; \Gamma \Vdash_{\text{tm}} e : \sigma}{P; \Gamma \Vdash_{\text{pgm}} e : \sigma} \text{PGMEXPR}$$

For brevity, if $P = \bullet$ and $\Gamma = \bullet$ we denote program typing as $\Vdash_{\text{pgm}} \text{pgm} : \sigma$.

A.3 Unification Algorithm

The unification algorithm takes the form $\text{unify}(\bar{a}; E) = \theta_{\perp}$ and is given by the following equations:

$$\begin{aligned} \text{unify}(\bar{a}; \bullet) &= \bullet \\ \text{unify}(\bar{a}; E, b \sim b) &= \text{unify}(\bar{a}; E) \\ \text{unify}(\bar{a}; E, b \sim \tau) &= \text{unify}(\bar{a}; \theta(E)) \cdot \theta \\ &\quad \text{where } b \notin \bar{a} \wedge b \notin \text{fv}(\tau) \wedge \theta = [\tau/b] \\ \text{unify}(\bar{a}; E, \tau \sim b) &= \text{unify}(\bar{a}; \theta(E)) \cdot \theta \\ &\quad \text{where } b \notin \bar{a} \wedge b \notin \text{fv}(\tau) \wedge \theta = [\tau/b] \\ \text{unify}(\bar{a}; E, (\tau_1 \rightarrow \tau_2) \sim (\tau_3 \rightarrow \tau_4)) &= \text{unify}(\bar{a}; E, \tau_1 \sim \tau_3, \tau_2 \sim \tau_4) \end{aligned}$$

Function *unify* is a straightforward extension of the standard first-order unification algorithm [5]. The only difference between the two lies in the additional argument: the *untouchable* variables \bar{a} .

¹⁰See GHC feature request #5927.

These variables are treated by the algorithm as skolem constants and therefore can not be substituted (they can be unified with themselves though).

A.4 Elaboration of Programs

Elaboration of programs is given by judgment $\mathcal{P}; \Gamma \vdash_{\text{pgm}} \text{pgm} : \sigma \rightsquigarrow \text{fpgm}$:

$\mathcal{P}; \Gamma \vdash_{\text{pgm}} \text{pgm} : \sigma \rightsquigarrow \text{fpgm}$ Program Elaboration

$$\frac{\Gamma \vdash_{\text{cls}} \text{cls} : \mathcal{A}_S; \Gamma_C \rightsquigarrow \text{fdata}; \overline{\text{fval}}}{\mathcal{P}; \mathcal{A}_S; \Gamma; \Gamma_C \vdash_{\text{pgm}} \text{pgm} : \sigma \rightsquigarrow \text{fpgm}} \text{PCLs}$$

$$\frac{\mathcal{P}; \Gamma \vdash_{\text{inst}} \text{inst} : \mathcal{A}_I \rightsquigarrow \text{fval}}{\mathcal{P}; \mathcal{A}_I; \Gamma \vdash_{\text{pgm}} \text{pgm} : \sigma \rightsquigarrow \text{fpgm}} \text{PINS}$$

$$\frac{\Gamma \vdash_{\text{tm}} e : \tau \rightsquigarrow t \mid \mathcal{A}; E \quad \theta = \text{unify}(\bullet; E) \quad \bar{a} = \text{fv}(\theta(\mathcal{A})) \cup \text{fv}(\theta(\tau))}{\bar{a}; \langle \bullet, \mathcal{A}_I, \mathcal{A}_L \rangle \models \theta(\mathcal{A}) \rightsquigarrow \bar{d} : C; \eta \quad \vdash_{\text{ct}} C_i \rightsquigarrow v_i} \text{PEXP}$$

$$\langle \mathcal{A}_S, \mathcal{A}_I, \mathcal{A}_L \rangle; \Gamma \vdash_{\text{pgm}} e : \sqrt{a}. \bar{C} \Rightarrow \theta(\tau) \rightsquigarrow \Delta \bar{a}. \lambda(\bar{d} : v). \eta(\theta(t))$$

Rules PCLs and PINS handle class and instance declarations, respectively, and they are entirely standard. Rule PEXP performs standard type-inference, simplification [17] and generalization for a top-level expression e . For simplicity, we do not utilize *interaction rules* (e.g. we do not simplify the constraints $\{Eq\ a, Ord\ a\}$ to $\{Ord\ a\}$), but it is straightforward to do so. Finally, observe that superclass axioms \mathcal{A}_S are not used for the simplification of wanted constraints. This is standard practice for Haskell but our distinction between the axioms within the program theory allows us to express this explicitly.

