

Visible Type Application

Richard A. Eisenberg, Stephanie Weirich, and Hamidhasan G. Ahmed

University of Pennsylvania
{eir,sweirich}@cis.upenn.edu
hamidhasan14@gmail.com

Abstract. The Hindley-Milner (HM) type system automatically infers the types at which polymorphic functions are used. In HM, the inferred types are unambiguous, and every expression has a principal type. Type annotations make HM compatible with extensions where complete type inference is impossible, such as higher-rank polymorphism and type-level functions. However, programmers cannot use annotations to explicitly provide type arguments to polymorphic functions, as HM requires type instantiations to be inferred.

We describe an extension to HM that allows visible type application. Our extension requires a novel type inference algorithm, yet its declarative presentation is a simple extension to HM. We prove that our extended system is a conservative extension of HM and admits principal types. We then extend our approach to a higher-rank type system with bidirectional type-checking. We have implemented this system in the Glasgow Haskell Compiler and show how our approach scales in the presence of complex type system features.

1 Introduction

The Hindley-Milner (HM) type system [7, 13, 18] achieves remarkable concision. While allowing a strong typing discipline, a program written in HM need not mention a single type. The brevity of HM comes at a cost, however: HM programs *must* not mention a single type. While this rule has long been relaxed by allowing visible type annotations (and even requiring them for various type system extensions), it remains impossible for languages based on HM, such as OCaml and Haskell, to use *visible type application* when calling a polymorphic function.¹

This restriction makes sense in the HM type system, where visible type application is unnecessary, as all type instantiations can be determined via unification. Suppose the function *id* has type $\forall a. a \rightarrow a$. If we wished to visibly instantiate the type variable *a* (in a version of HM extended with type annotations),

¹ Syntax elements appearing in a programmer’s source code are often called *explicit*, in contrast to *implicit* terms, which are inferred by the compiler. However, the implicit/explicit distinction also can indicate whether terms are computationally significant [19]. Our work applies only to the inferred vs. specified distinction, so we use *visible* to refer to syntax elements appearing in source code.

Declarative	Syntax-directed	
HM (§4.1)	C (§5.1)	from Damas and Milner [8] and Clément et al. [5]
HMV (§4.2)	V (§5.2)	HM types with visible type application
B (§6.2)	SB (§6.1)	Higher-rank types with visible type application

Fig. 1. The type systems studied in this paper

we could write the expression $(id :: Int \rightarrow Int)$. This annotation forces the type checker to unify the provided type $Int \rightarrow Int$ with the type $a \rightarrow a$, concluding that type a should be instantiated with Int .

However, this annotation is a roundabout way of providing information to the type checker. It would be much more direct if programmers could provide type arguments explicitly, writing the expression $id @Int$ instead.

Why do we want visible type application? In the Glasgow Haskell Compiler (GHC) – which is based on HM but extends it significantly – there are two main reasons:

First, type instantiation cannot always be determined by unification. Some Haskell features, such as type classes [28] and GHC’s type families [3, 4, 11], do not allow the type checker to unambiguously determine type arguments from an annotation. The current workaround for this issue is the *Proxy* type which clutters implementations and requires careful foresight by library designers. Visible type application improves such code. (See Sect. 2.)

Second, even when type arguments *can* be determined from an annotation, this mechanism is not always friendly to developers. For example, the variable to instantiate could appear multiple times in the type, leading to a long annotation. Partial type signatures help [29], but they do not completely solve the problem.²

Although the idea seems straightforward, adding visible type applications to the HM type system requires care, as we describe in Sect. 3. In particular, we observe that we can allow visible type application only at certain types: those with *specified type quantification*, known to the programmer via type annotation. Such types may be instantiated either visibly by the programmer, or when possible, invisibly through inference.

This paper presents a systematic study of the integration of visible type application within the HM typing system. In particular, the contributions of this paper are the four novel type systems (HMV, V, SB, B), summarized in Fig. 1. These systems come in pairs: a declarative version that justifies the compositionality of our extensions and a syntax-directed system that explains the structure of our type inference algorithm.

- System HMV extends the declarative version of the HM type system with a single, straightforward new rule for visible type application. In support of this feature, it also includes two other extensions: scoped type variables and

² The extended version of this paper [12] contains an example of this issue.

a distinction between specified and generalized type quantification. The importance of this system is that it demonstrates that visible type application can be added orthogonally to the HM type system, an observation that we found obvious only in hindsight.

- System V is a syntax-directed version of HMV. This type system directly corresponds to a type inference algorithm, called \mathcal{V} . Although Algorithm \mathcal{V} works differently than Algorithm \mathcal{W} [8], it retains the ability to calculate principal types. The key insight is that we can *delay* the instantiation of type variables until necessary. We prove that System V is sound and complete with respect to HMV, and that Algorithm \mathcal{V} is sound and complete with respect to System V. These results show the principal types property for HMV.
- System SB is a syntax-directed bidirectional type system with higher-rank types, *i.e.* higher-rank types. In extending GHC with visible type application, we were required to consider the interactions of System V with all of the many type system extensions featured in GHC. Most interactions are orthogonal, as expected from the design of V. However, GHC’s extension to support higher-rank types [23] changes its type inference algorithm to be bidirectional. System SB shows that our approach in designing System V straightforwardly extends to a bidirectional system. System SB’s role in this paper is twofold: to show how our approach to visible type application meshes well with type system extensions, and to be the basis for our implementation in GHC.
- System B is a novel, simple declarative specification of System SB. We prove that System SB is sound and complete with respect to System B. A similar declarative specification was not present in prior work [23]; this paper shows that an HM-style presentation is possible even in the case of higher-rank systems.

Our visible type application extension is part of GHC 8.0. The extended version [12] describes this implementation and elaborates on interactions between our algorithm and other features of GHC.³

2 Why visible type application?

Before we discuss how to extend HM type systems with visible type application, we first elaborate on why we would like this feature in the first place.

When a Haskell library author wishes to give a client the ability to control type variable instantiation, the current workaround is the standard library’s *Proxy* type.

```
data Proxy a = Proxy
```

³ Although the type inference algorithm in this paper uses unification to determine type instantiations, following Damas’s Algorithm \mathcal{W} [7], the results in this paper are applicable to GHC’s implementation based on constraint-solving. What matters here is how constraints are generated (specified by the syntax-directed versions of type systems) not how they are solved.

However, as we shall see, programming with *Proxy* is noisy and painfully indirect. With built-in visible type application, these examples are streamlined and easier to work with.⁴ In the following example and throughout this paper, unadorned code blocks are accepted by GHC 7.10, blocks with a solid gray bar at the left are ill-typed, and blocks with a gray background are accepted only by our implementation of visible type application.

Resolving type class ambiguity Suppose a programmer wished to normalize the representation of expression text by running it through a parser and then pretty printer. The *normalize* function below maps the string "7 - 1 * 0 + 3 / 3" to "((7 - (1 * 0)) + (3 / 3))", resolving precedence and making the meaning clear.⁵

```
normalize :: String → String
normalize x = show ((read :: String → Expr) x)
```

However, the designer of this function cannot make it polymorphic in a straightforward way. Adding a polymorphic type signature results in an ambiguous type, which GHC rightly rejects, as it cannot infer the instantiation for *a* at call sites.

```
normalizePoly :: ∀ a. (Show a, Read a) ⇒ String → String
normalizePoly x = show ((read :: String → a) x)
```

Instead, the programmer must add a *Proxy* argument, which is never evaluated, to allow clients of this polymorphic function to specify the parser and pretty-printer to use

```
normalizeProxy :: ∀ a. (Show a, Read a)
                ⇒ Proxy a → String → String
normalizeProxy _ x = show ((read :: String → a) x)
normalizeExpr :: String → String
normalizeExpr = normalizeProxy (Proxy :: Proxy Expr)
```

With visible type application, we can write these two functions more directly:⁶

⁴ Visible type application has been a GHC feature request since 2011. See <https://ghc.haskell.org/trac/ghc/ticket/5296>.

⁵ This example uses the following functions from the standard library,

```
show :: ∀ a. Show a ⇒ a → String
read :: ∀ a. Read a ⇒ String → a
```

as well as user-defined instances of the *Show* and *Read* classes for the type *Expr*.

⁶ Our new extension `TypeApplications` implies the extension `AllowAmbiguousTypes`, which allows our updated *normalize* definition to be accepted.

```

normalize :: ∀ a. (Show a, Read a) ⇒ String → String
normalize x = show (read @a x)
normalizeExpr :: String → String
normalizeExpr = normalize @Expr

```

Although the *show/read* ambiguity is somewhat contrived, proxies are indeed useful in more sophisticated APIs. For example, suppose a library designer would like to allow users to choose the representation of an internal data structure to best meet the needs of their application. If the type of that data structure is not included in the input and output types of the API, then a *Proxy* argument is a way to give this control to clients.⁷

Other examples More practical examples of the need for visible type application require a fair amount of build-up to motivate the need for the intricate types involved. We have included two larger examples in the extended version [12]. One builds from recent work on deferring constraints until runtime [2] and the other on translating a dependently typed program from Agda [16] into Haskell.

3 Our approach to visible type application

Visible type application seems like a straightforward extension, but adding this feature – both to GHC and to the HM type system that it is based on – turned out to be more difficult and interesting than we first anticipated. In particular, we encountered two significant questions.

3.1 Just *what* are the type parameters?

The first problem is that it is not always clear what the type parameters to a polymorphic function are!

