

Guarded Impredicative Polymorphism

Extended Version

Alejandro Serrano
Jurriaan Hage
A.SerranoMena@uu.nl
J.Hage@uu.nl
Utrecht University
Utrecht, The Netherlands

Dimitrios Vytiniotis
Simon Peyton Jones
dimitris@microsoft.com
simonpj@microsoft.com
Microsoft Research
Cambridge, United Kingdom

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.

CCS Concepts • Theory of computation → Type structures;

Keywords Type systems, impredicative polymorphism, constraint-based inference

This is an extended version of Alejandro Serrano, Jurriaan Hage, Dimitrios Vytiniotis, and Simon Peyton Jones. 2018. Guarded Impredicative Polymorphism. In PLDI '18: ACM SIGPLAN Conference on Programming Language Design and Implementation, June 18 – 22, 2018, Philadelphia, PA, USA.

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 3.4), `let` bindings (Section 3.5), 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 3.6). 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.
- 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.

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! In *head*'s type the variable p appears under a list type constructor. Given the type of *ids*, $\forall a. a \rightarrow a$, it is plain as a pikestaff that we must instantiate p with $\forall a. a \rightarrow a$. The difficulty comes when we have a “naked” or “un-guarded” type variable in one of the arguments, such as p in *single* $:: \forall p. p \rightarrow [p]$. Now if we examine (*single id*), it is not clear whether we should instantiate p with $\forall a. a \rightarrow a$, or with $Int \rightarrow Int$, or some other monomorphic type.

In fact, (*single id*) does not have a most general type. It has both of these two incomparable types $\forall a. [a \rightarrow a]$ and $[\forall a. a \rightarrow a]$. 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 ($f e_1 \dots e_n$). 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. *Given a n -argument call $f e_1 \dots e_n$ to a function $f :: \forall a_1 \dots a_p. \sigma_1 \rightarrow \dots \rightarrow \sigma_m \rightarrow \phi$, where $m \leq n$, a type variable a_i may be instantiated to:*

1. A polymorphic type ϕ , if a_i appears under a type constructor in one of the σ_j (note that we only take in consideration as many types as arguments given). In this case we say that a_i is guarded in the type of f .

2. A top-level monomorphic type μ (see Figure 3), if a_i appears in any of the σ_j at all.
3. A fully monomorphic type τ , otherwise.

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

- (*map poly*) is OK because in the type of *map*: $\forall a b. (a \rightarrow b) \rightarrow [a] \rightarrow [b]$, both type variables appear under the type constructor (\rightarrow) in the first argument. We can instantiate both to $(\forall a. a \rightarrow a)$, by case (1) of the Rule, as required to match the type of *poly*.
- (*single ids*) is OK because we can instantiate *single* $:: \forall a. a \rightarrow [a]$ with the top-level-monomorphic type $\forall b. [b]$, using case (2) of the Instantiation Rule.
- (*(:) id*) is a *partial application* of $(:)$. So although a appears guarded in the second argument of $(:)$ $:: \forall a. a \rightarrow [a] \rightarrow [a]$, in this call we can only take advantage of the first argument (see $m \leq n$ in the rule.). Hence a can only be instantiated by a top-level monomorphic type, by clause (2) of the rule. If we add a second argument in the call, such as $(:)$ *id ids*, instantiation may be polymorphic since the second argument is now taken into consideration.
- (*id poly* $(\lambda x. x)$) is a tricky one. Here *id* is applied to two arguments although its type, $\forall a. a \rightarrow a$, 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: the type of *poly* is $(\forall a. a \rightarrow a) \rightarrow (Int, Bool)$, a top-level monomorphic type. Thus the instantiation of *id* to that type is allowed by clause (2) of the Instantiation Rule.

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 Instantiation Rule *takes no account of the context of the call*. For example, consider *ids # single id*. We know that the result of (*single id*) must be $[\forall a. a \rightarrow a]$, given that *ids* has the same type, 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 GI, the types for *head* and *single ids* are instantiated independently, so we never need to consider the *interaction* between the 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.

| | | |
|--|---|--|
| $head :: \forall p. [p] \rightarrow p$ | $id :: \forall a. a \rightarrow a$ | $ids :: [\forall a. a \rightarrow a]$ |
| $tail :: \forall p. [p] \rightarrow [p]$ | $inc :: Int \rightarrow Int$ | $map :: \forall p q. (p \rightarrow q) \rightarrow [p] \rightarrow [q]$ |
| $[] :: \forall p. [p]$ | $choose :: \forall a. a \rightarrow a \rightarrow a$ | $app :: \forall a b. (a \rightarrow b) \rightarrow a \rightarrow b$ |
| $(:) :: \forall p. p \rightarrow [p] \rightarrow [p]$ | $poly :: (\forall a. a \rightarrow a) \rightarrow (Int, Bool)$ | $revapp :: \forall a b. a \rightarrow (a \rightarrow b) \rightarrow b$ |
| $single :: \forall p. p \rightarrow [p]$ | $auto :: (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$ | $flip :: \forall a b c. (a \rightarrow b \rightarrow c) \rightarrow b \rightarrow a \rightarrow c$ |
| $(\#) :: \forall p. [p] \rightarrow [p] \rightarrow [p]$ | $auto' :: (\forall a. a \rightarrow a) \rightarrow b \rightarrow b$ | $runST :: \forall v. (\forall s. ST s v) \rightarrow v$ |
| $length :: \forall p. [p] \rightarrow Int$ | | $argST :: \forall s. ST s Int$ |

Figure 1. Type signatures for functions used in the text

Thus, the empty list constructor $[] :: \forall a. [a]$ cannot be assigned a type $[\forall a. a \rightarrow a]$. We describe how to loosen this restriction in Section 3.3.

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 $g :: ((\forall a. a \rightarrow a) \rightarrow (Char, Bool)) \rightarrow Int$. 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. In common with many other systems, 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. Otherwise, the bound variable can only have a fully monomorphic type.*

By a “fully monomorphic type” we mean “no forall anywhere”. Nothing about guardedness here! In contrast, MLF requires an annotation only when the argument is used more than once in the body with different polymorphic types.

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 $\lambda(x :: \forall a. a \rightarrow a). x x$, GI infers the type $(\forall a. a \rightarrow a) \rightarrow b \rightarrow b$, and not $(\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$. To get the latter type, we have to write $\lambda(x :: \forall a. a \rightarrow a). (x x :: \forall a. a \rightarrow a)$ 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 $(runST \$ do \{ \dots \})$ (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 in GI, 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, GI does not accept the application $(g\ xs)$, 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 unpredictable 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 language 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

| | GI | MLF [1] | HMF [8] | FPH [23] | HML [9] |
|---|---|------------|------------|-------------|------------|
| A POLYMORPHIC INSTANTIATION | | | | | |
| A1 | $const2 = \lambda x y. y$ | ✓ | ✓ | ✓ | ✓ |
| | MLF infers $(b \geq \forall c. c \rightarrow c) \Rightarrow a \rightarrow b$, GI infers $a \rightarrow b \rightarrow b$. | | | | |
| A2 | $choose\ id$ | ✓ | ✓ | ✓ | ✓ |
| | MLF and HML infer $(a \geq \forall b. b \rightarrow b) \Rightarrow a \rightarrow a$, FPH, HMF, and GI infer $(a \rightarrow a) \rightarrow a \rightarrow a$. | | | | |
| A3 | $choose\ []\ ids$ | ✓ | ✓ | ✓ | ✓ |
| A4 | $\lambda(x :: \forall a. a \rightarrow a). x\ x$ | ✓ | ✓ | ✓ | ✓ |
| | MLF infers $(\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$, GI infers $(\forall a. a \rightarrow a) \rightarrow b \rightarrow b$. | | | | |
| A5 | $id\ auto$ | ✓ | ✓ | ✓ | ✓ |
| A6 | $id\ auto'$ | ✓ | ✓ | ✓ | ✓ |
| A7 | $choose\ id\ auto$ | ✓ | ✓ | No | No |
| A8 | $choose\ id\ auto'$ | No | ✓ | No | No |
| A9 | $f\ (choose\ id)\ ids$ | No | ✓ | No | ✓ |
| | where $f :: \forall a. (a \rightarrow a) \rightarrow [a] \rightarrow a$ GI needs an annotation on $id :: (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$ in the previous two examples. | | | | |
| A10 | $poly\ id$ | ✓ | ✓ | ✓ | ✓ |
| A11 | $poly\ (\lambda x. x)$ | ✓ | ✓ | ✓ | ✓ |
| A12 | $id\ poly\ (\lambda x. x)$ | ✓ | ✓ | ✓ | ✓ |
| B INFERENCE OF POLYMORPHIC ARGUMENTS | | | | | |
| B1 | $\lambda f. (f\ 1, f\ True)$ | No | No | No | No |
| | All systems require an annotation on $f :: \forall a. a \rightarrow a$. | | | | |
| B2 | $\lambda xs. poly\ (head\ xs)$ | No | ✓ | No | No |
| | All systems except for MLF require annotated $xs :: [\forall a. a \rightarrow a]$. | | | | |
| C FUNCTIONS ON POLYMORPHIC LISTS | | | | | |
| C1 | $length\ ids$ | ✓ | ✓ | ✓ | ✓ |
| C2 | $tail\ ids$ | ✓ | ✓ | ✓ | ✓ |
| C3 | $head\ ids$ | ✓ | ✓ | ✓ | ✓ |
| C4 | $single\ id$ | ✓ | ✓ | ✓ | ✓ |
| C5 | $id : ids$ | ✓ | ✓ | No | ✓ |
| C6 | $(\lambda x. x) : ids$ | ✓ | ✓ | No | ✓ |
| C7 | $single\ inc \# single\ id$ | ✓ | ✓ | ✓ | ✓ |
| C8 | $g\ (single\ id)\ ids$ | No | ✓ | No | ✓ |
| | where $g :: \forall a. [a] \rightarrow [a] \rightarrow a$ | | | | |
| C9 | $map\ poly\ (single\ id)$ | No | ✓ | ✓ | ✓ |
| | GI needs an annotation on $single\ id :: [\forall a. a \rightarrow a]$ in the previous two examples. | | | | |
| C10 | $map\ head\ (single\ ids)$ | ✓ | ✓ | ✓ | ✓ |
| D APPLICATION FUNCTIONS | | | | | |
| D1 | $app\ poly\ id$ | ✓ | ✓ | ✓ | ✓ |
| D2 | $revapp\ id\ poly$ | ✓ | ✓ | ✓ | ✓ |
| D3 | $runST\ argST$ | ✓ | ✓ | ✓ | ✓ |
| D4 | $app\ runST\ argST$ | ✓ | ✓ | ✓ | ✓ |
| D5 | $revapp\ argST\ runST$ | ✓ | ✓ | ✓ | ✓ |
| E η-EXPANSION | | | | | |
| E1 | $k\ h\ lst$ | No | No | No | No |
| E2 | $k\ (\lambda x. h\ x)\ lst$ | ✓ | ✓ | No | ✓ |
| | where $h :: Int \rightarrow \forall a. a \rightarrow a$, $k :: \forall a. a \rightarrow [a] \rightarrow a$, and $lst :: [\forall a. Int \rightarrow a \rightarrow a]$ | | | | |
| E3 | $r\ (\lambda x\ y \rightarrow y)$ | No | ✓ | No | No |
| | where $r :: (\forall a. a \rightarrow \forall b. b \rightarrow b) \rightarrow Int$ | | | | |

Figure 2. Comparison of type systems

| | | | |
|--------------------------|------------------------|-------|---|
| Skolem / rigid variables | \forall | \ni | a, b, c, \dots |
| Free type variables | $ftv(\sigma)$ | | |
| Type constructors | \mathbb{T} | \ni | $\rightarrow, \top, S, \dots$ |
| Fully mono. types | τ | $::=$ | $a \mid \top \bar{\tau}$ |
| Top-level mono. types | μ, η | $::=$ | $a \mid \top \bar{\sigma}$ |
| Polymorphic types | σ, ϕ | $::=$ | $\forall \bar{a}. \mu$ |
| | | | \bar{a} may be empty, $\bar{a} \subseteq ftv(\mu)$ |
| Term variables | | \ni | x, y |
| Expressions / terms | e | $::=$ | x |
| | | | $ e_0 e_1 \dots e_n \quad n \geq 0$ |
| | | | $ \lambda x. e$ |
| | | | $ \lambda(x :: \sigma). e$ |
| | | | $ (e_0 e_1 \dots e_n :: \sigma)$ |
| | | | $ \mathbf{let}\ x = e_1\ \mathbf{in}\ e_2$ |
| Environments | Γ | $::=$ | $\epsilon \mid \Gamma, x : \sigma$ |
| Substitutions | θ, φ, π | $::=$ | $[\bar{a} \mapsto \bar{\sigma}]$ |
| Sorts of variables | s | $::=$ | $u \mid t \mid m$ |
| | | | Sorts form a lattice, $m \sqsubseteq t \sqsubseteq u$ |
| Sort assignment | Δ | $::=$ | $[\bar{a} \mapsto \bar{s}]$ |
| Bit | ω | $::=$ | $\bullet \mid \star$ |
| Vector of bits | $\bar{\omega}$ | $::=$ | $\epsilon \mid \omega, \bar{\omega}$ |

Figure 3. Syntax of the language

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. Annotations and lets are treated in their own sections, we first focus on the core language with variables, applications, and lambdas.

Types (Figure 3) are classified by three “sorts”, u , t , and m . *Polymorphic* types σ, ϕ , of sort u , have unrestricted polymorphism. *Top-level monomorphic* types, μ, η , of sort t , have no polymorphism at the top level, but permit arbitrary nested polymorphic types under a type constructor. Finally, *fully monomorphic* types, τ , of sort 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 θ , the image of a type variable a is denoted by $\theta(a)$, and similarly for sort assignments.

