# Guarded impredicative polymorphism

Anonymous Author(s)

## Abstract

The design space for type systems that support impredicative instantiation is extremely complicated. One needs to strike a balance between expressiveness, simplicity for both the end programmer and the type system implementor, and how easily the system can be integrated with other advanced type system concepts. In this paper, we propose a new point in the design space, which we call guarded impredicativity. Its key idea is that impredicative instantiation in an application is allowed for type variables that occur under a type constructor. The resulting type system has a clean declarative specification — making it easy for programmers to predict what will type and what will not —, allows for a smooth integration with GHC's OUTSIDEIN(X) constraint solving framework, while giving up very little in terms of expressiveness compared to systems like HMF, HML, FPH and MLF. We give a sound and complete inference algorithm, and prove a principal type property for our system.

## 1 Introduction

Type inference for impredicative polymorphism is a deep, deep swamp. There is a dense literature of papers presenting type systems for impredicative polymorphism [1, 8, 9, 22, 23]. Alas none of them quite worked well enough to be deployed in a production compiler: either the system was too complicated for users to predict its behaviour, or it was too complicated to implement, or sometimes both.

Yet it is tantalising: what is wrong with the type $[\forall a.\, a \rightarrow a]$, a list of polymorphic functions? If we have $xs :: [\forall a.\, a \rightarrow a]$, why can't we write ($head\ xs$)? Not every programmer wants such types but, when they do, it is very annoying that they are disallowed, apparently for obscure technical reasons. That is why we keep trying.

So what is the problem? The difficulty is that to accept ($head\ xs$) we must instantiate the type variable of $head$'s type with a polymorphic type. More precisely, since $head :: \forall p.\, [p] \rightarrow p$, we must instantiate $p$ with ($\forall a.\, a \rightarrow a$). This instantiation seems deceptively simple, but in practice it is extremely hard to combine with type inference. We respond to this challenge by making the following contributions:

- Every attempt to combine type inference with impredicativity involves a design trade-off between complexity, expressiveness, and annotation burden. Our key contribution is a new trade-off, which we call *guarded instantiation* or GI (Section 2).

GI is simple: simple for the programmer to understand (Section 2.1-2.3), simple in its declarative specification and metatheory (Section 3); and simple in its implementation (Section 4). We do not extend the syntax of (System F) types in order to provide a specification of the type system (unlike previous work [1, 9, 23]), nor do we introduce new forms of annotations [19] or side-conditions that require principal types [8].

- Despite GI's relative simplicity, it accepts without annotation particularly celebrated and practically-important examples, such as $runST\ \$\ e$ (Figure 2).

- We give a declarative type system for GI for a small core language, highlighting the key ideas of our system (Section 3). Then we show simple extensions to handle a more full-fledged language, including type annotations (Section 3.3), **let** bindings (Section 3.4), and pattern matching (Appendix A). The system has a notion of *principal type* akin to Hindley-Milner type systems, that is, existence of a monomorphic substitution mediating between types. In particular, impredicativity is never guessed in GI (Section 3.5). The resulting system can express any System F program.

- We provide a sound and complete *inference algorithm* (Section 4) for GI, based on *constraints*. The type inference algorithm is a modest extension of the constraint-based algorithm already used by GHC. Type-correct programs can readily be elaborated into System FC, GHC's intermediate language, without extensions. Our inference algorithm scales readily to handle GADTs, type classes, higher kinds, type-level functions and other type system features.
Moreover, as we discuss in Section 5, we can reduce the annotation burden by relaxing the guardedness conditions during solving.

- We provide a prototype implementation of the whole system, integrated with Haskell's type classes.

Type inference for impredicativity is dense with related work, as we discuss in Section 6. A small but useful contribution to a dense field is Figure 2, which presents key examples from the literature and shows how each major system behaves on that example.

## 2 The key idea: intuition and examples

We begin with an informal introduction to GI, which we make fully precise in Section 3. In this discussion we make use of functions defined in Figure 1.

## 2.1 Exploiting the easy case

What is hard about typing (*head ids*)? Nothing! Since the type variable $p$ appears under a list type constructor in *head*'s type, and we know the type of *ids*, it is plain as a pikestaff that we must instantiate $p := \forall a. a \to a$. The difficulty comes when we have a "naked" or "un-guarded" type variable in the function type, such as *single* :: $\forall p. p \to [p]$. Now if we examine (*single id*), it is not clear whether we should instantiate $p$ with $\forall a. a \to a$, or with $Int \to Int$, or some other monomorphic type.

In fact, (*single id*) does not have a most general type. It has both of these two incomparable types $\forall a. [a \to a]$ and $[\forall a. a \to a]$. To make things worse (*single id*) is a perfectly typeable Hindley-Milner program (with the former type) so we must allow this type. But to support impredicativity, we must also allow the latter. But under which conditions?

Our approach is to exploit the common case. We focus on $n$-ary applications $(f\, e_1...e_n)$. It is more conventional to deal with binary applications, but in fact $n$-ary applications (unencumbered with intervening **let** or **case** constructs) are wildly dominant in practice, and we can get much better typing by treating the application all at once. Then we adopt this rule to type such $n$-ary applications:

> The **Instantiation Rule**. When instantiating an $n$-ary call of $f$ :: $\forall a_1 ... a_p. \tau$, a type variable $a_i$ may only be assigned a polymorphic type $\sigma$ if, *in the instantiated type*, $\sigma$ appears under a type constructor in argument position $1 . . . n$. In this case we say that $a$ appears *guarded*.

This rule is carefully crafted. To illustrate, consider these examples (consult Figure 1 for the types):

- (*map poly*) is OK because, after instantiation, *map*'s type becomes $((\forall a. a \to a) \to (Int, Bool)) \to [\forall a. a \to a] \to [(Int, Bool)]$. The instantiating type $(\forall a. a \to a)$ appears under a type constructor of one of the supplied arguments (*poly* in this case).
- (*single ids*) is OK because, after instantiation, *single*'s type becomes $[\forall a. a \to a] \to [[\forall a. a \to a]]$. The instantiating type appears under a list type constructor in one of the supplied arguments, namely *ids*.
- ((:) *id*): here we cannot instantiate (:) at the polytype $\sigma_{id} = \forall a. a \to a$ because, because there is only one supplied argument, and $\sigma_{id}$ does not appear guarded in the first argument of the instantiated type of (:). On the other hand ((:) *id ids*) would be fine.
- (*id poly* $(\lambda x. x)$) is a tricky one. Here *id* is applied to two arguments although its type, $\forall a. a \to a$, apparently only has one; moreover the type of *id*'s second argument must be polymorphic, and the $(\lambda x. x)$ must be generalised. But the Instantiation Rule says that this application is OK: in the instantiated type of *id*, the $\forall$ appears under an arrow type, in the type of the first argument.

The Instantiation Rule is still informal (which we remedy in Section 3) but it is very helpful to have a rule of thumb to explain to a programmer what will and will not work.

## 2.2 Ignoring the context of a call

Notice that the Intantiation Rule *takes no account of the context of the call*. For example, consider *ids* ++ *single id*. Here we append two lists, so we know that the result of (*single id*) must be $[\forall a. a \to a]$, and you might think that would be enough to fix the instantiation of *single*. But not in GI! The swamp beckons, and we stay on dry land.

Moreover, the Instantiation Rule allows the programmer to understand impredicativity in a simple bottom-up way. For example, consider the expression (*map head* (*single ids*)) (Figure 2). In our type system, the types for *head* and *single ids* are instantiated independently, so we never need to consider the *interaction* between the two arguments. This modularity pays off in the metatheory too.

There is a price to pay, however. As a degenerate case, a function application without any arguments – that is, a variable – may only instantiate *fully* monomorphically – no polymorphism, even if it appears under a type constructor. Thus, the empty list constructor [] :: $\forall a. [a]$ cannot be assigned a type $[\forall a. a \to a]$. We explore how to allow more programs in Section 5, and we also introduce annotations to cover the remaining cases.

## 2.3 Lambdas

In common with many other approaches to impredicativity, we take a conservative position on lambda-bound variables. Consider $g\,(\lambda f. (f\,'x', f\,True))$, where $g$ :: $((\forall a. a \to a) \to (Char, Bool)) \to Int$. Since $g$ can only be applied to a function whose argument is itself polymorphic, you could imagine that information being propagated to $f$ and so the program could be accepted. But, in common with every other system we know, we reject all programs that require a lambda-bound variable to be polymorphic, unless it is explicitly annotated:

> The **Lambda Rule**. *Every lambda abstraction whose argument is polymorphic must be annotated. If it is not annotated, the bound variable can only have a fully monomorphic type.*

By a "fully monomorphic type" we mean "no foralls anywhere". Nothing about guardedness here! Nevertheless, in Section 5 we explore ways to alleviate some of the burden for the programmer, for cases where the argument type, albeit polymorphic, can be inferred from its usage in the body.

While the Lambda Rule deals with the arguments to lambdas, it says nothing about the return type. To get a polymorphic return type, an annotation needs to be provided. For example, for $\lambda(x :: \forall a. a \to a). x\,x$, GI infers the type $(\forall a. a \to a) \to b \to b$, and not $(\forall a. a \to a) \to (\forall a. a \to a)$. To get the latter type, we have to write $\lambda(x :: \forall a. a \to a). (x\,x :: \forall a. a \to a)$ instead.

## 2.4 Expressiveness

By treating *n*-ary applications as a whole, and taking guardedness from both the function type and the argument types, we can infer impredicative instantiations in many practically-useful situations. We summarise a collection of examples culled from the literature in Figure 2. This table also compares our system with others, but we defer discussion of related work to Section 6.

A celebrated example is the function ($) :: $(a \rightarrow b) \rightarrow a \rightarrow b$. Haskellers use this function all the time to remove parentheses in their code, as in (*runST* $ **do** { … }) (the type of *runST* is given in Figure 2). This call absolutely requires impredicative instantiation of the variable *a* in the type of ($). It is so annoying to reject this program that GHC implements a special, built-in typing rule for $f \$ x$. Of course, that is horribly non-modular: if the programmer re-defines another version of ($), even with the same type, some programs cease to type check. In GI both type variables appear under the $(\rightarrow)$ constructor, so impredicative instantiation is allowed.

The lack of support for impredicative types is painful. For example, consider the following from Haskell's `lens` library:

**type** *Lens s t a b* $= \forall f.\ Functor\ f \Rightarrow (a \rightarrow f\ b) \rightarrow s \rightarrow f\ t$

Programmers think of a lens as a first-class value, and are perplexed when they cannot put a lens into a list or other data structure. With GI, many more lens-manipulating programs become well-typed.