One aspect of the HM type system is that it permits expressions to be given multiple types. For example, the identity function for pairs,

$$pid\ (x, y) = (x, y)$$

can be assigned any of the following most general types:

- (1) $\forall a\ b. (a, b) \rightarrow (a, b)$
- (2) $\forall a\ b. (b, a) \rightarrow (b, a)$
- (3) $\forall c\ a\ b. (a, b) \rightarrow (a, b)$

All of these types are principal; no type above is more general than any other. However, the type of the expression,

⁷ See <http://stackoverflow.com/questions/27044209/haskell-why-use-proxy>

Class constraints do not have a fixed ordering in types, and it is possible that a type variable is mentioned *only* in a constraint. Which of the following is preferred?

$$\begin{aligned} &\forall r m w a. (\text{MonadReader } r m, \text{MonadWriter } w m) \Rightarrow a \rightarrow m a \\ &\forall w m r a. (\text{MonadWriter } w m, \text{MonadReader } r m) \Rightarrow a \rightarrow m a \end{aligned}$$

Equality constraints and GADTs can add new quantified variables. Should we prefer the type $\forall a. a \rightarrow a$ or the equivalent type $\forall a b. (a \sim b) \Rightarrow a \rightarrow b$?

Type abbreviations mean that quantifying variables as they appear can be ambiguous without also specifying how type abbreviations are used and when they are expanded. Suppose

```
type Phantom a = Int
type Swap a b = (b, a)
```

Should we prefer $\forall a b. \text{Swap } a b \rightarrow \text{Int}$ or $\forall b a. \text{Swap } a b \rightarrow \text{Int}$? Similarly, should we prefer $\forall a. \text{Phantom } a \rightarrow \text{Int}$ or $\text{Int} \rightarrow \text{Int}$?

Fig. 2. Why specified polytypes?

```
pid @Int @Bool
```

is very different depending on which “equivalent” type is chosen for *pid*:

$$\begin{aligned} (\text{Int}, \text{Bool}) &\rightarrow (\text{Int}, \text{Bool}) && \text{-- } pid \text{ has type (1)} \\ (\text{Bool}, \text{Int}) &\rightarrow (\text{Bool}, \text{Int}) && \text{-- } pid \text{ has type (2)} \\ \forall b. (\text{Bool}, b) &\rightarrow (\text{Bool}, b) && \text{-- } pid \text{ has type (3)} \end{aligned}$$

Of course, there are ad hoc mechanisms for resolving this ambiguity. We could try to designate one of the above types (1–3) as the real principal type for *pid*, perhaps by disallowing the quantification of unused variables (ruling out type 3 above) and by enforcing an ordering on how variables are quantified (preferring type 1 over type 2 above). Our goal would be to make sure that each expression has a *unique* principal type, with respect to its quantified type variables. However, in the context of the full Haskell language, this strategy fails. There are just too many ways that types that are not α -equivalent can be considered equivalent by HM. See Fig. 2 for a list of language features that cause difficulties.

In the end, although it may be possible to resolve all of these ambiguities, we prefer not to. That approach leads to a system that is fragile (a new extension could break the requirement that principal types are unique up to α -equivalence), difficult to explain to programmers (who must be able to determine which type is principal) and difficult to reason about.

Our solution: specified polytypes Our system is designed around the following principle:

Only user-specified type parameters can be instantiated via explicit type applications.

In other words, we allow visible type application to instantiate a polytype only when that type is given by a user annotation. This restriction follows in a long line of work requiring user annotations to support advanced type system features [14, 22, 23]. We refer to variables quantified in type annotations as *specified variables*, distinct from compiler-generated quantified variables, which we call *generalized variables*.

There is one nuance to this rule in practice. Haskell allows programming to omit variable quantification, allowing a type signature like

$$\text{const} :: a \rightarrow b \rightarrow a \quad \text{-- NB: no } \forall$$

Are these variables specified? We have decided that they are. There is a very easy rule at work here: just order the variables left-to-right as the user wrote them. We thus consider variables from type signatures to be specified, even when not bound by an explicit \forall .

3.2 Is our extension compatible with the rest of the type system?

We do not want to extend just the type inference algorithm that GHC uses. We would also like to extend its *specification*, which is rooted in HM. This way, we will have a concise description (and better understanding) of what programs type check, and a simple way to reason about the properties of the type system.

Our first attempt to add type application to GHC was based on our understanding of Algorithm \mathcal{W} , the standard algorithm for HM type inference. This algorithm instantiates polymorphic functions only at occurrences of variables. So, it seems that the only new form we need to allow is a visible type right after variable occurrences:

$$x @\tau_1 \dots @\tau_n$$

However, this extension is not very robust to code refactoring. For example, it is not closed under substitution. If type application is only allowed at variables, then we cannot substitute for this variable and expect the code to still type check. Therefore our algorithm should allow visible type applications at other expression forms. But where else makes sense?

For example, it seems sensible to allow a type instantiation is after a polymorphic type annotation (such an annotation certainly specifies the type of the expression):

$$(\lambda x \rightarrow x :: \forall a b. (a, b) \rightarrow (a, b)) @Int$$

Likewise, we should also allow a visible instantiation after a **let** to enable refactoring:⁸

$$(\mathbf{let} \ y = ((\lambda x \rightarrow x) :: \forall a b. (a, b) \rightarrow (a, b)) \ \mathbf{in} \ y) @Int$$

⁸ In fact, the Haskell 2010 Report [15] *defines* type annotations by expanding to a **let**-declaration with a signature.

However, how do we know that we have identified all sites where visible type applications should be allowed? Furthermore, we may have identified them all for core HM, but what happens when we go to the full language of GHC, which includes features that may expose new potential sites?

One way to think about this issue in a principled way is to develop a compositional specification of the type system, which allows type application for *any* expression that can be assigned a polytype. Then, if our algorithm is complete with respect to this specification, we will know that we have allowed type applications in all of the appropriate places. This specification is itself useful in its own right, as we will have a concise description (and better understanding) of what programs type check and a simple way to reason about the properties of the type system.

Once we started thinking of specifications, we found that Algorithm \mathcal{W} could not be matched up with the compositional specification that we wanted. That led us to reconsider our algorithm and develop a new approach to HM type inference.

Our solution: lazy instantiation for specified polytypes Our new type inference algorithm, which we call Algorithm \mathcal{V} , is based on the following design principle:

Delay instantiation of “specified” type parameters until absolutely necessary.

Although Algorithm \mathcal{W} instantiates all polytypes immediately, it need not do so. In fact, it is possible to develop a sound and complete alternative implementation of the HM type system that does not do this immediate instantiation. Instead, instantiation is done only on demand, such as when a polymorphic function is applied to arguments. Lazy instantiation has been used in (non-HM) type inference before [10] and may be folklore; however this work contains the first proof that it can be used to implement the HM type system.

In the next section, we give this algorithm a simple specification, presented as a small extension of HM’s existing declarative specification. We then continue with a syntax-directed account of the type system, characterizing where lazy instantiations actually must occur during type checking.

4 HM with visible type application

To make our ideas precise, we next review the declarative specification of the HM type system [7, 13, 18] (which we call System HM), and then show how to extend this specification with visible type arguments.

4.1 System HM

The grammar of System HM is shown in Fig. 3. The expression language comprises the Curry-style typed λ -calculus with the addition of numeric literals (of

HM	HMV
Metavariables: x, y term variables a, b, c type variables n numeric literals	This grammar extends HM: $e ::= \dots \mid e @ \tau \mid (\Lambda \bar{a}. e : v)$ expressions $\tau ::= \dots$ monotypes $v ::= \tau \mid \forall a. v$ spec. polytypes $\sigma ::= v \mid \forall \{a\}. \sigma$ type schemes $\Gamma ::= \cdot \mid \Gamma, x : \sigma \mid \Gamma, a$ contexts
$e ::= x \mid \lambda x. e \mid e_1 e_2$ expressions $\mid n \mid \mathbf{let} x = e_1 \mathbf{in} e_2$ $\tau ::= a \mid \tau_1 \rightarrow \tau_2 \mid \mathit{Int}$ monotypes $\sigma ::= \forall \{\bar{a}\}. \tau$ type schemes $\Gamma ::= \cdot \mid \Gamma, x : \sigma$ contexts	

We write $(e : v)$ to mean $(\Lambda \cdot . e : v)$, $ftv(\sigma)$ to be the set of type variables free in σ , and lift ftv to contexts with $ftv(\bar{x} : \bar{\sigma}, \bar{a})$ to be $\bar{a} \cup \bigcup_i ftv(\sigma_i)$.

Fig. 3. Grammars for Systems HM and HMV

type Int) and **let**-expressions. Monotypes are as usual, but we diverge from standard notation in type schemes as they quantify over a possibly-empty *set* of type variables. Here, we write these type variables in braces to emphasize that they should be considered order-independent. We sometimes write τ for the type scheme $\forall \{\cdot\}. \tau$ with an empty set of quantified variables and write $\forall \{a\}. \forall \{\bar{b}\}. \tau$ to mean $\forall \{a, \bar{b}\}. \tau$. Here – and throughout this paper – we liberally use the Barendregt convention that bound variables are always distinct from free variables.

The declarative typing rules for System HM appear in Fig. 4. (The figure also includes the definition for our extended system, called System HMV, described in Sect. 4.2.) System HM is not syntax-directed; rules HM_GEN and HM_SUB can apply anywhere.

So that we can better compare this system with others in the paper, we make two changes to the standard HM rules. Neither of these changes are substantial; our version types the same programs as the original.⁹ First, in HM_LET, we allow the type of a **let** expression to be a polytype σ , instead of restricting it to be a monotype τ . We discuss this change further in Sect. 5.2. Second, we replace the usual instantiation rule with HM_SUB. This rule allows the type of any expression to be converted to any less general type in one step (as determined by the subsumption relation $\sigma_1 \leq_{\text{hm}} \sigma_2$). Note that in rule HM_INSTG the lists of variables \bar{a}_1 and \bar{a}_2 need not be the same length.