3.2 Typing rules

The typing judgment $\Gamma \vdash e : \sigma$ is given in Figure 4, 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 ϕ : we

$$\begin{array}{c}
\boxed{\sigma \text{ respects } s} \qquad \boxed{\theta \text{ respects } \Delta} \\
\frac{}{\sigma \text{ respects } u} \quad \frac{}{\mu \text{ respects } t} \quad \frac{}{\tau \text{ respects } m} \quad \frac{\forall a \in \text{dom}(\theta), \theta(a) \text{ respects } \Delta(a)}{\theta \text{ respects } \Delta} \\
\frac{\mu \triangleright^g \Delta}{\forall \bar{a}. \mu \triangleright^g \Delta \setminus \bar{a}} \text{ARGPOLY} \quad \frac{\boxed{\sigma \triangleright^g \Delta}}{\top \bar{\phi} \triangleright^g [\text{ftv}(\bar{\phi}) \mapsto u]} \text{ARGGUARD} \quad \frac{}{a \triangleright^g [a \mapsto t]} \text{ARGTYVAR} \\
\frac{}{\mu \triangleright_{\epsilon}^s [\text{ftv}(\mu) \mapsto s]} \text{ARGGRES} \quad \frac{}{a \triangleright_{\omega, \bar{\omega}}^s [a \mapsto m]} \text{ARGSTYVAR} \\
\frac{\mu \triangleright_{\omega}^s \Delta}{\forall \bar{a}. \mu \triangleright_{\omega}^s \Delta \setminus \bar{a}} \text{ARGSPOLY} \quad \frac{\sigma_1 \triangleright^g \Delta_1 \quad \sigma_2 \triangleright_{\omega}^s \Delta_2}{\sigma_1 \rightarrow \sigma_2 \triangleright_{\bullet, \bar{\omega}}^s \Delta_1 \sqcup \Delta_2} \text{ARGSAW} \\
\boxed{\phi \leq_{\omega}^s \bar{\sigma}; \mu} \\
\frac{}{\mu \leq_{\epsilon}^s \epsilon; \mu} \text{INSTMONO} \quad \frac{\phi_2 \leq_{\omega}^s \bar{\sigma}; \mu}{\phi_1 \rightarrow \phi_2 \leq_{\omega, \bar{\omega}}^s \phi_1, \bar{\sigma}; \mu} \text{INSTARROW} \quad \frac{\mu \triangleright_{\omega}^s \Delta \quad \theta \text{ respects } \Delta \quad \theta \mu \leq_{\omega}^s \bar{\sigma}; \mu}{\forall \bar{a}. \mu \leq_{\omega}^s \bar{\sigma}; \mu} \text{INSTPOLY} \\
\boxed{\Gamma \vdash^{\text{fun}} e : \sigma} \qquad \boxed{\Gamma \vdash_{\omega}^{\text{arg}} e : \sigma} \\
\frac{x : \sigma \in \Gamma}{\Gamma \vdash^{\text{fun}} x : \sigma} \text{VARHEAD} \quad \frac{\Gamma \vdash e : \sigma}{\Gamma \vdash^{\text{fun}} e : \sigma} \text{EXPRHEAD} \quad \frac{\Gamma \vdash e : \forall \bar{a}. \mu \quad \bar{b} \notin \Gamma}{\Gamma \vdash_{\bullet}^{\text{arg}} e : \forall \bar{b}. [\bar{a} \mapsto \bar{\tau}] \mu} \text{ARGGEN} \\
\boxed{\Gamma \vdash e : \sigma} \\
\frac{\Gamma, x : \tau \vdash e : \sigma}{\Gamma \vdash \lambda x. e : \tau \rightarrow \sigma} \text{ABS} \quad \frac{\Gamma, x : \phi \vdash e : \sigma}{\Gamma \vdash \lambda(x :: \phi). e : \phi \rightarrow \sigma} \text{ANNABS} \quad \frac{\Gamma \vdash e_1 : \phi \quad \Gamma, x : \phi \vdash e_2 : \sigma}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \sigma} \text{LET} \\
\frac{n \geq 0 \quad \Gamma \vdash^{\text{fun}} e_0 : \phi \quad \phi \leq_{\omega_1, \dots, \omega_n}^m \sigma_1, \dots, \sigma_n; \mu \quad \Gamma \vdash_{\omega_1}^{\text{arg}} e_1 : \sigma_1 \quad \dots \quad \Gamma \vdash_{\omega_n}^{\text{arg}} e_n : \sigma_n}{\Gamma \vdash e_0 e_1 \dots e_n : \mu} \text{APP} \\
\frac{n \geq 0 \quad \Gamma \vdash^{\text{fun}} e_0 : \phi \quad \phi \leq_{\omega_1, \dots, \omega_n}^u \sigma_1, \dots, \sigma_n; \eta \quad \Gamma \vdash_{\omega_1}^{\text{arg}} e_1 : \sigma_1 \quad \dots \quad \Gamma \vdash_{\omega_n}^{\text{arg}} e_n : \sigma_n}{\Gamma \vdash (e_0 e_1 \dots e_n :: \forall \bar{b}. \eta) : \forall \bar{b}. \eta} \text{ANNAPP}
\end{array}$$

Figure 4. Declarative type system

simply bring x into scope with type ϕ . As discussed in Section 2.3, where there is no annotation (rule ABS) we insist that x has a fully-monomorphic type τ .

All the action is in rule APP for n -ary applications ($e_0 e_1 \dots e_n$). First, note that a lone variable is treated as a nullary application. Second, the typing rules allow us to break an n -ary application in different ways, because e_0 could itself be an application. In practice we always choose e_0 not to be an application, so we get as many arguments as possible, and thereby maximise the opportunities for guardedness.

The first step in APP is typing the head of the application e_0 . The corresponding judgment \vdash^{fun} either looks for a variable in the environment or uses the normal typing judgment if the head is another kind of term. After typing the head, rule APP instantiates the type of the head with the *instantiation* judgement \leq_{ω}^s , yielding a list of argument types $\sigma_1 \dots \sigma_n$. APP uses \vdash^{arg} to check each argument e_i against the corresponding σ_i . This argument-checking judgement \vdash^{arg} can *generalise* the inferred type of the argument to match the type σ expected by the function, to support higher-rank polymorphism. Without such a rule, we would not be able

to type check *poly* ($\lambda x. x$), which requires the argument to be of type $\forall a. a \rightarrow a$.

For example, an application of a function to three arguments would give rise to the following instantiation:

$$\forall ab. a \rightarrow (\forall c. [a] \rightarrow b \rightarrow c \rightarrow [b]) \\ \leq_{\omega_1, \omega_2, \omega_3}^m \sigma_a, [\sigma_a], \tau_b; (\tau_c \rightarrow [\tau_b])$$

This example illustrates several points. First, the bit-vector $\bar{\omega}$ corresponds 1-1 with the arguments in the application. For now, the bit-values are irrelevant, but we will use them in Section 3.3. Second, the instantiation judgement returns a list of argument types that matches the length of the vector $\bar{\omega}$ (rule INSTARROW); in this case, there are three arguments, and the instantiated argument types are σ_a , $[\sigma_a]$, and τ_b . Third, the judgement can instantiate nested forall.

Fourth, and core to our contribution, each type variable is instantiated, by rule INSTPOLY, with a type whose sort reflects the way the type variable appears in the function type. This analysis is performed by the judgement $\sigma \triangleright_{\bar{\omega}}^s \Delta$, which returns a classification Δ of the free type variables of σ . This classification directly implements the three cases of the Instantiation Rule (Section 2.1):

1. If a appears under a type constructor (i.e. guarded) in any of the first n arguments of σ , then a may be instantiated by an unrestricted type ϕ (rule ARGGUARD).
2. If a appears in one of the first n arguments of σ , a may be instantiated by a top-level monomorphic type μ (rule ARGTYVAR).
3. If a appears only in the result type of σ , after stripping off n arguments, then a may be instantiated by a type of sort s (rule ARGRES). In the invocation of \leq in rule APP the sort s is always fully monomorphic m , but we need the extra generality for annotations (Section 3.4).

Once we have Δ to classify each type variable rule, INSTPOLY instantiates the function type with a Δ -respecting substitution θ , and recurses. Figure 4 also defines what it means for a substitution θ to “respect” a classification Δ .

Unlike some other systems (see [13] for a comprehensive account), all type constructors in GI are invariant, *including functions*. This means that neither $[\forall a. a \rightarrow a] \not\leq_{\bar{\omega}}^s [Int \rightarrow Int]$ nor $Int \rightarrow (\forall a. a \rightarrow a) \not\leq_{\bar{\omega}}^s Int \rightarrow Int \rightarrow Int$. However, because \leq handles forall nested to the right of arrows (rule INSTPOLY), we can often work around the lack of covariance using η -expansion. For example, suppose we have functions $f :: \forall a. a \rightarrow (\forall b. b \rightarrow a)$, and $g :: (Int \rightarrow Int \rightarrow Int) \rightarrow Bool$. Then $(g f)$ is ill-typed; but $(g (\lambda xy. f x y))$ is well typed.

3.3 Single variables

The type system described up to now propagates polymorphism between arguments and from the arguments to the return type of the function – this is the essence of guardedness. Alas, this creates a problem for single variables, which are treated as nullary applications: since there are no arguments, no type variable is considered guarded, and thus

$$\frac{x : \forall \bar{p}. \tau \in \Gamma \quad \bar{b} \notin \Gamma}{\Gamma \vdash_{\star}^{\text{arg}} x : \forall \bar{b}. [a \mapsto \sigma] \tau} \text{VARGEN} \\ \frac{\Delta_1 = [\text{ftv}(\sigma_1) \mapsto m] \quad \sigma_2 \triangleright_{\bar{\omega}}^s \Delta_2}{\sigma_1 \rightarrow \sigma_2 \triangleright_{\star, \bar{\omega}}^s \Delta_1 \sqcup \Delta_2} \text{ARGSTAR}$$

Figure 5. Decl. type system with single variables

impredicativity is forbidden. We can see this by expanding the derivation of $\Gamma \vdash_{\bullet}^{\text{arg}} x : \sigma$ for the case of a variable x .

$$\frac{x : \phi \in \Gamma}{\Gamma \vdash^{\text{fun}} x : \phi} \text{VARHEAD} \quad \frac{\phi \leq_{\epsilon}^m \epsilon; \mu}{\bar{b} \notin \Gamma \quad \Gamma \vdash x : \mu} \text{APP} \\ \frac{}{\Gamma \vdash_{\bullet}^{\text{arg}} x : \forall \bar{b}. \mu} \text{ARGGEN}$$

The highlighted premise, $\phi \leq_{\epsilon}^m \epsilon; \mu$, forces instantiation to use only fully monomorphic types.

That is embarrassing, because we cannot typecheck, say *choose* $[\]$ *ids* from Figure 2. From the type of *ids* it is obvious that $[\]$ should be given a type $[\forall a. a \rightarrow a]$. But we cannot do so, because such instantiation is not fully monomorphic. We get back into swampy waters.

For the case of single variables with a rank-1 polymorphic type (quantifiers appear only at top-level) we can get out of the swamp. The VARGEN rule in Figure 5 formalizes this idea: if a single variable x appears as the argument of an application – hence the use of the judgement \vdash^{arg} – and it has a rank-1 type $\forall \bar{p}. \tau$, we may instantiate impredicatively before generalizing. This is enough to cover the case of *choose* $[\]$ *ids*, because the type of $[\] :: \forall a. [a]$ has the right shape. Note, however, that it is still the case that the justification for instantiating the type of the first argument $[\]$ with top-level monomorphic type has to come from the second argument *ids*. If such justification is not forthcoming, as it should be in the case of *choose* $[\]$ $[\]$, then both arguments should be instantiated fully monomorphically.

We manage the bookkeeping to distinguish these situations by keeping track of a vector of bits $\bar{\omega}$. Its elements correspond to the arguments in an application: a \star means that the rule VARGEN was applied to type the corresponding argument; application of ARGGEN is represented by \bullet . The rule ARGSTAR ensures that whenever the rule VARGEN was used to type a given argument, we reset the sorts of the free variables to m so that however these type variables were instantiated, this information cannot be used to justify the impredicative instantiation of other arguments.

3.4 Annotations

The VARGEN rule makes it possible to accept more programs, but it is restricted to a very special class of expressions. For the general case, we provide *annotations* as part of the syntax. Since annotations fully specify the types, we do not

need to impose guardedness restrictions on those variables appearing in the result. 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 τ . If we want instead to obtain $[\forall a. a \rightarrow a]$, we can just annotate the result, thus *single* $(\lambda x. x) :: [\forall a. a \rightarrow a]$.

Rule ANNAPP is almost identical to APP, except for the choice of parameter to the instantiation judgment, which is 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.

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.

3.5 let bindings

The simplest way to type **let** $x = e_1$ **in** e_2 is to see it as a shorthand for $(\lambda x. e_2) e_1$. Alas, such a translation imposes an important restriction on the type of x : it must 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 4 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 .

3.6 Metatheory

Impredicativity is a great tool, but we do not want to lose those programs which only require top-level polymorphism. The following theorem states this fact, except for the different take on **let** generalization we discussed in Section 3.5.

Theorem 3.1 (Compatibility with rank-1 polymorphism). *Let e be an expression in the syntax of the lambda-calculus with predicative rank-1 polymorphism. If $\Gamma \vdash e : \tau$ in that type system, then $\Gamma \vdash e : \tau$ in GI.*

One key property of GI is that all impredicative instantiations are settled by the shape of the expression and the types in the environment, 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 Γ be an environment and e an expression. For every pair of fully monomorphic substitutions θ_1 and θ_2 , if $\theta_1 \Gamma \vdash e : \sigma_1$ and $\theta_2 \Gamma \vdash e : \sigma_2$, then there exists a polymorphic type σ^* and fully monomorphic substitutions φ_1 and φ_2 such that $\sigma_i = \varphi_i \sigma^*$.*

Corollary 3.3. *Let Γ be a closed environment (that is, no type in Γ contains a free variable). If $\Gamma \vdash e : \sigma_1$ and $\Gamma \vdash e : \sigma_2$, then there is a polymorphic type σ^* and fully monomorphic substitutions φ_1 and φ_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 σ^* for which any other type assignment ϕ to e is equal to $\theta \sigma^*$ for a fully monomorphic substitution θ . 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. Note also that this theorem only promises that if GI accepts an expression, there was no guessing involved. But there are expressions for which only one choice of polymorphism is possible, yet GI cannot find it and an annotation is required.

In the remainder of this section we look at some properties of GI concerning derivations and stability under transformations. We use the notation $e_1[e_2]$ refers to a context in which e_2 appears. Proofs are given in Appendix D.1.

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

The converse result does not hold: we cannot in general abstract over part of an expression. In practice, that means that we cannot always introduce a **let**, as in changing $e_0 e_1 e_2$ to **let** $x = e_0 e_1$ **in** $x e_2$. Our APP rule is responsible for propagating information between arguments, if we introduce a **let** for part of an expression, this bound is lost.

Theorem 3.5. *Let $app :: \forall a b. (a \rightarrow b) \rightarrow a \rightarrow b$ and $revapp :: \forall a b. a \rightarrow (a \rightarrow b) \rightarrow b$ be the application and reverse application functions, respectively. Given two expressions f and e such that $\Gamma \vdash^{fun} f : \sigma_0$, and $\sigma_0 \leq_e^m \epsilon; \sigma_1 \rightarrow \phi$ then:*

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

The hypothesis $\sigma_0 \leq_e^m \epsilon; \sigma_1 \rightarrow \phi$ means that type variables in f may only be instantiated with fully monomorphic variables. Thus, this transformation only respects well-typedness for the predicative fragment, but not in general when impredicative instantiation is allowed. One example of this restriction is the application *f ids*, where f has type $\forall a. [a] \rightarrow [a]$, which cannot be turned into *app f ids*. The reason is that in the original application the type variable a in f has to be instantiated with a polymorphic type $\forall b. b \rightarrow b$. However, this is not allowed in the form with *app*.

One property which does not hold in GI is *full subject reduction*. If $\Gamma \vdash e : \sigma$ and e β -reduces to e' , it may not hold

that $\Gamma \vdash e' : \sigma$. For example, *app auto* is typeable, but not its reduced form $\lambda x. \text{auto } x$, since it requires an annotation on x . Subject reduction holds in a milder form: if $\Gamma \vdash e : \sigma$, e β -reduces to e' and $\Gamma \vdash e' : \phi$, then σ and ϕ coincide.

