## ANONYMOUS AUTHOR(S)

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.

#### ACM Reference format:

Anonymous Author(s). 2017. Guarded impredicative polymorphism. *Proc. ACM Program. Lang.* 1, 1, Article 1 (January 2017), 39 pages.

https://doi.org/10.1145/nnnnnnnnnnnn

#### **1 INTRODUCTION**

Type inference for impredicative polymorphism is a deep, deep swamp. The literature is crowded with corpses (including several of our own): type systems for impredicative polymorphism that never quite worked well enough to be deployed in a production compiler [Botlan and Rémy 2003; Leijen 2008, 2009; Vytiniotis et al. 2006, 2008]. 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 2.2-2.4), 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 [Botlan and Rémy 2003; Leijen 2009; Vytiniotis et al. 2008]),
  nor do we introduce new forms of annotations [Russo and Vytiniotis 2009] or side-conditions
  that require principal types [Leijen 2008].

<sup>47</sup> 2017. 2475-1421/2017/1-ART1 \$15.00

<sup>48</sup> https://doi.org/10.1145/nnnnnnnnnn

Anon.

50	head $:: \forall p. [p] \rightarrow p$	$id \qquad :: \forall a. a \to a$
51	tail :: $\forall p. [p] \rightarrow [p]$	inc ::: Int $\rightarrow$ Int
52	$\begin{bmatrix} 1 \\ \vdots \\ \forall p \end{bmatrix} \begin{bmatrix} p \end{bmatrix}$	twice :: $\forall a \ (a \rightarrow a) \rightarrow a \rightarrow a$
53		
	$(:) \qquad :: \forall p. p \to [p] \to [p]$	$choose :: \forall a. a \to a \to a$
54	single V p p [p]	nolu (Vaa so) (Int Dool)
55	single :: $\forall p. p \rightarrow [p]$	$poly  :: (\forall a. a \to a) \to (Inl, bool)$
56	$(\texttt{+})  :: \forall p. [p] \to [p] \to [p]$	<i>auto</i> ::: $(\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$
50	longth V p [p] . Int	
57	$lengin :: \lor p. [p] \rightarrow Ini$	$app ::: \forall a \ b. \ (a \to b) \to a \to b$
58	map :: $\forall p a (p \rightarrow a) \rightarrow [p] \rightarrow [a]$	$\mathbf{x} = \mathbf{x} + $
30		$revapp :: \forall a \ b. \ a \to (a \to b) \to b$
59	$ids  \cdots [\forall a \ a \rightarrow a]$	CT () ( CT )
		$runS1 :: (\forall s. S1 s v) \rightarrow v$
60		$argST \dots \forall s ST s Int$
61		urg51 v 5.51 5 mi
62		

Fig. 1. Type signatures for functions used in the text

- Despite GI's relative simplicity, it accepts without annotation particularly celebrated and practically-important examples, such as *runST* \$ *e* (Figure 2).
- We give a declarative type system for GI for a small core language, highlighting the key ideas of our system (Section 3). Then we show simple extensions to handle a more full-fledged language, including type annotations (Section 3.3), let bindings, do notation, and pattern matching (Section 3.4). 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, as we show in Figure 8.
- 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, which is already based in constraints. Type-correct programs can easily be elaborated into System FC, GHC's intermediate language, without extensions. Our inference algorithm scales readily to handle GADTs and type classes, as we show in Section 5. Higher kinds, type-level functions and other type system features pose no additional problem.
  - Moreover, as we discuss in Section 6, we can reduce the annotation burden by relaxing the guardedness conditions during solving.
  - We provide a prototype implementation of the whole system, integrated with Haskell's type classes.

Type inference for impredicativity is dense with related work, as we discuss in Section 7. A small but useful contribution to a dense field is Figure 2, which presents key examples from the literature and shows how each major system behaves on that example.

## 2 THE KEY IDEA: INTUITION AND EXAMPLES

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

#### 2.1 Exploiting the easy case

What is hard about typing (*head ids*)? Nothing! Since the type variable p appears under a list 94 95 type constructor in *head*'s type, and we know the type of *ids*, it is plain as a pikestaff that we must instantiate  $p := \forall a. a \rightarrow a$ . The difficulty comes when we have a "naked" or "un-guarded" 96 type variable in the function type, such as *single* ::  $\forall p. p \rightarrow [p]$ . Now if we examine (*single id*), it 97

Proc. ACM Program. Lang., Vol. 1, No. 1, Article 1. Publication date: January 2017.

63 64 65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87 88

89

90

91 92

99 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 100 incomparable types  $\forall a. [a \rightarrow a]$  and  $[\forall a. a \rightarrow a]$ . To make things worse (*single id*) is a perfectly 101 typeable Hindley-Milner program (with the former type) so we must allow this type. But to support 102 impredicativity, we must also allow the latter. But under which conditions? This is clearly a hard 103 case for type inference because it requires information from the context. In GI our approach is to 104 exploit the easy case and turn the hard case to easy by requiring type annotations; finally when 105 106 annotations are not present we will fall back to monomorphic instantiation. We hope thereby to get 90% of the joy for 10% of the pain. Specifically: 107

- We focus on n-ary applications  $(fe_1...e_n)$ . It is more conventional to deal with binary applications, but in fact n-ary applications (unencumbered with intervening let or case constructs) are wildly dominant in practice, and we can get much better typing by treating the application all at once.
- For the choices of impredicative instantiation, we accept a particular information flow in an n-ary call: start with a function with a known type; guide its instantiation by the types of its explicitly-supplied arguments; and do not attempt to use the context of the call to guide instantiation.
  We instantiation.
  - We instantiate a function with a polymorphic type (i.e. impredicatively) only if the polymorphism is "guarded", either in the type of the function (as we saw with (*head xs*) above), or the type of the argument.

What does this last bullet mean? Consider the call (:) id ids (see Figure 1 for type signatures). Here the variable p in (:)'s type appears under a list type constructor in the second argument, so its instantiation is fixed by the second argument ids. All is good.

But what if we had a partially-applied call (:) *id*? Now *p* does not appear guarded (under a type constructor) in the first argument of (:)'s type. And indeed, it is not immediately clear whether to instantiate *p* with  $\forall a. a \rightarrow a$ , or *Int*  $\rightarrow$  *Int*, etc; it is like the (*single id*) case. We do not reject this program, of course; rather we just insist that *p* is instantiated with a monotype. So for us, the notion of "guardedness" clearly involves the number of arguments supplied in an n-ary call.

Similarly, consider the call *map poly ids*. Here, *p* (in *map*'s type) appears under a type constructor in the second argument of *map*, so this application is fine despite the polymorphism in *poly*'s type. What about the partially-applied call (*map poly*)? In GI that too is OK, because we treat  $(\rightarrow)$ as invariant, like any other type constructor. This occasionally means that you need to do some eta-expansion (Section 3.2), but it makes more type variables guarded, and that is very worthwhile. But a call of *map* applied to no arguments (for example, it might be passed as an argument to another function) cannot be instantiated impredicatively.

<sup>134</sup> Up to now we have spoken of guardedness as a property of the function's type, given the number <sup>135</sup> of arguments to which it is applied. However, in GI guardedness is not *only* a property of the <sup>136</sup> type of the function; it *also* involves the types of its arguments. Consider the call *single xs*, where <sup>138</sup>  $xs :: [\forall a. a \rightarrow a]$ . Now there is no doubt at what type to instantiate *single* ::  $\forall p. p \rightarrow [p]$ ! It must be <sup>139</sup> instantiated with  $p := [\forall a. a \rightarrow a]$ . Why is there no doubt? Again because the scary polymorphic <sup>140</sup> type appears under a list constructor.

# <sup>141</sup> 2.2 A general rule for n-ary applications

#### 143 Here is the general rule:

- Rule 1: When instantiating a type at an n-ary call, a type variable may only be assigned a polymorphic type  $\sigma$  if, in the instantiated type,  $\sigma$  appears under a type constructor in one or more of the supplied arguments.
- 147

108

109

110

111

117

118

123

124

125

126

127

128

129

130

131

132

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

- $(map \ poly)$  is OK because, after instantiation, map's type becomes  $((\forall \ a. \ a \rightarrow \ a) \rightarrow (Int, Bool)) \rightarrow [\forall \ a. \ a \rightarrow \ a] \rightarrow [(Int, Bool)]$ . The instantiating type  $(\forall \ a. \ a \rightarrow \ a)$  appears under a type constructor of one of the supplied arguments (poly in this case).
  - (single ids) is OK because, after instantiation, single's type becomes [∀ a. a → a] → [[∀ a. a → a]]. The instantiating type appears under a list type constructor in one of the supplied arguments, namely ids.
  - (*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 to match. But Rule 1 says that this application is OK: in the instantiated type of *id*, the forall appears under an arrow type, in the type of the first argument.

Rule 1 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.

#### 2.3 Ignoring the context of a call

Notice that Rule 1 *takes no account of the context of the call.* For example, consider *ids* + *single id.* Here we append two lists, so we know that the result of (*single id*) must be  $[\forall a. a \rightarrow a]$ , and you might think that would be enough to fix the instantiation of *single*. But not in GI! The swamp beckons, and we stay on dry land.

Moreover, Rule 1 allows the programmer to understand impredicativity in a simple bottom-up way. For example, consider the well-typed expression (*map head* (*single ids*)) (Figure 2)? In our type system, the types for *head* and *single ids* are determined, in fact, instantiated, independently, so we never need to consider the *interaction* between the two arguments. This modularity pays off in the metatheory too.

There is a price to pay, however. As a degenerate case, a function application without any arguments – that is, a variable – may only instantiate *fully* monomorphically – no polymorphism, even if it appears under a type constructor. Thus, the empty list constructor  $[] :: \forall a. [a]$  cannot be assigned a type  $[\forall a. a \rightarrow a]$ . We explore how to allow more programs in Section 6, and we also introduce annotations to cover the remaining cases.

### 2.4 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. But, in common with every other system we know, we reject all programs that require a lambda-bound variable to be polymorphic, unless it is explicitly annotated:

Rule 2: Every lambda abstraction whose argument is polymorphic must be annotated. If is it not annotated, the bound variable can only have a fully monomorphic type.

By a "fully monomorphic type" we mean "no foralls anywhere". Nothing about guardedness here!
Nevertheless, in Section 6 we explore a way to alleviate some of the burden for the programmer,
for cases where the argument type, albeit polymorphic, can be inferred from its usage in the body.

While Rule 2 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  $(\forall a. a \rightarrow a). x x$ , 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

154

155

156

157

158

159

160

161

162

163 164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179 180

181

182

183

184

185

186

187

188

189

#### 2.5 Expressiveness 200

199

212

213

214 215 216

217

218

219 220

221

222

223

224 225

226

227

232

233

240

241

242

243

245

201 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 202 situations. We summarise a collection of examples culled from the literature in Figure 2. This table 203 204 also compares our system with others, but we defer discussion of related work to Section 7.

One example is the function (\$) ::  $(a \rightarrow b) \rightarrow a \rightarrow b$ . Haskellers use this function all the time to 205 remove parentheses in their code, as in (runST  $do \{\dots\}$ ) (the type of runST is given in Figure 1). 206 This call absolutely requires impredicative instantiation of the variable *a* in the type of (\$). It is so 207 208 annoying to reject this program that GHC implements a special, built-in typing rule for  $f \ x$ . Of 209 course, that is horribly non-modular: if the programmer re-defines another version of (\$), even 210 with the same type, some programs cease to type check. In GI both type variables appear under the 211  $(\rightarrow)$  constructor, so impredicative instantiation is allowed.