B System F Semantics

Both the typing rules and call-by-name operational semantics for System F are entirely standard and can be found elsewhere, we include them here to keep the presentation self-contained. In the following, we denote System F typing environments by Δ :

$\Delta ::= \bullet \mid \Delta, T \mid \Delta, K : v \mid \Delta, a \mid \Delta, x : v$ typing environment

B.1 Term Typing

$\Delta \vdash_{\text{tm}}^F t : v$ Term Typing

$$\frac{(x : v) \in \Delta \quad \Delta \vdash_{\text{tm}}^F x : v}{\Delta \vdash_{\text{tm}}^F x : v} \text{TMVAR} \quad \frac{x \notin \text{dom}(\Gamma) \quad \Delta, x : v_1 \vdash_{\text{tm}}^F t : v_2 \quad \Delta \vdash_{\text{ty}}^F v_1}{\Delta \vdash_{\text{tm}}^F \lambda(x : v_1). t : v_1 \rightarrow v_2} (\rightarrow\text{I})$$

$$\frac{(K : v) \in \Delta \quad \Delta \vdash_{\text{tm}}^F K : v}{\Delta \vdash_{\text{tm}}^F K : v} \text{TMCON} \quad \frac{\Delta \vdash_{\text{tm}}^F t_1 : v_1 \rightarrow v_2 \quad \Delta \vdash_{\text{tm}}^F t_2 : v_1}{\Delta \vdash_{\text{tm}}^F t_1 t_2 : v_2} (\rightarrow\text{E})$$

$$\frac{a \notin \Gamma \quad \Delta, a \vdash_{\text{tm}}^F t : v}{\Delta \vdash_{\text{tm}}^F \Lambda a. t : \forall a. v} (\forall\text{I}) \quad \frac{\Delta \vdash_{\text{tm}}^F t : \forall a. v \quad \Delta \vdash_{\text{ty}}^F v_1}{\Delta \vdash_{\text{tm}}^F t v_1 : [v_1/a]v} (\forall\text{E})$$

$$\frac{x \notin \text{dom}(\Gamma) \quad \Delta, x : v_1 \vdash_{\text{tm}}^F t_1 : v_1 \quad \Delta \vdash_{\text{ty}}^F v_1 \quad \Delta, x : v_1 \vdash_{\text{tm}}^F t_2 : v_2}{\Delta \vdash_{\text{tm}}^F (\text{let } x : v_1 = t_1 \text{ in } t_2) : v_2} \text{TMLET}$$

$$\frac{\Delta \vdash_{\text{tm}}^F t_1 : T v \quad \bar{x} \notin \text{dom}(\Gamma) \quad (K : \forall a. \bar{v} \rightarrow T a) \in \Delta \quad \Delta, x : [v/a]v \vdash_{\text{tm}}^F t_2 : v_2}{\Delta \vdash_{\text{tm}}^F (\text{case } t_1 \text{ of } K \bar{x} \rightarrow t_2) : v_2} \text{TMCASE}$$

B.2 Well-formedness of Types

$\Delta \vdash_{\text{ty}}^F v$ Type Well-formedness

$$\frac{a \in \Delta \quad \Delta \vdash_{\text{ty}}^F a}{\Delta \vdash_{\text{ty}}^F a} \text{TYVAR} \quad \frac{T \in \Delta \quad \Delta \vdash_{\text{ty}}^F T}{\Delta \vdash_{\text{ty}}^F T} \text{TYCON} \quad \frac{\Delta \vdash_{\text{ty}}^F v_1 \quad \Delta \vdash_{\text{ty}}^F v_2}{\Delta \vdash_{\text{ty}}^F v_1 \rightarrow v_2} \text{TYARR}$$

$$\frac{a \notin \Gamma \quad \Delta, a \vdash_{\text{ty}}^F v}{\Delta \vdash_{\text{ty}}^F \forall a. v} \text{TYALL} \quad \frac{\Delta \vdash_{\text{ty}}^F v_1 \quad \Delta \vdash_{\text{ty}}^F v_2}{\Delta \vdash_{\text{ty}}^F v_1 v_2} \text{TYAPP}$$

B.3 Program Typing

$\Delta \vdash_{\text{pgm}}^F \text{fpgm} : v$ Program Typing

$$\frac{\Delta \vdash_{\text{tm}}^F t : v}{\Delta \vdash_{\text{pgm}}^F t : v} \text{PGMEXPR}$$

$$\frac{\Delta \vdash_{\text{val}}^F \text{fval} : \Delta_v \quad \Delta, \Delta_v \vdash_{\text{pgm}}^F \text{fpgm} : v}{\Delta \vdash_{\text{pgm}}^F (\text{fval}; \text{fpgm}) : v} \text{PGMVAL}$$

$$\frac{\Delta \vdash_{\text{data}}^F \text{fdata} : \Delta_d \quad \Delta, \Delta_d \vdash_{\text{pgm}}^F \text{fpgm} : v}{\Delta \vdash_{\text{pgm}}^F (\text{fdata}; \text{fpgm}) : v} \text{PGMDATA}$$

For brevity, if $\Delta = \bullet$ we denote System F program typing as $\vdash_{\text{pgm}}^F \text{fpgm} : v$.