One might worry about the order of quantifiers. Take:

$f\ \ :: (\forall a\ b.\ a \rightarrow b \rightarrow b) \rightarrow Int$
$x\ \ :: \forall b\ a.\ a \rightarrow b \rightarrow b$
$g\ \ :: [\forall a\ b.\ a \rightarrow b \rightarrow b] \rightarrow Int$
$xs :: [\forall b\ a.\ a \rightarrow b \rightarrow b]$

The application $(f\ x)$ is well-typed, despite the differing quantifier ordering, because we compare $f$'s argument type and $x$'s actual type using *subsumption*; effectively we instantiate and re-generalise. In contrast, the application $(g\ xs)$ is ill-typed, because under a list constructor we compare the types using *equality*. Happily, while *top-level* quantifiers (such as those for $x$) are invisibly inferred (with unpredicatable ordering), *nested* quantifiers, such as those in $g$ and $xs$'s type, are never inferred but rather declared through a type signature. This makes accidental incompatibility vanishingly rare in practice, as we verify in Section 5.

## 3 Declarative specification

We first present a systematic description of the declarative specification of GI. We use the term declarative in the sense of not syntax-directed. After we have proven soundness and completeness for the constraint-based variant, the programmer can take this declarative specification – easier to understand but without a direct inference algorithm – as a basis to understand when and why annotations are needed.

## 3.1 Syntax

The syntax of the core term language, for our initial discussion, is given in Figure 3. The language has some distinctive features. First, as discussed in Section 2, we deal with *n*-ary applications instead of binary ones. A lone term variable is treated as nullary application. Because of the Lambda Rule (Section 2.3), we provide explicitly-annotated lambda abstractions, $\lambda(x :: \sigma).\ e$, to support lambdas whose bound variable must have a polymorphic type.

Types (Figure 3) are classified by three "sorts", $\mathfrak{u}, \mathfrak{t}$, and $\mathfrak{m}$. *Polymorphic* types $\sigma, \phi$, of sort $\mathfrak{u}$, have unrestricted polymorphism. *Top-level monomorphic* types, $\mu, \eta$, of sort $\mathfrak{t}$, have no polymorphism at the top level, but permit arbitrary nested polymorphic types under a type constructor. Finally, *fully-monomorphic types*, $\tau$, of sort $\mathfrak{m}$, have no trace of polymorphism. Fully monomorphic types correspond to monotypes in the Hindley-Milner tradition. These are the only types which can be assigned to un-annotated lambda-bound variables. We extend this notion to substitutions, and sometimes speak of fully monomorphic substitution to mean that the image of the substitution contains only types of that sort.

For a substitution $\theta$, the image of a type variable $a$ is denoted by $\theta(a)$, and similarly for sort assignments. We denote the application of a substitution to, e.g., a type scheme, or an environment, by juxtaposition.

The definitions related to the guardedness requirement are given in Figure 4. Given an *n*-ary application, the judgment $\sigma \rhd_{\mathfrak{s}}^{n} \Delta$ classifies the free type variables of $\sigma$. Specifically, $\Delta$ maps each variable $a$ to the sort of types that are allowed to instantiate $a$. We combine information from multiple occurrences with $\sqcup$, the least upper bound of the (trivial) sort lattice in Figure 3; for example if a variable appear both unrestricted and monomorphically we can treat it as unrestricted. Using the $\rhd_{\mathfrak{s}}^{n}$ judgment, Figure 4 also defines what it means for a substitution $\theta$ to "respect the guardedness of a type", written $\theta$ respects $\Delta$.

$$\frac{\text{for all } a \in \text{dom}(\theta), \vdash \theta(a) : \Delta(a)}{\theta \text{ respects } \Delta}$$

## 3.2 Typing rules

The typing judgment $\Gamma \vdash e : \sigma$ is given in Figure 5, along with some auxiliary judgments.

Rules ABS and ANNABS concern lambda abstractions, and are straightforward. ANNABS deals with a lambda ($\lambda(x :: \phi).\ e$) where the user has supplied a type annotation $\phi$: we simply bring $x$ into scope with type $\phi$. As discussed in Section 2.3, where there is no annotation (rule ABS) we insist that $x$ has a fully-monomorphic type $\tau$.

All the action is in rule APP for *n*-ary applications ($h\ e_1 \dots e_n$). The first step is typing the head of the application $h$. The corresponding judgment $\vdash^h$ either looks for a variable in the environment or uses the normal typing judgment if the head is another kind of term.

$$head \quad :: \forall p. [p] \rightarrow p \qquad\qquad id \quad :: \forall a. a \rightarrow a \qquad\qquad ids \quad :: [\forall a. a \rightarrow a]$$
$$tail \quad :: \forall p. [p] \rightarrow [p] \qquad\qquad inc \quad :: Int \rightarrow Int \qquad\qquad map \quad :: \forall p\, q. (p \rightarrow q) \rightarrow [p] \rightarrow [q]$$
$$[\,] \quad :: \forall p. [p] \qquad\qquad twice :: \forall a. (a \rightarrow a) \rightarrow a \rightarrow a \qquad app \quad :: \forall a\, b. (a \rightarrow b) \rightarrow a \rightarrow b$$
$$(:) \quad :: \forall p.\, p \rightarrow [p] \rightarrow [p] \qquad choose :: \forall a. a \rightarrow a \rightarrow a \qquad revapp :: \forall a\, b.\, a \rightarrow (a \rightarrow b) \rightarrow b$$
$$single :: \forall p.\, p \rightarrow [p] \qquad\qquad poly \quad :: (\forall a. a \rightarrow a) \rightarrow (Int, Bool) \qquad flip \quad :: \forall a\, b\, c. (a \rightarrow b \rightarrow c) \rightarrow b \rightarrow a \rightarrow c$$
$$(\mathbin{+\!\!+}) \quad :: \forall p. [p] \rightarrow [p] \rightarrow [p] \qquad auto \quad :: (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a) \qquad runST :: \forall v. (\forall s.\, ST\, s\, v) \rightarrow v$$
$$length :: \forall p. [p] \rightarrow Int \qquad\qquad auto' \quad :: (\forall a. a \rightarrow a) \rightarrow b \rightarrow b \qquad argST :: \forall s.\, ST\, s\, Int$$

**Figure 1.** Type signatures for functions used in the text

After typing the head, we instantiate the type of the head by means of the $\leqslant_{\mathfrak{s}}^{n}$ judgment. Note that this *instantiation* judgment is parametrized by a sort $\mathfrak{s}$ (this is needed to support applications with and without annotations).

- All type constructors are *invariant*, even functions. This means that neither $[\forall a. a \rightarrow a] \not\leqslant_{\mathfrak{s}}^{n} [Int \rightarrow Int]$ nor $Int \rightarrow (\forall a. a \rightarrow a) \not\leqslant_{\mathfrak{s}}^{n} Int \rightarrow Int \rightarrow Int$.
- For function types the judgment $\leqslant_{\mathfrak{s}}^{n}$ embodies a bit of their covariance. Given a function $f$ of the type $\forall a.\, a \rightarrow (\forall b.\, b \rightarrow a)$ we are allowed to first instantiate $a$ with some type $\tau_1$, which returns a result $\forall b.\, b \rightarrow \tau_1$, and then instantiate $b$ if a second argument is supplied.

  As a consequence, we can work around the lack of variance of function types by $\eta$-expanding. In the previous example, if we need an expression of type $\tau_1 \rightarrow \tau_2$ for some hole, we cannot directly use $f$. But we can use $\lambda x\, y.\, f\, x\, y$ instead, which allows the $\leqslant_{\mathfrak{s}}^{2}$ judgment to kick in and instantiate the two type variables.
- The status of the variables in the result type depends on the parameter $\mathfrak{s}$ to the judgment. In the case of App, this parameter is set to $\mathfrak{m}$. As a consequence, those variables appearing only in the result of the function – the part of the type past the given number of arguments – have to be instantiated with fully monomorphic types.

Whereas the type of the head of the application may only be instantiated, the arguments may also ongo *generalisation*, that is, abstraction of some of the type variables. Generalisation is embodied by the ArgGen rule. Without such a rule, we would not be able to type check *twice* $(\lambda x.\, x)$. Here the type of $\lambda x.\, x$ must be of the form $\tau \rightarrow \tau$ for some fully monomorphic $\tau$. If we choose $\tau$ to be a fresh Skolem variable, we can derive $\Gamma \vdash^{arg} \lambda x.\, x : \forall a.\, a \rightarrow a$, as needed.

**Theorem 3.1** (Compatibility with Hindley-Milner). *Let $e$ be an expression in the syntax of Hindley-Milner without* **let***. If $\Gamma \vdash^{HM} e : \tau$, then $\Gamma \vdash e : \tau$.*

### 3.3 Annotations in applications

With the rules in Figure 5, a single variable (not applied to any arguments) is treated as a nullary application. Rule App cannot get any guardedness information from its arguments, and its type can only be instantiated with fully monomorphic types. As a result, *we cannot assign the polymorphic type* $[\forall a.\, a \rightarrow a]$ *to the empty list* $[\,]$. That is embarrassing, because we cannot typecheck, say $([\,] \mathbin{+\!\!+} ids)$. Figure 2 has other examples.

One solution is to allow the programmer to give a type annotation for an application; see the syntax in Figure 6. Now, since annotations fully specify the types, we do not need to impose guardedness restrictions on those variables appearing in the result. Thus, the expression $[\,] :: [\forall a.\, a \rightarrow a]$ is accepted, even though the type variable in $a$ is not guarded by a type constructor. GI is quite symmetrical in terms of syntax: both abstractions and applications may be annotated.

Annotations also free us from having a different judgment for declarations and expressions. For every combination $f :: \sigma; f = e$ in the source code, we just need to pose the problem of checking $f = e :: \sigma$ for well-typedness.

The extensions required for annotated applications are described in Figure 6. Rule AnnApp is almost identical to App, except for the choice of parameter to the instantiation judgment, which is $\mathfrak{u}$. This implies that in contrast to non-annotated applications, variables in the result type of the function might be substituted by any type, polymorphic or not. This is sensible, the annotation tells us exactly what the types are that those variables should be instantiated with.