The lack of support for impredicative types sometimes forces programmers to define new types as a work around. Consider the following from Haskell's lens library<sup>1</sup>:

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

Because *Lens* is a polymorphic type, forming a list of lenses [*Lens s t a b*] is illegal in Haskell, because you cannot instantiate the list constructor (:) at a polymorphic type. So programmers have to resort to additional wrapping and unwrapping

**newtype** ReifiedLens s t a  $b = Lens \{ runLens :: Lens s t a b \}$ 

Now you can build a list of reified lenses, of type [*ReifiedLens s t a b*], but doing so is tiresome, because it requires many explicit conversions between Lens and ReifiedLens. It would be much better if *Lens* can be treated as a first-class type, and that is exactly what impredicativity allows.

#### **DECLARATIVE SPECIFICATION** 3

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 228 for the constraint-based variant, the programmer can take this declarative specification - easier 229 to understand but without a direct inference algorithm – as a base to understand when and why 230 annotations are needed in their programs. 231

## 3.1 Syntax

The syntax of the core term language, for our initial discussion, is given in Figure 3. The language 234 235 has some distinctive features. First, as discussed in Section 2, we deal with n-ary applications instead 236 of binary ones. A lone term variable is treated as 0-ary application. Because of Rule 2 (Section 2.4), 237 we provide explicitly-annotated lambda abstractions,  $\lambda(x::\sigma)$ . *e*, to support lambdas whose bound variable must have a polymorphic type. 238 239

- Types (Figure 3) are classified by three "sorts", u, t, and m:
  - Polymorphic types  $\sigma$ ,  $\phi$ , of sort  $\mathfrak{u}$ , have unrestricted polymorphism.
  - Top-level monomorphic types,  $\mu, \eta$ , of sort t, have no polymorphism at the top level, but permit arbitrary nested polymorphic types under a type constructor.

Proc. ACM Program. Lang., Vol. 1, No. 1, Article 1. Publication date: January 2017.

<sup>&</sup>lt;sup>1</sup>https://hackage.haskell.org/package/lens-4.15.3/docs/Control-Lens-Reified.html 244

246		GI	GI	MIE	пие	FDU	имі		
247		decl. spec.	relaxed	WILI	111/11	1.1.11	TIML		
248	POLYMORPHIC INSTANTIATION	$:$ given $f :: \forall$	$a. (a \rightarrow a)$	$\rightarrow [a]$	$\rightarrow a$				
249	$const2 = \lambda x \ y. \ y$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
250	MLF infers $(b \ge \forall c. c \rightarrow c) \Rightarrow a \rightarrow b$	b, whereas GI i	nfers $a \rightarrow l$	$b \rightarrow b.$					
251	choose id	$\checkmark$	$\checkmark$	$\checkmark$		$\checkmark$	$\checkmark$		
252	MLF/HML infer $(a \ge \forall b. b \rightarrow b) \Rightarrow$	$a \rightarrow a$ , wherea	s FPH/GI in	ifer ( $a \rightarrow$	$a) \rightarrow a -$	<i>→ a</i> .			
253	choose [ ] ids	No	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
254	GI needs an annotation on [] :: [ $\forall a$ .	$a \rightarrow a$ ].							
255	$auto = \lambda(x :: \forall a. a \rightarrow a). x x$	$\checkmark$	$\checkmark$	$\checkmark$	No	$\checkmark$	$\checkmark$		
256	MLF infers $(\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$	a), whereas G	I infers $(\forall a$	$a \rightarrow a$	$\rightarrow b \rightarrow b$	•			
257	id auto	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
258	choose id auto	No	No	$\checkmark$	No	No	$\checkmark$		
259	Relaxed GI accepts the version with a	auto :: (∀ a. a —	$(\forall a) \rightarrow (\forall a)$	$a \rightarrow a$ ).					
260	f (choose id) ids	No	$\checkmark$	$\checkmark$	No	$\checkmark$	$\checkmark$		
261	GI needs an annotation on $id :: (\forall a. d)$	$a \to a) \to (\forall a)$	$a \rightarrow a$ in t	the previ	ous two e	xamples.			
262	poly id	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
263	poly $(\lambda x. x)$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
264	id poly $(\lambda x. x)$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
265	INFERENCE OF POLYMORPHIC A	RGUMENTS	1		1				
266	$\lambda f.(f 1, f True)$	No	No	No	No	No	No		
267	All systems require an annotation on	$f :: \forall a. a \rightarrow a$		I	1				
268	$\lambda xs. poly (head xs)$	No	No	$\checkmark$	No	No	No		
269	All systems except for MLF require an annotation on $xs :: [\forall a. a \rightarrow a]$ .								
270	FUNCTIONS OPERATING ON PO	LYMORPHIC	LISTS: giv	ren g :: ∀	/ a. [ a] –	$\rightarrow [a] \rightarrow a$			
271	length ids	$\checkmark$	<ul> <li>✓</li> </ul>	$\checkmark$	$\checkmark$	✓	$\checkmark$		
272	tail ids	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
273	head ids	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
274	single id	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
275	id : ids	$\checkmark$	$\checkmark$	$\checkmark$	No	$\checkmark$	$\checkmark$		
276	$(\lambda x. x): ids$	$\checkmark$	$\checkmark$	$\checkmark$	No	$\checkmark$	$\checkmark$		
277	single inc + single id	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
278	g (single id) ids	No	No	$\checkmark$	No	$\checkmark$	$\checkmark$		
279	map poly (single id)	No	No	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
280	GI needs an annotation on single id ::	$[\forall a. a \rightarrow a]$ i	n the two p	revious e	examples.		•		
281	map head (single ids)	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		
282	APPLICATION FUNCTIONS						-		
283	app poly id	1	1	$\checkmark$	1	1	$\checkmark$		
284	revapp poly id		· ·	•	·	·			
285	runST aroST			• •	·		·		
286	abb runST argST	·	· ·	, ,	· ·	· ·			
287	revabb runST argST	• √	·	<b>√</b>	·	·	$\checkmark$		
288	n-EXPANSION: given h. Int	$a \rightarrow a k$	$\forall a a = \rangle$	$[a] \rightarrow c$	a let [ \-	$a  Int \rightarrow c$	$ \rightarrow a^{1} $		
289	h h let	No. No.	No	$\begin{bmatrix} u \end{bmatrix} \rightarrow t$	No	$u. m \rightarrow t$	No		
290	$k (\lambda x + x)$ let	INU	INU .	INU .	No	INO	110		
291	л (лл. 11 л) isi	V	<b></b>	V	110	<b>v</b>	v		
292	Fig. 2. Comparison of type syste	ms (Figure 1	gives the	types of	the fund	tions men	tioned)		

293 294 Fig. 2. Comparison of type systems (Figure 1 gives the types of the functions mentioned)

Proc. ACM Program. Lang., Vol. 1, No. 1, Article 1. Publication date: January 2017.



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

341

Anon.

$$\begin{array}{c}
\begin{array}{c}
\begin{array}{c}
\Gamma \vdash h \ e: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array} \\
\begin{array}{c}
\hline \Gamma \vdash h \ x: \sigma \\
\end{array}$$



The definitions related to the guardedness requirement are given in Figure 4. Given an *n*-ary application, the judgment  $\sigma \triangleright_{\mathfrak{s}}^{\mathfrak{n}} \Delta$  classifies the free type variables of  $\sigma$ . Specifically,  $\Delta$  maps each variable *a* to the sort of types that are allowed to instantiate *a*:

- If *a* appears under a type constructor (i.e. guarded) in at least one of the first *n* arguments of  $\sigma$ , then *a* may be instantiated by an unrestricted type  $\phi$ .
- If a appears in one of the first n arguments of  $\sigma$ , a may be instantiated by a top-level monomorphic type  $\mu$ .
- If a appears in the result type of  $\sigma$ , after stripping off n arguments, then a may be instantiated by a type of sort 5.

We combine information from multiple occurrences with  $\sqcup$ , the least upper bound of the (trivial) sort lattice in Figure 3; for example if a variable appear both unrestricted and monomorphically we can treat it as unrestricted. Using the  $\triangleright_{s}^{n}$  judgment, Figure 4 also defines what it means for a substitution  $\theta$  to "respect the guardedness of a type", written  $\theta$  respects  $\Delta$ . 

### 3.2 Typing rules

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

Rules ABs and ANNABs concern lambda abstractions, and are straightforward. ANNABs deals with a lambda ( $\lambda(x::\phi)$ . e) where the user has supplied a type annotation  $\phi$ : we simply bring x into

1:8

Expressions / terms  $e ::= \dots | (h e_1 \dots e_n :: \sigma)$ 

$$\frac{\Gamma \vdash^{\mathsf{h}} h : \phi \qquad \phi \leqslant_{\mathfrak{u}}^{n} \sigma_{1}, \dots, \sigma_{n}, \eta \qquad \Gamma \vdash^{\mathsf{arg}} e_{1} : \sigma_{1} \qquad \dots \qquad \Gamma \vdash^{\mathsf{arg}} e_{n} : \sigma_{n}}{\Gamma \vdash (h e_{1} \dots e_{n} :: \forall \overline{b}, \eta) : \forall \overline{b}, \eta} \qquad \text{AnnApp}$$

 $\Gamma \vdash e : \sigma$ 

409

411

412

413

415

417

418

419

393 394

#### Fig. 6. Declarative type system with annotations in applications

scope with type  $\phi$ . As discussed in Section 2.4, where there is no annotation (rule ABS) we insist 403 that *x* has a fully-monomorphic type  $\tau$ . 404

All the action is in rule APP for n-ary applications ( $h e_1 \dots e_n$ ). The first step is typing the head of 405 the application h. The corresponding judgment  $\vdash^{h}$  either looks for a variable in the environment or 406 uses the normal typing judgment if the head is another kind of term. 407

After typing the head, we instantiate the type of the head by means of the  $\leq_{\alpha}^{n}$  judgment. Note 408 that this instantiation judgment is parametrized by a sort \$ (this is needed to support applications with and without annotations). 410

- All type constructors are considered *invariant*, even functions. This means that neither  $[\forall a. a \rightarrow a] \leq _{5}^{n} [Int \rightarrow Int] \text{ nor } Int \rightarrow (\forall a. a \rightarrow a) \leq _{5}^{n} Int \rightarrow Int \rightarrow Int.$
- For function types the judgment  $\leq_{5}^{n}$  embodies a bit of their covariance. Given a function f of the type  $\forall a. a \rightarrow (\forall b. b \rightarrow a)$  we are allowed to first instantiate a with some type  $\tau_1$ , which 414 returns a result  $\forall b. b \rightarrow \tau_1$ , and then instantiate b if a second argument is supplied. 416
  - As a consequence, we can work around the lack of variance of function types by  $\eta$ -expanding. Continuing the previous example, if we need an expression of type  $\tau_1 \rightarrow \tau_2$  for some hole, we cannot directly use f. But we can use  $\lambda x \ y$ . f x y instead, which allows the  $\leq_{s}^{2}$  judgment to kick in and instantiate the two type variables.
- The status of the variables in the result type depends on the parameter 5 to the judgment. 420 In the case of APP, this parameter is set to m. As a consequence, those variables appearing 421 only in the result of the function – the part of the type past the given number of arguments – 422 have to be instantiated with fully monomorphic types. 423

424 Whereas the type of the head of the application may only be instantiated, the arguments may also 425 ongo generalisation, that is, abstraction of some of the type variables. Generalisation is embodied 426 by the ArgGen rule. Without such a rule, we would not be able to type check *twice* ( $\lambda x$ . x). Here 427 the type of  $\lambda x$ . x must be of the form  $\tau \to \tau$  for some fully monomorphic  $\tau$ . If we choose  $\tau$  to be a 428 fresh Skolem variable, we can derive  $\Gamma \vdash^{\operatorname{arg}} \lambda x. x : \forall a. a \rightarrow a$ , as needed.