B.4 Value Binding Typing

$\Delta \vdash_{\text{val}}^F \text{fval} : \Delta_{\text{fval}}$ Value Binding Typing

$$\frac{x \notin \text{dom}(\Gamma) \quad \Delta, x : v \vdash_{\text{tm}}^F t : v \quad \Delta \vdash_{\text{ty}}^F v}{\Delta \vdash_{\text{val}}^F (\text{let } x : v = t) : [x : v]} \text{VAL}$$

B.5 Datatype Declaration Typing

$\Delta \vdash_{\text{data}}^F \text{fdata} : \Delta_{\text{fdata}}$ Datatype Declaration Typing

$$\frac{\Delta, a \vdash_{\text{ty}}^F v}{\Delta \vdash_{\text{val}}^F (\text{data } T a = K \bar{v}) : [T, K : \forall a. \bar{v} \rightarrow T a]} \text{DATA}$$

B.6 Call-by-name Operational Semantics

The small-step, call-by-name operational semantics of System F are presented below:

$t \rightarrow t'$ Operational Semantics (Small-step)

$$\begin{array}{c}
\frac{}{(\lambda a.t) v \longrightarrow [v/a]t} \text{TyBETA} \quad \frac{}{(\lambda(x : v).t) t' \longrightarrow [t'/x]t} \text{TmBETA} \\
\frac{t_1 \longrightarrow t'_1}{(\text{case } t_1 \text{ of } K \bar{x} \rightarrow t_2) \longrightarrow (\text{case } t'_1 \text{ of } K \bar{x} \rightarrow t_2)} \text{CASESTEP} \\
\frac{}{(\text{case } K \bar{t} \text{ of } K \bar{x} \rightarrow t) \longrightarrow [\bar{t}/\bar{x}]t} \text{CASEBETA} \\
\frac{}{(\text{let } x : v = t_1 \text{ in } t_2) \longrightarrow [\text{let } x : v = t_1 \text{ in } t_1/x]t_2} \text{LETBETA}
\end{array}$$

Acknowledgments

We would like to thank Steven Keuchel, Thomas Winant and Leonidas Lampropoulos, for their constructive feedback. This research was partially supported by the Flemish Fund for Scientific Research (FWO). This work was partially supported by project GRACeFUL, which has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 640954.