Of course, annotated applications serve for arbitrary applications, not just nullary ones. Take the expression *single* $(\lambda x.\, x)$. Due to the type of *single* being $\forall p.\, p \rightarrow [p]$, the type of the expression must be $[\tau \rightarrow \tau]$ for a monomorphic $\tau$. If we want instead to obtain $[\forall a.\, a \rightarrow a]$, we can just annotate the result, thus *single* $(\lambda x.\, x) :: [\forall a.\, a \rightarrow a]$. Since the variable $a$ appears in the result type of *single*, this is perfectly fine by rule AnnApp.

There is some room for choice when introducing annotations in the language. In particular, the annotations can be taken as *rigid* – the type of the decorated expression is exactly the one appearing as annotation – or *soft* – the resulting type might be an instance of the one stated in the annotation. We take the former route, which is shared by most of the previous work in impredicative polymorphism.

### 3.4 let bindings

Following the long-established tradition, we could translate **let** $x = e_1$ **in** $e_2$ as $(\lambda x.\, e_2)\ e_1$. Alas, such a translation imposes an important restriction on the type of $x$: it must

| | GI decl. | GI relax. | MLF | HMF | FPH | HML |
|---|---|---|---|---|---|---|
| **POLYMORPHIC INSTANTIATION** | | | | | | |
| $const2 = \lambda x\, y.\, y$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| MLF infers $(b \geqslant \forall c.\, c \to c) \Rightarrow a \to b$, GI infers $a \to b \to b$. | | | | | | |
| $choose\ id$ | ✓ | ✓ | ✓ | | ✓ | ✓ |
| MLF and HML infer $(a \geqslant \forall b.\, b \to b) \Rightarrow a \to a$, FPH and GI infer $(a \to a) \to a \to a$. | | | | | | |
| $choose\ [\,]\ ids$ | No | ✓ | ✓ | ✓ | ✓ | ✓ |
| GI needs an annotation on $[\,] :: [\forall a.\, a \to a]$. | | | | | | |
| $\lambda(x :: \forall a.\, a \to a).\, x\, x$ | ✓ | ✓ | ✓ | No | ✓ | ✓ |
| MLF infers $(\forall a.\, a \to a) \to (\forall a.\, a \to a)$, GI infers $(\forall a.\, a \to a) \to b \to b$. | | | | | | |
| $id\ auto$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $id\ auto'$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $choose\ id\ auto$ | No | ✓ | ✓ | No | No | ✓ |
| $choose\ id\ auto'$ | No | No | ✓ | No | No | ✓ |
| $f\ (choose\ id)\ ids$ | No | ✓ | ✓ | No | ✓ | ✓ |
| where $f :: \forall a.\, (a \to a) \to [a] \to a$ | | | | | | |
| GI needs an annotation on $id :: (\forall a.\, a \to a) \to (\forall a.\, a \to a)$ in the previous two examples. | | | | | | |
| $poly\ id$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $poly\ (\lambda x.\, x)$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $id\ poly\ (\lambda x.\, x)$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| **INFERENCE OF POLYMORPHIC ARGUMENTS** | | | | | | |
| $\lambda f.\, (f\ 1, f\ True)$ | No | No | No | No | No | No |
| All systems require an annotation on $f :: \forall a.\, a \to a$. | | | | | | |
| $\lambda xs.\, poly\ (head\ xs)$ | No | No | No | No | No | No |
| All systems except for MLF require an annotation on $xs :: [\forall a.\, a \to a]$. | | | | | | |
| **FUNCTIONS ON POLYMORPHIC LISTS** | | | | | | |
| $length\ ids$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $tail\ ids$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $head\ ids$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $single\ id$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $id : ids$ | ✓ | ✓ | ✓ | No | ✓ | ✓ |
| $(\lambda x.\, x) : ids$ | ✓ | ✓ | ✓ | No | ✓ | ✓ |
| $single\ inc \mathbin{+\!\!+} single\ id$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $g\ (single\ id)\ ids$ | No | No | ✓ | No | ✓ | ✓ |
| where $g :: \forall a.\, [a] \to [a] \to a$ | | | | | | |
| $map\ poly\ (single\ id)$ | No | No | ✓ | ✓ | ✓ | ✓ |
| GI needs an annotation on $single\ id :: [\forall a.\, a \to a]$ in the previous two examples. | | | | | | |
| $map\ head\ (single\ ids)$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| **APPLICATION FUNCTIONS** | | | | | | |
| $app\ poly\ id$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $revapp\ poly\ id$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $runST\ argST$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $app\ runST\ argST$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| $revapp\ runST\ argST$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| **η-EXPANSION** | | | | | | |
| $k\ h\ lst$ | No | No | No | No | No | No |
| $k\ (\lambda x.\, h\ x)\ lst$ | ✓ | ✓ | ✓ | No | ✓ | ✓ |
| where $h :: Int \to \forall a.\, a \to a$, $k :: \forall a.\, a \to [a] \to a$, and $lst :: [\forall a.\, Int \to a \to a]$ | | | | | | |

**Figure 2.** Comparison of type systems

| Skolem / rigid variables | $\mathbb{V}$ | $\ni$ | $a, b, c, \ldots$ |
|---|---|---|---|
| Type constructors | $\mathbb{T}$ | $\ni$ | $\to, \mathsf{T}, \mathsf{S}, \ldots$ |
| Fully mono. types | $\tau$ | $::=$ | $a \mid \mathsf{T}\ \overline{\tau}$ |
| Top-level mono. types | $\mu, \eta$ | $::=$ | $a \mid \mathsf{T}\ \overline{\sigma}$ |
| Polymorphic types | $\sigma, \phi$ | $::=$ | $\forall \overline{a}.\, \mu$ |
| | | | $\overline{a}$ may be empty, $\overline{a} \subseteq \mathrm{ftv}(\mu)$ |
| Term variables | | $\ni$ | $x, y$ |
| Heads | $h$ | $::=$ | $x \mid e$ |
| Expressions / terms | $e$ | $::=$ | $h\ e_1\ \ldots\ e_n$ |
| | | $\mid$ | $\lambda x.e$ |
| | | $\mid$ | $\lambda(x :: \sigma).e$ |
| Environments | $\Gamma$ | $::=$ | $\epsilon \mid \Gamma, x : \sigma$ |
| Substitutions | $\theta, \varphi, \pi$ | $::=$ | $[\overline{a \mapsto \sigma}]$ |
| Sorts of variables | $\mathfrak{s}$ | $::=$ | $\mathfrak{u} \mid \mathfrak{t} \mid \mathfrak{m}$ |
| | | | Sorts form a lattice, $\mathfrak{m} \sqsubset \mathfrak{t} \sqsubset \mathfrak{u}$ |
| Sort assignment | $\Delta$ | $::=$ | $[\overline{a \mapsto \mathfrak{s}}]$ |
| Free type variables | $\mathrm{ftv}(\sigma)$ | | |

**Figure 3.** Syntax of the language

be fully monomorphic even though the type of $e_1$ might be more general. The reason is that we try to guess the type of $e_1$ by looking at the way it is used in $e_2$, instead of looking at $e_1$ itself. But there is no need to be so restrictive! The rule LET in Figure 6 works in the other direction: the type obtained from typing $e_1$ is put in the environment as the one for $x$, allowing the type of $x$ to be fully polymorphic instead.

One difference between **let** bindings in GI and Hindley-Milner is that the latter always generalizes the type of a let-bound identifier before passing it to the body of the let. Vytiniotis et al. [21] argue however that **let**-generalisation is not so important in practice and that in complex type systems how to generalize is not completely clear. If desired, generalisation can be obtained by annotating the bound expression, **let** $x = (e_1 :: \phi)$ **in** $e_2$.

### 3.5 Metatheory

One key property of GI is that all impredicative instantiations are settled by the shape of the expression and the types in the environment. Formally, every possible type derivation for an expression $e$ results in the same type, modulo some monomorphic substitution. For that reason, we say that impredicative polymorphism is *not guessed* in GI.

**Theorem 3.2** (Impredicative instantiation is not guessed)**.** *Let $\Gamma$ be a (possibly open) environment and $e$ an expression. For every pair of fully monomorphic substitutions $\theta_1$ and $\theta_2$,*

1. *If $\theta_1\Gamma \vdash^h e : \sigma_1$ and $\theta_2\Gamma \vdash^h e : \sigma_2$, then there exists a polymorphic type $\sigma^*$ and fully monomorphic substitutions $\varphi_1$ and $\varphi_2$ such that $\sigma_i = \varphi_i\sigma^*$.*
2. *If $\theta_1\Gamma \vdash e : \sigma_1$ and $\theta_2\Gamma \vdash e : \sigma_2$, then there exists a polymorphic type $\sigma^*$ and fully monomorphic substitutions $\varphi_1$ and $\varphi_2$ such that $\sigma_i = \varphi_i\sigma^*$.*

$$\boxed{\sigma \vdash a \text{ guarded}}$$

$$\boxed{\vdash \sigma : \mathfrak{s}}$$

$$\frac{a \notin \overline{b} \qquad \mu \vdash a \text{ guarded}}{\forall \overline{b}. \mu \vdash a \text{ guarded}} \qquad \frac{a \in \text{ftv}(\overline{\phi})}{\top \overline{\phi} \vdash a \text{ guarded}}$$

$$\frac{}{\vdash \sigma : \mathfrak{u}} \qquad \frac{}{\vdash \mu : \mathfrak{t}} \qquad \frac{}{\vdash \tau : \mathfrak{m}}$$

$$\boxed{\sigma \rhd_{\mathfrak{s}}^{n} \Delta}$$

$$\frac{}{\mu \rhd_{\mathfrak{s}}^{0} [\text{ftv}(\mu) \mapsto \mathfrak{s}]} \qquad \frac{\mu \rhd_{\mathfrak{s}}^{n} \Delta}{\forall a. \mu \rhd_{\mathfrak{s}}^{n} \Delta \backslash \overline{a}} \qquad \frac{\sigma_2 \rhd_{\mathfrak{s}}^{n-1} \Delta \qquad \hat{\Delta} = [a \mapsto \text{ if } \sigma_1 \vdash a \text{ guarded}, \mathfrak{u} \text{ else } \mathfrak{t} \mid a \in \text{ftv}(\sigma_1)]}{\sigma_1 \to \sigma_2 \rhd_{\mathfrak{s}}^{n} \Delta \sqcup \hat{\Delta}}$$