429 THEOREM 3.1 (COMPATIBILITY WITH HINDLEY-MILNER). Let e be an expression in the syntax of 430 *Hindley-Milner without* let. *If*  $\Gamma \vdash^{\mathsf{HM}} e : \tau$ *, then*  $\Gamma \vdash e : \tau$ *.* 431

432 PROOF. By inspection of the rules for the syntax-directed version of Hindley-Milner. Notice that all instantiations done in Hindley-Milner are admissible in GI, since the fully monomorphic 433 types form a subset of the top-level monomorphic types which again form a subset of the fully 434 polymorphic types. So even if a type variable is sorted as t or u, an instance of sort m can be 435 chosen. 436 

#### 3.3 Annotations in applications 438

With the rules in Figure 5, a single variable (not applied to any arguments) is treated as a 0-ary 439 application. So rule APP cannot get any guardedness information from its arguments, and so its type 440

441

can only be instantiated with fully monomorphic types. As a result, we cannot assign the polymorphic 442 type  $[\forall a. a \rightarrow a]$  to the empty list []. That is embarrassing, because we cannot typecheck, say 443 444 ([] + ids). Figure 2 has other examples.

One solution is to allow the programmer to give a type annotation for an application; see the 445 syntax in Figure 6. Now, since annotations fully specify the types, we do not need to impose 446 guardedness restrictions on those variables appearing in the result. Thus, the expression []: 447  $[\forall a. a \rightarrow a]$  is accepted by the system, even though the type variable in a is not guarded by a type 448 449 constructor. GI is quite symmetrical in terms of syntax: both abstractions and applications may be annotated. 450

Annotations also free us from having a different judgment for declarations and expressions. For 451 every combination  $f :: \sigma; f = e$  in the source code, we just need to pose the problem of checking 452  $f = e :: \sigma$  to guarantee that the declaration is well-typed. 453

The extensions required for annotated applications are described in Figure 6. Rule ANNAPP is 455 almost identical to APP, except for the choice of parameter to the instantiation judgment, which 456 is u. This implies that in contrast to non-annotated applications, variables in the result type of 457 the function might be substituted by any type, polymorphic or not. This is sensible, since the 458 annotation tells us exactly what the types are those variables should be instantiated with. 459

Of course, annotated applications serve for arbitrary applications, not just 0-ary ones. Take the 460 expression single ( $\lambda x. x$ ). Due to the type of single being  $\forall p. p \rightarrow [p]$ , the type of the expression 461 must be  $[\tau \to \tau]$  for a monomorphic  $\tau$ . If we want instead to obtain  $[\forall a. a \to a]$ , we can just 462 annotate the result, thus single  $(\lambda x. x) :: [\forall a. a \to a]$ . Since the variable a appears in the result type 463 of *single*, this is perfectly fine by rule ANNAPP. 464

465 There is some room for choice when introducing annotations in the language. In particular, the annotations can be taken as *rigid* – the type of the decorated expression is exactly the one 466 appearing as annotation - or *soft* - the resulting type might be an instance of the one stated in the annotation. We take the former route, which is shared by most of the previous work in 468 impredicative polymorphism. One good consequence of this choice is that we can embed System F, 469 as described in Section 3.5. 470

#### 3.4 let bindings, do notation and pattern matching

473 A realistic language, such as Haskell [Marlow et al. 2010], includes other syntactic forms apart 474 from those inherited from the  $\lambda$ -calculus. In some cases, those forms are defined as translations to 475 the core language; let bindings and do notation are some important examples. As we shall see in 476 this section, adopting the standard translation for these language constructs to derive their typing 477 rule is not desirable: we can do better.

Let us start by looking at let-bindings. Following the long-established tradition, we could translate let  $x = e_1$  in  $e_2$  as  $(\lambda x. e_2) e_1$ , leading to the following derivation.

$$\frac{\Gamma, x: \tau \vdash e_2: \sigma}{\Gamma \vdash \lambda x. e_2: \tau \to \sigma} \quad \tau \to \sigma \leqslant^1_{\mathfrak{m}} \tau, \mu \qquad \frac{\Gamma \vdash e_1: \phi \qquad \phi \leqslant^0_{\mathfrak{m}} \tau}{\Gamma \vdash^{\operatorname{arg}} e_1: \tau}$$

$$\frac{\Gamma \vdash \operatorname{let} x = e_1 \text{ in } e_2: \mu}{\Gamma \vdash^{\operatorname{arg}} e_1: \tau}$$

But this translation imposes an important restriction on the type of x: it must be fully monomorphic 485 even though the type of  $e_1$  might be more general. The reason is that we try to guess the type of  $e_1$ 486 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 487 restrictive! The rule LET in Figure 7 works in the other direction: the type obtained from typing  $e_1$ 488 is put in the environment as the one for *x*, allowing the type of *x* to be fully polymorphic instead. 489

454

467

471

472

478

479

480 481 482

Data constructor

Expressions / terms

do-statements

49	1
49	2
49	3

$$\frac{\Gamma \vdash e: \sigma}{\Gamma \vdash \operatorname{let} x = e_{1} \text{ in } e_{2}: \sigma} \operatorname{Ler} \qquad \frac{\Gamma \vdash_{M}^{\operatorname{do}} \overline{s}: M \phi}{\Gamma \vdash \operatorname{do} \overline{s}: M \phi} \operatorname{Do}$$
for each branch  $K_{i} \overline{x_{i}} \to e_{i}$  with
$$\Gamma \vdash e_{0}: \sigma_{0} \qquad \sigma_{0} \leqslant_{\mathfrak{m}}^{0} \operatorname{T} \overline{\phi_{0}} \qquad K_{i} : \forall \overline{a} \overline{b} . \overline{\sigma_{i}} \to \operatorname{T} \overline{a} \in \Gamma$$

$$\Gamma, x_{i}: [\overline{a} \mapsto \overline{\phi_{0}}] \sigma_{i} \vdash e_{i}: \phi_{\star}$$

$$\Gamma \vdash \operatorname{case} e_{0} \operatorname{of} \{\overline{K \overline{x} \to e}\}: \phi_{\star}$$

$$\Gamma \vdash_{M}^{\operatorname{do}} e: \mu$$

$$\Gamma \vdash_{M}^{\operatorname{arg}} e: M \phi \qquad \Gamma, x: \phi \vdash_{M}^{\operatorname{do}} \overline{s}: \eta \qquad \Gamma \vdash e: \phi \qquad \Gamma, x: \phi \vdash_{M}^{\operatorname{do}} \overline{s}$$

 $e \mid x \leftarrow e; \overline{s} \mid \text{let } x = e; \overline{s}$ 

... | let  $x = e_1$  in  $e_2$  | do  $\overline{s}$  | case  $e_0$  of { $K \overline{x} \rightarrow e$ }

$$\frac{\Gamma \vdash^{\operatorname{arg}} e : \mu}{\Gamma \vdash^{\operatorname{do}}_{M} e : \mu} \qquad \frac{\Gamma \vdash^{\operatorname{arg}} e : M \phi \qquad \Gamma, x : \phi \vdash^{\operatorname{do}}_{M} s : \eta}{\Gamma \vdash^{\operatorname{do}}_{M} x \leftarrow e ; \overline{s} : \eta} \qquad \frac{\Gamma \vdash e : \phi \qquad \Gamma, x : \phi \vdash^{\operatorname{do}}_{M} s : \eta}{\Gamma \vdash^{\operatorname{do}}_{M} \operatorname{let} x = e ; \overline{s} : \eta}$$

Κ

Э

::=

e ::=

S



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. [2011] 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  $x = (e_1 :: \lambda \phi)$  in  $e_2$ .

Translating **do** notation to a sequence of applications of  $\gg$  raises the same problems as let. The Haskell 2010 Language Report [Marlow et al. 2010] defines **do** { $x \leftarrow e_1; e_2$ } to be translated to  $e_1 \gg \lambda x \rightarrow e_2$ . As before, this means that a fully monomorphic type has to be inferred for x, neglecting all information obtained from  $e_1$ . This restriction has been reported to forbid some classes of domain-specific languages [Augustsson 2011]. We give less restrictive rules in Figure 7, that again allow the type of x to be fully polymorphic.

Introducing pattern matching in the language gives no surprises, as witnessed by rule CASE. 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 T, 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 T, namely  $\overline{\phi_0}$ . The return type of every branch should be equal,  $\phi_{\star}$  (note that we do *not* allow the types of the branches to be instantiated).

#### <sup>533</sup> 3.5 Embedding of System F into GI

Our system is as expressive as System F, the gold standard for a fully-expressive impredicative type system:

THEOREM 3.2 (EMBEDDING OF SYSTEM F). Let e be a System F expression. If  $\Gamma \vdash^{\mathsf{F}} e_{\mathsf{F}} : \sigma \rightsquigarrow e'$ , as defined in Figure 8, then  $\Gamma \vdash e' : \sigma$  in GI.

System F terms  $e_{F} := x \mid \lambda(x :: \sigma). e_{F} \mid Aa. e_{F} \mid e_{F} e_{F} \mid e_{F} \sigma$  $\boxed{\Gamma \vdash^{F} e_{F} : \sigma \rightsquigarrow e'}$   $\frac{x : \mu \in \Gamma}{\Gamma \vdash^{F} x : \mu \rightsquigarrow x}$   $\frac{x : \sigma \in \Gamma}{\Gamma \vdash^{F} x : \sigma \rightsquigarrow (x :: \sigma)}$   $\frac{\Gamma, x : \tau \vdash^{F} e_{F} : \phi \rightsquigarrow e'}{\Gamma \vdash^{F} \lambda(x :: \tau). e_{F} : \tau \rightarrow \phi \rightsquigarrow \lambda x. e'}$   $\frac{\Gamma, x : \sigma \vdash^{F} e_{F} : \phi \rightsquigarrow e'}{\Gamma \vdash^{F} \lambda(x :: \sigma). e_{F} : \tau \rightarrow \phi \rightsquigarrow \lambda x. e'}$   $\frac{\Gamma \vdash^{F} e_{F} : \phi \restriction^{F} e_{F} : \phi \rightsquigarrow e'}{\Gamma \vdash^{F} \lambda(x :: \sigma). e_{F} : \sigma \rightarrow \phi \rightsquigarrow \lambda(x :: \sigma). e'}$   $\frac{\Gamma \vdash^{F} e_{0} : \sigma \rightsquigarrow e'_{0} \qquad \sigma_{0} = \forall \overline{a}_{1}. \sigma_{1} \rightarrow \forall \overline{a}_{2}. \sigma_{2} \rightarrow \dots \forall \overline{a}_{n}. \sigma_{n} \rightarrow \forall \overline{a}_{r}. \mu_{r}}{\sigma_{0} \leqslant^{m}_{u}} [\overline{a}_{1} \mapsto \phi_{1}]\sigma_{1}, [\overline{a}_{1} \mapsto \phi_{1}, \overline{a}_{2} \mapsto \phi_{2}]\sigma_{2}, \dots, [\overline{a}_{1} \mapsto \phi_{1}, \dots, \overline{a}_{n} \mapsto \phi_{n}]\sigma_{n} \rightsquigarrow e'_{n}}$   $\Gamma \vdash^{F} e_{0} : \overline{\phi}_{n} e_{n} \overline{\phi}_{r} : \forall \overline{a}_{r}. \mu_{r} \sim \begin{cases} e'_{0} e'_{1} \dots e'_{n} \qquad \text{if } \overline{a}_{p} \text{ is empty} \\ and \phi_{r} \text{ fully mono.} \\ (e'_{0} e'_{1} \dots e'_{n} : : \forall \overline{a}_{r}. \mu_{r}) \qquad \text{otherwise} \end{cases}$ 

Fig. 8. Translation from System F

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*<sub>0</sub> 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.