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:

$$\text{head } (\forall a. a \rightarrow a) \text{ ids } \text{Bool True}$$

That is, we instantiate *head* at type $(\forall a. a \rightarrow a)$, then apply it to *ids*, to produce a result of type $(\forall a. a \rightarrow a)$. Now we must in turn instantiate that type with *Bool* to get a function of type $(\text{Bool} \rightarrow \text{Bool})$ 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 α for some as-yet-unknown type α . So we want to defer the instantiation decision.

Sometimes we must defer generalisation decisions too. Consider the function application $(\text{:}) (\lambda x. x) \text{ ids}$. In System F terms, we want to infer the following elaborated program:

$$(\text{:}) (\forall a. a \rightarrow a) (\Lambda a. \lambda(x : a). x) \text{ ids}$$

in which (:) is instantiated at type $(\forall a. a \rightarrow a)$, 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, α say. Ultimately the type of *ids* forces α to be $\forall a. a \rightarrow a$, 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

| | | | |
|----------------------------|----------------------|-------|---|
| Unif. var. names | \bar{U} | \ni | $\alpha, \beta, \gamma, \delta, \dots$ |
| Unif. var. | v | $::=$ | α^s |
| Free unification variables | $\text{fuv}(\sigma)$ | | |
| Fully mono. types | τ | $::=$ | $\alpha^m \mid a \mid \bar{\tau}$ |
| Top-level mono. types | μ, η | $::=$ | $\alpha^\dagger \mid a \mid \bar{\tau}$ |
| Polymorphic types | σ, ϕ | $::=$ | $\alpha^n \mid \forall \bar{a}. \mu$ |
| | | | \bar{a} may be empty, $\bar{a} \subseteq \text{fuv}(\mu)$ |
| Types with generalisation | g | $::=$ | $\forall \{\bar{v}\}. C \Rightarrow \sigma$ |
| | | | \bar{v} may be empty |
| Constraints | C | $::=$ | \top |
| Conjunction | | | $C_1 \wedge C_2$ |
| Equality | | | $\sigma \sim \phi$ |
| Instantiation | | | $\sigma \leq_{\omega}^s \bar{\phi}; \mu$ |
| Generalisation | | | $g \leq \sigma$ |
| Quantification | | | $\forall \bar{a}. \exists \bar{v}. C$ |

Figure 6. Extended syntax

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 6 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 6, a key idea is that unification variables are drawn from three distinct ‘‘alphabets’’: α^s for each of the three sorts 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 s may only be unified with types belonging to that sort.

The syntax presents several kinds of constraints C . These constraints do not form part of the source language; they are internal to the solver. 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.

4.2 Constraint generation

Constraint *generation* is described in Figure 7 as a four-element judgment $\Gamma \vdash e : \sigma \rightsquigarrow C$. The first two elements are inputs: the environment Γ and the expression e for which to generate constraints. The output of the process is a type σ

$$\begin{array}{c}
\boxed{\Gamma \vdash^{\text{fun}} e : \sigma \rightsquigarrow C} \\
\frac{x : \sigma \in \Gamma}{\Gamma \vdash^{\text{fun}} x : \sigma \rightsquigarrow \epsilon} \text{VARHEAD} \qquad \frac{e \text{ not var. or app.} \quad \Gamma \vdash e : \sigma \rightsquigarrow C}{\Gamma \vdash^{\text{fun}} e : \sigma \rightsquigarrow C} \text{EXPRHEAD} \\
\boxed{\Gamma ; \bar{v} \vdash_{\omega}^{\text{arg}} e : \sigma \rightsquigarrow C} \\
\frac{x : \forall \bar{p}. \tau \in \Gamma \quad \forall \bar{p}. \tau \text{ closed} \quad \bar{\alpha} \text{ fresh}}{\Gamma ; \bar{v} \vdash_{\star}^{\text{arg}} x : \sigma \rightsquigarrow [\bar{p} \mapsto \alpha^{\text{II}}] \tau \leq \sigma} \text{VARGEN} \qquad \frac{\Gamma \vdash e : \phi \rightsquigarrow C \quad \bar{v}' = \text{fuv}(\phi, C) - \bar{v}}{\Gamma ; \bar{v} \vdash_{\bullet}^{\text{arg}} e : \sigma \rightsquigarrow \forall \{\bar{v}'\}. C \Rightarrow \phi \leq \sigma} \text{ARGGEN} \\
\boxed{\Gamma \vdash e : \sigma \rightsquigarrow C} \\
\frac{\alpha \text{ fresh} \quad \Gamma, x : \alpha^{\text{III}} \vdash e : \sigma \rightsquigarrow C}{\Gamma \vdash \lambda x. e : \alpha^{\text{III}} \rightarrow \sigma \rightsquigarrow C} \text{ABS} \qquad \frac{\Gamma, x : \phi \vdash e : \sigma \rightsquigarrow C}{\Gamma \vdash \lambda(x :: \phi). e : \phi \rightarrow \sigma \rightsquigarrow C} \text{ANNABS} \\
\frac{\bar{\alpha}, \beta \text{ fresh} \quad \Gamma \vdash^{\text{fun}} e_0 : \phi \rightsquigarrow C \quad \Gamma ; \text{fuv}(\Gamma, \phi, C) \vdash_{\omega_i}^{\text{arg}} e_i : \alpha_i^{\text{II}} \rightsquigarrow C_i}{\Gamma \vdash e_0 e_1 \dots e_n : \beta^{\text{t}} \rightsquigarrow C \wedge \phi \leq_{\omega_1, \dots, \omega_n}^{\text{m}} \alpha_1^{\text{II}}, \dots, \alpha_n^{\text{II}}; \beta^{\text{t}} \wedge C_1 \wedge \dots \wedge C_n} \text{APP} \\
\frac{\bar{\alpha} \text{ fresh} \quad \Gamma \vdash^{\text{fun}} e_0 : \phi \rightsquigarrow C \quad \bar{v}' = \text{fuv}(\phi, C) - \text{fuv}(\Gamma) \quad \Gamma ; \text{fuv}(\Gamma, \phi, C) \vdash_{\omega_i}^{\text{arg}} e_i : \alpha_i^{\text{II}} \rightsquigarrow C_i}{\Gamma \vdash (e_0 e_1 \dots e_n :: \forall \bar{b}. \eta) : \forall \bar{b}. \eta \rightsquigarrow \forall \bar{b}. \exists \bar{\alpha}^{\text{II}} \bar{v}'. (C \wedge \phi \leq_{\omega_1, \dots, \omega_n}^{\text{II}} \alpha_1^{\text{II}}, \dots, \alpha_n^{\text{II}}; \eta \wedge C_1 \wedge \dots \wedge C_n)} \text{ANNAPP} \\
\frac{\Gamma \vdash e_1 : \phi \rightsquigarrow C_1 \quad \Gamma, x : \phi \vdash e_2 : \sigma \rightsquigarrow C_2}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \sigma \rightsquigarrow C_1 \wedge C_2} \text{LET}
\end{array}$$

Figure 7. Constraint generation

assigned to the expression, possibly including some unification variables, and the set of extended constraints C that the types must satisfy.

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 4), the head of the application is typed using an ancillary judgment $\Gamma \vdash^{\text{fun}} e : \phi \rightsquigarrow C$, which either looks up a variable in the environment or threads the information to the normal gathering process.

Another ancillary judgment $\Gamma ; \bar{v} \vdash_b^{\text{arg}} e : \sigma \rightsquigarrow C$ is used to generate constraints on the arguments. In this judgement we use introduce a completely new constraint form, $g \leq \sigma$, 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 $(\forall \{\bar{v}\}. C \Rightarrow \phi) \leq \sigma$ should be read “a term of type ϕ with constraints C and unification variables \bar{v} can be instantiated and/or generalised to have type σ ”. 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 α_i from the function, into a generalisation constraint for the solver to deal with later. The reason for the additional

argument \bar{v} to the judgement is that we need to forbid generalisation over variables which are visible in the environment or the whole application.

Rule ANNAPP deals with a type-annotated application $(e_0 e_1 \dots e_n :: \sigma)$. It is similar to APP, but it is the first rule to introduce a *quantification constraint* $\forall \bar{b}. \exists \bar{v}. \bar{C}$. This binds the Skolem variables \bar{b} from the type signature, and existentially quantifies the unification variables free in the constraint but not used outside it. The other difference with APP is the use II as a parameter to \leq , rather than m , as done in the declarative specification.

4.3 Constraint solving

The solver takes the generated constraint C and its free unification variables $\bar{v} = \text{fuv}(C)$, and repeatedly applies the solver rules in Figure 8, 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. $\text{Int} \sim \text{Bool}$) represent type errors that can be reported to the user. We concentrate first on the solver rules that incrementally solve the constraint.

Each of the rules in Figure 8 rewrites a configuration $C; \bar{v}$ to another configuration. The unification variables \bar{v} are existentially quantified, so you can think of a configuration as representing $\exists \bar{v}. C$. 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

$$\boxed{\text{freshen}_{\omega}^s(\sigma) \Longrightarrow \langle \bar{v}, \mu \rangle}$$

$$\frac{\mu \triangleright_{\omega}^s \Delta \quad \bar{\alpha} \text{ fresh} \quad v_i = \alpha_i^{\Delta(a_i)}}{\text{freshen}_{\omega}^s(\forall \bar{a}. \mu) \Longrightarrow \langle \bar{v}, [\bar{a} \mapsto \bar{v}] \mu \rangle}$$

$$\boxed{C ; \bar{v} \Longrightarrow C' ; \bar{v}'}$$

| | | | | |
|-----------------------|--|--|----------|---|
| [CONJ] | $\frac{C_1 ; \bar{v} \Longrightarrow C'_1 ; \bar{v}'}{C_1 \wedge C_2 ; \bar{v} \Longrightarrow C'_1 \wedge C_2 ; \bar{v}'}$ | modulo assoc. and comm. of \wedge | [FORALL] | $\frac{C ; \bar{v}_{\text{in}} \Longrightarrow C' ; \bar{v}'_{\text{in}}}{\forall \bar{a}. \exists \bar{v}_{\text{in}}. C ; \bar{v} \Longrightarrow \forall \bar{a}. \exists \bar{v}'_{\text{in}}. C' ; \bar{v}}$ |
| [TIDENT] | $\top \wedge C ; \bar{v} \Longrightarrow C ; \bar{v}$ | | | |
| [TFORALL] | $\forall \bar{a}. \exists \bar{v}_{\text{in}}. \top ; \bar{v} \Longrightarrow \top ; \bar{v}$ | | | |
| [EQREFL] | $\sigma \sim \sigma ; \bar{v} \Longrightarrow \top ; \bar{v}$ | | | |
| [EQMONO] | $\top \sigma_1 \dots \sigma_n \sim \top \phi_1 \dots \phi_n ; \bar{v} \Longrightarrow (\sigma_1 \sim \phi_1) \wedge \dots \wedge (\sigma_n \sim \phi_n) ; \bar{v}$ | | | |
| [EQSUBST] | $(\alpha^s \sim \sigma) \wedge C ; \bar{v} \Longrightarrow (\alpha^s \sim \sigma) \wedge [\alpha^s \mapsto \sigma] C ; \bar{v}$ | | | if σ respects s , and $\alpha \notin \text{ftv}(\sigma)$ |
| [EQVAR] | $\alpha^{s_1} \sim \beta^{s_2} ; \bar{v} \Longrightarrow \beta^{s_2} \sim \alpha^{s_1} ; \bar{v}$ | | | if $s_1 \sqsubset s_2$ |
| [EQFULLY] | $\alpha^m \sim \sigma ; \bar{v} \Longrightarrow \{\beta^s \sim \gamma^m \mid \beta^s \in \text{fv}(\sigma), s \neq m\} ; \bar{v}, \gamma^m$ | | | \bar{v}, γ^m disjoint |
| [INST ϵ] | $\mu \leq_{\epsilon}^s \epsilon ; \eta ; \bar{v} \Longrightarrow \mu \sim \eta ; \bar{v}$ | | | |
| [INST \rightarrow] | $\mu \leq_{\omega, \bar{\omega}}^s \sigma, \bar{\phi} ; \eta ; \bar{v} \Longrightarrow (\mu \sim \sigma \rightarrow \beta^u) \wedge (\beta^u \leq_{\bar{\omega}}^s \bar{\phi} ; \eta) ; \bar{v}$ | | | |
| [INST \forall L] | $(\forall \bar{a}. \mu) \leq_{\bar{\omega}}^s \bar{\phi} ; \eta ; \bar{v} \Longrightarrow \mu' \leq_{\bar{\omega}}^s \bar{\phi} ; \eta ; \bar{v}, \bar{v}'$ | | | where $\text{freshen}_{\omega}^s(\forall \bar{a}. \mu) \Longrightarrow \langle \bar{v}', \mu' \rangle$ |
| [INST \forall L] | $(\forall \{\bar{v}'\}. \bar{C} \Rightarrow \sigma) \leq \eta ; \bar{v} \Longrightarrow \bar{C} \wedge (\sigma \leq_{\epsilon}^m \epsilon ; \eta) ; \bar{v}, \bar{v}'$ | | | \bar{v}, \bar{v}' disjoint |
| [INST \forall R] | $g \leq (\forall \bar{a}. \mu) ; \bar{v} \Longrightarrow \forall \bar{a}. (g \leq \mu) ; \bar{v}$ | | | |

Figure 8. Solving rules

the rules are read modulo commutativity and associativity of \wedge ; that is why CONJ only has to handle the left conjunct. The rules TIDENT and TFORALL remove \top constraints from the set, since they are identities for the conjunction.

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 the substitution does not lead to an infinite type (hence the 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.

Given the different behaviour embodied by 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. For example, when we have an unrestricted α^u and a top-level monomorphic β^t , then α^u

should be replaced by β^t , and not the other way around. Full monomorphism goes deeper: EQFULLY ensures that if a type σ is equated with a fully monomorphic variable α^m , all the variables in σ become fully monomorphic too.

One difference between these rules and other presentations is that we do not rewrite an unsatisfiable constraint, such as $\text{Int} \sim \text{Bool}$, to \perp . 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 α -equality) to match under the EQREFL rule. This means that $(\forall a b. a \rightarrow b \rightarrow b) \sim (\forall b a. a \rightarrow b \rightarrow b)$ does *not* hold in our system. As we discuss in Section 2.4, this is not problematic, since on application type variables are instantiated and regeneralized using the \leq relation.

4.3.2 Instantiation and generalisation constraints

For instantiation constraints $\leq_{\bar{\omega}}^s$, we follow closely the judgement with the same name in the declarative specification. Rule INST ϵ encodes the fact that once all arguments are processed, and thus $\bar{\omega}$ is an empty vector of bits, the remaining types must be equal. This is a consequence of the invariance of type constructors. On the other hand, if we have a top-level monomorphic type μ and the vector of bits

is not empty, the only possibility of for μ to be a function type. This is the goal of the $\text{INST}\rightarrow$ rule.