**Figure 4.** Definitions related to guardedness

$$\boxed{\Gamma \vdash^{h} e : \sigma}$$

$$\boxed{\Gamma \vdash^{\text{arg}} e : \sigma}$$

$$\frac{x : \sigma \in \Gamma}{\Gamma \vdash^{h} x : \sigma} \text{VarHead} \qquad \frac{e \text{ not var. or app.} \quad \Gamma \vdash e : \sigma}{\Gamma \vdash^{h} e : \sigma} \text{ExprHead} \qquad \frac{\Gamma \vdash e : \sigma \quad \overline{b} \text{ fresh} \quad \sigma \leqslant_{\mathfrak{m}}^{0} [a \mapsto \overline{b}]\mu}{\Gamma \vdash^{\text{arg}} e : \forall a. \mu} \text{ArgGen}$$

$$\boxed{\sigma \leqslant_{\mathfrak{s}}^{n} \phi_1, \ldots, \phi_n, \mu}$$

$$\frac{}{\mu \leqslant_{\mathfrak{s}}^{0} \mu} \text{InstMono} \qquad \frac{n \geqslant 1 \quad \phi_2 \leqslant_{\mathfrak{s}}^{n-1} \sigma_2, \ldots, \sigma_n, \mu}{\phi_1 \to \phi_2 \leqslant_{\mathfrak{s}}^{n} \phi_1, \sigma_2, \ldots, \sigma_n, \mu} \text{InstArrow} \qquad \frac{\mu \rhd_{\mathfrak{s}}^{n} \Delta \quad \theta \text{ respects } \Delta \quad \theta\mu \leqslant_{\mathfrak{s}}^{n} \overline{\sigma}, \eta}{\forall \overline{a}. \mu \leqslant_{\mathfrak{s}}^{n} \overline{\sigma}, \eta} \text{InstPoly}$$

$$\boxed{\Gamma \vdash e : \sigma}$$

$$\frac{\Gamma, x : \tau \vdash e : \sigma}{\Gamma \vdash \lambda x. e : \tau \to \sigma} \text{Abs} \qquad \frac{\Gamma, x : \phi \vdash e : \sigma}{\Gamma \vdash \lambda(x :: \phi). e : \phi \to \sigma} \text{AnnAbs}$$

$$\frac{\Gamma \vdash^{h} h : \phi \quad \phi \leqslant_{\mathfrak{m}}^{n} \sigma_1, \ldots, \sigma_n, \mu \quad \Gamma \vdash^{\text{arg}} e_1 : \sigma_1 \ldots \Gamma \vdash^{\text{arg}} e_n : \sigma_n}{\Gamma \vdash h \, e_1 \, \ldots \, e_n : \mu} \text{App}$$

**Figure 5.** Declarative type system

Expressions / terms $\quad e \quad ::= \quad \ldots \mid (h \, e_1 \, \ldots \, e_n :: \sigma)$
$$\mid \quad \text{let } x = e_1 \text{ in } e_2$$

$$\frac{\Gamma \vdash^{h} h : \phi \qquad \boxed{\phi \leqslant_{\mathfrak{u}}^{n} \sigma_1, \ldots, \sigma_n, \eta}}{\Gamma \vdash^{\text{arg}} e_1 : \sigma_1 \qquad \ldots \qquad \Gamma \vdash^{\text{arg}} e_n : \sigma_n}{\Gamma \vdash (h \, e_1 \, \ldots \, e_n :: \forall \overline{b}. \eta) : \forall \overline{b}. \eta} \text{AnnApp}$$

$$\frac{\Gamma \vdash e_1 : \phi \qquad \Gamma, x : \phi \vdash e_2 : \sigma}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \sigma} \text{Let}$$

**Figure 6.** Decl. system with annotations in apps. and **let**

**Corollary 3.3.** *Let $\Gamma$ be a closed environment. If $\Gamma \vdash e : \sigma_1$ and $\Gamma \vdash e : \sigma_2$, then there is a polymorphic type $\sigma^*$ and fully monomorphic substitutions $\varphi_1$ and $\varphi_2$ such that $\sigma_i = \varphi_i \sigma^*$.*

This property suggests a notion of *principal type* similar to the one found in Hindley-Milner. A principal type for an expression $e$ is defined as a type $\sigma^*$ for which any other type assignment $\phi$ to $e$ is equal to $\theta\sigma^*$ for a fully monomorphic substitution $\theta$. The fact that we only need to consider *fully monomorphic* substitutions here is a direct consequence of Theorem 3.2. The proof of the principal types property, however, is a corollary of other properties of the inference process, which we describe in Section 4.4.

In the remainder of this section we look at some properties of GI concerning derivations and stability under transformations. The latter are important as they provide a basis for the compiler to optimize the code while respecting the typing semantics. Proofs are given in Appendix E.1.

**Theorem 3.4.** *If $\Gamma \vdash u : \sigma$ and $\Gamma, x : \sigma \vdash e[x] : \phi$, then $\Gamma \vdash e[u] : \phi$.*

**Theorem 3.5.** *Let app $:: \forall a \, b. (a \to b) \to a \to b$ and revapp $:: \forall a \, b. a \to (a \to b) \to b$ be the application and reverse application functions, respectively. Given two expressions*

| Unif. var. names | $\mathbb{U}$ | $\ni$ | $\alpha, \beta, \gamma, \delta, \ldots$ |
|---|---|---|---|
| Unif. var. | $\upsilon$ | $::=$ | $\alpha^{\mathfrak{s}}$ |
| Fully mono. types | $\tau$ | $::=$ | $\boxed{\alpha^{\mathfrak{m}}} \mid a \mid \mathsf{T}\,\overline{\tau}$ |
| Top-level mono. types | $\mu, \eta$ | $::=$ | $\boxed{\alpha^{\mathfrak{t}}} \mid a \mid \mathsf{T}\,\overline{\sigma}$ |
| Polymorphic types | $\sigma, \phi$ | $::=$ | $\boxed{\alpha^{\mathfrak{u}}} \mid \forall \overline{a}.\,\mu$ |
| | | | $\overline{a}$ may be empty, $\overline{a} \subseteq \mathrm{ftv}(\mu)$ |
| Types with generalisation | $\mathfrak{g}$ | $::=$ | $\forall \{\overline{\upsilon}\}.\,C \Rightarrow \sigma$ |
| | | | $\overline{\upsilon}$ may be empty |
| Constraints | $C$ | $::=$ | $\top$ |
| Conjunction | | $\mid$ | $C_1 \wedge C_2$ |
| Equality | | $\mid$ | $\sigma \sim \phi$ |
| Instantiation | | $\mid$ | $\sigma \leqslant_{\mathfrak{s}}^n \mu$ |
| Generalisation | | $\mid$ | $\mathfrak{g} \leq \sigma$ |
| Quantification | | $\mid$ | $\forall \overline{a}.\,\exists \overline{\upsilon}.\,C$ |
| Free unification variables | $\mathrm{fuv}(\sigma)$ | | |

**Figure 7.** Extended syntax

*f and e such that* $\Gamma \vdash^{\mathsf{h}} f : \sigma_0$, *and* $\sigma_0 \leqslant_{\mathfrak{m}}^0 \sigma_1 \to \phi$ *then:*

$$\Gamma \vdash f\ e : \phi \iff \Gamma \vdash app\ f\ e : \phi \iff \Gamma \vdash revapp\ e\ f : \phi$$

The hypothesis $\sigma_0 \leqslant_{\mathfrak{m}}^0 \sigma_1 \to \phi$ means that type variables may only be instantiated with fully monomorphic variables. Thus, this transformation only respects well-typedness for the whole predicative and higher-rank fragment, but not in the fully impredicative system in general.

## 4 Type inference using constraints

In the previous section we described GI from a declarative perspective and now we turn to describing an efficient type inference algorithm for it.

Following Pottier and Rémy [15], we first walk over the syntax tree of the source program and generate *typing constraints*, a process that typically introduces many *unification variables* that stand for as-yet-unknown types. Next, we solve those constraints producing a *type substitution* for these unification variables. By separating type inference in two simpler problems, the implementation and conceptual overhead with new source language and type system features remains low. For example, earlier work dubbed OUTSIDEIN(X) applies these ideas to a language like Haskell, with type classes, type-level functions, GADTs, and the like [21]. Another advantage is more sophisticated type-error diagnosis [6, 20, 25].

### 4.1 Constraints

The main challenge of type inference for impredicativity concern instantiation and generalisation of terms with polymorphic type. Consider the call (*head ids True*), when fully elaborated we want to generate this System F term:

*head* $(\forall a.\,a \to a)$ *ids Bool True*

That is, we instantiate *head* at type $(\forall a.\,a \to a)$, then apply it to *ids*, to produce a result of type $(\forall a.\,a \to a)$. Now we must in turn instantiate that type with *Bool* to get a function of type ($Bool \to Bool$) which we can apply to *True*. This second instantiation is problematic because, at constraint generation time, we do not yet know what type we are going to instantiate *head* at; all we know is that (*head ids*) has type $\alpha$ for some as-yet-unknown type $\alpha$. So we want to defer the instantiation decision.

Sometimes we must defer generalisation decisions too. For example, consider the function application $((:)\ (\lambda x.\ x)\ ids)$. In System F terms, we are trying to infer the following elaborated program:

$(:)\ (\forall a.\,a \to a)\ (\Lambda a.\,\lambda(x:a).\,x)\ ids$

in which $(:)$ is instantiated at type $(\forall a.\,a \to a)$, and $(:)$'s first argument is generalised to have that polymorphic type. Now consider constraint generation for this expression. We may instantiate the type of $(:)$ with a fresh unification variable, $\alpha$ say. Ultimately the type of *ids* forces $\alpha$ to be $\forall a.\,a \to a$, but we don't know that yet. Moreover, in the final program we will need to generalise the type of $(\lambda x.\,x)$, but again at constraint generation time we don't know that type either.

When we don't know something at constraint generation time, the solution is to *defer the choice, by generating a constraint that represents that choice*. This is the key idea of the constraint solving approach. The game is to develop a constraint language that neatly embodies the choices that we want to defer, and a solver that can subsequently make those choices. With that in mind, Figure 7 gives the syntax of our constraint language.

As mentioned earlier, constraint generation produces many *unification variables*, each of which stands for an as-yet unknown type. Looking at Figure 7, a key idea is that unification variables are drawn from three distinct "alphabets": $\alpha^{\mathfrak{s}}$ for each of the threes sorts $\mathfrak{s}$. (Sorts were introduced in Figure 3.) The sort of a unification variable specifies the possible types that the unification variable can stand for; operationally, a unification variable of sort $\mathfrak{s}$ may only be unified with types belonging to that sort.