#### 3.6 Metatheory

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

THEOREM 3.3 (IMPREDICATIVE INSTANTIATION IS NOT GUESSED). Let  $\Gamma$  be a (possibly open) environment and e an expression. For every pair of fully monomorphic substitutions  $\theta_1$  and  $\theta_2$ ,

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

PROOF. See Appendix C.1.

COROLLARY 3.4. Let  $\Gamma$  be a closed environment. If  $\Gamma \vdash e : \sigma_1$  and  $\Gamma \vdash e : \sigma_2$ , then there exists a polymorphic type  $\sigma^*$  and fully monomorphic substitutions  $\varphi_1$  and  $\varphi_2$  such that  $\sigma_i = \varphi_i \sigma^*$ .

This property suggests a notion of *principal type* similar to the one found in Hindley-Milner. A principal type for an expression *e* is defined as a type  $\sigma^*$  for which any other type assignment  $\phi$  to *e* is equal to  $\theta\sigma^*$  for a fully monomorphic substitution  $\theta$ . The fact that we only need to consider *fully monomorphic* substitutions here is a direct consequence of Theorem 3.3. The proof of the principal types property, however, is a corollary of other properties of the inference process, which we describe in Section 4.5.

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

THEOREM 3.5. 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,  $u = h u_1 \dots u_n$ . Note that in that case the assigned types are always top-level monomorphic, so  $\sigma = \mu$  and  $\phi = \eta$ .

The derivations involved look like:

and the question is whether we can derive:

$$\frac{\Gamma \vdash^{\mathsf{h}} h : \sigma_{u} \qquad \sigma_{u} \leqslant_{\mathfrak{m}}^{m+n} \sigma_{1}, \dots, \sigma_{m}, \phi_{1}, \dots, \phi_{n}, \eta \qquad \overline{\Gamma} \vdash^{\mathsf{arg}} u_{i} : \sigma_{i} \qquad \Gamma \vdash^{\mathsf{arg}} e_{j} : \phi_{j}}{\Gamma \vdash h u_{1} \dots u_{m} e_{1} \dots e_{n} : \eta}$$

All the premises on this last rule come directly from those in the hypotheses, except for the instantiation  $\sigma_u \leq_m^{m+n} \sigma_1, \ldots, \sigma_m, \phi_1, \ldots, \phi_n, \eta$ . 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.

THEOREM 3.6. Let  $app :: \forall a \ b. (a \to b) \to a \to b$  and  $revapp :: \forall a \ b. a \to (a \to b) \to b$  be the application and reverse application functions, respectively. Given two expressions f and e such that  $\Gamma \vdash^{h} f : \sigma_{0}$ , and  $\sigma_{0} \leq_{m}^{0} \sigma_{1} \to \phi$  then:

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

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

PROOF. Let us compare the derivations of the first two expressions.

$$\frac{\Gamma \vdash^{\mathsf{h}} f : \sigma_0}{\Gamma \vdash f e : \mu} \qquad \Gamma \vdash^{\mathsf{arg}} e : \sigma_1$$

Proc. ACM Program. Lang., Vol. 1, No. 1, Article 1. Publication date: January 2017.

 $\Gamma \vdash^{\operatorname{arg}} e : \sigma_1$ 

1:14

643

644 645

647 648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

670

671

672

673

674 675

680

681

# $\Gamma \mapsto f: \sigma_0 \qquad \sigma_0 \leq_{\mathfrak{m}}^0 \sigma_1 \to \mu$ $\Gamma \mapsto^{\operatorname{arg}} f: \sigma_1 \to \mu$

$$\Gamma \vdash^{\mathsf{h}} app : \forall a \ b. \ (a \to b) \to a \to b$$
$$\forall a \ b. \ (a \to b) \to a \to$$

$$\xrightarrow{} a \xrightarrow{} b \leq_{\mathfrak{m}}^{2} (\sigma_{1} \xrightarrow{} \mu), \sigma_{1}, \mu$$
$$\Gamma \vdash f \ e : \mu$$

The application of *app* can be instantiated with any type, given that both variables *a* and *b* are guarded. The premises coming from *f* and *e* are the same, except for the highlighted ones. We know that if  $\sigma_0 \leq_m^0 \sigma_1 \rightarrow \phi$ , then  $\sigma_0 \leq_m^1 \sigma_1, \phi$ , by inspection of the rule relating instantiation and function types. The case for *revapp* is similar.

## 4 TYPE INFERENCE USING CONSTRAINTS

In the previous section we described GI from a declarative perspective. But the utility of the system depends critically on whether we can formulate an efficient type inference algorithm for GI, a question that is the focus of this section.

Our approach to type inference is based on [Pottier and Rémy 2005], which works in two steps:

- Constraint generation: walk over the syntax tree of the source program, generating typing constraints. Constraint-generation is a simple algorithm over a large source-language syntax; i.e. there are many simple cases to consider. Constraint generation typically generates many unification variables that stand for as-yet-unknown types.
- (2) Constraint solving: solve the typing constraints generated in step 1. Constraint solving is a complex algorithm that works over a small language of constraints; i.e. there are few cases to consider but constraint solving is subtle. Constraint solving not only determines satisfiability (i.e. whether the program is typeable at all) but also a *substitution* that specifies the type that each unification variable stands for.

An attractive feature of this two-step approach is that it scales well with the complexity of the type system. For example, earlier work dubbed OUTSIDEIN(X) explains how to apply the Pottier and Rémy [2005] approach to a language like Haskell, with type classes, type-level functions, GADTs, and the like [Vytiniotis et al. 2011]. Another advantage is that is supports more sophisticated type-error diagnosis [Heeren et al. 2003; Serrano and Hage 2016; Zhang et al. 2015].

Performing type inference via constraints may carry a performance cost, due to the increased book-keeping. GHC alleviates this by solving "easy" constraints on-the-fly using ordinary unification, rather than emitting them to be solved later. Doing so might threaten the constraint-based error diagnosis; however the on-the-fly solver can easily be disabled when one wants to prioritise error messages over performance.

## 4.1 Constraints

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

head  $(\forall a. a \rightarrow a)$  ids 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 (*Bool*  $\rightarrow$  *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

686

Proc. ACM Program. Lang., Vol. 1, No. 1, Article 1. Publication date: January 2017.

687	Unification variable names	U	Э	$\alpha, \beta, \gamma, \delta, \ldots$	
688	Unification variables	υ	::=	$\alpha^{5}$	
689	Fully monomorphic types	τ	::=	$\alpha^{\mathfrak{m}} \mid a \mid T \ \overline{\tau}$	
690 691	Top-level monomorphic types	$\mu, \eta$	::=	$\alpha^{t} \mid a \mid T \overline{\sigma}$	
692	Polymorphic types	$\sigma, \phi$	::=	$\alpha^{\mathfrak{u}} \mid \forall \overline{a}. \mu$	$\overline{a}$ may be empty, $\overline{a} \subseteq \operatorname{ftv}(\mu)$
693	Types with generalisation	g	::=	$\mathbb{V}\{\overline{v}\}. C \Rightarrow \sigma$	$\overline{v}$ may be empty
694	Constraints	С	::=	Т	Trivial
695				$C_1 \wedge C_2$	Conjunction
696				$\sigma \sim \phi$	Equality
697				$\sigma \leqslant_{\mathfrak{s}}^n \mu$	Instantiation
698				$\mathfrak{g} \preceq \sigma$	Generalisation
699				$\forall \overline{a}. \exists \overline{v}. C$	Quantification
700	Free unification variables	$fuv(\sigma)$	=	free unification	n variables of $\sigma$

Fig. 9. Extended syntax

we know is that (*head ids*) has type  $\alpha$  for some as-yet-unknown type  $\alpha$ . So we want to defer the instantiation decision.

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

701 702

703 704

705

706

707

708

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

<sup>712</sup> in which (:) is instantiated at type ( $\forall a. a \rightarrow a$ ), and (:)'s first argument is generalised to have that <sup>713</sup> polymorphic type. Now consider constraint generation for this expression. We may instantiate <sup>714</sup> the type of (:) with a fresh unification variable,  $\alpha$  say. Ultimately the type of *ids* forces  $\alpha$  to be <sup>715</sup>  $\forall a. a \rightarrow a$ , but we don't know that yet. Moreover, in the final program we will need to generalise <sup>716</sup> the type of ( $\lambda x. x$ ), but again at constraint generation time we don't know that type either.

When we don't know something at constraint generation time, the solution is to *defer the choice*, *by generating a constraint that represents that choice*. This is the key idea of the constraint solving approach. The game is to develop a constraint language that neatly embodies the choices that we want to defer, and a solver that can subsequently make those choices. With that in mind, Figure 9 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 9, a key idea is that unification variables are drawn from three distinct "alphabets":  $\alpha^{s}$  for each of the threes 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*. Equality constraints are self explanatory. Instantiation constraints arise from the occurrence of a polymorphic variable, whose type must be instantiated – but that decision must be deferred (embodied in a constraint). Quantification constraints arise from explicit user type signatures, and pattern matching on data types involving existentials and GADTs. Both are fairly conventional. However *generalisation constraints* are new; they precisely embody the deferred decision about generalisation that we mention above. We will elaborate on all these forms in what follows.

Anon.