References

- [1] Jean-marc Andreoli. 1992. Logic Programming with Focusing Proofs in Linear Logic. *Journal of Logic and Computation* 2 (1992), 297–347.
- [2] Richard S. Bird and Lambert G. L. T. Meertens. 1998. Nested Datatypes. In *MPC '98*. Springer, 52–67.
- [3] Manuel M. T. Chakravarty, Gabriele Keller, and Simon Peyton Jones. 2005. Associated Type Synonyms. *SIGPLAN Not.* 40, 9 (2005), 241–253.
- [4] Satvik Chauhan, Piyush P. Kurur, and Brent A. Yorgey. 2016. How to Twist Pointers Without Breaking Them. In *Haskell 2016*. ACM, 51–61.
- [5] Luis Damas and Robin Milner. 1982. Principal Type-schemes for Functional Programs. In *POPL '82*. ACM, 207–212.
- [6] Peng Fu, Ekaterina Komendantskaya, Tom Schrijvers, and Andrew Pond. 2016. Proof Relevant Corecursive Resolution. In *FLOPS 2016*. Springer, 126–143.
- [7] Jean-Yves Girard, Paul Taylor, and Yves Lafont. 1989. *Proofs and Types*. Cambridge University Press.
- [8] Cordelia V. Hall, Kevin Hammond, Simon L. Peyton Jones, and Philip L. Wadler. 1996. Type Classes in Haskell. *ACM TOPLAS* 18, 2 (1996), 109–138.
- [9] Ronald Harrop. 1956. On disjunctions and existential statements in intuitionistic systems of logic. *Math. Ann.* 132, 4 (1956), 347–361.
- [10] Ralf Hinze. 2000. Perfect trees and bit-reversal permutations. *JFP* 10, 3 (2000), 305–317.
- [11] Ralf Hinze. 2010. Adjoint Folds and Unfolds: Or: Scything Through the Thicket of Morphisms. In *MPC '10*. Springer, 195–228.
- [12] Ralf Hinze and Simon Peyton Jones. 2000. Derivable Type Classes. In *Proceedings of the Fourth Haskell Workshop*. Elsevier Science, 227–236.
- [13] Mauro Jaskelioff. 2011. Monatron: an extensible monad transformer library. In *IFL '08*. Springer, Berlin, Heidelberg, 233–248.
- [14] Mark P. Jones. 1992. A theory of qualified types. In *ESOP '92*, Bernd Krieg-Brückner (Ed.). LNCS, Vol. 582. Springer Berlin Heidelberg, 287–306.
- [15] Mark P. Jones. 1995. Functional Programming with Overloading and Higher-Order Polymorphism. In *Advanced Functional Programming*. Springer, 97–136.
- [16] Mark P. Jones. 1995. *Qualified Types: Theory and Practice*. Cambridge University Press.
- [17] Mark P. Jones. 1995. Simplifying and Improving Qualified Types. In *FPCA '95*. ACM, 160–169.
- [18] Mark P. Jones. 2000. Type Classes with Functional Dependencies. In *Programming Languages and Systems*. LNCS, Vol. 1782. Springer Berlin Heidelberg, 230–244.
- [19] Simon Peyton Jones, Mark Jones, and Erik Meijer. 1997. Type classes: an exploration of the design space. In *Proceedings of the 1997 Haskell Workshop*. ACM.
- [20] Edward A. Kmett. 2017. The constraint package. (2017). <https://hackage.haskell.org/package/constraints-0.9.1>.
- [21] Ralf Lämmel and Simon Peyton Jones. 2003. Scrap Your Boilerplate: A Practical Design Pattern for Generic Programming. *SIGPLAN Not.* 38, 3 (2003), 26–37.
- [22] Ralf Lämmel and Simon Peyton Jones. 2005. Scrap Your Boilerplate with Class: Extensible Generic Functions. *SIGPLAN Not.* 40, 9 (2005), 204–215.
- [23] Chuck Liang and Dale Miller. 2009. Focusing and Polarization in Linear, Intuitionistic, and Classical Logics. *Theor. Comput. Sci.* 410, 46 (2009), 4747–4768.
- [24] Dale Miller, Gopalan Nadathur, Frank Pfenning, and Andre Seedorf. 1989. *Uniform Proofs As a Foundation for Logic Programming*. Technical Report.
- [25] J. Garrett Morris. 2014. A Simple Semantics for Haskell Overloading. *SIGPLAN Not.* 49, 12 (2014), 107–118.
- [26] Bruno C.d.S. Oliveira, Adriaan Moors, and Martin Odersky. 2010. Type Classes As Objects and Implicits. *SIGPLAN Not.* 45, 10 (2010), 341–360.
- [27] Bruno C.d.S. Oliveira, Tom Schrijvers, Wontae Choi, Wonchan Lee, and Kwangkeun Yi. 2012. The Implicit Calculus: A New Foundation for Generic Programming. *SIGPLAN Not.* 47, 6 (2012), 35–44.
- [28] Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Mark Shields. 2007. Practical Type Inference for Arbitrary-rank Types. *JFP* 17, 1 (2007).
- [29] Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Geoffrey Washburn. 2006. Simple Unification-based Type Inference for GADTs. *SIGPLAN Not.* 41, 9 (2006), 50–61.
- [30] Frank Pfenning. 2010. Lecture Notes on Focusing. (2010). <https://www.cs.cmu.edu/~fp/courses/oregon-m10/04-focusing.pdf>.
- [31] Tom Schrijvers and Bruno C.d.S. Oliveira. 2011. Monads, Zippers and Views: Virtualizing the Monad Stack. *SIGPLAN Not.* 46, 9 (2011), 32–44.
- [32] Tom Schrijvers, Bruno C. d. S. Oliveira, and Philip Wadler. 2017. *Cochis: Deterministic and Coherent Implicits*. Report CW 705. KU Leuven.
- [33] Matthieu Sozeau and Nicolas Oury. 2008. First-Class Type Classes. In *TPHOLS 2008 (LNCS)*, Vol. 5170. Springer, 278–293.
- [34] Mike Spivey. 2017. Faster Coroutine Pipelines. In *ICFP 2017*. accepted.
- [35] Martin Sulzmann, Gregory J. Duck, Simon Peyton-Jones, and Peter J. Stuckey. 2007. Understanding Functional Dependencies via Constraint Handling Rules. *JFP* 17, 1 (2007), 83–129.
- [36] Valery Trifonov. 2003. Simulating Quantified Class Constraints. In *Haskell '03*. ACM, 98–102.
- [37] Dimitrios Vytiniotis, Simon Peyton Jones, and Tom Schrijvers. 2010. Let Should Not Be Generalized. In *TLDI '10*. ACM, 39–50.
- [38] P. Wadler and S. Blott. 1989. How to Make Ad-hoc Polymorphism Less Ad Hoc. In *POPL '89*. ACM, 60–76.