$\text{INST}\forall\text{L}$ instantiates a polytype σ 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 $\text{freshen}_{\overline{\alpha}}^s$, which in turn uses the already-introduced classification judgment $\triangleright_{\overline{\alpha}}^s$ (Figure 4). Finally, the new variables enter the set of existentially quantified variables.

Finally, we come to generalisation constraints, which (recall Section 4.1) express a deferred generalisation decision. Rule $\text{INST}\forall\text{L}$ is simple: if the right hand side has no top-level forall (it is of the form η) 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.

$\text{INST}\forall\text{R}$ is where actual generalisation takes place. In order to forge some intuition, let us look at the constraint

$$(\forall\{\alpha\beta\}. (\alpha \leq_{\epsilon}^m \epsilon; \beta) \Rightarrow \alpha \rightarrow \beta) \leq (\forall p. p \rightarrow p)$$

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

$$(\alpha \leq_{\epsilon}^m \epsilon; \beta) \wedge (\alpha \rightarrow \beta \sim p \rightarrow p)$$

We obtain a solution by making $\alpha \sim \beta \sim p$. 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

$$\forall p. \exists \alpha \beta. (\alpha \leq_{\epsilon}^m \epsilon; \beta) \wedge (\alpha \rightarrow \beta \sim a \rightarrow a)$$

Rule $\text{INST}\forall\text{R}$ achieves this rather neatly simply by doing skolemisation and pushing the \forall inside; then FORALL and $\text{INST}\forall\text{L}$ 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 $\alpha^{\text{u}} \leq_{\epsilon}^s \epsilon; \mu$ we cannot turn it directly into $\alpha^{\text{u}} \sim \mu$, because α^{u} might be unified later to a polymorphic type and we need instantiation.
2. Similarly, if we have $(\forall\{\overline{\alpha}\}. C \Rightarrow \mu) \leq \alpha^{\text{u}}$, and α^{u} is later substituted by a polytype, we must skolemise.

The guardedness restrictions are carefully crafted to ensure that the solver is *confluent* and that it is never completely *stuck*, unless the constraint set as a whole is inconsistent. A single constraint can be stuck for some time, but if the whole set is consistent, by steps applied to other constraints, it will eventually 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.*

$$\begin{array}{c} \frac{}{\overline{a}; \overline{\alpha}; \emptyset \vdash \top \text{ solved}} \text{SOLVEDT} \\ \frac{\sigma \text{ respects } s \quad \text{ftv}(\sigma) \subseteq \overline{a} \cup \overline{\alpha}}{\overline{a}; \overline{\alpha}; \{\beta\} \vdash \beta^s \sim \sigma \text{ solved}} \text{SOLVEDVAR} \\ \frac{\overline{a}; \overline{\alpha}; \overline{\beta}_1 \vdash C_1 \text{ solved} \quad \overline{a}; \overline{\alpha}; \overline{\beta}_2 \vdash C_2 \text{ solved}}{\overline{a}; \overline{\alpha}; \overline{\beta}_1 \uplus \overline{\beta}_2 \vdash C_1 \wedge C_2 \text{ solved}} \text{SOLVEDCONJ} \\ \frac{\overline{v} = \overline{\gamma}_1 \uplus \overline{\gamma}_2 \quad \overline{a} \cup \overline{b}; \overline{\alpha} \cup \overline{\gamma}_1; \overline{\gamma}_2 \vdash C \text{ solved}}{\overline{a}; \overline{\alpha}; \emptyset \vdash \forall \overline{b}. \exists \overline{v}. C \text{ solved}} \text{SOLVEDQUANT} \end{array}$$

Figure 9. Definition of solved set of constraints

Theorem 4.1 tells us that the process of solving can be divided into two phases. In the first phase all constraints, including instantiation and generalization constraints, are turned into a set of equalities, possibly with different quantification levels. This is an instance of the problem of first-order unification under a mixed prefix [3], for which a complete solving algorithm is described by Pottier and Rémy [15].

4.4 Soundness, principality and completeness

The inference algorithm presented here – gathering the constraints from an expression followed by solving them – satisfies the usual properties of soundness, principality, and completeness with respect to the declarative specification.

In order to state the theorems we need some ancillary notions. A constraint is in *solved form* if it consists only of quantification and equality constraints ($v \sim \sigma$); and the equalities constitute a well-sorted idempotent substitution of its unification variables. For example

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

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

$$\exists \alpha. (\alpha \sim \text{Int}) \wedge (\alpha \sim \text{Bool}) \quad \exists \alpha. \dots (\forall b. \alpha \sim [b]) \dots$$

In the first there are two equalities for α (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 9 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 β , checks well-sortedness, and also checks that σ 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{v} into the unconstrained sets, γ_1 and γ_2 resp. Rule SOLVEDT simply states that \top is a solved form without consuming any variables.

A second auxiliary notion which we need to state the results is *substitution induced by a solved form*.

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

Theorem 4.2 (Soundness). *Let Γ be a closed environment and e an expression. If $\Gamma \vdash e : \sigma \rightsquigarrow C$ and C_s is a solution for C with an induced substitution \widehat{C}_s , then we have $\Gamma \vdash e : \widehat{C}_s(\sigma)$.*

Theorem 4.3 (Principality). *Suppose $\Gamma \vdash e : \sigma$. Then there exists a type ϕ such that $\Gamma \vdash e : \phi$, and for every other $\Gamma \vdash e : \sigma'$, there is a fully monomorphic substitution π such that $\sigma' = \pi\phi$.*

Theorem 4.4 (Completeness). *Let Γ 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.*

Proofs of these results are given in Appendix D.2.

4.5 Alternative solver for equalities

Unfortunately, once we extend the language of types, by introducing type classes and local assumptions, 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 these extensions. The approach is rather simple – a single rule `FLOAT` in Figure 10 – and is directly inspired by how GHC handles constraints: by floating constraints out from inside a quantification constraint. When can we do that? Precisely when the constraint does not mention the skolems. But what about the existentials? Consider

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

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

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

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

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

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

Conjecture 4.5. *The solver presented in Figure 8 is complete for unification problems under a mixed prefix.*

5 Practical matters

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 column.

GI does not support any co- or contra-variance in function types. For example, the constraint $Int \rightarrow (\forall a. a \rightarrow a) \leq Int \rightarrow Int \rightarrow Int$ does not hold. In contrast, GHC with the *RankNTypes* extension supports some amount of variance. Libraries such as Scrap Your Boilerplate use this fact very often in definitions with a \forall to the right of an arrow:

$$f :: \forall a. a \rightarrow (\forall b. b \rightarrow b)$$

$$f \ x \ y = y$$

Other uses of variance can be worked around by η -expansion. 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 η -expand the argument, thus (*flip* $(\lambda x \rightarrow f \ x)$). GHC does this automatically at the moment, but in fact this η -expansion is unsound in general, since a change in the laziness behavior can be observed.

One might worry that if GI is integrated in GHC 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 to support function definition in the style of Scrap Your Boilerplate.

With this done, very few packages required modifications, and modifications were always η -expansions. In particular, of the 2,400 packages in Stackage, 609 use *RankNTypes*; of these, only 75 required manual changes, all of which were simple η -expansions. One (*singletons*) would require larger changes, because it uses Template Haskell to *generate* Haskell code; so it needs to generate η -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 η -expansion strategy is unsound in the first place.

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 higher-order 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 *predicative* instantiation 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.

$$\text{[FLOAT]} \frac{\text{ftv}(F) \cap \bar{a} = \emptyset \quad \bar{\alpha}^s = \text{fuv}(F) \cap v_{\text{in}} \quad \bar{\gamma}^s \text{ fresh} \quad \bar{E} = \bigwedge_{\alpha^s \in \bar{\alpha}} \alpha^s \sim \gamma^s}{\forall \bar{a}. \exists \bar{v}_{\text{in}}. (C \wedge F) ; \bar{v} \implies [\alpha^s \mapsto \gamma^s] F \wedge \forall \bar{a}. \exists \bar{v}_{\text{in}}. (C \wedge E) ; \bar{v}, \bar{\gamma}^s}$$

Figure 10. Solving rule for quantification constraints

MLF [1, 2, 18] is an *extension* of System F based on quantification with instance and equality 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], also known as HML, 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 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 by Eisenberg et al. [5] also proposes a distinction between an implicitly and an explicitly instantiable \forall ; only the latter is impredicative.

We return now to the table in Figure 2, where we present a collection of examples appearing in selected related works on type inference for impredicativity. We have selected those systems because they all strike a good balance between expressivity and requiring 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 (as in QML). We also show how to recover the typeability of a program in cases where a valid type exists.

In the table we see that MLF can type all programs that do not require implicit η -expansion (*k h lst*) or the use of a polymorphic function argument at two types ($\lambda f. (f \ 1, f \ \text{True})$), but uses constrained types and a unification algorithm that is not straightforward to integrate in a type inference engine with elaboration to System F. HML – a simplification of MLF with only “flexible” type bounds in constraints – is only slightly less powerful than MLF. HML cannot type $\lambda xs. \text{poly} \ (\text{head } xs)$ because xs would have to be assigned a polymorphic type, even though it’s only used at this one type. FPH – which is based on System F types – is equally expressive to these systems for the fragment of System F consisting only of variables and applications. However, the FPH treatment of λ -abstractions requires the returned type to be a fully-resolved top-level monomorphic type. Therefore, FPH fails to type check *choose id auto* because *auto* will be

assigned the same type as GI infers ($(\forall a. a \rightarrow a) \rightarrow b \rightarrow b$). GI – modulo the quirk about not generalizing the bodies of lambda abstractions that only MLF and HML can tackle – type checks almost the same programs as those systems. 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* to $(\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$ and no argument has a type with a top-level constructor.

HMF is 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*), though it does generalize in argument positions (e.g. *snoc id ids* is accepted, for *snoc :: \forall a. [a] \rightarrow a \rightarrow [a]*). Leijen [8] proposes an extension of the basic algorithm to n -ary applications that makes these examples type check: after the function type is instantiated enough to cover all the arguments, we proceed to type check the arguments in *a computed order*, instead of left-to-right. The arguments that must be type checked against a naked type variable coming from the instantiated function type are postponed and checked last, in the hope that the rest of the arguments will by that time determine any impredicative instantiations. The procedure is iterated, possibly uncovering impredicative instantiations in each round. Under that extension *id : ids* is accepted. Alas, the new algorithm is not accompanied with a declarative specification. Our system achieves a similar effect (and we conjecture is equally expressive), thanks to the delaying we get from the use of constraints. In the end, the main important difference of HMF with GI is that GI provides a declarative specification and an algorithm that easily integrates in a pre-existing constraint-based type inference engine.

7 Further work

GI seeks a sweet spot that balances *simplicity* with *expressiveness*. We are also exploring some nearby variants. For example, extending VARGEN to handle larger expressions would allow us to accept examples A9, C8, C9 in Figure 2; and by improving skolemisation we could accept E3. It remains to be seen whether the extra expressiveness justifies the extra complexity, but GI seems to be an encouragingly robust base camp.

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 Löb. 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://crystal.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 ML^F : 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>

Data constructor $\ni K$
 Expressions / terms $e ::= \dots \mid \text{case } e_0 \text{ of } \{K \bar{x} \rightarrow e\}$

$$\frac{\boxed{\Gamma \vdash e : \sigma} \quad \Gamma \vdash e_0 : \sigma_0 \quad \sigma_0 \leq_m^\epsilon \epsilon; \top \bar{\phi}_0 \quad \text{for each branch } K_i \bar{x}_i \rightarrow e_i \text{ with } K_i : \forall \bar{a} \bar{b}. \bar{\sigma}_i \rightarrow \top \bar{a} \in \Gamma \quad \Gamma, x_i : [a \mapsto \bar{\phi}_0] \sigma_i \vdash e_i : \phi_\star}{\Gamma \vdash \text{case } e_0 \text{ of } \{K \bar{x} \rightarrow e\} : \phi_\star} \text{CASE}$$

Figure 11. Decl. type system with pattern matching

A Pattern matching

Introducing pattern matching in the language gives no surprises, as witnessed by the rule CASE in Figure 11. We first need to check that the scrutinized expression e_0 can be given a type compatible with the indicated data constructor. It is important to note that the mere presence of the data constructors is enough to know the type constructor \top , there is no inference at that point. Then we introduce new term variables in each branch, whose type is obtained by combining the type of the constructor with the inferred values for the type variables in \top , namely $\bar{\phi}_0$. The return type of every branch should be equal, ϕ_\star (note that we do *not* allow the types of the branches to be instantiated).

The corresponding CASE rule for constraint gathering is given in Figure 12. In this case we need a quantification constraint to introduce new skolem constants, for the existentially quantified type variables \bar{b}_i in the data constructor.

B Integration with other language features

One of the main advantages of organising a type checker around the concept of constraints is that extensions to the type system can be accommodated for quite easily. Indeed, this is one of the main advantages of our system over existing work. In this section we describe how to deal with other forms of constraint beyond the standard equality constraints we now support.

For example, Haskell supports *type classes*, which restrict the scope of a polymorphic abstraction to a subset of types. The archetypal example is *Eq*, which describes the types with support for decidable equality. Such a type class constraint is visible in the type of the equality operator (\equiv) :: $\forall a. Eq a \Rightarrow a \rightarrow a \rightarrow Bool$. Similarly, languages like OCaml and PureScript feature *row* (or *record*) *types* like $\{x :: Point, y :: Point\}$ in addition to usual ADTs.

The good news is that a constraint-based formulation of typing makes it easy to cope with new concepts if they can also be described in terms of constraints. Several examples in the literature, like $HM(X)$ [15] and $OUTSIDEIN(X)$ [21], are

actually frameworks which can be parametrized by different constraint systems, hence the X in their names.

The modifications needed to accommodate the new kinds of constraints in GI are given in Figure 13. First, we split up the syntax of constraints into so-called simple and extended constraints. The former consists of constraints that can be used by programmers in their programs, while the second category consists of additional constraints that are internal to the solver. Syntax for any new kind of constraints can be added to the syntax of simple constraints Q ; we have in fact done so for the specific example of type classes. We also modify the syntax of polymorphic types so that they may contain a number of simple constraints.

Now, suppose that we need to check that $\forall a. Ord a \Rightarrow a \rightarrow Bool$ is an instance of $\forall a. Eq a \Rightarrow a \rightarrow Bool$, in other words $(\forall \{\alpha\}. Eq \alpha \Rightarrow \alpha \rightarrow Bool) \leq (\forall a. Ord a \Rightarrow a \rightarrow Bool)$. During this process we are allowed to *assume* that *Ord a* holds in order to discharge *Eq a*. To be able to store this information in our extended constraints, we modify the syntax of quantification constraints to include the assumed information as part of an implication, in this case

$$\forall a. (Ord a \supset \forall \{\alpha\}. Eq \alpha \Rightarrow \alpha \rightarrow Bool \leq a \rightarrow Bool)$$

Implication constraints are also introduced by the updated rules for constraint gathering. An annotated application – rule ANNAAPP – may also mention constraints, which are assumed while checking the enclosed expression. Rule CASE allows some constraints to be locally valid, which means that the system gains support for generalized algebraic data types (GADTs).