4.2 System HMV: HM with visible types

System HMV is an extension of System HM, adding visible type application. A key detail in its design is its separation of specified type variables from those arising from generalization, as initially explored in Sect. 3.1. Types may be generalized at any time in HMV, quantifying over a variable free in a type but

⁹ Both ways of the equivalence proof proceed by induction, liberally using instantiation, generalization, and subsumption to bridge the gap between the two systems.

$\Gamma \vdash_{\text{hm}} e : \sigma$	Typing rules for System HM
$\frac{x:\sigma \in \Gamma}{\Gamma \vdash_{\text{hm}} x : \sigma} \text{HM_VAR} \quad \frac{\Gamma, x:\tau_1 \vdash_{\text{hm}} e : \tau_2}{\Gamma \vdash_{\text{hm}} \lambda x. e : \tau_1 \rightarrow \tau_2} \text{HM_ABS}$ $\frac{\Gamma \vdash_{\text{hm}} e_1 : \tau_1 \rightarrow \tau_2 \quad \Gamma \vdash_{\text{hm}} e_2 : \tau_1}{\Gamma \vdash_{\text{hm}} e_1 e_2 : \tau_2} \text{HM_APP} \quad \frac{}{\Gamma \vdash_{\text{hm}} n : \text{Int}} \text{HM_INT}$ $\frac{\Gamma \vdash_{\text{hm}} e_1 : \sigma_1 \quad \Gamma, x:\sigma_1 \vdash_{\text{hm}} e_2 : \sigma_2}{\Gamma \vdash_{\text{hm}} \text{let } x = e_1 \text{ in } e_2 : \sigma_2} \text{HM_LET}$ $\frac{\Gamma \vdash_{\text{hm}} e : \sigma \quad a \notin \text{ftv}(\Gamma)}{\Gamma \vdash_{\text{hm}} e : \forall\{a\}.\sigma} \text{HM_GEN} \quad \frac{\Gamma \vdash_{\text{hm}} e : \sigma_1 \quad \sigma_1 \leq_{\text{hm}} \sigma_2}{\Gamma \vdash_{\text{hm}} e : \sigma_2} \text{HM_SUB}$	
$\sigma_1 \leq_{\text{hm}} \sigma_2$	HM subsumption
$\frac{\tau_1[\bar{\tau}/\bar{a}_1] = \tau_2 \quad \bar{a}_2 \notin \text{ftv}(\forall\{\bar{a}_1\}.\tau_1)}{\forall\{\bar{a}_1\}.\tau_1 \leq_{\text{hm}} \forall\{\bar{a}_2\}.\tau_2} \text{HM_INSTG}$	

$\Gamma \vdash_{\text{hmv}} e : \sigma$	Extra typing rules for System HMV
$\frac{\Gamma \vdash \tau}{\Gamma \vdash_{\text{hmv}} e : \forall a. v} \text{HMV_TAPP} \quad \frac{\Gamma \vdash v \quad v = \forall \bar{a}, \bar{b}. \tau}{\Gamma, \bar{a} \vdash_{\text{hmv}} e : \tau \quad \bar{a}, \bar{b} \notin \text{ftv}(\Gamma)} \text{HMV_ANNOT}$ $\frac{\Gamma \vdash_{\text{hmv}} e @ \tau : v[\tau/\bar{a}]}{\Gamma \vdash_{\text{hmv}} (\Lambda \bar{a}. e : v) : v}$	
$\sigma_1 \leq_{\text{hmv}} \sigma_2$	HMV subsumption
$\frac{\tau_1[\bar{\tau}/\bar{b}] = \tau_2}{\forall \bar{a}, \bar{b}. \tau_1 \leq_{\text{hmv}} \forall \bar{a}. \tau_2} \text{HMV_INSTS} \quad \frac{v_1[\bar{\tau}/\bar{a}_1] \leq_{\text{hmv}} v_2 \quad \bar{a}_2 \notin \text{ftv}(\forall\{\bar{a}_1\}.\tau_1)}{\forall\{\bar{a}_1\}.\tau_1 \leq_{\text{hmv}} \forall\{\bar{a}_2\}.\tau_2} \text{HMV_INSTG}$	
$\Gamma \vdash v$	Type well-formedness
$\frac{\text{ftv}(v) \subseteq \Gamma}{\Gamma \vdash v} \text{TY_SCOPED}$	

Fig. 4. Typing rules for Systems HM and HMV

not free in the typing context. The type variable generalized in this manner is *not* specified, as the generalization takes place absent any direction from the programmer. By contrast, a type variable mentioned in a type annotation *is* specified, precisely because it is written in the program text.

The grammar of System HMV appears in Fig. 3. The type language is enhanced with a new intermediate form v that quantifies over an ordered list of type variables. (We sometimes write $\forall a. \forall b. \tau$ as $\forall a, b. \tau$.) This form sits between

$\forall\{a, b\}. a \rightarrow b \leq_{\text{hmv}} \forall\{a\}. a \rightarrow a$	Works the same as \leq_{hm} for type schemes
$\forall a, b. a \rightarrow b \leq_{\text{hmv}} \text{Int} \rightarrow \text{Int}$	Can instantiate specified vars
$\forall a, b. a \rightarrow b \leq_{\text{hmv}} \forall a. a \rightarrow \text{Int}$	Can instantiate only a <i>tail</i> of the specified vars
$\forall a, b. a \rightarrow b \leq_{\text{hmv}} \forall\{a, b\}. a \rightarrow b$	Variables can be regeneralized
$\forall a, b. a \rightarrow b \leq_{\text{hmv}} \forall\{b\}. \text{Int} \rightarrow b$	Right-to-left nature of HMV_INSTS forces regeneralization
$\forall a, b. a \rightarrow b \not\leq_{\text{hmv}} \forall b. \text{Int} \rightarrow b$	Known vars are instantiated from the right, never the left
$\forall\{a\}. a \rightarrow a \not\leq_{\text{hmv}} \forall a. a \rightarrow a$	Specified quantification is more general than generalized quantification

Fig. 5. Examples of HMV subsumption relation

type schemes and monotypes; σ s contain ν s, which then contain τ s.¹⁰ Thus the full form of a type scheme σ can be written as $\forall\{\bar{a}\}, \bar{b}. \tau$, including both a set of generalized variables $\{\bar{a}\}$ and a list of specified variables \bar{b} . Note that order never matters for generalized variables (they are in a set) while order does certainly matter for specified variables (the list specifies their order). We say that ν is the metavariable for *specified polytypes*, distinct from *type schemes* σ .

Expressions in HMV include two new forms: $e@_\tau$ instantiates a specified type variable with a monotype τ , while $(\Lambda\bar{a}. e : \nu)$ allows us to annotate an expression with its type, potentially binding scoped type variables if the type is polymorphic. Requiring a type annotation in concert with scoped type variable binding ensures that the order of quantification is specified: the type annotation is a specified polytype ν . We do not allow annotation by type schemes σ : if the user writes the type, all quantified variables are considered specified.

Typing contexts Γ in HMV are enhanced with the ability to store type variables. This feature is used to implement scoped type variables, where the type variables \bar{a} , bound in $\Lambda\bar{a}. e$, are available for use in types occurring within e .

Typing rules The type system of HMV includes all of the rules of HM plus the new rules and relation shown at the bottom of Fig. 4. The HMV rules inherited from System HM are modified to recur back to System HMV relations: in effect, replace all hm subscripts with hmv subscripts. Note, in particular, rule HM_SUB; in System HMV, this rule refers to the $\sigma_1 \leq_{\text{hmv}} \sigma_2$ relation, described below.

The most important addition to this type system is HMV_TAPP, which enables visible type application when the type of the expression is quantified over a specified type variable.

¹⁰ The grammar for System HMV redefines several metavariables. These metavariables then have (slightly) different meanings in different sections of this paper, but disambiguation should be clear from context. In analysis relating systems with different grammars (for example, in Lemma 1), the more restrictive grammar takes precedence.

A type annotation $(\Lambda \bar{a}. e : v)$, typed with `HMV_ANNOT`, allows an expression to be assigned a specified polytype. We require the specified polytype to have the form $\forall \bar{a}, \bar{b}. \tau$; that is, a prefix of the specified polytype’s quantified variables must be the type variables scoped in the expression.¹¹ The inner expression e is then checked at type τ , with the type variables \bar{a} (but not the \bar{b}) in scope. Types that appear in expressions (such as in type annotations and explicit type applications) may mention only type variables that are currently in scope.

Of course, in the $\Gamma, \bar{a} \vdash_{\text{hmV}} e : \tau$ premise, the variables \bar{a} and \bar{b} may appear in τ . We call such variables *skolems* and say that *skolemizing* v yields τ . In effect, these variables form new type constants when type-checking e . When the expression e has type τ , we know that e cannot make any assumptions about the skolems \bar{a}, \bar{b} , so we can assign e the type $\forall \bar{a}, \bar{b}. \tau$. This is, in effect, *specified generalization*.

The relation $\sigma_1 \leq_{\text{hmV}} \sigma_2$ (Fig. 4) implements subsumption for System HMV. The intuition is that, if $\sigma_1 \leq_{\text{hmV}} \sigma_2$, then an expression of type σ_1 can be used wherever one of type σ_2 is expected. For type schemes, the standard notion of σ_1 being a more general type than σ_2 is sufficient. However for specified polytypes, we must be more cautious.