The syntax presents several kinds of constraints $C$. Equality constraints are self explanatory. Instantiation constraints arise from the occurrence of a polymorphic variable, whose type must be instantiated – but that decision must be deferred (embodied in a constraint). Quantification constraints arise from explicit user type signatures, and pattern matching on data types involving existentials and GADTs. Both are fairly conventional. However *generalisation constraints* are new; they precisely embody the deferred decision about generalisation that we mention above. We will elaborate on all these forms in what follows.

In GI, constraints do not form part of the source language; they are internal to the solver. But once we extend this approach to work with type system extensions such as type

classes in Appendix B, we shall impose a distinction between simple constraints – those which the programmer can type – and extended ones – which are internal.

### 4.2 Constraint generation

Constraint *generation* is described in Figure 8 as a four-element judgment $\Gamma \vdash e : \sigma \rightsquigarrow C$. The first two elements are inputs: the environment $\Gamma$ and the expression $e$ for which to generate constraints. The output of the process is a type $\sigma$ assigned to the expression, possibly including some unification variables, and the set of extended constraints $C$ that the types must satisfy. We first focus on the generation process for the core calculus, which we extend later to cover the rest of the language.

Rules Abs and AbsAnn are not surprising: they just extend the environment with a new unification variable or a given polymorphic type, respectively, and then proceed to generate constraints for the body of the abstraction. The usage of a fully monomorphic variable in Abs mimics the restriction imposed by the declarative specification.

Rule App is where most of the work is done. Just like the declarative specification (Figure 5), the head of the application is typed using an ancillary judgment $\Gamma \vdash^h h : \phi \rightsquigarrow C$, which either looks up a variable in the environment or threads the information to the normal gathering process.

The first column of constraints, of form $\beta_i^u \leqslant_m^{n-i} \alpha_{i+1}^u \to \beta_{i+1}^u$, successively decomposes the function type $\phi$ to expose its arguments. At each argument we may need to instantiate the top-level foralls of the function type to expose the arrow; hence the use of subsumption $\leqslant_s^n$ rather than type equality. Note that what we express declaratively in a single relation $\sigma_0 \leqslant_s^n \sigma_1, \ldots, \sigma_n, \mu$ is expressed as a sequence of constraints

$$\sigma_0 \leqslant_s^n \sigma_1 \to \phi_1 \qquad \phi_1 \leqslant_s^{n-1} \sigma_2 \to \phi_2$$
$$\cdots \qquad \phi^{n-1} \leqslant_s^1 \sigma_n \to \phi_n \qquad \phi_n \leqslant_s^0 \mu$$

The reason will become apparent once we describe how solving proceeds.

The second column of constraints introduces a completely new constraint form, $\mathfrak{g} \leq \sigma$, which we call *generalisation* constraint. These constraints allow us to defer the generalisation decisions to the solver, as sketched in Section 4.1. The constraint $(\forall\{\overline{v}\}. C \Rightarrow \phi) \leq \sigma$ should be read *"a term of type $\phi$ with constraints $C$ and unification variables $\overline{v}$ can be instantiated and/or generalised to have type $\sigma$"*. Even looking at the syntax alone, you can see that the fruits of constraint generation for each argument $e_i$ are wrapped up, along with the expected argument type $\alpha_i$ from the function, into a generalisation constraint for the solver to deal with later.

Rule AnnApp deals with a type-annotated application $(h\, e_1 \, \ldots \, e_n :: \sigma)$. It is similar to App, but it is the first rule to introduce a *quantification constraint* $\forall \overline{b}.\, \exists \overline{v}.\, \overline{C}$. This binds the Skolem variables $\overline{b}$ from the type signature, and existentially quantifies the unification variables free in the constraint but not used outside it. The other significant difference is the use

of $\leqslant_u^n$ in the column of instantiation constraints, with suffix u, rather than m in rule App, exactly following the difference between AnnApp and App in Figure 6 and 5 resp.

### 4.3 Constraint solving

The solver takes the generated constraint $C$ and its free unification variables $\overline{v} = \text{fuv}(C)$, and repeatedly applies the solver rules in Figure 9, until no rule applies. The result is a *residual constraint*. If the residual constraint is in *solved form*, then the program is well typed; if not, the unsolved constraints (e.g. Int ~ Bool) represent type errors that can be reported to the user. We will discuss solved form shortly, in Section 4.3.3, but first we concentrate on the solver rules that incrementally solve the constraint.

Each of the rules in Figure 9 rewrites a configuration $C; \overline{v}$ to another configuration. The unification variables $\overline{v}$ are existentially quantified, so you can think of a configuration as representing $\exists \overline{v}.C$. Rule conj and forall are structural rules: the former allows a rule to be applied to one part of a conjunction, while the latter allows a rule to be applied under a quantification. To avoid clutter we implicitly assume that the rules are read modulo commutativity and associativity of $\wedge$; that is why conj only has to handle the left conjunct.

#### 4.3.1 Basic rules

Rule eqrefl removes trivial equality constraints $\sigma \sim \sigma$. Rule eqmono indicates that two types headed by constructors are equal if and only if their heads coincide and all the arguments are equal. eqsubst is the only rule that involves the interaction of two constraints. It applies the substitution of a unification variable to any other constraints conjoined with it (remember the implicit associativity and commutativity of $\wedge$), provided sorts are respected, and there is no occurs-check. Notice that the equality constraint is not discarded; it remains in case it is needed again; indeed, these equality constraints remain in a solved constraint (Section 4.3.3).

Given this different behaviour of the different sorts of variables, the solver has to propagate this information. eqvar ensures that whenever we have two variables with different sorts, the least restrictive one is substituted by the most restrictive one. For example, when we have an unrestricted $\alpha^u$ and a top-level monomorphic $\beta^t$, then $\alpha^u$ should be replaced by $\beta^t$, and not the other way around. Full monomorphism goes deeper: eqfully ensures that if a type $\sigma$ is equated with a fully monomorphic variable $\alpha^m$, all the variables in $\sigma$ become fully monomorphic too.

One difference between these rules and other presentations is that we do not rewrite an unsatisfiable constraint, such as Int ~ Bool, to $\bot$. Instead, that constraint is simply stuck, and we can report it at the end.

Note that two polymorphic types need to be *syntactically equal* (modulo $\alpha$-equality) to match under the eqrefl rule. This means that $(\forall\, a\, b.\, a \to b \to b) \sim (\forall\, b\, a.\, a \to b \to b)$ does *not* hold in our system. As we discuss in Section 2.4,

$$\boxed{\Gamma \vdash^{\mathsf{h}} e : \sigma \rightsquigarrow C}$$

$$\frac{x : \sigma \in \Gamma}{\Gamma \vdash^{\mathsf{h}} x : \sigma \rightsquigarrow \epsilon}\ \textsc{VarHead} \qquad\qquad \frac{e \text{ not var. or app.} \quad \Gamma \vdash e : \sigma \rightsquigarrow C}{\Gamma \vdash^{\mathsf{h}} e : \sigma \rightsquigarrow C}\ \textsc{ExprHead}$$

$$\boxed{\Gamma \vdash e : \sigma \rightsquigarrow C}$$

$$\frac{\alpha \text{ fresh} \quad \Gamma, x : \alpha^{\mathsf{m}} \vdash e : \sigma \rightsquigarrow C}{\Gamma \vdash \lambda x.\, e : \alpha^{\mathsf{m}} \rightarrow \sigma \rightsquigarrow C}\ \textsc{Abs} \qquad\qquad \frac{\Gamma, x : \phi \vdash e : \sigma \rightsquigarrow C}{\Gamma \vdash \lambda(x::\phi).\, e : \phi \rightarrow \sigma \rightsquigarrow C}\ \textsc{AnnAbs}$$

$$\frac{\Gamma \vdash^{\mathsf{h}} h : \phi \rightsquigarrow C \qquad \begin{array}{c}\Gamma \vdash e_i : \sigma_i \rightsquigarrow C_i \\ \overline{v}_i = \mathsf{fuv}(\sigma_i, C_i) - \mathsf{fuv}(\Gamma, \phi, C)\end{array} \qquad \overline{\alpha}, \overline{\beta}, \delta \text{ fresh}}{\begin{array}{l}\Gamma \vdash h\, e_1\, \ldots\, e_n : \delta_{\mathsf{t}} \rightsquigarrow C \quad \wedge \quad \phi \leqslant^n_{\mathsf{m}} \alpha_1^{\mathsf{u}} \rightarrow \beta_1^{\mathsf{u}} \quad \wedge \quad \mathbb{V}\{\overline{v}_1\}.\, C_1 \Rightarrow \sigma_1 \leq \alpha_1^{\mathsf{u}} \\ \qquad\qquad\qquad\qquad\quad \wedge \quad \beta_1^{\mathsf{u}} \leqslant^{n-1}_{\mathsf{m}} \alpha_2^{\mathsf{u}} \rightarrow \beta_2^{\mathsf{u}} \quad \wedge \quad \mathbb{V}\{\overline{v}_2\}.\, C_2 \Rightarrow \sigma_2 \leq \alpha_2^{\mathsf{u}} \\ \qquad\qquad\qquad\qquad\quad \wedge \qquad\qquad \vdots \qquad\qquad\quad \wedge \qquad\qquad \vdots \qquad\qquad\quad \wedge \quad \beta_n^{\mathsf{u}} \leqslant^0_{\mathsf{m}} \delta^{\mathsf{t}}\end{array}}\ \textsc{App}$$

$$\frac{\Gamma \vdash^{\mathsf{h}} h : \phi \rightsquigarrow C \qquad \begin{array}{c}\Gamma \vdash e_i : \sigma_i \rightsquigarrow C_i \\ \overline{v}' = \mathsf{fuv}(\phi, C) - \mathsf{fuv}(\Gamma) \qquad \overline{v_i} = \mathsf{fuv}(\sigma_i, C_i) - \mathsf{fuv}(\Gamma, \phi, C)\end{array} \qquad \overline{\alpha}, \overline{\beta} \text{ fresh}}{\begin{array}{l}\Gamma \vdash (h\, e_1\, \ldots\, e_n :: \forall \overline{b}.\, \eta) : \forall \overline{b}.\, \eta \\ \quad \rightsquigarrow \forall \overline{b}.\, \exists \overline{\alpha^{\mathsf{u}}}\, \overline{\beta^{\mathsf{u}}}\, \overline{v}'.\, \left( \begin{array}{llll} C & \wedge & \phi \leqslant^n_{\mathsf{u}} \alpha_1^{\mathsf{u}} \rightarrow \beta_1^{\mathsf{u}} & \wedge \quad \mathbb{V}\{v_1\}.\, C_1 \Rightarrow \sigma_1 \leq \alpha_1^{\mathsf{u}} \\ & \wedge & \beta_1^{\mathsf{u}} \leqslant^{n-1}_{\mathsf{u}} \alpha_2^{\mathsf{u}} \rightarrow \beta_2^{\mathsf{u}} & \wedge \quad \mathbb{V}\{v_2\}.\, C_2 \Rightarrow \sigma_2 \leq \alpha_2^{\mathsf{u}} \\ & \wedge & \vdots \qquad \wedge \qquad \vdots & \wedge \quad \beta_n^{\mathsf{u}} \leqslant^0_{\mathsf{u}} \eta \end{array} \right)\end{array}}\ \textsc{AnnApp}$$

$$\frac{\Gamma \vdash e_1 : \phi \rightsquigarrow C_1 \qquad \Gamma, x : \phi \vdash e_2 : \sigma \rightsquigarrow C_2}{\Gamma \vdash \mathbf{let}\ x = e_1\ \mathbf{in}\ e_2 : \sigma \rightsquigarrow C_1 \wedge C_2}\ \textsc{Let}$$