Figure 14 presents the updates needed for the solver to cope with this new form of constraints. The solver now relates two sets of constraints: the *assumptions* and the *wanted* constraints. Whenever we go under an implication, we introduce the constraints that are part of the antecedent as additional assumptions. The solver is allowed to rewrite in both sets, and more importantly, to use an assumed constraint to rewrite a wanted one. The modifications to the rules are very much in line with Vytiniotis et al. [21].

The rules that dealt with polymorphic types, namely INSTVL and INSTVR need to be updated. In the former case, a constraint in the type becomes an obligation to prove that it holds. In the latter case the constraints become assumptions that we are allowed to use while solving; these are stored in an implication constraint. Note that the guardedness restrictions do not change by the introduction of constraints.

Floating of constraints and promotion of variables also require changes, as described in rule FLOAT. In particular, we are only allowed to float equality constraints, which in the vanilla system were the only kind of simple constraints. The reason is that other kinds of constraints may incorporate information of the assumptions while they are solved, and these assumptions only hold in their respective branches. Equality constraints do not pose this problem.

$$\begin{array}{c}
\boxed{\Gamma \vdash e : \sigma \rightsquigarrow C} \\
\text{for each branch } K_i \bar{x}_i \rightarrow e_i \\
K_i : \forall \bar{a} \bar{b}_i. \bar{\sigma}_i \rightarrow \top \bar{a} \in \Gamma \\
\Gamma, x_i : [a \mapsto \alpha^u] \sigma_i \vdash e_i : \phi_i \rightsquigarrow C_i \\
\bar{v}_i = \text{fuv}(\phi_i, C_i) - \text{fuv}(\Gamma) - \bar{\alpha} \\
\hline
\Gamma \vdash e_0 \text{ of } \{\overline{K \bar{x}} \rightarrow e\} : \beta^u \rightsquigarrow C_0 \wedge (\sigma_0 \leq_{\epsilon}^m \epsilon; \top \bar{\alpha}^u) \wedge \forall \bar{b}_i. \exists \bar{v}_i. (C_i \wedge \beta^u \sim \phi_i) \quad \text{CASE}
\end{array}$$

Figure 12. Constraint generation for pattern matching

| | | |
|--------------------------------|---|---------------------------|
| Polymorphic types | $\sigma, \phi ::= \alpha^u \mid \forall \bar{a}. \boxed{Q} \Rightarrow \mu$ | |
| Simple constraints | $Q ::= \top \mid Q_1 \wedge Q_2$ | |
| | $\mid \sigma \sim \phi$ | Equality |
| | $\mid \dots$ | Open for extension |
| | $\mid C \sigma_1 \dots \sigma_n$ | Type classes, for example |
| Extended constraints | $C ::= Q \mid C_1 \wedge C_2$ | |
| Instantiation | $\mid \sigma \leq_{\omega}^s \bar{\phi}; \mu$ | |
| Generalisation | $\mid g \leq \sigma$ | |
| Quantification and implication | $\mid \forall \bar{a}. \exists \bar{v}. (\boxed{Q} \supset C)$ | |

$$\begin{array}{c}
\boxed{\Gamma \vdash e : \sigma \rightsquigarrow C} \\
\vdots \quad \vdots \\
\hline
\Gamma \vdash (h e_1 \dots e_n :: \forall \bar{b}. Q \Rightarrow \eta) : \forall \bar{b}. Q \Rightarrow \eta \rightsquigarrow \forall \bar{b}. \exists \bar{v}. (\boxed{Q} \supset C \wedge \dots) \quad \text{ANNAPP} \\
\text{for each branch } K_i \bar{x}_i \rightarrow e_i \\
K_i : \forall \bar{a} \bar{b}_i. \boxed{Q}_i \Rightarrow \bar{\sigma}_i \rightarrow \top \bar{a} \in \Gamma \\
\Gamma, x_i : [a \mapsto \alpha^u] \sigma_i \vdash e_i : \phi_i \rightsquigarrow C_i \\
\bar{v}_i = \text{fuv}(\phi_i, C_i) - \text{fuv}(\Gamma) - \bar{\alpha} \\
\hline
\Gamma \vdash e_0 \text{ of } \{\overline{K \bar{x}} \rightarrow e\} : \beta^u \rightsquigarrow C_0 \wedge (\sigma_0 \leq_{\epsilon}^m \epsilon; \top \bar{\alpha}^u) \wedge \forall \bar{b}_i. \exists \bar{v}_i. (\boxed{Q}_i \supset C_i \wedge \beta^u \sim \phi_i) \quad \text{CASE}
\end{array}$$

Figure 13. Extensions for integration with other constraints

The $C ; \bar{v} \Longrightarrow C' ; \bar{v}'$ relation only deals with one set of constraints. But as implications enter the game, we need to consider the interaction between two sets: the assumed and the wanted ones. For that matter we introduce a new rewriting judgment $Q_a ; C ; \bar{v} \Longrightarrow Q'_a ; C' ; \bar{v}'$, where assumed constraints Q_a are rewritten to Q'_a and wanted constraints C to C' . Note that set of assumed constraints consists only of simple constraints.

Rules ASSUMED and WANTED allow us to apply any rule to any of the two sets involved in solving. Rule INTERACT deals with the interaction of an assumption and a wanted constraint; the result is put on the wanted set. For example, if you have $\alpha \sim \text{Int} ; \beta \sim [\alpha] ; \alpha, \beta$, the rule moves to $\alpha \sim \text{Int} ; \beta \sim [\text{Int}] ; \alpha, \beta$. Rule DUPL is a simple case of interaction, in which a constraint in the wanted set is “crossed out” if it is already in the assumed set. Note that

information may only flow from assumptions to wanted constraints, and never the other way around.

The last rule, an updated version of FORALL, is responsible for dealing with implication constraints. The constraints C inside of the implication are rewritten in an environment where the set of assumed constraints is enlarged with the antecedent Q from the implication.

C From and to System F

Our system is as expressive as System F, the gold standard for a fully-expressive impredicative type system. We show so by describing a translation from every System F expression to GI in Figure 15. We also give the converse translation from GI terms to System F terms in Figure 16 (as usual, ϵ denotes an empty list). In particular, this proves that GI is sound as a type system. For the sake of conciseness we do not include the translations of pattern matching.

$$\begin{array}{c}
\boxed{C; \bar{v} \Longrightarrow C'; \bar{v}'} \\
\text{[INST}\forall\text{L]} \quad (\forall \bar{a}. \boxed{Q \Rightarrow \mu}) \leq_{\omega}^s \bar{\phi}; \eta; \bar{v} \Longrightarrow \boxed{[a \mapsto v']Q} \wedge (\mu' \leq_{\omega}^s \bar{\phi}; \eta); \bar{v}, \bar{v}' \\
\text{where } \text{freshen}_{\omega}^s(\forall \bar{a}. \mu) \Longrightarrow \langle \bar{v}', \mu' \rangle \\
\text{[INST}\forall\text{R]} \quad g \leq (\forall \bar{a}. \boxed{Q \Rightarrow \mu}); \bar{v} \Longrightarrow \forall \bar{a}. \boxed{Q \supset g \leq \mu}; \bar{v} \\
\text{[FLOAT]} \quad \frac{\begin{array}{c} F \text{ is an equality} \quad \text{ftv}(F) \cap \bar{a} = \emptyset \\ \bar{\alpha}^s = \text{fuv}(F) \cap v_{\text{in}} \quad \bar{\gamma}^s \text{ fresh} \quad \bar{E} = \bigwedge_{\alpha^s \in \bar{\alpha}} \alpha^s \sim \gamma^s \end{array}}{\forall \bar{a}. \exists \bar{v}_{\text{in}}. (Q \supset C \wedge F); \bar{v} \Longrightarrow [\bar{\alpha}^s \mapsto \bar{\gamma}^s]F \wedge \forall \bar{a}. \exists \bar{v}_{\text{in}}. (Q \supset C \wedge E); \bar{v}, \bar{\gamma}^s} \\
\boxed{Q_a; C; \bar{v} \Longrightarrow Q'_a; C'; \bar{v}'} \\
\text{[ASSUMED]} \quad \frac{Q_a; \bar{v} \Longrightarrow Q'_a; \bar{v}_{\text{forget}}}{Q_a; C; \bar{v} \Longrightarrow Q'_a; C'; \bar{v}} \\
\text{[WANTED]} \quad \frac{C; \bar{v} \Longrightarrow C'; \bar{v}'}{Q_a; C; \bar{v} \Longrightarrow Q_a; C'; \bar{v}'} \\
\text{[INTERACT]} \quad \frac{Q \wedge C^*; \bar{v} \Longrightarrow C'; \bar{v}'}{Q_a \wedge Q; C \wedge C^*; \bar{v} \Longrightarrow Q_a \wedge Q; C \wedge C'; \bar{v}'} \\
\text{[DUPL]} \quad \frac{Q_a \wedge Q; C \wedge Q; \bar{v} \Longrightarrow Q_a \wedge Q; C; \bar{v}}{Q_a \wedge Q; C; \bar{v}_{\text{in}} \Longrightarrow Q'_a; C'; \bar{v}'_{\text{in}}} \\
\text{[FORALL]} \quad \frac{Q_a \wedge Q; C; \bar{v}_{\text{in}} \Longrightarrow Q'_a; C'; \bar{v}'_{\text{in}}}{Q_a; \forall \bar{a}. \exists \bar{v}_{\text{in}}. (Q \supset C); \bar{v} \Longrightarrow Q_a; \forall \bar{a}. \exists \bar{v}'_{\text{in}}. (Q \supset C'); \bar{v}}
\end{array}$$

Figure 14. Solving rules for implication constraints

Theorem C.1 (Embedding of System F). *Let e be a System F expression. If $\Gamma \vdash^F e_F : \sigma \rightsquigarrow e'$, as defined in Figure 15, then $\Gamma \vdash e' : \sigma$ in GI.*

The main difference between GI and System F is that in the former impredicative instantiation is restricted by guardedness. The presented translation relies on annotations to work around them.

- Following GI we present n -ary application – which in the case of System F also includes type application. The given application rule only works if the guardedness restrictions are satisfied, otherwise an annotation on e_0 needs to be added before applying the rule.
- In GI the result of a non-annotated application always gets a top-level monomorphic type. This is not the case in System F, and thus an annotation may be required to further generalize.
- Completely unrestricted instantiation is only available via annotations. In order to apply this translation, we need to split applications so that guardedness guarantees are always met. Every time we split, we introduce a new annotation guiding the type checking process.

Proof. By induction on the typing derivation in System F, $\Gamma \vdash^F e : \sigma$.

Variable. We need to derive $\Gamma \vdash x : \sigma$, for a general $\sigma = \forall \bar{a}. \mu$. We recall that single variables are treated like 0-ary application, so some amount of instantiation must take place.

$$\frac{\frac{x : \sigma \in \Gamma}{\Gamma \vdash^{\text{fun}} x : \sigma} \text{VARHEAD} \quad \sigma \leq_{\epsilon}^m \epsilon; \mu}{\Gamma \vdash x : \mu} \text{APP}$$

If the set of quantified variables \bar{a} is empty, then this derivation is all we need. If it is not, we need an annotation to re-generalise those.

$$\frac{\frac{x : \sigma \in \Gamma}{\Gamma \vdash^{\text{fun}} x : \sigma} \text{VARHEAD} \quad \sigma \leq_{\epsilon}^u \epsilon; \mu}{\Gamma \vdash (x :: \forall \bar{a}. \mu) : \forall \bar{a}. \mu} \text{ANNAPP}$$

We could get a smaller translation in some cases by using the **VARGEN** rule, but the simpler type system with annotations is enough for our goals.

Abstraction. By induction hypothesis we are able to type check the body of the abstraction. The distinction between fully monomorphic and unrestricted types in the translation ensures that we can use the right **ABS** or **ANNABS** rule from the declarative specification.

$$\begin{array}{c}
\text{System F terms } e_F ::= x \mid \lambda(x :: \sigma). e_F \mid \Lambda a. e_F \mid e_F e_F \mid e_F \sigma \\
\boxed{\Gamma \vdash^F e_F : \sigma \rightsquigarrow e'} \\
\frac{x : \mu \in \Gamma}{\Gamma \vdash^F x : \mu \rightsquigarrow x} \qquad \frac{x : \sigma \in \Gamma}{\Gamma \vdash^F x : \sigma \rightsquigarrow (x :: \sigma)} \\
\frac{\Gamma, x : \tau \vdash^F e_F : \phi \rightsquigarrow e'}{\Gamma \vdash^F \lambda(x :: \tau). e_F : \tau \rightarrow \phi \rightsquigarrow \lambda x. e'} \qquad \frac{\Gamma, x : \sigma \vdash^F e_F : \phi \rightsquigarrow e'}{\Gamma \vdash^F \lambda(x :: \sigma). e_F : \sigma \rightarrow \phi \rightsquigarrow \lambda(x :: \sigma). e'} \\
\frac{\Gamma \vdash^F e_F : \phi[a_1, \dots, a_n] \rightsquigarrow e'}{\Gamma \vdash^F \Lambda a_1 \dots \Lambda a_n. e_F : \forall a_1. \dots \forall a_n. \phi \rightsquigarrow ([e'] :: \forall a_1 \dots a_n. \phi)} \\
\Gamma \vdash^F e_0 : \sigma \rightsquigarrow e'_0 \quad \sigma_0 = \forall \bar{a}_1. \sigma_1 \rightarrow \forall \bar{a}_2. \sigma_2 \rightarrow \dots \forall \bar{a}_n. \sigma_n \rightarrow \forall \bar{a}_r \bar{a}_p. \mu_r \\
\sigma_0 \leq_{\epsilon}^{\text{u}} \underbrace{\bullet, \dots, \bullet}_{n \text{ times}} \quad \underbrace{[a_1 \mapsto \bar{\phi}_1] \sigma_1, [a_1 \mapsto \bar{\phi}_1, a_2 \mapsto \bar{\phi}_2] \sigma_2, \dots, [a_1 \mapsto \bar{\phi}_1, \dots, a_n \mapsto \bar{\phi}_n; a_r \mapsto \bar{\phi}_r] \mu_r}_{n \text{ times}} \\
\frac{\Gamma \vdash^F e_1 : [a_1 \mapsto \bar{\phi}_1] \sigma_1 \rightsquigarrow e'_1 \quad \dots \quad \Gamma \vdash^F e_n : [a_1 \mapsto \bar{\phi}_1, \dots, a_n \mapsto \bar{\phi}_n] \sigma_n \rightsquigarrow e'_n}{\Gamma \vdash^F e_0 \bar{\phi}_1 e_1 \bar{\phi}_2 e_2 \dots \bar{\phi}_n e_n \bar{\phi}_r : \forall \bar{a}_r. \mu_r \rightsquigarrow \begin{cases} e'_0 e'_1 \dots e'_n & \text{if } \bar{a}_p \text{ is empty} \\ & \text{and } \bar{\phi}_r \text{ fully mono.} \\ (e'_0 e'_1 \dots e'_n :: \forall \bar{a}_r. \mu_r) & \text{otherwise} \end{cases}} \\
\text{where } \begin{cases} [e :: \sigma] = [e] \\ [e] = e & \text{otherwise} \end{cases}
\end{array}$$

Figure 15. Translation from System F

Type abstraction. There are two cases to consider:

- $[e']$ is not an application. In this case the annotation is treated as an annotated application with zero arguments. By induction hypothesis, we know that $\Gamma \vdash e' : \phi[\bar{a}]$. We assume without loss of generality that $\phi = \forall \bar{b}. \mu$, and thus $\Gamma \vdash e' : \forall \bar{b}. \mu[\bar{a}]$. Now we can build the following derivation:

$$\frac{\frac{\Gamma \vdash e' : \forall \bar{b}. \mu[\bar{a}]}{\Gamma \vdash^{\text{fun}} e' : \forall \bar{b}. \mu[\bar{a}]} \quad \forall \bar{b}. \mu[\bar{a}] \leq_{\epsilon}^{\text{u}} \epsilon; \mu[\bar{a}]}{\Gamma \vdash (e' :: \forall \bar{a} \bar{b}. \mu) : \forall \bar{a} \bar{b}. \mu} \text{ANNAPP}$$

- $[e']$ is an application. In this case e' results from the application of either APP or ANNAPP in the declarative specification. By inspection of the rule ANNAPP, we can see that we can always choose to quantify over more variables via an annotation.