Suppose an expression $x @_{\tau_1} @_{\tau_2}$ type checks, where x has type $\forall a, b. v_1$. The subsumption rule means that we can arbitrarily change the type of x to some v , as long as $v \leq_{\text{hmV}} \forall a, b. v_1$. Therefore, v must be of the form $\forall a, b. v_2$ so that $x @_{\tau_1} @_{\tau_2}$ will continue to instantiate a with τ_1 and b with τ_2 . Accordingly, we cannot, say, allow subsumption to reorder specified variables.

However, it is safe to allow *some* instantiation of specified variables as part of subsumption, as in rule `HMV_INSTS`. Examine this rule closely: it instantiates variables *from the right*. This odd-looking design choice is critical. Continuing the example above, v could also be of the form $\forall a, b, c. v_3$. In this case, the additional specified variable c causes no trouble – it need not be instantiated by a visible application. But we cannot allow instantiation *left-to-right* as that would allow the visible type arguments to skip instantiating a or b .

Further examples illustrating \leq_{hmV} appear in Fig. 5.

4.3 Properties of System HMV

We wish System HMV to be a conservative extension of System HM. That is, any expression that is well-typed in HM should remain well-typed in HMV, and any expression not well-typed in HM (but written in the HM subset of HMV) should also not be well-typed in HMV.

Lemma 1 (Conservative Extension for HMV). *Suppose Γ and e are both expressible in HM; that is, they do not include any type instantiations, type annotations, scoped type variables, or specified polytypes. Then, $\Gamma \vdash_{\text{hm}} e : \sigma$ if and only if $\Gamma \vdash_{\text{hmV}} e : \sigma$.*

¹¹ Note that the Barendregt convention allows bound variables to α -vary, so only the number of scoped type variables is important.

This property follows directly from the definition of HMV as an extension of HM. Note, in particular, that no HM typing rule is changed in HMV and that the \leq_{hmv} relation contains \leq_{hm} ; furthermore, the new rules all require constructs not found in HM.

We also wish to know that making generalized variables into specified variables does not disrupt types:

Lemma 2 (Extra knowledge is harmless). *If $\Gamma, x:\forall\{\bar{a}\}.\tau \vdash_{\text{hmv}} e : \sigma$, then $\Gamma, x:\forall\bar{a}.\tau \vdash_{\text{hmv}} e : \sigma$.*

This property follows directly from the context generalization lemma below, noting that $\forall\bar{a}.\tau \leq_{\text{hmv}} \forall\{\bar{a}\}.\tau$.

Lemma 3 (Context generalization for HMV). *If $\Gamma \vdash_{\text{hmv}} e : \sigma$ and $\Gamma' \leq_{\text{hmv}} \Gamma$, then $\Gamma' \vdash_{\text{hmv}} e : \sigma$.*

This lemma is proved in the extended version [12].

In practical terms, Lemma 2 means that if an expression contains **let** $x = e_1$ **in** e_2 , and the programmer figures out the type assigned to x (say, $\forall\{\bar{a}\}.\tau$) and then includes that type in an annotation (as **let** $x = (e_1 : \forall\bar{a}.\tau)$ **in** e_2), the outer expression's type does not then change.

However, note that, by design, context generalization is not as flexible for specified polytypes as it is for type schemes. In other words, suppose the following expression type-checks.

let $x = ((\lambda y \rightarrow y) :: \forall a b. (a, b) \rightarrow (a, b))$ **in** ...

The programmer cannot then replace the type annotation with the type $\forall a. a \rightarrow a$, because x may be used with visible type applications. This behavior may be surprising, but it follows directly from the fact that $\forall a. a \rightarrow a \not\leq_{\text{hmv}} \forall a b. (a, b) \rightarrow (a, b)$.

Finally, we would also like to show that HMV retains the principal types property, defined with respect to the enhanced subsumption relation $\sigma_1 \leq_{\text{hmv}} \sigma_2$.

Theorem 4 (Principal types for HMV). *For all terms e well-typed in a context Γ , there exists a type scheme σ_p such that $\Gamma \vdash_{\text{hmv}} e : \sigma_p$ and, for all σ such that $\Gamma \vdash_{\text{hmv}} e : \sigma$, $\sigma_p \leq_{\text{hmv}} \sigma$.*

Before we can prove this, we first must show how to extend HM's type inference algorithm (Algorithm \mathcal{W} [8]) to include visible type application. Once we do so, we can prove that this new algorithm always computes principal types.

5 Syntax-directed versions of HM and HMV

The type systems in the previous section declare when programs are well-formed, but they are fairly far removed from an algorithm. In particular, the rules HM_GEN and HM_SUB can appear at any point in a typing derivation.

$\Gamma \vdash_{\mathbb{C}} e : \tau$

 Typing rules for System C
$$\frac{x:\forall\{\bar{a}\}. \tau \in \Gamma}{\Gamma \vdash_{\mathbb{C}} x : \tau[\bar{\tau}/\bar{a}]} \text{C_VAR} \quad \frac{\Gamma, x:\tau_1 \vdash_{\mathbb{C}} e : \tau_2}{\Gamma \vdash_{\mathbb{C}} \lambda x. e : \tau_1 \rightarrow \tau_2} \text{C_ABS}$$

$$\frac{\Gamma \vdash_{\mathbb{C}} e_1 : \tau_1 \rightarrow \tau_2 \quad \Gamma \vdash_{\mathbb{C}} e_2 : \tau_1}{\Gamma \vdash_{\mathbb{C}} e_1 e_2 : \tau_2} \text{C_APP} \quad \frac{}{\Gamma \vdash_{\mathbb{C}} n : \text{Int}} \text{C_INT}$$

$$\frac{\Gamma \stackrel{\text{gen}}{\vdash}_{\mathbb{C}} e_1 : \sigma \quad \Gamma, x:\sigma \vdash_{\mathbb{C}} e_2 : \tau_2}{\Gamma \vdash_{\mathbb{C}} \text{let } x = e_1 \text{ in } e_2 : \tau_2} \text{C_LET}$$

$\Gamma \stackrel{\text{gen}}{\vdash}_{\mathbb{C}} e : \sigma$

 Generalization for System C
$$\frac{\bar{a} = \text{ftv}(\tau) \setminus \text{ftv}(\Gamma) \quad \Gamma \vdash_{\mathbb{C}} e : \tau}{\Gamma \stackrel{\text{gen}}{\vdash}_{\mathbb{C}} e : \forall\{\bar{a}\}. \tau} \text{C_GEN}$$

Fig. 6. Syntax-directed version of the HM type system

5.1 System C

We can explain the HM type system in a more algorithmic manner by using a syntax-directed specification, called System C, in Fig. 6. This version of the type system, derived from Clément et al. [5], clarifies exactly where generalization and instantiation occur during type checking. Notably, instantiation occurs only at the usage of a variable, and generalization occurs only at a **let**-binding. These rules are syntax-directed because the conclusion of each rule in the main judgment $\Gamma \vdash_{\mathbb{C}} e : \tau$ is syntactically distinct. Thus, from the shape of an expression, we can determine the shape of its typing derivation.

However, the judgment $\Gamma \vdash_{\mathbb{C}} e : \tau$ is still not quite an algorithm: it makes non-deterministic guesses. For example, in the rule C_ABS, the type τ_1 is guessed; there is no indication in the expression what the choice for τ_1 should be. The advantage of studying a syntax-directed system such as System C is that doing so separates concerns: System C fixes the structure of the typing derivation (and of any implementation) while leaving monotype-guessing as a separate problem. Algorithm \mathcal{W} deduces the monotypes via unification, but a constraint-based approach [25, 27] would also work.

5.2 System V: Syntax-directed visible types

Just as System C is a syntax-directed version of HM, we can also define System V, a syntax-directed version of H MV (Fig. 7). However, although we could define H MV by a small addition to HM (two new rules, plus subsumption), the difference between System C and System V is more significant.

Like System C, System V uses multiple judgments to restrict where generalization and instantiation can occur. In particular, the system allows an expres-

$\Gamma \vdash e : \tau$	Monotype checking for System V
$\frac{\Gamma, x:\tau_1 \vdash e : \tau_2}{\Gamma \vdash \lambda x. e : \tau_1 \rightarrow \tau_2} \text{V_ABS} \quad \frac{\Gamma \vdash e_1 : \tau_1 \rightarrow \tau_2 \quad \Gamma \vdash e_2 : \tau_1}{\Gamma \vdash e_1 e_2 : \tau_2} \text{V_APP}$	
$\frac{}{\Gamma \vdash n : \text{Int}} \text{V_INT} \quad \frac{\Gamma \vdash^* e : \forall \bar{a}. \tau \quad \text{no other rule matches}}{\Gamma \vdash e : \tau[\bar{\tau}/\bar{a}]} \text{V_INSTS}$	
$\Gamma \vdash^* e : v$	Specified polytype checking for System V
$\frac{x:\forall\{\bar{a}\}. v \in \Gamma}{\Gamma \vdash^* x : v[\bar{\tau}/\bar{a}]} \text{V_VAR} \quad \frac{\Gamma \vdash^{gen} e_1 : \sigma_1 \quad \Gamma, x:\sigma_1 \vdash^* e_2 : v_2}{\Gamma \vdash^* \text{let } x = e_1 \text{ in } e_2 : v_2} \text{V_LET}$	
$\frac{\Gamma \vdash \tau \quad \Gamma \vdash^* e : \forall a. v}{\Gamma \vdash^* e @\tau : v[\tau/a]} \text{V_TAPP} \quad \frac{\Gamma \vdash v \quad v = \forall \bar{a}, \bar{b}. \tau \quad \Gamma, \bar{a} \vdash e : \tau \quad \bar{a}, \bar{b} \notin \text{ftv}(\Gamma)}{\Gamma \vdash^* (\Lambda \bar{a}. e : v) : v} \text{V_ANNOT}$	
$\frac{\Gamma \vdash e : \tau \quad \text{no other rule matches}}{\Gamma \vdash^* e : \tau} \text{V_MONO}$	
$\Gamma \vdash^{gen} e : \sigma$	Generalization for System V
$\frac{\bar{a} = \text{ftv}(v) \setminus \text{ftv}(\Gamma) \quad \Gamma \vdash^* e : v}{\Gamma \vdash^{gen} e : \forall\{\bar{a}\}. v} \text{V_GEN}$	

Fig. 7. Typing rules for System V

sion to have a type scheme only as a result of generalization (using the judgment $\Gamma \vdash^{gen} e : \sigma$). Generalization is, once again, available only in **let**-expressions.