ANNAPP

1:16

$$\frac{x:\sigma \in \Gamma}{\Gamma \vdash^{h} x: \sigma \rightsquigarrow \epsilon} \text{VarHead} \qquad \frac{e \text{ not var. or app.} \qquad \Gamma \vdash e: \sigma \rightsquigarrow C}{\Gamma \vdash^{h} e: \sigma \rightsquigarrow C} \text{ ExprHead} \\
\frac{e \text{ not var. or app.} \qquad \Gamma \vdash e: \sigma \rightsquigarrow C}{\Gamma \vdash^{h} e: \sigma \rightsquigarrow C} \text{ AnnAbs} \\
\frac{\sigma \text{ fresh}}{\Gamma \vdash \lambda x. e: \alpha^{m} \rightarrow \sigma \rightsquigarrow C} \text{ Abs} \qquad \frac{\Gamma, x: \phi \vdash e: \sigma \rightsquigarrow C}{\Gamma \vdash \lambda (x:: \phi). e: \phi \rightarrow \sigma \rightsquigarrow C} \text{ AnnAbs} \\
\frac{\Gamma \vdash^{h} h: \phi \rightsquigarrow C}{\Gamma \vdash h: \phi \rightsquigarrow C} \qquad \frac{\Gamma \vdash^{h} e: \sigma_{i} \rightsquigarrow C_{i}}{\Gamma \vdash h(e_{1} \dots e_{n}: \delta_{t} \rightsquigarrow C)} \text{ Abs} \qquad \frac{\Gamma, x: \phi \vdash^{h} e: \sigma \rightsquigarrow C}{\Gamma \vdash^{h} (\sigma_{i}, C_{i}) - (\sigma_{i}, C_{i}) + \sigma_{i} \land^{h} (\sigma_{i}, C_{i}) - \sigma_{i} \land^{h} (\sigma_{i}, C_{i}) + \sigma_{i} \land^{h} (\sigma_{i}, C_{i})$$

Fig. 10. Constraint generation, part 1

In GI constraints do not form part of the source language; they are internal to the solver. But once we extend this approach to work with type system extensions such as type classes in in Section 5, we shall impose a distinction between simple constraints – those which the programmer can type in a program - and extended ones - which are internal to the algorithm.

#### **Constraint generation** 4.2

Constraint generation is described in Figure 10 as a four-element judgment  $\Gamma \vdash e : \sigma \rightsquigarrow C$ . The first two elements are inputs: the environment  $\Gamma$  and the expression *e* for which to generate constraints. The output of the process is a type  $\sigma$  assigned to the expression, possibly including some unification variables, and the set of constraints C that the types must satisfy. We first focus on the generation process for the core calculus, which we extend later to cover the rest of the language: annotations, let bindings and pattern matching. 

Rules ABS and ABSANN are not surprising: they just extend the environment with a new unifica-tion 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. 

Proc. ACM Program. Lang., Vol. 1, No. 1, Article 1. Publication date: January 2017.

Rule APP is where most of the work is done. Just like the declarative specification (Figure 5), the 785 head of the application is typed using an ancillary judgment  $\Gamma \vdash^{h} h : \phi \rightsquigarrow C$ , which either looks up 786 a variable in the environment or threads the information to the normal gathering process. 787

Then the first column of constraints, each of form  $\beta_i^{\mathfrak{u}} \leq_{\mathfrak{m}}^{n-i} \alpha_{i+1}^{\mathfrak{u}} \to \beta_{i+1}^{\mathfrak{u}}$ , successively decompose the function type  $\phi$  to expose its arguments. At each argument we may need to instantiate the 788 top-level foralls of the function type to expose the arrow; hence the use of subsumption  $\leq_{\alpha}^{n}$  rather than type equality. Note that what we express declaratively in a single relation  $\sigma_0 \leq_s^n \sigma_1, \ldots, \sigma_n, \mu$ is expressed as a sequence of constraints

 $\sigma_0 \leqslant_{\mathfrak{s}}^{n} \sigma_1 \to \phi_1 \quad \phi_1 \leqslant_{\mathfrak{s}}^{n-1} \sigma_2 \to \phi_2 \quad \dots \quad \phi^{n-1} \leqslant_{\mathfrak{s}}^{1} \sigma_n \to \phi_n \quad \phi_n \leqslant_{\mathfrak{s}}^{0} \mu$ 

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

796 The second column of constraints introduces a completely new constraint form,  $g \leq \sigma$ , which 797 we call generalisation constraint. These constraints allow us to defer the generalisation decisions to 798 the solver, as sketched in Section 4.1. The constraint  $(\mathbb{V}\{\overline{v}\}, C \Rightarrow \phi) \leq \sigma$  should be read "*a term of* type  $\phi$  with constraints C and unification variables  $\overline{v}$  can be instantiated and/or generalised to have 799 800 type  $\sigma$ ". Even looking at the syntax alone, you can see that the fruits of constraint generation for each argument  $e_i$  are wrapped up, along with the expected argument type  $\alpha_i$  from the function, 801 802 into a generalisation constraint for the solver to deal with later.

Rule ANNAPP deals with a type-annotated application  $(h e_1 \dots e_n :: \sigma)$ . It is similar to APP, 803 804 but it is the first rule to introduce a *quantification constraint*  $\forall \overline{b}$ .  $\exists \overline{v}$ .  $\overline{C}$ . This binds the Skolem 805 variables  $\overline{b}$  from the type signature, and existentially quantifies the unification variables free in the 806 constraint but not used outside it. The other significant difference is the use of  $\leq_{\mu}^{n}$  in the column of 807 instantiation constraints, with suffix u, rather than m in rule APP, exactly following the difference 808 between ANNAPP and APP in Figure 6 and 5 resp. 809

#### **Constraint solving** 4.3

The solver takes the generated constraint C and its free unification variables  $\overline{v} = fuv(C)$ , and repeatedly applies the solver rules in Figure 11, until no rule applies. The result is a residual constraint. If the residual constraint is in solved form, then the program is well typed; if not, the unsolved constraints (e.g. Int  $\sim$  Bool) represent type errors that can be reported to the user. We will discuss solved form shortly, in Section 4.3.3, but first we concentrate on the solver rules that incrementally solve the constraint.

Each of the rules in Figure 11 rewrites a configuration  $C; \overline{v}$  to another configuration. The unification variables  $\overline{v}$  are existentially quantified, so you can think of a configuration as representing  $\exists \overline{v}.C$ . Rule conj and FORALL are structural rules: the former allows a rule to be applied to one part of a conjunction, while the latter allows a rule to be applied under a quantification constraint. To avoid clutter we implicitly assume that the rules are read modulo commutativity and associativity of  $\wedge$ ; that is why CONJ only has to handle the left conjunct.

*Basic rules.* Rule EQREL removes trivial equality constraints of the form  $\sigma \sim \sigma$ . Note that 824 4.3.1 two polymorphic types need to be *syntactically equal* (modulo  $\alpha$ -equality) to match this rule. Rule 825 EQMONO indicates that two types headed by constructors are equal if and only if their heads coincide 826 and all the arguments are equal. 827

EQSUBST is the only rule that involves the interaction of two constraints. It applies the substi-828 tution of a unification variable to any other constraints conjoined with it (remember the implicit 829 associativity and commutativity of  $\wedge$ ), provided sorts are respected, and there is no occurs-check. 830 Notice that the equality constraint is not discarded; it remains in case it is needed again; indeed, 831 these equality constraints remain in a solved constraint (Section 4.3.3). 832

789

790

791 792

793

794

795

810

811

812

813

814

815

816

817

818

819

820

821

822

Anon.

 $\begin{array}{c} \hline \text{freshen}_{s}^{n}(\sigma) \implies \langle \overline{v}, \mu \rangle \\ \\ \mu \rhd_{s}^{n} \Delta \quad \overline{\alpha} \text{ fresh} \quad v_{i} = \alpha_{i}^{\Delta(a_{i})} \\ \\ \hline \text{freshen}_{s}^{n}(\forall \overline{a}. \mu) \implies \langle \overline{v}, [\overline{a \mapsto v}] \mu \rangle \end{array}$ 834 835 836 837 838 839  $C ; \overline{v} \implies C' ; \overline{v}'$ 840  $[\text{CONJ}] \xrightarrow{\hline C_1; \overline{v} \implies} C'_1; \overline{v}'$   $\overline{C_1 \land C_2; \overline{v} \implies} C'_1 \land C_2; \overline{v}'$ 841 842 843  $[\text{FORALL}] \xrightarrow{C ; \overline{v}_{\text{in}} \implies C' ; \overline{v}'_{\text{in}}}_{\forall \overline{a}. \exists \overline{v}_{\text{in}}. C ; \overline{v} \implies \forall \overline{a}. \exists \overline{v}'_{\text{in}}. C' ; \overline{v}}$ 844 845 846  $\begin{array}{cccc} \sigma \sim \sigma \; ; \; \overline{v} & \Longrightarrow & \top \; ; \; \overline{v} \\ \mathsf{T} \; \sigma_1 \; \dots \; \sigma_n \sim \mathsf{T} \; \phi_1 \; \dots \; \phi_n \; ; \; \overline{v} & \Longrightarrow & \sigma_1 \sim \phi_1, \dots, \sigma_n \sim \phi_n \; ; \; \overline{v} \\ (\alpha^{\mathfrak{s}} \sim \sigma) \land C \; ; \; \overline{v} & \Longrightarrow & (\alpha^{\mathfrak{s}} \sim \sigma) \land [\alpha^{\mathfrak{s}} \mapsto \sigma]C \; ; \; \overline{v} \end{array}$ 847 [EOREFL] 848 [EOMONO] 849 [EOSUBST] 850 if  $\vdash \sigma : \mathfrak{s}$ , and  $\alpha \notin \operatorname{ftv}(\sigma)$  $\begin{array}{cccc} \alpha^{\mathfrak{s}_{1}} \sim \beta^{\mathfrak{s}_{2}} ; \overline{\upsilon} & \Longrightarrow & \beta^{\mathfrak{s}_{2}} \sim \alpha^{\mathfrak{s}_{1}} ; \overline{\upsilon} & \text{if } \mathfrak{s}_{1} \sqsubset \mathfrak{s}_{2} \\ \alpha^{\mathfrak{m}} \sim \sigma ; \overline{\upsilon} & \Longrightarrow & \{\beta^{\mathfrak{s}} \sim \gamma^{\mathfrak{m}} \mid \beta^{\mathfrak{s}} \in \mathsf{fuv}(\sigma), \mathfrak{s} \neq \mathfrak{m}\} ; \overline{\upsilon}, \overline{\gamma^{\mathfrak{m}}} \end{array}$ 851 [EQVAR] 852 [EQFULLY] 853  $\begin{array}{cccc} \mu \leqslant_{\mathfrak{s}}^{n} \eta ; \overline{v} & \Longrightarrow & \mu \sim \eta ; \overline{v} \\ (\forall \overline{a}. \, \mu) \leqslant_{\mathfrak{s}}^{n} \eta ; \overline{v} & \Longrightarrow & \mu' \sim \eta ; \overline{v}. \overline{v}' \\ & & \text{where freshen}_{\mathfrak{s}}^{n} (\forall \overline{a}. \, \mu) \implies \langle \overline{v}', \mu' \rangle \end{array}$ [INSTMONO] 854 [INST∀L] 855 856  $(\mathbb{V}\{\overline{v}'\},\overline{C} \Rightarrow \sigma) \leq \eta ; \overline{v} \implies \overline{C} \land \left(\sigma \leqslant^{0}_{\mathfrak{m}} \eta\right) ; \overline{v}, \overline{v}'$ 857 [INST\U] 858  $\mathfrak{g} \leq (\forall \overline{a}, \mu); \overline{v} \implies \forall \overline{a}. (\mathfrak{g} \leq \mu); \overline{v}$ [inst∀r] 859 860

Fig. 11. Solving rules

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

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

4.3.2 Instantiation and generalisation constraints. The last set of rules take care of instantiation
 and generalisation constraints.