(Type) application. The premises about guardedness ensure that we can apply the rule APP or ANNAPP in the declarative specification. As discussed in the main text, if at a certain application we cannot apply this rule, we can always split the application, annotate the head and then use the application rule again. \square

D Proofs

D.1 Declarative specification

Lemma D.1. *Let Γ be an environment and e an expression. For every pair of fully monomorphic substitutions θ_1 and θ_2 , if $\theta_1 \Gamma \vdash^{\text{fun}} e : \sigma_1$ and $\theta_2 \Gamma \vdash^{\text{fun}} e : \sigma_2$, then there exists a polymorphic type σ^* and fully monomorphic substitutions φ_1 and φ_2 such that $\sigma_i = \varphi_i \sigma^*$.*

Theorem 3.2 (Impredicative instantiation is not guessed). *Let Γ be an environment and e an expression. For every pair of fully monomorphic substitutions θ_1 and θ_2 , if $\theta_1 \Gamma \vdash e : \sigma_1$ and $\theta_2 \Gamma \vdash e : \sigma_2$, then there exists a polymorphic type σ^* and fully monomorphic substitutions φ_1 and φ_2 such that $\sigma_i = \varphi_i \sigma^*$.*

Proof. We prove Lemma D.1 and Theorem 3.2 by mutual induction over the typing derivation of e .

Proof of Lemma D.1. We distinguish two cases:

- *Case VARHEAD.* The two derivations to consider are:

$$\frac{x : \theta_1 \sigma \in \theta_1 \Gamma}{\theta_1 \Gamma \vdash^{\text{fun}} x : \theta_1 \sigma} \qquad \frac{x : \theta_2 \sigma \in \theta_2 \Gamma}{\theta_2 \Gamma \vdash^{\text{fun}} x : \theta_2 \sigma}$$

Thus we can take $\sigma^* = \sigma$ and as substitutions those applied to Γ , which are fully monomorphic by our assumptions.

- *Case EXPRHEAD.* Follows by induction over the premise.

$$\boxed{\sigma \leq_{\omega}^s \bar{\phi}; \mu \rightsquigarrow \bar{\psi}_1, \dots, \bar{\psi}_n, \bar{\psi}_r}$$

$$\frac{}{\mu \leq_{\epsilon}^s \epsilon; \mu \rightsquigarrow \epsilon} \text{INSTMONO}$$

$$\frac{\phi_2 \leq_{\omega}^s \sigma_2, \dots, \sigma_n; \mu \rightsquigarrow \bar{\psi}_2, \dots, \bar{\psi}_n, \bar{\psi}_r}{\phi_1 \rightarrow \phi_2 \leq_{\omega, \bar{\omega}}^s \phi_1, \sigma_2, \dots, \sigma_n; \mu \rightsquigarrow \epsilon, \bar{\psi}_2, \dots, \bar{\psi}_n, \bar{\psi}_r} \text{INSTARROW}$$

$$\frac{\forall \bar{a}. \mu \triangleright_{\omega}^s \Delta \quad \theta \text{ respects } \Delta \quad \theta \mu \leq_{\omega}^s \sigma_1, \dots, \sigma_n; \eta \rightsquigarrow \epsilon, \bar{\psi}_2, \dots, \bar{\psi}_n, \bar{\psi}_r}{\forall \bar{a}. \mu \leq_{\omega}^s \sigma_1, \dots, \sigma_n; \eta \rightsquigarrow \theta(\bar{a}), \bar{\psi}_2, \dots, \bar{\psi}_n, \bar{\psi}_r} \text{INSTPOLY}$$

$$\boxed{\Gamma \vdash^{\text{fun}} e : \sigma \rightsquigarrow e_{\text{F}}}$$

$$\frac{x : \sigma \in \Gamma}{\Gamma \vdash^{\text{fun}} x : \sigma \rightsquigarrow x} \text{VARHEAD}$$

$$\frac{\Gamma \vdash e : \sigma \rightsquigarrow e_{\text{F}}}{\Gamma \vdash^{\text{fun}} e : \sigma \rightsquigarrow e_{\text{F}}} \text{EXPRHEAD}$$

$$\boxed{\Gamma \vdash_{\omega}^{\text{arg}} e : \sigma \rightsquigarrow e_{\text{F}}}$$

$$\frac{\Gamma \vdash e : \forall \bar{a}. \mu \rightsquigarrow e_{\text{F}} \quad \bar{b} \notin \Gamma}{\Gamma \vdash_{\bullet}^{\text{arg}} e : \forall \bar{b}. [\bar{a} \mapsto \bar{\tau}] \mu \rightsquigarrow \Lambda \bar{b}. e_{\text{F}} \bar{\tau}} \text{ARGGEN}$$

$$\frac{x : \forall \bar{p}. \tau \in \Gamma \quad \bar{b} \notin \Gamma}{\Gamma \vdash_{\star}^{\text{arg}} x : \forall \bar{b}. [\bar{a} \mapsto \bar{\sigma}] \tau \rightsquigarrow \Lambda \bar{b}. x \bar{\sigma}} \text{VARGEN}$$

$$\boxed{\Gamma \vdash e : \sigma \rightsquigarrow e_{\text{F}}}$$

$$\frac{\Gamma, x : \tau \vdash e : \sigma \rightsquigarrow e_{\text{F}}}{\Gamma \vdash \lambda x. e : \tau \rightarrow \sigma \rightsquigarrow \lambda(x :: \tau). e_{\text{F}}} \text{ABS}$$

$$\frac{\Gamma, x : \phi \vdash e : \sigma \rightsquigarrow e_{\text{F}}}{\Gamma \vdash \lambda(x :: \phi). e : \phi \rightarrow \sigma \rightsquigarrow \lambda(x :: \sigma). e_{\text{F}}} \text{ANNABS}$$

$$\frac{\Gamma \vdash^{\text{fun}} e_0 : \phi \rightsquigarrow e_{0, \text{F}} \quad \phi \leq_{\omega_1, \dots, \omega_n}^m \sigma_1, \dots, \sigma_n; \mu \rightsquigarrow \bar{\psi}_1, \dots, \bar{\psi}_n, \bar{\psi}_r \quad \Gamma \vdash_{\omega_1}^{\text{arg}} e_1 : \sigma_1 \rightsquigarrow e_{\text{F}, 1} \quad \dots \quad \Gamma \vdash_{\omega_n}^{\text{arg}} e_n : \sigma_n \rightsquigarrow e_{\text{F}, n}}{\Gamma \vdash e_0 e_1 \dots e_n : \mu \rightsquigarrow e_{0, \text{F}} \bar{\psi}_1 e_{\text{F}, 1} \dots \bar{\psi}_n e_{\text{F}, n} \bar{\psi}_r} \text{APP}$$

$$\frac{\Gamma \vdash^{\text{fun}} e_0 : \phi \rightsquigarrow e_{0, \text{F}} \quad \phi \leq_{\omega_1, \dots, \omega_n}^u \sigma_1, \dots, \sigma_n; \eta \rightsquigarrow \bar{\psi}_1, \dots, \bar{\psi}_n, \bar{\psi}_r \quad \Gamma \vdash_{\omega_1}^{\text{arg}} e_1 : \sigma_1 \rightsquigarrow e_{\text{F}, 1} \quad \dots \quad \Gamma \vdash_{\omega_n}^{\text{arg}} e_n : \sigma_n \rightsquigarrow e_{\text{F}, n}}{\Gamma \vdash (e_0 e_1 \dots e_n :: \forall \bar{b}. \eta) : \forall \bar{b}. \eta \rightsquigarrow \Lambda \bar{b}. h_{\text{F}} \bar{\psi}_1 e_{\text{F}, 1} \dots \bar{\psi}_n e_{\text{F}, n} \bar{\psi}_r} \text{ANNAPP}$$

$$\frac{\Gamma \vdash e_1 : \phi \rightsquigarrow e_{\text{F}, 1} \quad \Gamma, x : \phi \vdash e_2 : \sigma \rightsquigarrow e_{\text{F}, 2}}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \sigma \rightsquigarrow (\lambda(x :: \phi). e_{\text{F}, 2}) e_{\text{F}, 1}} \text{LET}$$

Figure 16. Translation to System F

Proof of Theorem 3.2. We distinguish seven cases:

- *Case ABS.* The derivations look like:

$$\frac{\theta_1\Gamma, x : \tau_1 \vdash e : \sigma_1}{\theta_1\Gamma \vdash \lambda x.e : \tau_1 \rightarrow \sigma_1} \quad \frac{\theta_2\Gamma, x : \tau_2 \vdash e : \sigma_2}{\theta_2\Gamma \vdash \lambda x.e : \tau_2 \rightarrow \sigma_2}$$

We cannot apply the induction hypothesis yet, since the environments in the premises are not of the right shape. Consider instead the environment $\Gamma' = \Gamma, x : \alpha$ for a fresh α , and the substitutions

$$\theta'_1 = [\alpha \mapsto \tau_1] \circ \theta_1 \quad \theta'_2 = [\alpha \mapsto \tau_2] \circ \theta_2$$

These substitutions are fully monomorphic, since all of θ_i and τ_i ($1 \leq i \leq 2$) are fully monomorphic by hypothesis. We have then the following equalities over the environments

$$\theta_1\Gamma, x : \tau_1 = \theta'_1\Gamma' \quad \theta_2\Gamma, x : \tau_2 = \theta'_2\Gamma'$$

and thus we can apply the induction hypothesis to e to obtain σ^* , φ_1 and φ_2 such that $\sigma_i = \varphi_i\sigma^*$. Consider now the extended substitutions:

$$\varphi'_1 = [\alpha \mapsto \tau_1] \circ \varphi \quad \varphi'_2 = [\alpha \mapsto \tau_2] \circ \varphi_2$$

Since τ_1 and τ_2 are fully monomorphic types, φ'_i are fully monomorphic substitutions. Take $\sigma' = \alpha \rightarrow \sigma^*$. We have, for each derivation, that

$$\varphi'_i\sigma' = \varphi'_i(\alpha \rightarrow \sigma^*) = \varphi'_i\alpha \rightarrow \varphi'_i\sigma^* = \tau_i \rightarrow \sigma_i$$

And we are done: σ' , φ'_1 and φ'_2 are the desired outputs of the theorem.

- *Case ANNAbs.* Similar to *ABS*.
- *Case APP.* In this case the derivations have the following shape for each $i \in \{1, 2\}$.

$$\frac{\theta_i\Gamma \vdash^{\text{fun}} e_0 : \theta_i\phi \quad \frac{\phi \leq^m_{\omega} \sigma_{i,1}, \dots, \sigma_{i,n}; \mu_i \quad \theta_i\Gamma \vdash^{\text{arg}}_{\omega_j} e_j : \sigma_{i,j}}{\theta_i\Gamma \vdash e_0 e_1 \dots e_n : \mu_i}}{\theta_i\Gamma \vdash e_0 e_1 \dots e_n : \mu_i}$$

We have two choices: either ϕ is a top-level monomorphic type η , or $\theta_i\phi = \sqrt{a}. \theta_i\eta$. This shall remain true during the derivation of \leq^m_{ω} . Intuitively, the substitutions do not alter the “polymorphism structure” of the type. In particular, as result type of the instantiation, we get μ_1 and μ_2 for which we know that there exists a common μ^* such that:

$$\mu_1 = \varphi_1^m \varphi_1^r \theta_1 \mu^* \quad \mu_2 = \varphi_2^m \varphi_2^r \theta_2 \mu^*$$

where φ_i^m correspond to the fully monomorphic substitutions applied to those type variables which do not appear in any of the arguments. But we should be careful here: φ_i^r are unrestricted, and for those we do not directly obtain the conclusion of the theorem.

Let us first consider the case in which $\omega_j = \bullet$ for all j . In other words, all arguments are typed using the *ARGGEN* rule for \vdash^{arg} . By the definition of the respects relation for substitutions, for each quantified type variable a in ϕ we have *at least one* $\sigma_{i,j}$ for which $\varphi_i^r \theta_i \sigma_{i,j}$

is of the form $\sqrt{a}. \top \bar{\psi}$ for both derivations. We can see this by distinguishing between unrestricted variables, for which we already have the type constructor in ϕ ; and top-level monomorphic variables for which the corresponding substitution must map to top-level monomorphic types.

Take the expression corresponding to this type variable a – let us call it e_a – and the derivations $\theta_1\Gamma \vdash e_a : \varphi_1^r \theta_1 \sigma_{1,a}$ and $\theta_2\Gamma \vdash e_a : \varphi_2^r \theta_2 \sigma_{2,a}$. By the induction hypothesis there exists a σ_j^* and fully monomorphic substitutions $\xi_{i,a}$ such that $\varphi_i^r \theta_i \sigma_{i,a} = \xi_{i,a} \sigma_a^*$.

We know that for each type variable a appearing in some argument of the application we have at least one ξ_a so that a is in its domain. It does not matter which one we choose: they all have to agree or otherwise the derivation is ill-formed. Take now the following substitutions:

$$\xi_i = \bigcirc_{a \in \text{args}^n(\phi)} [a \mapsto \xi_{i,a}(a)],$$

where \bigcirc denotes substitution composition. By construction, we know that $\varphi_i^r = \xi_i$: they have the same domain and they map to the same outputs. But now we are sure that ξ_i are fully monomorphic, since arise from a composition of fully monomorphic substitutions. In conclusion, we have that

$$\mu_1 = \varphi_1^m \xi_1 \theta_1 \mu^* \quad \mu_2 = \varphi_2^m \xi_2 \theta_2 \mu^*$$

and all of φ_i^m , ξ_i and θ_i are fully monomorphic.