However, the main difference that enables visible type annotation is the separation of the main typing judgment into two: $\Gamma \vdash e : \tau$ and $\Gamma \vdash^* e : v$. The key idea is that, sometimes, we need to be lazy about instantiating specified type variables so that the programmer has a chance to add a visible instantiation. Therefore, the system splits the rules into a judgment \vdash that requires e to have a monotype, and those in \vdash^* that can retain specified quantification.

The first set of rules in Fig. 7, as in System C, infers a monotype for the expression. The premises of the rule V_ABS uses this same judgment, for example, to require that the body of an abstraction have a monotype. All expressions can be assigned a monotype; if the first three rules do not apply, the last rule V_INSTS infers a polytype instead, then instantiates it to yield a monotype. Because implicit instantiation happens all at once in this rule, we do not need to worry about instantiating specified variables out of order, as we did in System HMV.

The second set of rules (the \vdash_v^* judgment) allows e to be assigned a specified polytype. Note that the premise of rule V_TAPP uses this judgment.

Rule V_VAR is like rule C_VAR : both look up a variable in the environment and instantiate its generalized quantified variables. The difference is that C_VAR 's types can contain *only* generalized variables; System V's types can have specified variables after the generalized ones. Yet we instantiate only the generalized ones in the V_VAR rule, lazily preserving the specified ones.

Rule V_LET is likewise similar to C_LET . The only difference is that the result type is not restricted to be a monotype. By putting V_LET in the \vdash_v^* judgment and returning a specified polytype, we allow the following judgment to hold:

$$\cdot \vdash_v (\mathbf{let} \ x = (\lambda y. y : \forall a. a \rightarrow a) \ \mathbf{in} \ x) @ \mathit{Int} : \mathit{Int} \rightarrow \mathit{Int}$$

The expression above would be ill-typed in a system that restricted the result of a \mathbf{let} -expression to be a monotype. It is for this reason that we altered System HM to include a polytype in its HM_LET rule, for consistency with HMV.

Rule V_ANNOT is identical to rule HMV_ANNOT . It uses the \vdash_v judgment in its premise to force instantiation of all quantified type variables before regeneralizing to the specified polytype v . In this way, the V_ANNOT rule is effectively able to reorder specified variables. Here, reordering is acceptable, precisely because it is user-directed.

Finally, if an expression form cannot yield a specified polytype, rule V_MONO delegates to \vdash_v to find a monotype for the expression.

5.3 Relating System V to System HMV

Systems HMV and V are equivalent; they type check the same set of expressions. We prove this correspondence using the following two theorems.

Theorem 5 (Soundness of V against HMV).

1. If $\Gamma \vdash_v e : \tau$, then $\Gamma \vdash_{\mathit{HMV}} e : \tau$.
2. If $\Gamma \vdash_v^* e : v$, then $\Gamma \vdash_{\mathit{HMV}} e : v$.
3. If $\Gamma \vdash_v^{gen} e : \sigma$, then $\Gamma \vdash_{\mathit{HMV}} e : \sigma$.

Theorem 6 (Completeness of V against HMV). *If $\Gamma \vdash_{\mathit{HMV}} e : \sigma$, then there exists σ' such that $\Gamma \vdash_v^{gen} e : \sigma'$ where $\sigma' \leq_{\mathit{HMV}} \sigma$.*

The proofs of these theorems appear in the extended version [12].

Having established the equivalence of System V with System HMV, we can note that Lemma 2 (“Extra knowledge is harmless”) carries over from HMV to V. This property is quite interesting in the context of System V. It says that a typing context where all type variables are specified admits all the same expressions as one where some type variables are generalized. In System V, however, specified and generalized variables are instantiated via different mechanisms, so this is a powerful theorem indeed.

It is mechanical to go from the statement of System V in Fig. 7 to an algorithm. In the extended version [12], we define Algorithm \mathcal{V} which implements

System V, analogous to Algorithm \mathcal{W} which implements System C. We then prove that Algorithm \mathcal{V} is sound and complete with respect to System V and that Algorithm \mathcal{V} finds principal types. Linking the pieces together gives us the proof of the principal types property for System HMV (Theorem 4). Furthermore, Algorithm \mathcal{V} is guaranteed to terminate, yielding this theorem:

Theorem 7. *Type-checking System V is decidable.*

6 Higher-rank type systems

We now extend the design of System HMV to include *higher-rank polymorphism* [17]. This extension allows function parameters to be used at multiple types. Incorporating this extension is actually quite straightforward. We include this extension to show that our framework for visible type application is indeed easy to extend – the syntax-directed system we study in this section is essentially a merge of System V and the bidirectional system from our previous work [23]. This system is also the basis for our implementation in GHC.

As an example, the following function does not type check in the vanilla Hindley-Milner type system, assuming id has type $\forall a. a \rightarrow a$.

```
let foo =  $\lambda f \rightarrow (f\ 3, f\ True)$  in foo id
```

Yet, with the `RankNTypes` language extension and the following type annotation, GHC is happy to accept

```
let foo :: ( $\forall a. a \rightarrow a$ )  $\rightarrow (Int, Bool)$ 
    foo =  $\lambda f \rightarrow (f\ 3, f\ True)$ 
in foo id
```

Visible type application means that higher-rank arguments can also be explicitly instantiated. For example, we can instantiate lambda-bound identifiers:

```
let foo :: ( $\forall a. a \rightarrow a$ )  $\rightarrow (Int \rightarrow Int, Bool)$ 
    foo =  $\lambda f \rightarrow (f\ @Int, f\ True)$ 
in foo id
```

Higher-rank types also mean that visible instantiations can occur after other arguments are passed to a function. For example, consider this alternative type for the `pair` function:

```
pair ::  $\forall a. a \rightarrow \forall b. b \rightarrow (a, b)$ 
pair =  $\lambda x\ y \rightarrow (x, y)$ 
```

If `pair` has this type, we can instantiate b after providing the first component for the pair, thus:

```
bar = pair 'x' @Bool
-- bar inferred to have type Bool  $\rightarrow (Char, Bool)$ 
```

In the rest of this section, we provide the technical details of these language features and discuss their interactions. In contrast to the presentation above, we present the syntax-directed higher-rank system first. We do so for two reasons: understanding a bidirectional system requires thinking about syntax, and thus the syntax-directed system seems easier to understand; and we view the declarative system as an expression of properties – or a set a metatheorems – about the higher-rank type system.

6.1 System SB: Syntax-directed Bidirectional Type Checking

Figures 8 and 9 show System SB, the higher-rank, bidirectional analogue of System V, supporting predicative higher-rank polymorphism and visible type application.

This system shares the same expression language of Systems HMV and V, retaining visible type application and type annotations. However, types in System SB may have non-prenex quantification. The body of a specified polytype v is now a *phi-type* ϕ : a type that has no top-level quantification but may have quantification to the left or to the right of arrows. Note also that these inner quantified types are vs , not σs . In other words, non-prenex quantification is over only *specified* variables, never generalized ones. As we will see, inner quantified types are introduced only by user annotation, and thus there is no way the system could produce an inner type scheme, even if the syntactic restriction were not in place.

The grammar also defines *rho-types* ρ , which also have no top-level quantification, but do allow inner quantification to the *left* of arrows. We convert specified polytypes (which may quantify to the right of arrows) to corresponding rho-types by means of the *prenex* operation, which appears in Fig. 9.

System SB is defined by five mutually recursive judgments: $\Gamma \vdash_{\text{sb}} e \Rightarrow \phi$, $\Gamma \vdash_{\text{sb}}^* e \Rightarrow v$, and $\Gamma \vdash_{\text{sb}}^{\text{gen}} e \Rightarrow \sigma$ are synthesis judgments, producing the type as an output; $\Gamma \vdash_{\text{sb}} e \Leftarrow \rho$ and $\Gamma \vdash_{\text{sb}}^* e \Leftarrow v$ are checking judgments, requiring the type as an input.

Type synthesis The synthesis judgments are very similar to the judgments from System V, ignoring direction arrows. The differences stem from the non-prenex quantification allowed in SB. The level of similarity is unsurprising, as the previous systems essentially all work only in synthesis mode; they derive a type given an expression. The novelty of a bidirectional system is its ability to propagate information about specified polytypes toward the leaves of an expression.