For instantiation constraints  $\leq_{s}^{n}$ , Rule INSTMONO encodes the fact that in our system if two toplevel monomorphic types  $\mu$ ,  $\eta$  are in an instance relation, they must be equal. This is a consequence of the invariance of type constructors. INST $\forall$ L instantiates a polytype  $\sigma$  with fresh unification variables, much as in the usual Damas-Milner algorithm, except that we must use a sort-respecting instantiation. This is done by freshen<sup>n</sup><sub>s</sub>, which in turn uses the already-introduced classification

882

861

862 863 864

865

866

867

868

869

870

judgment  $\triangleright_{\mathfrak{s}}^{n}$  (Figure 4). Finally, the new variables enter the set of existentially quantified vari-883 ables. Notice that the right hand side of an instantiation constraint  $\sigma \leq_{5}^{n} \mu$  is always a top-level 884 885 monomorphic type  $\mu$ , so we do not need to worry about having a polymorphic type on the right.

Finally, we come to generalisation constraints, which (recall Section 4.1) express a deferred 886 generalisation decision. Rule INSTVL is simple: if the right hand side has no top-level foralls (it is 887 of the form  $\eta$ ) then there is no generalisation to be done, so it suffices to release all the captured 888 constraints *C* and existentials v' into the current constraint. 889

INST WR is where actual generalisation takes place. In order to forge some intuition, let us look at 890 the constraint

$$(\mathbb{V}\{\alpha\,\beta\}.\,\alpha\leqslant^0_\mathfrak{m}\beta\Rightarrow\alpha\rightarrow\beta)\leq(\forall p.\,p\rightarrow p)$$

This generalisation constraint says that by performing some solving and possibly abstracting over some of the variables  $\alpha$  and  $\beta$ , 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.

$$\alpha \leqslant^0_{\mathfrak{m}} \beta \land \alpha \to \beta \sim p \to p$$

We obtain a solution by making  $\alpha \sim \beta \sim a$ . But in order for this solution to remain valid, we must guarantee that the skolem *p* does not escape to the outer world. We recall this restriction by means of a fresh quantification constraint  $\forall p. \exists \alpha \beta. (\alpha \leq_{\mathfrak{m}}^{0} \beta \land \alpha \rightarrow \beta \sim a \rightarrow a)$ . Rule INST  $\mathbb{V}_{\mathbb{R}}$  achieves this rather neatly simply by doing skolemisation and pushing the g inside; then INST VL 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  $\alpha^{\mu} \leq_{\mathfrak{s}}^{n} \mu$  we cannot turn it directly into  $\alpha^{\mu} \sim \mu$ , because  $\alpha^{\mu}$  might be unified later to a polymorphic type and we need instantiation.
- (2) Similarly, if we have  $(\forall \{\overline{\alpha}\}, C \Rightarrow \mu) \leq \alpha^{\mu}$ , and  $\alpha^{\mu}$  is later substituted by a polytype, we must skolemise.

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

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

The resulting constraint set is an instance of the problem of first-order unification under a mixed prefix. [Comon and Lescanne 1988] (which in our system is expressed by quantification constraints). A complete algorithm to solve this kind of problem is described by Pottier and Rémy [2005]. As a consequence, we have this result:

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

Unfortunately, once we extend the language of types, by introducing type classes and local assumptions, completeness is known not to hold [Vytiniotis et al. 2011]. Furthermore, the approach 928 by Pottier and Rémy [2005] is no longer applicable. With that in mind, Figure 12 presents one extra rule that, applied exhaustively like the others, also solves the problem of unification under a mixed 930

931

891 892 893

894

895

896

897 898

899

900

901

902

903 904

905

906

907

908

909

910

911

912

913

914

915 916

917

918

919

920

921

922

923 924

925

926

927

1:20

$$\begin{array}{ccc} & & & \\ \hline 932 \\ \hline 933 \\ \hline 934 \\ \hline 934 \\ \hline 934 \\ \hline 934 \\ \hline 935 \\ \hline \end{array} \begin{bmatrix} \operatorname{FLOAT} \end{bmatrix} & \frac{\operatorname{ftv}(F) \cap \overline{a} = \emptyset & \overline{\alpha^{\,\mathrm{s}}} = \operatorname{fuv}(F) \cap v_{\mathrm{in}} & \overline{\gamma^{\,\mathrm{s}}} & \operatorname{fresh} & \overline{E} = \bigwedge_{\alpha^{\,\mathrm{s}} \in \overline{\alpha}} \alpha^{\,\mathrm{s}} \sim \gamma^{\,\mathrm{s}} \\ \hline \forall \overline{a} . \exists \overline{v}_{\mathrm{in}} . (C \wedge F) ; \overline{v} \implies [\overline{\alpha^{\,\mathrm{s}}} \mapsto \gamma^{\,\mathrm{s}}] F \wedge \forall \overline{a} . \exists \overline{v}_{\mathrm{in}} . (C \wedge E) ; \overline{v}, \overline{\gamma^{\,\mathrm{s}}} \\ \hline 935 \\ \end{array}$$

Fig. 12. Solving rules for quantification constraints

$$\boxed{\begin{array}{c}\overline{a} ; \overline{\alpha} ; \overline{\beta} \vdash C \text{ solved}}\\ \hline \hline \overline{a} ; \overline{\alpha} ; \overline{\beta} \vdash \alpha^{s} \sim \sigma \text{ solved}} \end{array}}$$

$$\frac{\overline{a}; \overline{\alpha}; \overline{\beta}_{1} \vdash C_{1} \text{ solved } \overline{a}; \overline{\alpha}; \overline{\beta}_{2} \vdash C_{2} \text{ solved }}{\overline{a}; \overline{\alpha}; \overline{\beta}_{1} \uplus \overline{\beta}_{2} \vdash C_{1} \land C_{2} \text{ solved }}$$
SolvedConj

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

Fig. 13. Definition of solved set of constraints

prefix, and does scale to handle local assumptions. It is also rather simple, and is directly inspired by how GHC handles constraints.

Rule FLOAT simply floats constraints F from inside a quantification constraint to outside. When can we do that? Precisely when the constraint does not mention the skolems. But what about the existentials? For example, suppose we have

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

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

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

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

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

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

CONJECTURE 4.3. The solver presented in Figure 11 and Figure 12 is complete for unification problems under a mixed prefix.

4.3.3 Solved form. 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

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

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

 $\exists \alpha.(\alpha \sim \text{Int}) \land (\alpha \sim \text{Bool}) \qquad \exists \alpha...(\forall b.\alpha \sim [b])...$ 

Proc. ACM Program. Lang., Vol. 1, No. 1, Article 1. Publication date: January 2017.

981 982

- 985
- 986
- 987 988
- 989 990

991 992

1000 1001

1010

1011

1012

1020

1021

1022

1023

1024

1025 1026

1027

 $\frac{\overline{\Gamma \vdash e: \sigma \rightsquigarrow C}}{\overline{\Gamma \vdash e_1: \phi \rightsquigarrow C_1}} \xrightarrow{\Gamma \vdash e_1: \phi \rightsquigarrow C_1} \Gamma, x: \phi \vdash e_2: \sigma \rightsquigarrow C_2}_{\overline{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2: \sigma \rightsquigarrow C_1 \land C_2}} \text{ Let}$   $\frac{\Gamma \vdash e_0: \sigma_0 \rightsquigarrow C_0}{\overline{\alpha}, \beta \text{ fresh}} \xrightarrow{\text{for each branch } K_i \overline{x_i} \rightarrow e_i}_{\begin{array}{c}K_i: \forall \overline{a} \ \overline{b}_i. \overline{\sigma_i} \rightarrow T \ \overline{a} \in \Gamma \\ \overline{\Gamma, x_i: [\overline{a \mapsto \alpha^u}]} \sigma_i \vdash e_i: \phi_i \rightsquigarrow C_i \\ \overline{v_i} = \text{fuv}(\phi_i, C_i) - \text{fuv}(\Gamma) - \overline{\alpha}\end{array}}_{\overline{\Gamma} \vdash \text{case } e_0 \text{ of}\{\overline{K \overline{x} \rightarrow e}\}: \beta^u \rightsquigarrow C_0 \land (\sigma_0 \leqslant_{\mathfrak{m}}^0 T \overline{\alpha^u}) \land \forall \overline{b}_i. \exists \overline{v}_i. (C_i \land \beta^u \sim \phi_i)} \text{ Case}$ 

Fig. 14. Constraint generation, part 2

$$\widehat{C_{s}} = \theta$$

$$C_{s} = E \land R, \text{ where } E \text{ are all the equalities in } C_{s}$$

$$\widehat{C_{s}} = [\overline{\alpha \mapsto \sigma} \mid \alpha \sim \sigma \in E]$$

Fig. 15. Substitution induced by a solved form

<sup>1002</sup> In the first there are two equalities for  $\alpha$  (we should apply EQSUBST to make progress); in the second, <sup>1003</sup> there is a skolem-escape problem. However it is OK for a unification variable to have *no* equalities; <sup>1004</sup> it is simply unconstrained.

Figure 13 defines solved form precisely. We keep a set of variables  $\overline{\beta}$  for which we ensure that there is precisely one equality constraint, and another set  $\overline{\alpha}$  (the unconstrained variables) for which there are none. Rule SOLVEDVAR expects precisely one  $\beta$ , checks well-sortedness, and also checks that  $\sigma$  does not mention any variables other than the skolems and unconstrained unification variables – the latter check ensures idempotence.

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

### 1013 4.4 let bindings and pattern matching

The generation rules for the rest of the language are given in Figure 14. The rules are rather unsurprising. In rule CASE we need a quantification constraint to introduce new skolem constants, for the existentially quantified type variables  $\overline{b}_i$  in the data constructor. We refrain from describing the rules for **do** notation for the sake of conciseness; they do not pose any particular problem.

## <sup>1019</sup> 4.5 Soundness and principality

The presented gathering and solving process satisfies the usual properties of soundness and principality with respect to the declarative specification. In order to state these results, we need the auxiliary notion of substitution induced by a solved form, given in Figure 15.

THEOREM 4.4 (SOUDNESS). Let  $\Gamma$  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 4.5 (PRINCIPALITY). Suppose  $\Gamma \vdash e : \sigma$ . Then there exists a type  $\sigma^*$  such that  $\Gamma \vdash e : \sigma^*$  and  $\sigma = \pi \sigma^*$  where  $\pi$  is a fully monomorphic substitution.

Anon.

 $\sigma, \phi \quad ::= \quad \alpha^{\mathfrak{u}} \mid \forall \overline{a}. \ Q \Longrightarrow \mu$ Polymorphic types 1030 1031  $\begin{array}{ccc} Q & & \vdots = & \top \mid Q_1 \land Q_2 \\ & & \mid & \sigma \sim \phi \end{array}$ Simple constraints 1032 Equality 1033 Open for extension • • • 1034  $C \sigma_1 \dots \sigma_n$ Type classes, for example 1035  $:= Q \mid C_1 \land Q_2$ Extended constraints C 1036  $| \sigma \leqslant_{\mathfrak{s}}^{n} \mu$ Instantiation 37  $\mathfrak{g} \leq \sigma$ Generalisation 38  $\forall \overline{a}. \exists \overline{v}. (Q \supset C)$  Quantification and implication 39 40  $\Gamma \vdash e : \sigma \rightsquigarrow C$ 041 1042 1043 1045 for each branch  $K_i \overline{x_i} \rightarrow e_i$  $K_i: \forall \overline{a} \, \overline{b}_i. \ Q_i \implies \overline{\sigma_i} \to \mathsf{T} \, \overline{a} \in \Gamma$  $\Gamma \vdash e_0 : \sigma_0 \rightsquigarrow C_0 \qquad \overline{\alpha}, \beta \text{ fresh}$ 1048 q 1050 1051 1052 1053

Fig. 16. Extensions for integration with other constraints

The proofs of these results are given in Appendix C.2.