When we use *VARGEN* we need to work a bit more. If for each type variable there is one e_a where *ARGGEN* has been applied, the proof still works. On the other hand, if for a type variable a all expressions where it appears make use of *VARGEN*, we cannot guarantee the existence of the σ_j^* and fully monomorphic substitutions $\xi_{i,a}$ by induction hypothesis. However, in that case the definition of $\triangleright^s_{\omega}$ ensures that the substitution for the type variable a is fully monomorphic, giving us the desired result.

- *Case ANNApP.* This case is trivial, since the annotation ensures that both derivations yield the same type.
- *Case LET.* In this case the derivations are:

$$\frac{\theta_i\Gamma \vdash e_1 : \phi_i \quad \theta_i\Gamma, x : \phi_i \vdash e_2 : \sigma_i}{\theta_i\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \sigma_i} \text{ for } i \in \{1, 2\}$$

We cannot readily apply the induction hypothesis to the second premise – which would give us the desired conclusion –, since the environment may not be in the right shape. By the induction hypothesis on the first premise, there exists a type ϕ^* and fully monomorphic substitutions φ_1 and φ_2 such that $\phi_i = \varphi_i \phi^*$. We can build a new version of ϕ^* , ϕ' , where all the free variables are fresh – and thus disjoint from those in Γ – and corresponding φ'_i where the domain is replaced

by these free variables. Consider $\theta'_i = \theta_i \circ \varphi'_i$:

$$\begin{aligned} \theta'_i(\Gamma, x : \phi') &= \theta'_i(\Gamma), x : \theta'_i \phi' \\ &= (\text{substitutions have disjoint domains}) \\ &\quad \theta_i \Gamma, x : \varphi'_i \phi' \\ &= (\text{by definition of } \varphi'_i) \\ &\quad \theta_i \Gamma, x : \phi_i \end{aligned}$$

This means that we have found a fully monomorphic substitution θ'_i for the environment, which allows us to apply induction on the second premise and reach the desired conclusion.

- *Case CASE.* By induction on the first premise of CASE, we get σ_0^* and two substitutions φ_1 and φ_2 such that $\sigma_{0,i} = \varphi_i \sigma_0^*$ for each of the two derivations. Without loss of generality, we can assume that $\sigma_0^* = \forall \bar{a}. \mu_0^*$. Then $\sigma_{0,i} = \forall \bar{a}. \varphi_i \mu_0^*$.

The next step is to notice that the instantiation judgment is nullary. As a result, all the type variables in \bar{a} must be substituted by fully monomorphic types. Let π_i be the substitution for each of the derivations. Then we know that $\top \overline{\phi_{0,i}} = \pi_i \varphi_i \mu_0^* = \top \overline{\pi_i \varphi_i \phi_0^*}$. Playing the same game we did for LET, we can build substitutions θ'_i such that,

$$\theta'_i(\Gamma, x_i : [a \mapsto \phi_{0,i}] \sigma_j) = \theta_i \Gamma, x_i : [a \mapsto \pi_i \varphi_i \phi_0^*] \sigma_j$$

The last step is applying induction hypothesis to one of the branches and obtain the desired common polymorphic type and fully monomorphic instantiations. The branch we choose does not matter, since they all give the same type as result. \square

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

Proof. By induction over the expression $e[x]$. The only interesting case is when x is the head of an application, that is, $e = x e_1 \dots e_n$ and u is also an application, $u = u_0 u_1 \dots u_n$. Note that in that case the assigned types are always top-level monomorphic, so $\sigma = \mu$ and $\phi = \eta$.

The derivations involved look like:

$$\frac{\Gamma \vdash^{\text{fun}} u_0 : \sigma_u \quad \frac{\sigma_u \leq_{\omega_1, \dots, \omega_m}^m \sigma_1, \dots, \sigma_m; \mu \quad \Gamma \vdash_{\omega_i}^{\text{arg}} u_i : \sigma_i}{\Gamma \vdash u_0 u_1 \dots u_m : \mu}}{\Gamma \vdash^{\text{fun}} x : \mu} \quad \frac{\mu \leq_{\omega'_1, \dots, \omega'_n}^m \phi_1, \dots, \phi_n; \eta \quad \Gamma \vdash_{\omega'_j}^{\text{arg}} e_j : \phi_j}{\Gamma \vdash x e_1 \dots e_n : \eta}$$

and the question is whether we can derive:

$$\frac{\Gamma \vdash^{\text{fun}} u_0 : \sigma_u \quad \frac{\sigma_u \leq_{\omega_1, \dots, \omega_m, \omega'_1, \dots, \omega'_n}^m \sigma_1, \dots, \sigma_m, \phi_1, \dots, \phi_n; \eta \quad \Gamma \vdash_{\omega_i}^{\text{arg}} u_i : \sigma_i \quad \Gamma \vdash_{\omega'_j}^{\text{arg}} e_j : \phi_j}{\Gamma \vdash u_0 u_1 \dots u_m e_1 \dots e_n : \eta}}{\Gamma \vdash u_0 u_1 \dots u_m e_1 \dots e_n : \eta}$$

All the premises on this last rule come directly from those in the hypotheses, except for the instantiation

$$\sigma_u \leq_{\omega}^m \sigma_1, \dots, \sigma_m, \phi_1, \dots, \phi_n; \eta$$

For the n last components we can reuse the derivation in the second hypothesis. For the first m components, inspection on the rules of guardedness show that any instantiation with m arguments is admissible when $m + n$ are considered. \square

Theorem 3.5. *Let $\text{app}::\forall a b. (a \rightarrow b) \rightarrow a \rightarrow b$ and $\text{revapp}::\forall a b. a \rightarrow (a \rightarrow b) \rightarrow b$ be the application and reverse application functions, respectively. Given two expressions f and e such that $\Gamma \vdash^{\text{fun}} f : \sigma_0$, and $\sigma_0 \leq_{\epsilon}^m \epsilon; \sigma_1 \rightarrow \phi$ then:*

$$\Gamma \vdash f e : \phi \iff \Gamma \vdash \text{app } f e : \phi \iff \Gamma \vdash \text{revapp } e f : \phi$$

Proof. Let us compare the derivations of the first two.

$$\frac{\Gamma \vdash^{\text{fun}} f : \sigma_0 \quad \sigma_0 \leq_{\omega_e}^m \sigma_1; \mu \quad \Gamma \vdash_{\omega_e}^{\text{arg}} e : \sigma_1}{\Gamma \vdash f e : \mu} \quad \frac{\Gamma \vdash^{\text{fun}} \text{app} : \forall a b. (a \rightarrow b) \rightarrow a \rightarrow b \quad \Gamma \vdash_{\omega_f}^{\text{arg}} f : \sigma_1 \rightarrow \mu \quad \Gamma \vdash_{\omega_e}^{\text{arg}} e : \sigma_1 \quad \forall a b. (a \rightarrow b) \rightarrow a \rightarrow b \leq_{\omega_f, \omega_e}^m (\sigma_1 \rightarrow \mu), \sigma_1; \mu}{\Gamma \vdash f e : \mu}$$

The application of app can be instantiated with any type, given that both variables a and b are guarded. We need to consider the two different ways in which $\Gamma \vdash_{\omega_f}^{\text{arg}} f : \sigma_1 \rightarrow \mu$ can be derived. If it is derived using ARGGEN, then we have:

$$\frac{\Gamma \vdash f : \sigma_0 \quad \sigma_0 \leq_{\epsilon}^m \epsilon; \sigma_1 \rightarrow \phi}{\Gamma \vdash_{\bullet}^{\text{arg}} f : \sigma_1 \rightarrow \phi}$$

Then the premises for the whole derivation are the same, except for the highlighted one. We know that if $\sigma_0 \leq_{\epsilon}^m \epsilon; \sigma_1 \rightarrow \phi$, then $\sigma_0 \leq_{\bullet}^m \sigma_1; \phi$, by inspection of the rule relating instantiation and function types.

The other possibility is that f is a variable and the rule for single variables, VARGEN, applies:

$$\frac{f : \forall \bar{p}. \tau \in \Gamma \quad \forall \bar{p}. \tau \leq_{\epsilon}^n \epsilon; \sigma_1 \rightarrow \phi}{\Gamma \vdash_{\star}^{\text{arg}} f : \sigma_1 \rightarrow \phi}$$

In this case VARGEN does not rule our polymorphic instantiation. However, the \leq judgment applied to app will ignore this argument for guardedness purposed. That means that we can only instantiate impredicative those type variables coming for e , exactly the same as we could with $f e$.

The proof for revapp is similar to the one for app . \square

D.2 Constraint-based formulation

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.*

Lemma D.2. *Suppose $\Gamma \vdash^{\text{fun}} e : \sigma \rightsquigarrow C$. Then C is either inconsistent, or can be rewritten to a new set C' without instantiation or generalisation constraints which fixes the value of all unrestricted and top-level monomorphic variables.*

Proof. We prove both theorems by mutual induction. Note that Lemma D.2 follows simply from Theorem 4.1, as the head typing judgment either produces no constraints with the rule VARHEAD or refers to the normal typing judgment as a premise with the rule EXPRHEAD.

For Theorem 4.1, we proceed by induction on the constraint generation judgment $\Gamma \vdash e : \sigma \rightsquigarrow C$.

Case ABS and ANNABS. Follows directly from induction hypothesis, since the constraints from the body are copied as the output.

Case APP. We need to consider three sets of constraints: the set C coming from generating constraints for the head, the set of instantiation constraints C_{\leq} which extract the expected types of arguments from the function types, and the set of generalisation constraints C_{\leq} which check that the actual arguments fit into those expected types. The third set is obtained from the ancillary judgement \vdash^{arg} .

The desired properties for the first set of constraints, C , follow by induction. Since C fixes all the unrestricted and top-level monomorphic variables, this implies that ϕ is fixed up to fully monomorphic variables. As a result, we can rewrite all the constraints in C_{\leq} into equalities – but note that instantiations may introduce new variables to be fixed.

First consider the case in which all applications of \vdash^{arg} use the ARGGEN rule. In that case, the newly-introduced variables can be divided between unrestricted \bar{u} , top-level monomorphic \bar{t} , and fully monomorphic \bar{m} . Only the first two interact with arguments, by definition. For the first set we know that for each variable α there is at least one argument of the form $\forall \bar{a}. \top \bar{\phi}$, where $\alpha \in \text{ftv}(\bar{\phi})$. For that argument there is a corresponding expression e_j and a generalisation constraint,

$$(\forall \{\bar{y}_j\}. C_j \Rightarrow \sigma_j) \leq (\forall \bar{a}. \top \bar{\phi})$$

Rule INST \forall R and INST \forall L apply, rewriting this constraint to $C_j \wedge \sigma_j \leq_{\epsilon}^{\text{m}} \epsilon; \top \bar{\phi}$, possibly under an universal quantifier $\forall \bar{b}$. Now we can apply the induction hypothesis to C_j , which means that we can rewrite these constraints to a set of type equalities which fixes the unrestricted and top-level monomorphic variables.

We still have a remaining constraint $\sigma_j \leq_{\epsilon}^{\text{m}} \epsilon; \top \bar{\phi}$ and we have not proven yet that the variable $\alpha \in \text{ftv}(\bar{\phi})$ is fixed. But

since σ_j is fixed up to fully monomorphic components, we can decide which instantiation rule to apply; in either case a type equality is produced. In turn, this equality fixes the value for α . As a result, the generalisation constraint can be completely turned into equalities, as desired. Furthermore, this value of α is floated out of the universal quantification – if it was ever introduced – up to the point where the variable was introduced.

Once the value of α is fixed – because of its appearance under a type constructor –, we can deal with those cases in which the variable appears alone as an argument,

$$(\forall \{\bar{y}_k\}. C_k \Rightarrow \sigma_k) \leq \alpha$$

But now we already know the type for α up to its fully monomorphic components! Thus, we are able to decide which rule in the solver to apply, and then apply the induction hypothesis. The end result is again a set of equalities, maybe under a universal quantifier.

Arguments which feature a top-level monomorphic variable $\beta \in \bar{t}$ are dealt with in the same way; the fact that the variable is top-level monomorphic is enough to unwrap the generalisation constraint. Then, the solving proceeds as with α unrestricted.

Finally, no constraints are generated over the set of fully monomorphic types \bar{m} . This is OK, since these variables are already fixed up to fully monomorphic components, by definition. This case also applies when all the expressions featuring a variable α result from the application of VARGEN.

Case ANNAPP. This case is almost identical to APP. The only difference is that now there is a set of variables \bar{u}' which would have been classified as either top-level or fully monomorphic which are now classified as unrestricted.

As in the case of APP, we apply the induction hypothesis over the head to guarantee that the set of constraints C can be turned into a set of equalities under a mixed prefix. Since these equalities fix σ up to fully monomorphic components, we can still apply the same reasoning to the set of instantiation constraints C_{\leq} . In particular, the constraint $\beta_n \leq_{\epsilon}^{\text{u}} \epsilon; \eta$ is turned into an equality, too, and fixes the value of all those \bar{u}' variables – those are the ones appearing in the assignment to β_n and η is completely known since it is explicitly given by the programmer.

For the remaining variables we apply the same reasoning as before, distinguishing between unrestricted and top-level monomorphic variables and going through each of the arguments. \square

D.2.1 Solved form and solutions

In Section 4.4 we introduced the notion of a *solved form*. In the following results we also need the notion of when a solved form C_s is a *solution* for a set of constraints C . We give the corresponding judgment $C_s \models C$ in Figure 17.

$$\begin{array}{c}
\boxed{C_s \models C} \quad \boxed{C_s \models' C} \\
\frac{C_s \models' \widehat{C}_s(C)}{C_s \models C} \\
\frac{}{C_s \models' \sigma \sim \sigma} \quad \frac{C_s \models' C_1 \quad C_s \models' C_2}{C_s \models' C_1 \wedge C_2} \quad \frac{\forall \bar{b}. C'_s \in C_s \quad C'_s \models C}{C_s \models' \forall \bar{b}. C} \\
\frac{C_s \models \mu \sim \eta}{C_s \models' \leq_{\epsilon}^s \epsilon; \eta} \quad \frac{C_s \models \mu \sim \sigma \rightarrow \sigma' \quad C_s \models \sigma' \leq_{\omega}^s \bar{\phi}; \eta}{C_s \models' \mu \leq_{\omega, \bar{\omega}}^s \sigma, \bar{\phi}; \eta} \quad \frac{\bar{a} \text{ were instantiated to } \bar{\alpha} \quad C_s \models [\bar{a} \mapsto \bar{\alpha}] \mu \leq_{\omega}^s \bar{\phi}; \eta}{C_s \models' \forall \bar{a}. \mu \leq_{\omega}^s \bar{\phi}; \eta} \\
\frac{C_s \models C \quad C_s \models \sigma \leq_{\epsilon}^m \epsilon; \eta}{C_s \models' (\forall \{\bar{\alpha}\}. C \Rightarrow \sigma) \leq \eta} \quad \frac{C_s \models' \forall \bar{b}. (\forall \{\bar{\alpha}\}. C \Rightarrow \sigma \leq \eta)}{C_s \models' (\forall \{\bar{\alpha}\}. C \Rightarrow \sigma) \leq (\forall \bar{b}. \eta)}
\end{array}$$

Figure 17. Definition of solution

Given two solutions C_s and D_s which range over the same set of variables, we say that C_s is *more general* than D_s if there exists a fully monomorphic substitution π such that for every pair of assignments $\alpha \sim \sigma$ in C_s , and $\alpha \sim \phi$ in D_s , $\phi = \pi\sigma$.