Type checking Rule SB_DABS is what makes the system higher-rank. The checking judgment $\Gamma \vdash_{\text{sb}} e \Leftarrow \rho$ pushes in a rho-type, with no top-level quantification. Thus, SB_DABS can recognize an arrow type $v_1 \rightarrow \rho_2$. Propagating this type into an expression $\lambda x. e$, SB_DABS uses the type v_1 as x 's type when

The grammar for SB extends that for System HMV (Fig. 3):

$e ::= \dots$	expressions	$\rho ::= \tau \mid v_1 \rightarrow \rho_2$	rho-types
$\tau ::= \dots$	monotypes	$\phi ::= \tau \mid v_1 \rightarrow v_2$	phi-types
$\Gamma ::= \cdot \mid \Gamma, x:\sigma \mid \Gamma, a$	contexts	$v ::= \phi \mid \forall a. v$	specified polytypes
		$\sigma ::= v \mid \forall\{a\}. \sigma$	type schemes

$\boxed{\Gamma \vdash_{\text{sb}} e \Rightarrow \phi}$ Synthesis of types without top-level quantifiers

$$\frac{\Gamma, x:\tau \vdash_{\text{sb}}^* e \Rightarrow v}{\Gamma \vdash_{\text{sb}} \lambda x. e \Rightarrow \tau \rightarrow v} \text{SB_ABS} \quad \frac{\Gamma \vdash_{\text{sb}}^* e \Rightarrow \forall \bar{a}. \phi \quad \text{no other rule matches}}{\Gamma \vdash_{\text{sb}} e \Rightarrow \phi[\bar{\tau}/\bar{a}]} \text{SB_INSTS}$$

$$\frac{}{\Gamma \vdash_{\text{sb}} n \Rightarrow \text{Int}} \text{SB_INT}$$

$\boxed{\Gamma \vdash_{\text{sb}}^* e \Rightarrow v}$ Synthesis of specified polytypes

$$\frac{x:\forall\{\bar{a}\}. v \in \Gamma}{\Gamma \vdash_{\text{sb}}^* x \Rightarrow v[\bar{\tau}/\bar{a}]} \text{SB_VAR} \quad \frac{\Gamma \vdash_{\text{sb}} e_1 \Rightarrow v_1 \rightarrow v_2 \quad \Gamma \vdash_{\text{sb}}^* e_2 \Leftarrow v_1}{\Gamma \vdash_{\text{sb}}^* e_1 e_2 \Rightarrow v_2} \text{SB_APP}$$

$$\frac{\Gamma \vdash \tau}{\Gamma \vdash_{\text{sb}}^* e \Rightarrow \forall a. v} \text{SB_TAPP} \quad \frac{\Gamma \vdash v \quad v = \forall \bar{a}, \bar{b}. \phi \quad \Gamma, \bar{a} \vdash_{\text{sb}}^* e \Leftarrow \phi \quad \bar{a}, \bar{b} \notin \text{ftv}(\Gamma)}{\Gamma \vdash_{\text{sb}}^* (\Lambda \bar{a}. e : v) \Rightarrow v} \text{SB_ANNOT}$$

$$\frac{\Gamma \vdash_{\text{sb}}^{gen} e_1 \Rightarrow \sigma_1 \quad \Gamma, x:\sigma_1 \vdash_{\text{sb}}^* e_2 \Rightarrow v_2}{\Gamma \vdash_{\text{sb}}^* \text{let } x = e_1 \text{ in } e_2 \Rightarrow v_2} \text{SB_LET} \quad \frac{\Gamma \vdash_{\text{sb}} e \Rightarrow \phi \quad \text{no other rule matches}}{\Gamma \vdash_{\text{sb}}^* e \Rightarrow \phi} \text{SB_PHI}$$

$\boxed{\Gamma \vdash_{\text{sb}}^{gen} e \Rightarrow \sigma}$ Synthesis with generalization

$$\frac{\Gamma \vdash_{\text{sb}}^* e \Rightarrow v \quad \bar{a} = \text{ftv}(v) \setminus \text{ftv}(\Gamma)}{\Gamma \vdash_{\text{sb}}^{gen} e \Rightarrow \forall\{\bar{a}\}. v} \text{SB_GEN}$$

$\boxed{\Gamma \vdash_{\text{sb}} e \Leftarrow \rho}$ Checking against types without top-level quantifiers

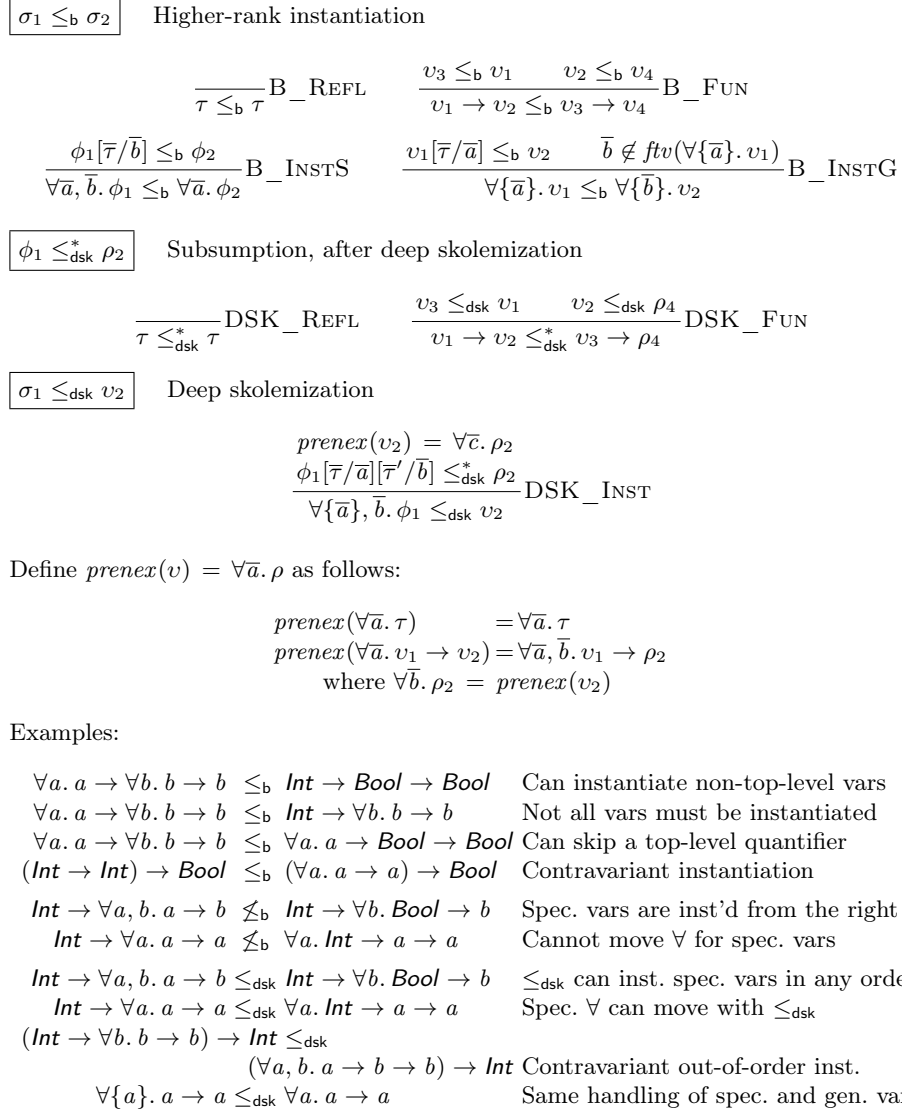
$$\frac{\Gamma \vdash_{\text{sb}}^{gen} e_1 \Rightarrow \sigma_1 \quad \Gamma, x:\sigma_1 \vdash_{\text{sb}} e_2 \Leftarrow \rho_2}{\Gamma \vdash_{\text{sb}} \text{let } x = e_1 \text{ in } e_2 \Leftarrow \rho_2} \text{SB_DLET}$$

$$\frac{\Gamma, x:v_1 \vdash_{\text{sb}}^* e \Leftarrow \rho_2}{\Gamma \vdash_{\text{sb}} \lambda x. e \Leftarrow v_1 \rightarrow \rho_2} \text{SB_DABS} \quad \frac{\Gamma \vdash_{\text{sb}}^* e \Rightarrow v_1 \quad v_1 \leq_{\text{dsk}} \rho_2 \quad \text{no other rule matches}}{\Gamma \vdash_{\text{sb}} e \Leftarrow \rho_2} \text{SB_INFER}$$

$\boxed{\Gamma \vdash_{\text{sb}}^* e \Leftarrow v}$ Checking against specified polytypes

$$\frac{\text{prenex}(v) = \forall \bar{a}. \rho \quad \bar{a} \notin \text{ftv}(\Gamma) \quad \Gamma \vdash_{\text{sb}} e \Leftarrow \rho}{\Gamma \vdash_{\text{sb}}^* e \Leftarrow v} \text{SB_DEEPSKOL}$$

Fig. 8. Syntax-directed bidirectional type system

**Fig. 9.** Higher-rank subsumption relations

checking e . This is the only place in system SB where a lambda-term can abstract over a variable with a polymorphic type. Note that the synthesis rule SB_ABS uses a monotype for the type of x .¹²

¹² Higher-rank systems can also include an “annotated abstraction” form, $\lambda x:v. e$. This form allows higher-rank types to be synthesized for lambda expressions as well as checked. However, this form is straightforward to add but is not part of GHC, which

Rule `SB_INFER` mediates between the checking and synthesis judgments. When no checking rule applies, we synthesize a type and then check it according to the \leq_{dsk} deep skolemization relation, taken directly from previous work and shown in Fig. 9. For brevity, we do not explain the details of this relation here, instead referring readers to Peyton Jones et al. [23, Sect. 4.6] for much deeper discussion. However, we note that there is a design choice to be made here; we could have also used Odersky–Läufer’s slightly less expressive higher-rank subsumption relation [21] instead. We present the system with deep skolemization for backwards compatibility with GHC. See the extended version [12] for a discussion of this alternative.