**Figure 8.** Constraint generation

$$\boxed{\mathsf{freshen}^n_{\mathsf{s}}(\sigma) \implies \langle \overline{v}, \mu \rangle}$$

$$\frac{\mu \rhd^n_{\mathsf{s}} \Delta \qquad \overline{\alpha} \text{ fresh} \qquad v_i = \alpha_i^{\Delta(a_i)}}{\mathsf{freshen}^n_{\mathsf{s}}(\forall \overline{a}.\, \mu) \implies \langle \overline{v}, [\overline{a \mapsto v}]\mu \rangle}$$

$$\boxed{C\ ;\ \overline{v} \implies C'\ ;\ \overline{v}'}$$

$$[\textsc{conj}]\ \frac{C_1\ ;\ \overline{v} \implies C_1'\ ;\ \overline{v}'}{C_1 \wedge C_2\ ;\ \overline{v} \implies C_1' \wedge C_2\ ;\ \overline{v}'} \qquad\qquad [\textsc{forall}]\ \frac{C\ ;\ \overline{v}_{\mathsf{in}} \implies C'\ ;\ \overline{v}'_{\mathsf{in}}}{\forall a.\, \exists \overline{v}_{\mathsf{in}}.\, C\ ;\ \overline{v} \implies \forall a.\, \exists \overline{v}'_{\mathsf{in}}.\, C'\ ;\ \overline{v}}$$

$$
\begin{array}{llll}
[\textsc{eqrefl}] & \sigma \sim \sigma\ ;\ \overline{v} & \implies & \top\ ;\ \overline{v} \\
[\textsc{eqmono}] & \mathsf{T}\, \sigma_1\, \ldots\, \sigma_n \sim \mathsf{T}\, \phi_1\, \ldots\, \phi_n\ ;\ \overline{v} & \implies & \sigma_1 \sim \phi_1, \ldots, \sigma_n \sim \phi_n\ ;\ \overline{v} \\
[\textsc{eqsubst}] & (\alpha^{\mathsf{s}} \sim \sigma) \wedge C\ ;\ \overline{v} & \implies & (\alpha^{\mathsf{s}} \sim \sigma) \wedge [\alpha^{\mathsf{s}} \mapsto \sigma]C\ ;\ \overline{v} & \text{if } \vdash \sigma : \mathsf{s}, \text{ and } \alpha \notin \mathsf{ftv}(\sigma) \\
[\textsc{eqvar}] & \alpha^{\mathsf{s}_1} \sim \beta^{\mathsf{s}_2}\ ;\ \overline{v} & \implies & \beta^{\mathsf{s}_2} \sim \alpha^{\mathsf{s}_1}\ ;\ \overline{v} & \text{if } \mathsf{s}_1 \sqsubset \mathsf{s}_2 \\
[\textsc{eqfully}] & \alpha^{\mathsf{m}} \sim \sigma\ ;\ \overline{v} & \implies & \{\beta^{\mathsf{s}} \sim \gamma^{\mathsf{m}} \mid \beta^{\mathsf{s}} \in \mathsf{fuv}(\sigma), \mathsf{s} \neq \mathsf{m}\}\ ;\ \overline{v}, \overline{\gamma^{\mathsf{m}}} \\
[\textsc{instmono}] & \mu \leqslant^n_{\mathsf{s}} \eta\ ;\ \overline{v} & \implies & \mu \sim \eta\ ;\ \overline{v} \\
[\textsc{inst}\forall\textsc{l}] & (\forall \overline{a}.\, \mu) \leqslant^n_{\mathsf{s}} \eta\ ;\ \overline{v} & \implies & \mu' \sim \eta\ ;\ \overline{v}, \overline{v}' & \text{where } \mathsf{freshen}^n_{\mathsf{s}}(\forall \overline{a}.\, \mu) \implies \langle \overline{v}', \mu' \rangle \\
[\textsc{inst}\forall\textsc{l}] & (\mathbb{V}\{\overline{v}'\}.\, \overline{C} \Rightarrow \sigma) \leq \eta\ ;\ \overline{v} & \implies & \overline{C} \wedge (\sigma \leqslant^0_{\mathsf{m}} \eta)\ ;\ \overline{v}, \overline{v}' \\
[\textsc{inst}\forall\textsc{r}] & \mathfrak{g} \leq (\forall \overline{a}.\, \mu)\ ;\ \overline{v} & \implies & \forall \overline{a}.\, (\mathfrak{g} \leq \mu)\ ;\ \overline{v}
\end{array}
$$

**Figure 9.** Solving rules

this is not problematic, since on application type variables are instantiated and regeneralized using the $\preceq$ relation.

### 4.3.2 Instantiation and generalisation constraints

For instantiation constraints $\leqslant_s^n$, Rule INSTMONO encodes the fact that in our system if two top-level monomorphic types $\mu, \eta$ are in an instance relation, they must be equal. This is a consequence of the invariance of type constructors. INST∀L instantiates a polytype $\sigma$ with fresh unification variables, much as in the usual Damas-Milner algorithm, except that we must use a sort-respecting instantiation. This is done by $\mathrm{freshen}_s^n$, which in turn uses the already-introduced classification judgment $\rhd_s^n$ (Figure 4). Finally, the new variables enter the set of existentially quantified variables. Notice that the right hand side of an instantiation constraint $\sigma \leqslant_s^n \mu$ is always a top-level monomorphic type $\mu$, so we do not need to worry about having a polymorphic type on the right.

Finally, we come to generalisation constraints, which (recall Section 4.1) express a deferred generalisation decision. Rule INST∀L is simple: if the right hand side has no top-level foralls (it is of the form $\eta$) then there is no generalisation to be done, so it suffices to release all the captured constraints $C$ and existentials $v'$ into the current constraint.

INST∀R is where actual generalisation takes place. In order to forge some intuition, let us look at the constraint

$$(\mathbb{V}\{\alpha\,\beta\}.\,\alpha \leqslant_m^0 \beta \Rightarrow \alpha \to \beta) \preceq (\forall p.\,p \to p)$$

This generalisation constraint says that by performing some solving and possibly abstracting over some of the variables $\alpha$ and $\beta$, we should get the polymorphic type $\forall p.\,p \to p$. Following standard practice, we skolemise the type on the right, introducing a fresh skolem or rigid variable $p$, which should not be unified.

$$\alpha \leqslant_m^0 \beta \quad \wedge \quad \alpha \to \beta \sim p \to p$$

We obtain a solution by making $\alpha \sim \beta \sim a$. But in order for this solution to remain valid, we must guarantee that the skolem $p$ *does not escape* to the outer world. We recall this restriction by means of a fresh quantification constraint $\forall p.\exists \alpha\,\beta.(\alpha \leqslant_m^0 \beta \quad \wedge \quad \alpha \to \beta \sim a \to a)$. Rule INST∀R achieves this rather neatly simply by doing skolemisation and pushing the g inside; then INST∀L will do the rest.

The rules applicable to instantiation and generalisation constraints do not handle every case. In particular, whenever an unrestricted variable appears in one of the sides of the constraint, there are good reasons to wait:

1. If we have $\alpha^u \leqslant_s^n \mu$ we cannot turn it directly into $\alpha^u \sim \mu$, because $\alpha^u$ might be unified later to a polymorphic type and we need instantiation.
2. Similarly, if we have $(\mathbb{V}\{\overline{\alpha}\}.\,C \Rightarrow \mu) \preceq \alpha^u$, and $\alpha^u$ is later substituted by a polytype, we must skolemise.

The guardedness restrictions are carefully crafted to ensure that the solver is never completely *stuck*, unless the constraint set as a whole is inconsistent. A single constraint can be stuck for some time, but if we are dealing with a consistent constraint set, the promise is that it will, by solving steps applied to other constraints in the set, become unstuck.

**Theorem 4.1.** *Suppose* $\Gamma \vdash e : \sigma \rightsquigarrow C$. *Then* $C$ *is either inconsistent, or can be rewritten to a new set* $C'$ *without instantiation and generalisation constraints which fixes the value of all unrestricted and top-level monomorphic variables.*

The resulting constraint set is an instance of the problem of first-order unification under a mixed prefix. [3] (which in our system is expressed by quantification constraints). A complete algorithm to solve this kind of problem is described by Pottier and Rémy [15]. As a consequence, we have:

**Corollary 4.2.** *Suppose* $\Gamma \vdash e : \sigma \rightsquigarrow C$. *Then* $C$ *is either inconsistent, or can be rewritten to a solved form.*

Unfortunately, once we extend the language of types, by introducing type classes and local assumptions, completeness is known not to hold [21]. Furthermore, the approach by Pottier and Rémy [15] is no longer applicable. With that in mind, we introduce a *different* approach to solve the problem of unification under a mixed prefix, which *does* scale to handle local assumptions. The approach is rather simple – a single rule FLOAT in Figure 10 – and is directly inspired by how GHC handles constraints: by floating constraints $F$ from inside a quantification constraint to outside. When can we do that? Precisely when the constraint does not mention the skolems. But what about the existentials? For example, suppose we have

$$\exists \alpha. \ldots. (\forall a.\exists \beta.(\alpha \sim [\beta]) \wedge C) \ldots$$

We would like to float the constraint $(\alpha \sim [\beta])$ out of the quantification constraint, but then $\beta$ would be out of scope. We can solve this by "promoting" $\beta$: producing a fresh $\beta'$ that lives in the outer scope, and making $\beta$ equal to it, thus:

$$\exists \alpha, \beta'. \ldots. (\alpha \sim [\beta'])) \wedge (\forall a.\exists \beta.(\beta \sim \beta' \wedge C)) \ldots$$

All this is expressed directly by rule FLOAT. If we cannot float, we have a skolem escape error; for example, consider:

$$\exists \alpha. \ldots. (\forall a.\exists \beta.(\alpha \sim [a]) \wedge C) \ldots$$

Here we cannot float $(\alpha \sim [a])$ because it mentions the skolem $a$, so an inner skolem has leaked into an outer scope ($\alpha$ is bound further out). Floating makes manifest that skolem escape has not happened, and brings the constraint nearer to *solved form*, which we treat next.