### 4.6 Implementation

A prototype implementation of the type inference process, optionally including the relaxed condi-1059 tions described in Section 6, is available at url.anonymized. The expressions in Figure 2 are accepted 1060 or rejected as described by the table. The prototype also includes support for Haskell's type classes 1061 by extending GI as described in Section 5. 1062

#### INTEGRATION WITH OTHER LANGUAGE FEATURES 5

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 1069 1070 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 1071  $(\equiv)$  ::  $\forall a. Eq a \Rightarrow a \rightarrow a \Rightarrow Bool$ . Similarly, languages like OCaml and PureScript feature row (or 1072 *record*) *types* like { *x* :: *Point*, *y* :: *Point* } in addition to usual ADTs. 1073

The good news is that a constraint-based formulation of typing makes it easy to cope with new 1074 1075 concepts if they can also be described in terms of constraints. Several examples in the literature, like HM(X) [Pottier and Rémy 2005] and OUTSIDEIN(X) [Vytiniotis et al. 2011], are actually frameworks 1076 which can be parametrized by different constraint systems, hence the X in their names. 1077 1078

1054 1055

1056 1057

1058

1063

1064

1065

1066

1067

1079

1092 1093

The modifications needed to accommodate the new kinds of constraints in GI are given in Figure 16. 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  $Eq \alpha$ . 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

$$\forall a. (Ord \ a \supset \mathbb{V}\{\alpha\}. Eq \ \alpha \Longrightarrow \alpha \longrightarrow Bool \le a \longrightarrow Bool$$

In turn, this requires changes to the solving rules (Figure 11 and Figure 12), which will then need to relate 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. 2011] so we do not discuss it further in the paper; we do provide the modified rules in Appendix A.

Implication constraints are also introduced by the updated rules for constraint gathering. An
 annotated application – rule ANNAPP – 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).

#### 1106 6 RELAXING GUARDEDNESS

In this section we show a very simple change to the constraint solver that allows it to solve more constraints. The downside is that, while it accepts all programs that satisfy the declarative specification, we do not yet have a simple declarative specification that precisely characterises all the programs it accepts.

We start with an example. Consider the expression (*choose* [] *ids*) from Figure 2. It fails to typecheck in GI, for the reason discussed in Section 3.3: Rule 1 says that we cannot instantiate the empty list [] polymorphically, forbidding the instance [ $\forall a. a \rightarrow a$ ] which is required. In this case the guardedness restriction seems to be too restrictive: all the information to infer the correct type seems to be there and we are not violating the "polymorphism is not guessed" policy.

What actually goes wrong in the solver? Here are the constraints generated for (choose [] ids):

1119

1105

$$\begin{array}{ll} \forall a. a \to a \to a \leqslant^2_{\mathfrak{m}} \alpha_1^{\mathfrak{u}} \to \beta_1^{\mathfrak{u}} & \beta_1^{\mathfrak{u}} \leqslant^2_{\mathfrak{m}} \alpha_2^{\mathfrak{u}} \to \beta_2^{\mathfrak{u}} & \beta_2^{\mathfrak{u}} \leqslant^0_{\mathfrak{m}} \delta_3^{\mathfrak{t}} \\ \mathbb{V}\{\delta_1^{\mathfrak{t}}\}. \forall a. [a] \leqslant^0_{\mathfrak{m}} \delta_1^{\mathfrak{t}} \Rightarrow \delta_1^{\mathfrak{t}} \le \alpha_1^{\mathfrak{u}} & \mathbb{V}\{\delta_2^{\mathfrak{t}}\}. [\forall a. a \to a] \leqslant^0_{\mathfrak{m}} \delta_2^{\mathfrak{t}} \Rightarrow \delta_2^{\mathfrak{t}} \le \alpha_2^{\mathfrak{u}} \end{array}$$

From rule APP we get one instantiation constraint and one generalisation constraint for each argument, plus one instantiation constraint for the result. The first instantiation is rewritten (rule INST  $\forall$ L) to  $\epsilon^t \rightarrow \epsilon^t \rightarrow \epsilon^t \leqslant_s^2 \alpha_1^u \rightarrow \beta_1^u$  where  $\epsilon^t$  is fresh.  $\epsilon$  is chosen to be top-level monomorphic since it appears in at least one given argument, but never under a constructor. This constraint in combination with the rest of the first row leads to a series of equalities  $\alpha_1^u \sim \alpha_2^u \sim \delta_3^t \sim \epsilon^t$ .

We replace  $\alpha_1^{\mu}$ ,  $\alpha_2^{\mu}$ , and  $\delta_1^{t}$  by the more restrictive  $\epsilon^{t}$  (rules EQVAR and EQSUBST). Now we can unwrap the first generalisation constraint (using INST VL), because  $\epsilon^{t}$  is guaranteed to be top-level 1136 1137

1128	[covr]		1 -	hafana	
1120	[CONJ]		As	belore	
1100	[FORALL]		As	before	
1130	[FLOAT]		As	before	
1131	[EQREFL]		As	before	
1132	[EOMONO]		As	before	
1133	[FOSUBST]	$(\alpha^5 \sim \sigma) \wedge C \cdot \overline{v}$	$\rightarrow$	$(\alpha^5 \sim \sigma) \wedge [\alpha^5 \mapsto \sigma]C : \overline{n}$	if $\vdash \sigma = \sigma$ and $\alpha \notin fty(\sigma)$
1134	[EQ30B31]	(u • 0) / C , 0		(u = 0) / [u = 70]e, 0	$11 \pm 0.5$ , and $11 \pm 1000$
1135					

Fig. 17. Solving rules for phase 2 of relaxed GI

<sup>1138</sup> monomorphic. That gives the constraint  $(\forall a.a \rightarrow a) \leq_{\mathfrak{m}}^{0} \epsilon^{\mathfrak{t}}$ , which we can instantiate (INST $\forall$ L) to <sup>1139</sup> get  $\epsilon^{\mathfrak{t}} \sim [\nu^{\mathfrak{m}}]$ . The key point here is that  $\nu^{\mathfrak{m}}$  is chosen to be *fully* monomorphic because the type <sup>1140</sup> of [] has no arguments. The final step is to unwrap the second generalisation constraint to get <sup>1141</sup>  $[\forall a. a \rightarrow a] \leq_{\mathfrak{m}}^{\mathfrak{m}} [\nu^{\mathfrak{m}}]$ , which quickly leads to the residual constraint  $\nu^{\mathfrak{m}} \sim \forall a. a \rightarrow a$ .

And at this point we are stuck! A monomorphic variable is equated with a quantified type, which is signaled as an error by the solver. But *why* is it an error? Remember, we are trying to find a substitution for the unification variables that solves the entire constraint. Well, making  $v^{m}$  stand for  $\forall a. a \rightarrow a$  will certainly solve this one! The reason for the sortedness constraints was to guide instantiation and generalisation. But since there are no instantiation or generalisation constraints, the presence of sorts only stops us from making progress. So at this point we simply ignore the sortedness restriction in EQSUBST and continue solving.

For example, suppose that at the moment that no further rules from Figure 11 can be applied, we are left with the following constraints:

$$C = \alpha^{\mathfrak{m}} \sim \forall a. a \rightarrow a \land \beta^{\mathfrak{t}} \sim [\alpha^{\mathfrak{m}}]$$

<sup>1153</sup> Note that this is not a solved form, because the equality over  $\alpha^{m}$  is ill-sorted, and the equality over <sup>1154</sup>  $\beta^{t}$  mentions another unification variable,  $\alpha^{m}$ , for which a further equality holds. Now we start <sup>1155</sup> phase 2 of the solver, using only the rules in Figure 17. Apart from the structural rules CONJ, FORALL <sup>1156</sup> and FLOAT, in this second phase we only handle equality constraints, *but now without respecting the* <sup>1157</sup> *sorts*. In other words, the  $\vdash \sigma : \mathfrak{s}$  condition is dropped from the updated rule Eqsubst, and the old <sup>1158</sup> rules EqvAR and EqFULLY – which were responsible for propagating sort restrictions – are gone. <sup>1159</sup> Applied to our set of constraints *C*, running through phase 2 leaves us with the following:

1160 1161

1171

1150

1151 1152

$$\alpha^{\mathfrak{m}} \sim \forall a. a \rightarrow a \land \beta^{\mathfrak{t}} \sim [\forall a. a \rightarrow a].$$

If we disregard the sortedness restrictions imposed by guardedness, these constraints are in solved form, and the program is accepted.

1164 A reasonable question to ask is: why do we need phase 2 at all? The reason is that until we have 1165 performed substitutions exhaustively, we cannot be sure that we have a solution free of occurs-1166 check errors. This happens, for example, if in the above example we replace the first constraint 1167 with  $\alpha^{m} \sim \forall a.a \rightarrow \beta^{t} \rightarrow a.$ 

We conclude this section by noting that we have also experimented with other, yet more powerful extensions of the solver, some involving iteration. It remains future work to decide their cost/benefit trade-off, and to seek better declarative specifications. The swamp is calling!

#### 1172 7 RELATED WORK

Early work has showed that full type inference for System F is an undecidable problem [Wells 1993]
– partial type inference with known positions of generalisation but unknown instantiations can be
reduced to higher-order unification [Pfenning 1988]. It is also folklore that simple specifications of

1:25

System F lack principal types, making thus modular type inference and the additions of ML-stylelet-bindings impossible.