The entry point into the type checking judgments is through the $\Gamma \vdash_{\text{sb}}^* e \Leftarrow v$ judgment. This judgment has just one rule, `SB_DEEPSKOL`. The rule skolemizes all type variables appearing at the top-level and to the right of arrows. Skolemizing here is necessary to expose a rho-type to the $\Gamma \vdash_{\text{sb}} e \Leftarrow \rho$ judgment, so that rule `SB_DABS` can fire.¹³ The reason this rule requires deep skolemization instead of top-level skolemization is subtle, but this choice is not due to visible type application or lazy instantiation; the same choice is made in prior work [23, rule `GEN2` of Fig. 8]. We refer readers to the extended version [12] for the details.

6.2 System B: Declarative specification

Figure 10 shows the typing rules of System B, a declarative system that accepts the same programs as System SB. This declarative type system itself is a novel contribution of this work. (The systems presented in related work [10, 21, 23] are more similar to SB than to B.)

Although System B is *bidirectional*, it is also *declarative*. In particular, the use of generalization (`B_GEN`), subsumption (`B_SUB`), skolemization (`B_SKOL`), and mode switching (`B_INFER`), can happen arbitrarily in a typing derivation. Understanding what expressions are well-typed does not require knowing precisely when these operations take place.

The subsumption rule (`B_SUB`) in the synthesis judgment corresponds to `HMV_SUB` from `HMV`. However, the novel subsumption relation \leq_{b} used by this rule, shown at the top of Fig. 9, is *different* from the \leq_{dsk} deep skolemization relation used in System SB. This $\sigma_1 \leq_{\text{b}} \sigma_2$ judgment extends the action of \leq_{hmv} to higher-rank types: in particular, it allows subsumption for generalized type variables (which can be quantified only at the top level) and instantiation (only) for specified type variables. We could say that this judgment enables *inner instantiation* because instantiations are not restricted to top level. See also the examples at the bottom of Fig. 9.

In contrast, rule `B_INFER` (in the checking judgment) uses the stronger of the two subsumption relations \leq_{dsk} . This rule appears at precisely the spot in

uses patterns (beyond the scope of this paper) to bind variables in abstractions. Therefore we omit the annotated abstraction form from our formalism.

¹³ Our choice to skolemize before `SB_DLET` is arbitrary, as `SB_DLET` does not interact with the propagated type.

$\Gamma \Vdash e \Rightarrow \sigma$	Synthesis rules for System B
$\frac{x:\sigma \in \Gamma}{\Gamma \Vdash x \Rightarrow \sigma} \text{B_VAR} \quad \frac{\Gamma, x:\tau \Vdash e \Rightarrow v}{\Gamma \Vdash \lambda x. e \Rightarrow \tau \rightarrow v} \text{B_ABS}$	
$\frac{\Gamma \Vdash e_1 \Rightarrow v_1 \rightarrow v_2 \quad \Gamma \Vdash e_2 \Leftarrow v_1}{\Gamma \Vdash e_1 e_2 \Rightarrow v_2} \text{B_APP} \quad \frac{}{\Gamma \Vdash n \Rightarrow \text{Int}} \text{B_INT}$	
$\frac{\Gamma \Vdash e_1 \Rightarrow \sigma_1 \quad \Gamma, x:\sigma_1 \Vdash e_2 \Rightarrow \sigma}{\Gamma \Vdash \text{let } x = e_1 \text{ in } e_2 \Rightarrow \sigma} \text{B_LET}$	
$\frac{\Gamma \Vdash e \Rightarrow \sigma \quad a \notin \text{ftv}(\Gamma)}{\Gamma \Vdash e \Rightarrow \forall \{a\}. \sigma} \text{B_GEN} \quad \frac{\Gamma \Vdash e \Rightarrow \sigma_1 \quad \sigma_1 \leq_b \sigma_2}{\Gamma \Vdash e \Rightarrow \sigma_2} \text{B_SUB}$	
$\frac{\Gamma \vdash \tau}{\Gamma \Vdash e \Rightarrow \forall a. v} \text{B_TAPP} \quad \frac{\Gamma \vdash v \quad v = \forall \bar{a}, \bar{b}. \phi \quad \Gamma, \bar{a} \Vdash e \Leftarrow \phi \quad \bar{a}, \bar{b} \notin \text{ftv}(\Gamma)}{\Gamma \Vdash (\Lambda \bar{a}. e : v) \Rightarrow v} \text{B_ANNOT}$	
$\Gamma \Vdash e \Leftarrow v$	Checking rules for System B
$\frac{\Gamma, x:v_1 \Vdash e \Leftarrow v_2}{\Gamma \Vdash \lambda x. e \Leftarrow v_1 \rightarrow v_2} \text{B_DABS} \quad \frac{\Gamma \Vdash e_1 \Rightarrow \sigma_1 \quad \Gamma, x:\sigma_1 \Vdash e_2 \Leftarrow v}{\Gamma \Vdash \text{let } x = e_1 \text{ in } e_2 \Leftarrow v} \text{B_DLET}$	
$\frac{\Gamma \Vdash e \Leftarrow v \quad a \notin \text{ftv}(\Gamma)}{\Gamma \Vdash e \Leftarrow \forall a. v} \text{B_SKOL} \quad \frac{\Gamma \Vdash e \Rightarrow \sigma_1 \quad \sigma_1 \leq_{\text{dsk}} v_2}{\Gamma \Vdash e \Leftarrow v_2} \text{B_INFER}$	

Fig. 10. System B

the derivation where a specified type from synthesis mode meets the specified type from checking mode. The relation \leq_{dsk} subsumes \leq_b ; that is, $\sigma_1 \leq_b v_2$ implies $\sigma_1 \leq_{\text{dsk}} v_2$.

Properties of System B and SB We can show that Systems SB and B admit the same expressions.

Lemma 8 (Soundness of System SB).

1. If $\Gamma \Vdash_{\text{sb}} e \Rightarrow \phi$ then $\Gamma \Vdash e \Rightarrow \phi$.
2. If $\Gamma \Vdash_{\text{sb}}^* e \Rightarrow v$ then $\Gamma \Vdash e \Rightarrow v$.
3. If $\Gamma \Vdash_{\text{sb}}^{\text{gen}} e \Rightarrow \sigma$ then $\Gamma \Vdash e \Rightarrow \sigma$.
4. If $\Gamma \Vdash_{\text{sb}}^* e \Leftarrow v$ then $\Gamma \Vdash e \Leftarrow v$.
5. If $\Gamma \Vdash_{\text{sb}} e \Leftarrow \rho$ then $\Gamma \Vdash e \Leftarrow \rho$.

Lemma 9 (Completeness of System SB).

1. If $\Gamma \Vdash e \Rightarrow \sigma$ then $\Gamma \Vdash_{\text{sb}}^{\text{gen}} e \Rightarrow \sigma'$ where $\sigma' \leq_b \sigma$.
2. If $\Gamma \Vdash e \Leftarrow v$ then $\Gamma \Vdash_{\text{sb}}^* e \Leftarrow v$.

What is the role of System B? In our experience, programmers tend to prefer the syntax-directed presentation of the system because that version is more algorithmic. As a result, it can be easier to understand why a program type checks (or doesn't) by reasoning about System SB.

However, the fact that System B is sound and complete with respect to System SB provides properties that we can use to reason about SB. The main difference between the two is that System B divides subsumption into two different relations. The weaker \leq_b can be used at any time during synthesis, but it can only instantiate (but not generalize) specified variables. The stronger \leq_{dsk} is used at only the check/synthesis boundary but can generalize and reorder specified variables.

The connection between the two systems tells us that B_SUB is *admissible* for SB. As a result, when refactoring code, we need not worry about precisely where a type is instantiated, as we see here that instantiation need not be determined syntactically.

Likewise, the proof also shows that System B (and System SB) is flexible with respect to the instantiation relation \leq_b in the context. As in System HMV, this result implies that making generalized variables into specified variables does not disrupt types.

Lemma 10 (Context Generalization). *Suppose $\Gamma' \leq_b \Gamma$.*

1. *If $\Gamma \Vdash_b e \Rightarrow \sigma$ then there exists $\sigma' \leq_b \sigma$ such that $\Gamma' \Vdash_b e \Rightarrow \sigma'$.*
2. *If $\Gamma \Vdash_b e \Leftarrow v$ and $v \leq_b v'$ then $\Gamma' \Vdash_b e \Leftarrow v'$.*

Proofs of these properties appear in the extended version [12].

6.3 Integrating visible type application with GHC

System SB is the direct inspiration for the type-checking algorithm used in our version of GHC enhanced with visible type application. It is remarkably straightforward to implement the system described here within GHC; accounting for the behavior around imported functions (Sect. 3.1) was the hardest part. The other interactions (the difference between this paper's scoped type variables and GHC's, how specified type variables work with type classes, etc.) are generally uninteresting; see the extended version [12] for further comments.

One pleasing synergy between visible type application and GHC concerns GHC's recent *partial type signature* feature [29]. This feature allows wildcards, written with an underscore, to appear in types; GHC infers the correct replacement for the wildcard. These work well in visible type applications, allowing the user to write $@_$ as a visible type argument where GHC can infer the argument. For example, if f has type $\forall a b. a \rightarrow b \rightarrow (a, b)$, then we can write $f @_ @[Int] True []$ to let GHC infer that a should be *Bool* but to visibly instantiate b to be *[Int]*. Getting partial type signatures to work in the new context of visible type applications required nothing more than hooking up the pieces.

7 Related work and Conclusions

Explicit type arguments The F# language [26] permits explicit type arguments in function applications and method invocations. These explicit arguments, typically mandatory, are used to resolve ambiguity in type-dependent operations. However, the properties of this feature have not been studied.