**Conjecture 4.3.** *The solver presented in Figure 9 is complete for unification problems under a mixed prefix.*

We stress that for constraints without local assumptions we can either choose the algorithm in Pottier and Rémy [15] – which is known to be complete and thus allows us to prove

$$[\textsc{float}] \; \frac{\mathsf{ftv}(F) \cap \overline{a} = \emptyset \qquad \overline{\alpha^{\mathfrak{s}}} = \mathsf{fuv}(F) \cap \upsilon_{\mathsf{in}} \qquad \overline{\gamma^{\mathfrak{s}}} \text{ fresh} \qquad \overline{E} = \bigwedge_{\alpha^{\mathfrak{s}} \in \overline{\alpha}} \alpha^{\mathfrak{s}} \sim \gamma^{\mathfrak{s}}}{\forall \overline{a}. \exists \overline{\upsilon}_{\mathsf{in}}. (C \wedge F) \, ; \, \overline{\upsilon} \implies [\overline{\alpha^{\mathfrak{s}} \mapsto \gamma^{\mathfrak{s}}}]F \wedge \forall \overline{a}. \exists \overline{\upsilon}_{\mathsf{in}}. (C \wedge E) \, ; \, \overline{\upsilon}, \overline{\gamma^{\mathfrak{s}}}}$$

**Figure 10.** Solving rule for quantification constraints

$$\frac{\vdash \sigma : \mathfrak{s} \qquad \mathsf{ftv}(\sigma) \subseteq \overline{a} \cup \overline{\alpha}}{\overline{a} \, ; \, \overline{\alpha} \, ; \, \{\beta\} \vdash \beta^{\mathfrak{s}} \sim \sigma \text{ solved}} \; \textsc{SolvedVar}$$

$$\frac{\overline{a} \, ; \, \overline{\alpha} \, ; \, \overline{\beta}_1 \vdash C_1 \text{ solved} \qquad \overline{a} \, ; \, \overline{\alpha} \, ; \, \overline{\beta}_2 \vdash C_2 \text{ solved}}{\overline{a} \, ; \, \overline{\alpha} \, ; \, \overline{\beta}_1 \uplus \overline{\beta}_2 \vdash C_1 \wedge C_2 \text{ solved}} \; \textsc{SolvedConj}$$

$$\frac{\overline{\upsilon} = \overline{\gamma}_1 \uplus \overline{\gamma}_2 \qquad \overline{a} \cup \overline{b} \, ; \, \overline{\alpha} \cup \overline{\gamma}_1 \, ; \, \overline{\gamma}_2 \vdash C \text{ solved}}{\overline{a} \, ; \, \overline{\alpha} \, ; \, \emptyset \vdash \forall \overline{b}. \exists \overline{\upsilon}.C \text{ solved}} \; \textsc{SolvedQuant}$$

**Figure 11.** Definition of solved set of constraints

a completeness theorem (Theorem 4.6) – or the one with the float rule – which we only conjecture as complete but scales when the constraint language is extended.

### 4.3.3 Solved form

A constraint is in *solved form* if it consists only of quantification and equality constraints ($\upsilon \sim \sigma$); and the equalities constitute a well-sorted idempotent substitution of its unification variables. For example

$$\exists \alpha^{\mathfrak{m}}.(\alpha^{\mathfrak{m}} \sim \mathsf{Int}) \wedge (\forall b. \exists \beta^{\mathfrak{u}}. \beta^{\mathfrak{u}} \sim (b \rightarrow \mathsf{Int}))$$

is in solved form. Being in solved form is more than just syntactic; here are two constraints that are not:

$$\exists \alpha.(\alpha \sim \mathsf{Int}) \wedge (\alpha \sim \mathsf{Bool}) \qquad \exists \alpha. \ldots (\forall b. \alpha \sim [b]) \ldots$$

In the first there are two equalities for $\alpha$ (we should apply eqsubst to make progress); in the second, there is a skolem-escape problem. However it is OK for a unification variable to have *no* equalities; it is simply unconstrained.

Figure 11 defines solved form precisely. We keep a set of variables $\overline{\beta}$ for which we ensure that there is precisely one equality constraint, and another set $\overline{\alpha}$ (the unconstrained variables) for which there are none. Rule SolvedVar expects precisely one $\beta$, checks well-sortedness, and also checks that $\sigma$ does not mention any variables other than the skolems and unconstrained unification variables – the latter check ensures idempotence.

Rule SolvedConj partitions the $\overline{\beta}$ between the two conjunctions. Rule SolvedQuant partitions the local existentials $\overline{\upsilon}$ into the unconstrained sets, $\gamma_1$ and $\gamma_2$ resp.

### 4.4 Soundness, principality and completeness

The inference algorithm presented here satisfies the usual properties of soundness, principality, and completeness with respect to the declarative specification. In order to state these

results, we need the auxiliary notion of substitution induced by a solved form.

$$\frac{C_{\mathsf{s}} = E \wedge R, \text{ where } E \text{ are all the equalities in } C_s}{\widehat{C}_{\mathsf{s}} = [\alpha \mapsto \sigma \mid \alpha \sim \sigma \in E]}$$

The proofs are given in Appendix E.2.

**Theorem 4.4** (Soundness). *Let $\Gamma$ be a closed environment and $e$ an expression. If $\Gamma \vdash e : \sigma \rightsquigarrow C$ and $C_{\mathsf{s}}$ is a solution for $C$ with an induced substitution $\widehat{C}_{\mathsf{s}}$, then we have $\Gamma \vdash e : \widehat{C}_{\mathsf{s}}(\sigma)$.*

**Theorem 4.5** (Principality). *Suppose $\Gamma \vdash e : \sigma$. Then there exists a type $\sigma^{\star}$ such that $\Gamma \vdash e : \sigma^{\star}$ and $\sigma = \pi\sigma^{\star}$ where $\pi$ is a fully monomorphic substitution.*

**Theorem 4.6** (Completeness). *Let $\Gamma$ be a closed environment and $e$ an expression. If $\Gamma \vdash e : \sigma$ then $\Gamma \vdash e : \phi \rightsquigarrow C$ and $C$ can reach a solved form.*

## 5 Practical matters

***Prototype.*** We have implemented a prototype of the type inference process described in this paper, including support for Haskell's type classes by extending GI as described in Appendix B. The expressions in Figure 2 are accepted or rejected as described by the table, in the GI decl. column.

***Backward compatibility.*** GI does not support any co- or contra-variance in function types, nor deep skolemization; but GHC does, at least with the *RankNTypes* extension. One might worry that many existing Haskell libraries would need to be modified. To quantify this impact, we modified GHC to impose those restrictions and rebuilt all the packages in Stackage which require the *RankNTypes* extension. In order to minimize the annotation burden, we added a simple special case that supports function definition with a forall to the right of an arrow:

$$f :: \forall a. \, a \rightarrow (\forall b. \, b \rightarrow b)$$
$$f \; x \; y = y$$

Such definitions are very common in libraries such as Scrap Your Boilerplate. With this done, very few packages required modifications, and modifications were always $\eta$-expansions. Consider (*flip f*), where *flip*'s type is in Figure 1. This is ill-typed because *flip* requires an argument of type $a \rightarrow b \rightarrow c$, but $f$'s type, after instantiation, looks like $\tau \rightarrow \forall b. \, b \rightarrow b$. The fix is simple: just $\eta$-expand the argument, thus (*flip* ($\lambda x \rightarrow f \; x$)). GHC does this automatically at the moment, but in fact this $\eta$-expansion is unsound in general,

because of *seq*, so requiring manual $\eta$-expansion is probably the right choice anyway.

Of the 2,400 packages in Stackage, 609 use *RankNTypes*; of these, only 75 required manual changes, all of which were simple $\eta$-expansions. One (*singletons*) would require larger changes, because it uses Template Haskell to *generate* Haskell code; so it needs to generate $\eta$-expanded code. Two more failed for reasons we have yet to investigate. Our conclusion is that the impact of our proposed changes is extremely minor, especially since GHC's current covert $\eta$-expansion strategy is unsound in the first place.

***Relaxed inference.*** Section 2.2 describes some limitations of our system, like the impossibility of inferring a polymorphic element type for the empty list [ ]. There are some scenarios, though, where one can relax the guardedness restrictions and still obtain a correct solution for the generated constraints. We describe that *relaxed* solver in Appendix C. There is a net gain on the amount of programs accepted by this approach, as Figure 2 shows under the GI relax. heading.

## 6 Related work

Full type inference for System F is undecidable [24] – partial type inference with known generalisation positions but unknown instantiations can be reduced to h.o. unification [14]. System F lacks principal types, making modular type inference and the addition of ML-style let-bindings impossible. Higher-rank type inference with instantiation restricted to *monomorphic types* has some successful solutions [4, 13, 17], exploiting a mix of annotation propagation and unification.

On the other hand, no solution for impredicativity with a good benefit-to-weight ratio has been presented to-date. MLF [1, 2, 18] is an *extension* of System F based on quantification with instance bounds. The resulting system is powerful, but also quite complex to implement; in return we get back principal types. There have been several attempts to simplify the user-facing part of MLF to System F types. FPH [23] exposes a "box" structure around inferred types (that would be hidden under a constraint in MLF). Flexible types [9] avoid quantification over equality constraints. Implementing these systems in a working compiler is a significant undertaking, and so is the integration with features like type classes[10].