1179 Type inference in the presence of higher-rank types but where instantiation of quantified variables is restricted to monomorphic types is a problem that has received successful and practical 1180 solutions [Dunfield and Krishnaswami 2013; Peyton Jones et al. 2007; Rémy 2005]. These systems 1181 exploit type annotations and a mixture of bidirectional propagation and unification and enjoy 1182 fairly simple declarative specifications while algorithmically being modest extensions of Algorithm 1183 W. Because of predicative instantiation they are typically more expressive than a predicative 1184 1185 restriction of System F and type programs in  $F\eta$ , a variant of System F with implicit  $\eta$ -conversions. Notably [Rémy 2005] present a modification that does allow impredicative - but explicit, based on 1186 annotations and quantified constraints - instantiation. 1187

On the other hand no satisfactory solution has been presented to-date that tackles impredica-1188 tivity with a good benefit-to-weight ratio. The MLF work [Botlan and Rémy 2003, 2009; Rémy 1189 1190 and Yakobowski 2008] presents a declarative specification and type inference algorithm for an (impredicative) extension of System F that relies on extending quantification with instance bounds. 1191 The resulting system is powerful, but also quite complex to implement; in return we get back 1192 principal types. Stemming from this work there have been several attempts to simplify the spec-1193 ification and restrict the user-facing part of MLF to System F types, while still relying on the 1194 1195 same powerful unification algorithm. For example FPH [Vytiniotis et al. 2008] exposes an extra "box" structure around inferred types (that would be hidden under a constraint in the context 1196 in MLF). Flexible types [Leijen 2009] accept quantification over instantiation constraints but not 1197 rigid equality constraints, thanks to a slightly different treatment of type annotations than MLF. 1198 Implementing any of these systems in a working compiler is a significant undertaking, and so is 1199 the integration with other forms of polymorphism such as type classes [Leijen and Löh 2005]. 1200

There have been also proposals of systems with simpler algorithms than MLF and hence more 1201 likely to be incorporated in existing type inference engines. Boxy Types [Vytiniotis et al. 2006] has 1202 been an early - largely unsuccessful due to its specification complexity - attempt to push bidirec-1203 tional type inference for higher-rank types to do a small amount of impredicative instantiation. 1204 HMF [Leijen 2008] on the other hand, stems from a generalisation of the Hindley-Milner algorithm, 1205 and keeps some of its nice properties such as principal types but imposes side conditions which 1206 range over all possible derivations, making it more difficult to grasp why a program fails to type 1207 check. QML [Russo and Vytiniotis 2009], inspired by earlier work on boxed polymorphism [O'Toole 1208 and Gifford 1989; Rémy 1994], proposes the introduction of a non-implicitly-instantiable univer-1209 sal quantifier. In QML, let-bindings can still be generalized a-la Hindley-Milner and implicitly 1210 instantiated. In addition first-class polymorphic values can be passed and returned, but can only 1211 be instantiated and constructed explicitly, through annotations. The presence of two forms of 1212 quantifiers as well as a rather inflexible form of annotation hindered the practicality of the system. 1213 Recent work [Eisenberg et al. 2016] also proposes the distinction between an implicitly instantiable 1214 and an explicitly instantiable universal quantifier; explicit type applications can be impredicative, 1215 whereas implicit ones have to be monomorphic. 1216

1217 We return now to Figure 2, where we present a collection of examples appearing in selected related works on type inference for impredicativity, namely MLF [Botlan and Rémy 2003], HMF [Leijen 1218 2008], FPH [Vytiniotis et al. 2008], HML [Leijen 2009]. We have selected those systems to compare 1219 because they strike extremely well in expressivity and require very few type annotations. The table 1220 shows the flexibility/expressivity price we pay in order to keep the implementation and specification 1221 1222 costs low, and avoid the introduction of new type system features (such as types with constraints or boxes) or new forms of annotations (such as QML-style instantiation annotations). We also show 1223 the annotation that can recover the typeability of a program in cases where a valid type exists. As 1224 1225

one observes, MLF can type all of the programs that do not require an implicit  $\eta$ -expansion (k h lst) 1226 or the use of a polymorphic function argument at two types  $\lambda f$ . (f 1, f True). Unsurprisingly, HML 1227 1228 is only minimally less powerful as it cannot type  $\lambda xs$ . poly (head xs) because xs would have to be assigned a polymorphic type, even though it's only used at this type. FPH is equally expressive for 1229 the applicative fragment of System F – however its treatment of  $\lambda$ -abstractions requires the returned 1230 type to be a fully-resolved top-level monomorphic type and hence fails to type check choose id auto 1231 because *auto* will be assigned the same type as GI infers (( $\forall a. a \rightarrow a) \rightarrow b \rightarrow b$ ). HMF on 1232 1233 the other hand is just based on local decisions about polymorphic instantiations and - without an extension to *n*-ary applications – fails to type check programs where the local instantiation 1234 has to be delayed to take more arguments into account (e.g. fails to type check id : ids). The 1235 relaxed GI system on the other hand can - modulo the quirk about not generalizing the bodies 1236 of lambda abstractions that only MLF and HML can tackle – type check the same programs as 1237 those systems. The vanilla (non-relaxed) GI for which a declarative specification exists is more 1238 1239 restrictive than those systems but incomparable to HMF. To determine why a program fails to type check (and how it should be fixed) it suffices to determine whether some function needs to 1240 be instantiated to a function type with top-level polymorphism in its arguments. For example, 1241 f (choose id) ids fails to type check because it requires the instantiation of choose ::  $\forall a. a \rightarrow a \rightarrow a$ 1242 to  $(\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$  and none of the arguments have a type with a 1243 top-level constructor. The fix is to add an annotation around the offending expression to determine 1244 its type: f (choose id ::  $(\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$ ) ids. The relaxed GI pretty much recovers some 1245 of the lost expressivity, and is by construction as expressive as the vanilla GI, but we have yet to 1246 determine whether it admits a simple declarative specification. 1247

#### 1249 8 CONCLUSION

1248

In this paper we have introduced the notion of guarded impredicativity, based on the key idea 1250 that in a function application, type variables that occur under a type constructor can be safely 1251 instantiated impredicatively. As Figure 2 shows, our system is somewhat less powerful than the 1252 most expressive impredicative type systems that can be found in the literature (HMF, HML, FPH 1253 and MLF). However, our system has the combined advantage of easy integration with existing 1254 constraint solving frameworks such as GHC's OUTSIDEIN(X) and a concise declarative specification 1255 making it easy for programmers to predict whether their program will compile. We have given a 1256 sound and complete inference algorithm, and proven a principal type property for our system. A 1257 prototype implementation, including integration with Haskell's type classes, is available. 1258

An obvious candidate for future work is to integrate our type system into GHC proper. Although it will certainly involve some work, e.g., the existing higher-ranked type system now implemented in GHC can be removed, our prototype was constructed in a way that is consistent with the architecture of the current OUTSIDEIN(X) implementation. We do not expect big complications to arise. Because we are then dealing with a single integrated contraint solving framework in GHC, we can investigate the use of heuristics for type error improving diagnosis and domain-specific type error diagnosis within the context of GHC [Hage and Heeren 2007; Serrano and Hage 2016].

As we hinted at the end of Section 6, we have already looked at yet more powerful extensions of the solver, some involving iteration, for which we need to come up with a suitable declarative specification.

1269 As the reader may have noticed, some expressions need to  $\eta$ -expanded in order to type check. 1270 This a direct consequence of the fact that all our type constructors are invariant. Seeing how far 1271 we remove the need for  $\eta$ -expansion is yet another interesting topic for further research.

1273 1274

#### 1275 **REFERENCES**

- Lennart Augustsson. 2011. Impredicative polymorphism, a use case. (2011). http://augustss.blogspot.nl/2011/07/
   impredicative-polymorphism-use-case-in.html.
- Didier Le Botlan and Didier Rémy. 2003. ML<sup>F</sup>: 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. DOI:http://dx.doi.org/10.1145/944705.944709
  - Didier Le Botlan and Didier Rémy. 2009. Recasting MLF. Inf. Comput. 207, 6 (2009), 726–785.
- Hubert Comon and Pierre Lescanne. 1988. Equational problems and disunification. Research Report RR-0904. INRIA.
   https://hal.inria.fr/inria-00075652
- 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. DOI:http://dx.doi.org/10.1145/2500365.
   2500582
- 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. DOI: http://dx.doi.org/10.1007/978-3-662-49498-1\_10
- Jurriaan Hage and Bastiaan Heeren. 2007. Heuristics for type error discovery and recovery. In *Implementation of Functional Languages IFL 2006*, Z. Horváth, V. Zsók, and A. Butterfield (Eds.), Vol. 4449. Springer Verlag, Heidelberg, 199 216.
- Bastiaan Heeren, Jurriaan Hage, and S. Doaitse Swierstra. 2003. Scripting the type inference process. *SIGPLAN Notices* 38, 9 (2003), 3–13. DOI:http://dx.doi.org/10.1145/944746.944707
- 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.
- Daan Leijen. 2008. HMF: simple type inference for first-class polymorphism, See [Hook and Thiemann 2008], 283–294. DOI: http://dx.doi.org/10.1145/1411204.1411245
- 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. DOI: http://dx.doi.org/10.1145/1480881.1480891
- Daan Leijen and Andres Löh. 2005. Qualified types for MLF. In Proceedings of the 10th ACM SIGPLAN International Conference on Functional Programming, ICFP 2005, Tallinn, Estonia, September 26-28, 2005, Olivier Danvy and Benjamin C. Pierce (Eds.). ACM, 144–155. DOI:http://dx.doi.org/10.1145/1086365.1086385
- Simon Marlow and others. 2010. Haskell 2010 Language Report. (2010). https://www.haskell.org/onlinereport/haskell2010/.
   Alberto Martelli and Ugo Montanari. 1982. An Efficient Unification Algorithm. ACM Trans. Program. Lang. Syst. 4, 2 (1982), 258–282. DOI: http://dx.doi.org/10.1145/357162.357169
- 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. DOI:http://dx.doi.org/10.1145/73141.74836
- 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.
- 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. DOI : http://dx.doi.org/10.
   1145/62678.62697
- François Pottier and Didier Rémy. 2005. The Essence of ML Type Inference. In Advanced Topics in Types and Programming Languages, Benjamin C. Pierce (Ed.). MIT Press, Chapter 10, 389–489. http://cristal.inria.fr/attapl/
- 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
- 1314Didier Rémy. 2005. Simple, Partial Type-inference for System F Based on Type-containment. In Proceedings of the Tenth1315ACM SIGPLAN International Conference on Functional Programming (ICFP '05). ACM, New York, NY, USA, 130–143. DOI:<br/>http://dx.doi.org/10.1145/1086365.1086383
- Didier Rémy and Boris Yakobowski. 2008. From ML to ML<sup>F</sup>: graphic type constraints with efficient type inference, See
   [Hook and Thiemann 2008], 63–74. DOI: http://dx.doi.org/10.1145/1411204.1411216
- 1318Claudio V. Russo and Dimitrios Vytiniotis. 2009. QML: Explicit First-class Polymorphism for ML. In Proceedings of the 20091319ACM SIGPLAN Workshop on ML (ML '09). ACM, New York, NY, USA, 3–14. DOI: http://dx.doi.org/10.1145/1596627.1596630
- 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
- 1323

- 2-8, 2016, Proceedings (Lecture Notes in Computer Science), Peter Thiemann (Ed.), Vol. 9632. Springer, 672-698. DOI: http://dx.doi.org/10.1007/978-3-662-49498-1\_26
- Dimitrios Vytiniotis, Simon L. Peyton Jones, Tom Schrijvers, and Martin Sulzmann. 2011. OUTSIDEIN(X): Modular type infer-ence with local assumptions. J. Funct. Program. 21, 4-5 (2011), 333-412. DOI: http://dx.doi.org/10.1017/S0956796811000098
- 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. DOI: http://dx.doi.org/10.1145/1159803.1159838
- Dimitrios Vytiniotis, Stephanie Weirich, and Simon L. Peyton Jones. 2008. FPH: first-class polymorphism for Haskell, See [Hook and Thiemann 2008], 295-306. DOI: http://dx.doi.org/10.1145/1411204.1411246
- J. B. Wells. 1993. Typability and Type Checking in the Second-Order Lambda-Calculus Are Equivalent and Undecidable. Technical Report. Boston, MA, USA.
- 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. DOI: http://dx.doi.org/10.1145/2737924.

Proc. ACM Program. Lang., Vol. 1, No. 1, Article 1. Publication date: January 2017.