Implicit arguments in dependently-typed languages Languages such as Coq [6], Agda [20], Idris [1] and Twelf [24] are not based on the HM type system, so their designs differ from Systems HMV and B. However, they do support invisible arguments. In these languages, an invisible argument is not necessarily a type; it could be any argument that can be inferred by the type checker.

Coq, Agda, and Idris require all quantification, including that for invisible arguments, to be specified by the user. These languages do not support generalization, i.e., automatically determining that an expression should quantify over an invisible argument (in addition to any visible ones). They differ in how they specify the visibility of arguments, yet all of them provide the ability to override an invisibility specification and provide such arguments visibly. These languages also provide a facility for *named* invisible arguments, allowing users to specify argument values by name instead of by position. This choice means that α -equivalent types are no longer fully interchangeable. Though we have not studied this possibility deeply, we conjecture that formally specifying a named-argument system would encounter many of the same subtleties (in particular, requiring two different subsumption relations in the metatheory) that we encountered with positional arguments.

Twelf, on the other hand, supports invisible arguments via generalization and visible arguments via specification. Although it is easy to convert between the two versions, there is no way to visibly provide an invisible argument as we have done. Instead, the user must rely on type annotations to control instantiations.

Specified vs. generalized variables Dreyer and Blume’s work on specifying ML’s type system and inference algorithm in the presence of modules [9] introduces a separation of (what we call) specified and generalized variables. Their work focuses on the type parameters to ML functors, finding inconsistencies between the ML language specification and implementations. They conclude that the ML specification as written is hard to implement and propose a new one. Their work includes a type system that allows functors to have invisible arguments alongside their visible ones. This specification is easier to implement, as they demonstrate.

Their work has similarities to ours in the separation of classes of variables and the need to alter the specification to make type inference reasonable. Interestingly, they come from the opposite direction from ours, adding invisible arguments in a place where arguments previously were all visible. However, despite these surface similarities, we have not found a deeper connection between our work and theirs.

Predicative, higher-rank type systems As we have already indicated, Systems B and SB are directly inspired by GHC’s design for higher-rank types [23]. However, in this work we have redesigned the algorithm to use lazy instantiation and have made a distinction between specified polytypes and generalized polytypes. Furthermore, we have pushed the design further, providing a declarative specification for the type system.

Our work is also closely related to recent work on using a bidirectional type system for higher-rank polymorphism by Dunfield and Krishnaswami [10], called DK below. The closest relationship is between their declarative system (Fig. 4 in their paper) and our System SB (Fig. 8). The most significant difference is that the DK system never generalizes. All polymorphic types in their system are specified; functions must have a type annotation to be polymorphic. Consequently, DK uses a different algorithm for type checking than the one proposed in this work. Nevertheless, it defers instantiations of specified polymorphism, like our algorithm.

Our relation \leq_{dsk} , which requires two specified polytypes, is similar to the DK subsumption relation. The DK version is slightly weaker as it does not use deep skolemization; but that difference is not important in this context. Another minor difference is that System SB uses the $\Gamma \vdash_{\text{sb}} e \Rightarrow \phi$ judgment to syntactically guide instantiation whereas the the DK system uses a separate application judgment form. System B – and the metatheory of System SB – also includes implicit subsumption \leq_{b} , which does not have an analogue in the DK system. A more extended comparison with the DK system appears in the extended version [12].

Conclusion This work extends the HM type system with visible type application, while maintaining important properties of that system that make it useful for functional programmers. Our extension is fully backwards compatible with previous versions of GHC. It retains the principal types property, leading to robustness during refactoring. At the same time, our new systems come with simple, compositional specifications.

While we have incorporated visible type application with all existing features of GHC, we do not plan to stop there. We hope that our mix of specified polytypes and type schemes will become a basis for additional type system extensions, such as impredicative types, type-level lambdas, and dependent types.

Acknowledgments Thanks to Simon Peyton Jones, Dimitrios Vytiniotis, Iavor Diatchki, Adam Gundry, Conor McBride, Neel Krishnaswami, and Didier Rémy for helpful discussion and feedback.

Bibliography

- [1] Edwin Brady. Idris, a general-purpose dependently typed programming language: Design and implementation. *J. Funct. Prog.*, 23, 2013.
- [2] Pablo Buiras, Dimitrios Vytiniotis, and Alejandro Russo. HLIO: Mixing static and dynamic typing for information-flow control in Haskell. In *International Conference on Functional Programming*, ICFP '15. ACM, 2015.
- [3] Manuel M. T. Chakravarty, Gabriele Keller, and Simon Peyton Jones. Associated type synonyms. In *International Conference on Functional Programming*, ICFP '05. ACM, 2005.
- [4] Manuel M. T. Chakravarty, Gabriele Keller, Simon Peyton Jones, and Simon Marlow. Associated types with class. In *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, 2005.
- [5] Dominique Clément, Thierry Despeyroux, Gilles Kahn, and Joëlle Despeyroux. A simple applicative language: Mini-ML. In *Conference on LISP and Functional Programming*, LFP '86. ACM, 1986.
- [6] Coq development team. *The Coq proof assistant reference manual*. LogiCal Project, 2004. URL <http://coq.inria.fr>. Version 8.0.
- [7] Luis Damas. *Type Assignment in Programming Languages*. PhD thesis, University of Edinburgh, 1985.
- [8] Luis Damas and Robin Milner. Principal type-schemes for functional programs. In *Symposium on Principles of Programming Languages*, POPL '82. ACM, 1982.
- [9] Derek Dreyer and Matthias Blume. Principal type schemes for modular programs. In *Proceedings of the 16th European Conference on Programming*, ESOP'07. Springer-Verlag, 2007.
- [10] Joshua Dunfield and Neelakantan R. Krishnaswami. Complete and easy bidirectional typechecking for higher-rank polymorphism. In *International Conference on Functional Programming*, ICFP '13. ACM, 2013.
- [11] Richard A. Eisenberg, Dimitrios Vytiniotis, Simon Peyton Jones, and Stephanie Weirich. Closed type families with overlapping equations. In *Principles of Programming Languages*, POPL '14. ACM, 2014.
- [12] Richard A. Eisenberg, Stephanie Weirich, and Hamidhasan Ahmed. Visible type application (extended version), 2015. URL <http://www.seas.upenn.edu/~sweirich/papers/type-app-extended.pdf>.
- [13] J. Roger Hindley. The principal type-scheme of an object in combinatory logic. *Transactions of the American Mathematical Society*, 146, 1969.
- [14] Didier Le Botlan and Didier Rémy. ML^F : Raising ML to the power of System F. In *International Conference on Functional Programming*. ACM, 2003.
- [15] Simon Marlow (editor). Haskell 2010 language report, 2010.
- [16] Conor McBride. Agda-curious? Keynote, ICFP'12, 2012.
- [17] Nancy McCracken. The typechecking of programs with implicit type structure. In Gilles Kahn, David B. MacQueen, and Gordon Plotkin, editors,

- Semantics of Data Types*, volume 173 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 1984.
- [18] Robin Milner. A Theory of Type Polymorphism in Programming. *Journal of Computer and System Sciences*, 17, 1978.
 - [19] Alexandre Miquel. The implicit calculus of constructions: Extending pure type systems with an intersection type binder and subtyping. In Samson Abramsky, editor, *Typed Lambda Calculi and Applications*, volume 2044 of *Lecture Notes in Computer Science*, pages 344–359. Springer Berlin Heidelberg, 2001.
 - [20] Ulf Norell. *Towards a practical programming language based on dependent type theory*. PhD thesis, Department of Computer Science and Engineering, Chalmers University of Technology, SE-412 96 Göteborg, Sweden, September 2007.
 - [21] Martin Odersky and Konstantin Läufer. Putting type annotations to work. In *Symposium on Principles of Programming Languages*, POPL '96. ACM, 1996.
 - [22] Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Geoffrey Washburn. Simple unification-based type inference for GADTs. In *International Conference on Functional Programming*, ICFP '06. ACM, 2006.
 - [23] Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Mark Shields. Practical type inference for arbitrary-rank types. *Journal of Functional Programming*, 17(1), January 2007.
 - [24] Frank Pfenning and Carsten Schürmann. System description: Twelf - a meta-logical framework for deductive systems. In H. Ganzinger, editor, *Proceedings of the 16th International Conference on Automated Deduction (CADE-16)*, number 1632 in LNAI, pages 202–206, Trento, Italy, July 1999. Springer-Verlag.
 - [25] François Pottier and Didier Rémy. *Advanced Topics in Types and Programming Languages*, chapter The Essence of ML Type Inference, pages 387–489. The MIT Press, 2005.
 - [26] Don Syme. *The F# 2.0 Language Specification*. Microsoft Research and the Microsoft Developer Division, 2012. Accessed from <http://fsharp.org/specs/language-spec/2.0/FSharpSpec-2.0-April-2012.pdf>.
 - [27] Dimitrios Vytiniotis, Simon Peyton Jones, Tom Schrijvers, and Martin Sulzmann. OutsideIn(X) modular type inference with local assumptions. *Journal of Functional Programming*, 21(4-5), September 2011.
 - [28] Philip Wadler and Stephen Blott. How to make ad-hoc polymorphism less ad-hoc. In *POPL*, pages 60–76. ACM, 1989.
 - [29] Thomas Winant, Dominique Devriese, Frank Piessens, and Tom Schrijvers. Partial type signatures for haskell. In *Practical Aspects of Declarative Languages*, volume 8324, pages 17–32. Springer International Publishing, January 2014.