For this reason, there are proposals for algorithms simpler than MLF. Boxy Types [22] is an early attempt to push bidirectional type inference to allow impredicativity, but resulted in a complex specification. HMF [8] imposes universal conditions on typing derivations to recover principal types. QML [19], inspired by boxed polymorphism [12, 16], introduces an only-explicitly-instantiable $\forall$. Recent work [5] also proposes a distinction between an implicitly and an explicitly instantiable $\forall$; only the latter is impredicative.

We return now to Figure 2, where we present a collection of examples appearing in selected related works on type inference for impredicativity, namely MLF [1], HMF [8], FPH [23], HML [9]. We have selected those systems to compare because they strike extremely well in expressivity and require very few type annotations. The table shows the flexibility/expressivity price we pay in order to keep the implementation and specification costs low, and avoid the introduction of new type system features (such as types with constraints or boxes) or new forms of annotations (such as QML-style instantiation annotations). We also show the annotation that can recover the typeability of a program in cases where a valid type exists. As one observes, MLF can type all of programs that do not require implicit $\eta$-expansion (*k h lst*) or the use of a polymorphic function argument at two types $\lambda f . (f\ 1, f\ True)$. Unsurprisingly, HML is only minimally less powerful as it cannot type $\lambda xs.\ poly\ (head\ xs)$ because *xs* would have to be assigned a polymorphic type, even though it's only used at this type. FPH is equally expressive for the applicative fragment of System F – however its treatment of $\lambda$-abstractions requires the returned type to be a fully-resolved top-level monomorphic type and hence fails to type check *choose id auto* because *auto* will be assigned the same type as GI infers $((\forall a.\ a \rightarrow a) \rightarrow b \rightarrow b)$. HMF on the other hand is just based on local decisions about polymorphic instantiations and – without an extension to *n*-ary applications – fails to type check programs where the local instantiation has to be delayed to take more arguments into account (e.g. fails to type check *id : ids*). The relaxed GI system on the other hand can – modulo the quirk about not generalizing the bodies of lambda abstractions that only MLF and HML can tackle – type check the same programs as those systems. GI (vanilla) is more restrictive than those systems but incomparable to HMF. To determine why a program fails to type check (and how it should be fixed) it suffices to determine whether some function has been instantiated to a type with top-level polymorphism in its arguments. For example, $f\ (choose\ id)\ ids$ fails to type check because it requires the instantiation of $choose :: \forall a.\ a \rightarrow a \rightarrow a$ to $(\forall a.\ a \rightarrow a) \rightarrow (\forall a.\ a \rightarrow a) \rightarrow (\forall a.\ a \rightarrow a)$ and none of the arguments have a type with a top-level constructor. The fix is to add an annotation around the offending expression to determine its type: $f\ (choose\ id :: (\forall a.\ a \rightarrow a) \rightarrow (\forall a.\ a \rightarrow a))\ ids$. The relaxed GI recovers pretty much of the lost expressivity, and is by construction as expressive as the vanilla GI, but we have yet to determine a simple declarative specification.

## 7 Further work

GI is relatively simple and predictable but, as mentioned earlier, we would like it to be just a bit more expressive. For example, the need to type-annotate occurrences of the empty list (Section 2.2) seems particularly tiresome – and it is hard to see why it should truly be necessary. Thus motivated, we are investigating some modest extensions, both of the declarative system and of the constraint solver, that would accept more programs. The swamp is calling!

# References

[1] Didier Le Botlan and Didier Rémy. 2003. ML$^F$: raising ML to the power of system F. In *Proceedings of the Eighth ACM SIGPLAN International Conference on Functional Programming, ICFP 2003, Uppsala, Sweden, August 25-29, 2003*, Colin Runciman and Olin Shivers (Eds.). ACM, 27–38. https://doi.org/10.1145/944705.944709

[2] Didier Le Botlan and Didier Rémy. 2009. Recasting MLF. *Inf. Comput.* 207, 6 (2009), 726–785.

[3] Hubert Comon and Pierre Lescanne. 1988. *Equational problems and disunification*. Research Report RR-0904. INRIA. https://hal.inria.fr/inria-00075652

[4] Joshua Dunfield and Neelakantan R. Krishnaswami. 2013. Complete and easy bidirectional typechecking for higher-rank polymorphism. In *ACM SIGPLAN International Conference on Functional Programming, ICFP'13, Boston, MA, USA - September 25 - 27, 2013*, Greg Morrisett and Tarmo Uustalu (Eds.). ACM, 429–442. https://doi.org/10.1145/2500365.2500582

[5] Richard A. Eisenberg, Stephanie Weirich, and Hamidhasan G. Ahmed. 2016. Visible Type Application. In *Proceedings of the 25th European Symposium on Programming Languages and Systems - Volume 9632*. Springer-Verlag New York, Inc., New York, NY, USA, 229–254. https://doi.org/10.1007/978-3-662-49498-1_10

[6] Bastiaan Heeren, Jurriaan Hage, and S. Doaitse Swierstra. 2003. Scripting the type inference process. *SIGPLAN Notices* 38, 9 (2003), 3–13. https://doi.org/10.1145/944746.944707

[7] James Hook and Peter Thiemann (Eds.). 2008. *Proceeding of the 13th ACM SIGPLAN international conference on Functional programming, ICFP 2008, Victoria, BC, Canada, September 20-28, 2008*. ACM.

[8] Daan Leijen. 2008. HMF: simple type inference for first-class polymorphism, See [7], 283–294. https://doi.org/10.1145/1411204.1411245

[9] Daan Leijen. 2009. Flexible types: robust type inference for first-class polymorphism. In *Proceedings of the 36th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2009, Savannah, GA, USA, January 21-23, 2009*, Zhong Shao and Benjamin C. Pierce (Eds.). ACM, 66–77. https://doi.org/10.1145/1480881.1480891

[10] Daan Leijen and Andres Löh. 2005. Qualified types for MLF. In *Proceedings of the 10th ACM SIGPLAN International Conference on Functional Programming, ICFP 2005, Tallinn, Estonia, September 26-28, 2005*, Olivier Danvy and Benjamin C. Pierce (Eds.). ACM, 144–155. https://doi.org/10.1145/1086365.1086385

[11] Alberto Martelli and Ugo Montanari. 1982. An Efficient Unification Algorithm. *ACM Trans. Program. Lang. Syst.* 4, 2 (1982), 258–282. https://doi.org/10.1145/357162.357169

[12] J. W. O'Toole, Jr. and D. K. Gifford. 1989. Type Reconstruction with First-class Polymorphic Values. In *Proceedings of the ACM SIGPLAN 1989 Conference on Programming Language Design and Implementation (PLDI '89)*. ACM, New York, NY, USA, 207–217. https://doi.org/10.1145/73141.74836

[13] Simon L. Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Mark Shields. 2007. Practical type inference for arbitrary-rank types. *Journal of Functional Programming* 17, 1 (2007), 1–82.

[14] Frank Pfenning. 1988. Partial Polymorphic Type Inference and Higher-order Unification. In *Proceedings of the 1988 ACM Conference on LISP and Functional Programming (LFP '88)*. ACM, New York, NY, USA, 153–163. https://doi.org/10.1145/62678.62697

[15] François Pottier and Didier Rémy. 2005. The Essence of ML Type Inference. In *Advanced Topics in Types and Programming Languages*, Benjamin C. Pierce (Ed.). MIT Press, Chapter 10, 389–489. http://cristal.inria.fr/attapl/

[16] Didier Rémy. 1994. Programming Objects with ML-ART, an Extension to ML with Abstract and Record Types. In *Proceedings of the International Conference on Theoretical Aspects of Computer Software (TACS '94)*. Springer-Verlag, London, UK, UK, 321–346. http://dl.acm.org/citation.cfm?id=645868.668492

[17] Didier Rémy. 2005. Simple, Partial Type-inference for System F Based on Type-containment. In *Proceedings of the Tenth ACM SIGPLAN International Conference on Functional Programming (ICFP '05)*. ACM, New York, NY, USA, 130–143. https://doi.org/10.1145/1086365.1086383

[18] Didier Rémy and Boris Yakobowski. 2008. From ML to MLF: graphic type constraints with efficient type inference, See [7], 63–74. https://doi.org/10.1145/1411204.1411216

[19] Claudio V. Russo and Dimitrios Vytiniotis. 2009. QML: Explicit First-class Polymorphism for ML. In *Proceedings of the 2009 ACM SIGPLAN Workshop on ML (ML '09)*. ACM, New York, NY, USA, 3–14. https://doi.org/10.1145/1596627.1596630

[20] Alejandro Serrano and Jurriaan Hage. 2016. Type Error Diagnosis for Embedded DSLs by Two-Stage Specialized Type Rules. In *Programming Languages and Systems - 25th European Symposium on Programming, ESOP 2016, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2016, Eindhoven, The Netherlands, April 2-8, 2016, Proceedings (Lecture Notes in Computer Science)*, Peter Thiemann (Ed.), Vol. 9632. Springer, 672–698. https://doi.org/10.1007/978-3-662-49498-1_26

[21] Dimitrios Vytiniotis, Simon L. Peyton Jones, Tom Schrijvers, and Martin Sulzmann. 2011. OUTSIDEIN(X): Modular type inference with local assumptions. *J. Funct. Program.* 21, 4-5 (2011), 333–412. https://doi.org/10.1017/S0956796811000098

[22] Dimitrios Vytiniotis, Stephanie Weirich, and Simon L. Peyton Jones. 2006. Boxy types: inference for higher-rank types and impredicativity. In *Proceedings of the 11th ACM SIGPLAN International Conference on Functional Programming, ICFP 2006, Portland, Oregon, USA, September 16-21, 2006*, John H. Reppy and Julia L. Lawall (Eds.). ACM, 251–262. https://doi.org/10.1145/1159803.1159838

[23] Dimitrios Vytiniotis, Stephanie Weirich, and Simon L. Peyton Jones. 2008. FPH: first-class polymorphism for Haskell, See [7], 295–306. https://doi.org/10.1145/1411204.1411246

[24] J. B. Wells. 1993. *Typability and Type Checking in the Second-Order Lambda-Calculus Are Equivalent and Undecidable*. Technical Report. Boston, MA, USA.

[25] Danfeng Zhang, Andrew C. Myers, Dimitrios Vytiniotis, and Simon L. Peyton Jones. 2015. Diagnosing type errors with class. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation, Portland, OR, USA, June 15-17, 2015*, David Grove and Steve Blackburn (Eds.). ACM, 12–21. https://doi.org/10.1145/2737924.2738009