The solver described in [15] returns most general solutions for a given set of equalities under a mixed prefix. This is a consequence of the fact that they reuse the first-order solver in [11], which always returns most general substitutions.

D.2.2 Soundness and principality

In the proofs of the following results we focus on the core language without **let**, **do** and pattern matching. Nevertheless, the results still hold when the language is extended: no surprises there.

Lemma D.3. *Suppose that C_s is a solution for the constraint:*

$$\sigma \leq_{\omega}^s \alpha_1^{\#}, \dots, \alpha_n^{\#}; \delta^{\dagger}$$

Then we have a derivation for:

$$\widehat{C}_s(\sigma) \leq_{\omega}^s \widehat{C}_s(\alpha_1^{\#}), \dots, \widehat{C}_s(\alpha_n^{\#}); \widehat{C}_s(\delta)$$

(as defined in the declarative specification)

Proof. The definition for \models correspond one to one to the rules for \leq in the declarative specification. \square

Theorem D.4 (Soundness, solution version). *Let Γ be an environment and e an expression. Then, for every set of constraints C' ,*

1. *If $\Gamma \vdash^{\text{fun}} e : \sigma \rightsquigarrow C$, and $C_s \models C \wedge C'$, then we can build a derivation for $\widehat{C}_s(\Gamma) \vdash^{\text{fun}} e : \widehat{C}_s(\sigma)$.*
2. *If $\Gamma \vdash e : \sigma \rightsquigarrow C$ and $C_s \models C \wedge C'$, then we can build a derivation for $\widehat{C}_s(\Gamma) \vdash e : \widehat{C}_s(\sigma)$.*
3. *If $\Gamma ; \bar{v} \vdash_{\omega}^{\text{arg}} e : \sigma \rightsquigarrow C$ and $C_s \models C \wedge C'$, then we have a derivation for $\widehat{C}_s(\Gamma) \vdash_{\omega}^{\text{arg}} e : \widehat{C}_s(\sigma)$.*

Proof. We prove this theorem by mutual induction over the expression e .

Proof of (1). We distinguish two cases by inversion:

- *Case VARHEAD.* We have in this case that $\Gamma \vdash^{\text{fun}} x : \sigma \rightsquigarrow \epsilon$ if $x : \sigma \in \Gamma$. In this case C_s is a solution of C' . By applying the substitution \widehat{C}_s to the environment we get that $x : \widehat{C}_s(\sigma) \in \widehat{C}_s(\Gamma)$, which allows us to conclude that $\widehat{C}_s(\Gamma) \vdash^{\text{fun}} x : \widehat{C}_s(\sigma)$.
- *Case EXPRHEAD.* Follows directly from the induction hypothesis and (2).

Proof of (2). We distinguish four cases by inversion:

- *Case ABS.* In this case the constraint gathering is:

$$\frac{\Gamma, x : \alpha^{\#} \vdash e : \mu \rightsquigarrow C}{\Gamma \vdash \lambda x. e : \alpha^{\#} \rightarrow \mu \rightsquigarrow C}$$

By induction hypothesis we know that for every C' , if $C \wedge C'$ is solved to C_s , then we have $\widehat{C}_s(\Gamma), x : \widehat{C}_s(\alpha^{\#}) \vdash e : \widehat{C}_s(\sigma)$. From this we can apply the corresponding rule from the declarative specification.

$$\frac{\widehat{C}_s(\Gamma), x : \widehat{C}_s(\alpha) \vdash e : \widehat{C}_s(\sigma)}{\widehat{C}_s(\Gamma) \vdash \lambda x. e : \widehat{C}_s(\alpha) \rightarrow \widehat{C}_s(\sigma)} \text{APP}$$

Now we just need to notice that $\widehat{C}_s(\alpha) \rightarrow \widehat{C}_s(\sigma) \sim \widehat{C}_s(\alpha \rightarrow \sigma)$ and we are done.

- *Case ANNABS.* Similar to the ABS case.
- *Case APP.* Suppose that the shape of the expression is $e_0 e_1 \dots e_n$. By applying (1) on e_0 we know that a solution C_s determines a derivation $\widehat{C}_s(\Gamma) \vdash^{\text{fun}} e_0 : \widehat{C}_s(\sigma)$. By Lemma D.3, a solution for the \leq constraint determines a derivation for $\widehat{C}_s(\sigma) \leq_{\dots}^m \widehat{C}_s(\alpha_1^{\#}), \dots, \widehat{C}_s(\alpha_n^{\#}); \widehat{C}_s(\delta^{\dagger})$. Finally, by applying (3) to each argument, we get the corresponding derivations for $\widehat{C}_s(\Gamma) \vdash_{b_i}^{\text{arg}} e_i : \widehat{C}_s(\alpha_i^{\#})$. As a result, we can apply the APP rule from the declarative specification.
- *Case ANNAPP.* Similar to the APP case.

Proof of (3). We distinguish two cases by inversion:

- *Case ARGGEN.* We can assume, without loss of generality, that the generalisation constraint has the form $(\forall \{\bar{v}\}. C \Rightarrow \sigma) \leq (\forall \bar{b}. \eta)$ for a possibly empty set of variables \bar{b} .

Let us first consider the simple case in which the set of quantified variables \bar{b} is empty. By the definition of solution, this implies that we have

$$C_s \models C \quad \text{and} \quad C_s \models \sigma \leq_{\epsilon}^m \epsilon; \eta$$

Now we can apply (2) to $C_s \models C$ and derive that $\widehat{C}_s(\Gamma) \vdash \widehat{C}_s(\sigma)$. By inspection of the rules, we know that a solution of the second constraint, $\sigma \leq_{\epsilon}^m \epsilon; \eta$, defines a fully monomorphic substitution in the declarative specification. As a result, we can apply rule ARGGEN and obtain the desired result. The case for a not empty set of constraints is similar, one just has to realize that the rule for \models' for generalization constraints introduce the Skolem variables required by ARGGEN.

- *Case VARGEN.* In this case the generated constraint is $[\overline{p \mapsto \alpha^u}] \tau \leq (\bar{b}. \eta)$. As in the last case of ARGGEN, the definition of solution tells us that there is a $\forall \bar{b}. C'_s \in C_s$ such that

$$C'_s \models C_s(C) \quad \text{and} \quad C'_s \models C_s([\overline{p \mapsto \alpha^u}] \tau) \leq_{\epsilon}^u \epsilon; C_s(\eta)$$

The second application of \models' is trivially equal to $C'_s \models C_s([\overline{p \mapsto \alpha^u}] \tau) \sim C_s(\eta)$, since the type in the left-hand side is by construction a top-level monomorphic one. This is enough to apply rule VARGEN from the declarative specification. \square

Theorem 4.2 (Soundness). *Let Γ be a closed environment and e an expression. If $\Gamma \vdash e : \sigma \rightsquigarrow C$ and C_s is a solution for C with an induced substitution \widehat{C}_s , then we have $\Gamma \vdash e : \widehat{C}_s(\sigma)$.*

Proof. By Theorem D.4 we know that the existence of such a solution C_s implies that $\widehat{C}_s(\Gamma) \vdash e : \widehat{C}_s(\sigma)$. Since Γ is closed, $\widehat{C}_s(\Gamma) = \Gamma$, as desired. \square

Theorem D.5 (Derivations provide solutions). *Let Γ be an environment and e an expression.*

1. *If $\theta \Gamma \vdash^{\text{fun}} e : \sigma$ and $\Gamma \vdash^{\text{fun}} e : \phi \rightsquigarrow C$, then there exists a solution C_s for C such that \widehat{C}_s is fully monomorphic and $\widehat{C}_s(\theta \phi) \sim \sigma$.*
2. *If $\theta \Gamma \vdash e : \sigma$ and $\Gamma \vdash e : \phi \rightsquigarrow C$, then there exists a solution C_s for C such that \widehat{C}_s is fully monomorphic and $\widehat{C}_s(\theta \phi) \sim \sigma$.*

Proof. We prove this theorem by mutual induction over the expression e .

Proof of (1). We distinguish two cases by inversion:

- *Case VARHEAD.* In this case the generated set of constraints C is empty, so we can just take an empty solved form as solution C_s . The corresponding substitution \widehat{C}_s is the identity function, which is trivially fully monomorphic.

Now, we just need to prove that $\widehat{C}_s(\theta \phi) \sim \theta \phi \sim \sigma$. By inversion of the rule VARHEAD in the constraint-based formulation we know that ϕ must come from an element $x : \phi \in \Gamma$. On the other hand, we know that in the declarative specification σ must come from $x : \sigma \in \theta \Gamma$. Since Γ is the same in both derivations, it must be the case that $\sigma \sim \theta \phi$.

- *Case EXPRHEAD.* Follows by induction on the premise.

Proof of (2). We distinguish four cases:

- *Case ABS.* The derivation in the declarative specification looks like:

$$\frac{\theta \Gamma, x : \tau \vdash e : \sigma}{\theta \Gamma \vdash \lambda x. e : \tau \rightarrow \sigma}$$

Let us first rewrite the premise as $\theta'(\Gamma, x : \alpha)$, where $\theta' = [a \mapsto \tau] \circ \theta$. Now we can apply the induction hypothesis to obtain a solution C'_s such that $\widehat{C}'_s(\theta' \phi) \sim \sigma$, where ϕ is the type assigned to the abstraction body during constraint gathering.

Now consider the problem $C'_s \wedge \alpha \sim \tau$. We know that a solution C_s exists, or otherwise the expression would be ill-typed. Furthermore, $\widehat{C}'_s \circ \theta' = \widehat{C}_s \circ \theta$ – we are just moving the α assignment from one place to the other. In conclusion, we have that:

$$C_s(\theta(\alpha \rightarrow \phi)) \sim C'_s(\theta'(\alpha \rightarrow \phi)) \sim \tau \rightarrow \sigma$$

- *Case ANNABS.* Similar to the ABS case.
- *Case APP.* The derivation in the specification is:

$$\frac{\theta \Gamma \vdash^{\text{fun}} e_0 : \sigma_0 \quad \sigma_0 \leq_{\omega}^m \sigma_1, \dots, \sigma_n; \mu \quad \theta \Gamma \vdash_{\omega_i}^{\text{arg}} e_i : \sigma_i}{\theta \Gamma \vdash e_0 e_1 \dots e_n : \phi} \text{APP}$$

The constraints generated by this sequence of APP rules are:

$$C_0 \wedge (\phi_0 \leq_{\omega}^m \alpha_1^u, \dots, \alpha_n^u; \delta^t) \wedge C_1 \wedge \dots \wedge C_n$$

By (1) applied to e_0 there exists a solution C_s^0 such that $\widehat{C}_s^0 \theta \phi_0 \sim \sigma_0$. From the derivation of $\sigma_0 \leq_{\omega}^m \sigma_1, \dots, \sigma_n; \mu$ in the declarative specification we obtain a list of equalities C_s^{inst} , in particular $\alpha_i^u \sim \sigma_i$ and $\delta^t \sim \mu$.

Let us call $C'_s = C_s^h \wedge C_s^{\text{inst}}$ and consider for each argument the derivation $\widehat{C}'_s(\theta \Gamma) \vdash_{\omega_i}^{\text{arg}} e_i : \widehat{C}'_s(\sigma_i)$; which holds by substitutivity in the typing judgment. In the case σ_i is polymorphic, the rules for \vdash^{arg} perform the job required to obtain a solution, namely introducing rigid variables. We are thus left with $\widehat{C}'_s(\theta \Gamma) \vdash e_i : \widehat{C}'_s(\eta_i)$, where we can apply induction hypothesis to

obtain $C_s^{\text{arg}_i}$. Define C_s^i to be exactly $C_s^{\text{arg}_i}$ if no rigid variables were introduced and $\forall b. C_s^{\text{arg}_i}$ otherwise. The final solution is the conjunction:

$$C_s^h \wedge C_s^{\text{inst}} \wedge C_s^1 \wedge \dots \wedge C_s^n$$

- *Case ANNAPP.* Similar to the APP case.

□

Corollary D.6. *Suppose $\Gamma \vdash e : \sigma$ and let $\Gamma \vdash e : \phi \rightsquigarrow C$. Then there exists a solution C_s for C such that \widehat{C}_s is monomorphic and $\widehat{C}_s(\phi) \sim \sigma$.*

Corollary D.7. *Suppose $\Gamma \vdash e : \sigma \rightsquigarrow C$. Then C is either inconsistent, or can be rewritten to a solved form.*

Proof. By Theorem 4.1 every consistent set of constraints can be turned into an instance of the problem of first-order unification under a mixed prefix. The solver in Pottier and Rémy [15] describes a complete algorithm for this problem, from which we obtain the desired solved form. □

Corollary D.8 (Principality, solution version). *Let $\Gamma \vdash e : \sigma$ and let $\Gamma \vdash e : \phi \rightsquigarrow C$. If C is solved to C_s , then $\Gamma \vdash e : \widehat{C}_s(\phi)$ and there exists a monomorphic substitution π such that $\pi \widehat{C}_s(\phi) \sim \sigma$.*

Proof. By soundness we know that if C is solved to C_s , then $\Gamma \vdash e : C_s(\phi)$. From the derivation of $\Gamma \vdash e : \sigma$ we obtain another solution C'_s such that $C'_s(\phi) \sim \sigma$. Since the solver produces most general solutions we know in particular that there exists a fully monomorphic substitution π such that $\pi \circ \widehat{C}_s = \widehat{C}'_s$. This gives the desired result. □

Theorem 4.3 (Principality). *Suppose $\Gamma \vdash e : \sigma$. Then there exists a type ϕ such that $\Gamma \vdash e : \phi$, and for every other $\Gamma \vdash e : \sigma'$, there is a fully monomorphic substitution π such that $\sigma' = \pi\phi$.*

Proof. Take $\Gamma \vdash e : \phi \rightsquigarrow C$. Given that the derivation $\Gamma \vdash e : \sigma$ exists, by Corollary D.7 we are guaranteed to obtain a solution C_s for C . Then by Corollary D.8 we know that $\Gamma \vdash e : \widehat{C}_s(\phi)$ and $\sigma \sim \pi \widehat{C}_s(\phi)$ for where π is a fully monomorphic substitution. Take $\sigma^\star = \widehat{C}_s(\phi)$ and we are done. □

Theorem 4.4 (Completeness). *Let Γ 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.*

Proof. Consider a closed environment Γ and an expression e . By Corollary D.6 we know that if we generate constraints $\Gamma \vdash e : \phi \rightsquigarrow C$, then there exists a solution for C . In particular, since there exists a solution the constraint set C is not inconsistent.

Corollary D.7 states that every constraint set is either inconsistent or can be rewritten to a solved form. Given that our constraint set C is not inconsistent, this implies that it can be rewritten to a solved form, as desired. □