



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

FreezeML

Citation for published version:

Emrich, F, Lindley, S, Stolarek, J, Cheney, J & Coates, J 2020, FreezeML: Complete and Easy Type Inference for First-Class Polymorphism. in *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation*. Association for Computing Machinery (ACM), pp. 423-437, 41st ACM SIGPLAN Conference on Programming Language Design and Implementation, London, United Kingdom, 15/06/20. <https://doi.org/10.1145/3385412.3386003>

Digital Object Identifier (DOI):

[10.1145/3385412.3386003](https://doi.org/10.1145/3385412.3386003)

Link:

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

Document Version:

Peer reviewed version

Published In:

Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation

General rights

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

Take down policy

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



FreezeML

Complete and Easy Type Inference for First-Class Polymorphism

Frank Emrich

The University of Edinburgh
frank.emrich@ed.ac.uk

Sam Lindley

The University of Edinburgh
Imperial College London
sam.lindley@ed.ac.uk

Jan Stolarek

The University of Edinburgh
Lodz University of Technology
jan.stolarek@ed.ac.uk

James Cheney

The University of Edinburgh
The Alan Turing Institute
jcheney@inf.ed.ac.uk

Jonathan Coates

The University of Edinburgh
s1627856@sms.ed.ac.uk

Abstract

ML is remarkable in providing statically typed polymorphism without the programmer ever having to write any type annotations. The cost of this parsimony is that the programmer is limited to a form of polymorphism in which quantifiers can occur only at the outermost level of a type and type variables can be instantiated only with monomorphic types.

Type inference for unrestricted System F-style polymorphism is undecidable in general. Nevertheless, the literature abounds with a range of proposals to bridge the gap between ML and System F.

We put forth a new proposal, FreezeML, a conservative extension of ML with two new features. First, let- and lambda-binders may be annotated with arbitrary System F types. Second, variable occurrences may be *frozen*, explicitly disabling instantiation. FreezeML is equipped with type-preserving translations back and forth between System F and admits a type inference algorithm, an extension of algorithm W, that is sound and complete and which yields principal types.

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

Keywords: first-class polymorphism, type inference, impredicative types

ACM Reference Format:

Frank Emrich, Sam Lindley, Jan Stolarek, James Cheney, and Jonathan Coates. 2020. FreezeML: Complete and Easy Type Inference for First-Class Polymorphism. In *Proceedings of the 41st ACM SIGPLAN*

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

PLDI '20, June 15–20, 2020, London, UK

© 2020 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-7613-6/20/06.

<https://doi.org/10.1145/3385412.3386003>

International Conference on Programming Language Design and Implementation (PLDI '20), June 15–20, 2020, London, UK. ACM, New York, NY, USA, 15 pages. <https://doi.org/10.1145/3385412.3386003>

1 Introduction

The design of ML [19] hits a sweet spot in providing statically typed polymorphism without the programmer ever having to write type annotations. The Hindley-Milner type inference algorithm on which ML relies is *sound* (it only yields correct types) and *complete* (if a program has a type then it will be inferred). Moreover, inferred types are *principal*, that is, most general. Alas, this sweet spot is rather narrow, depending on a delicate balance of features; it still appears to be an open question how best to extend ML type inference to support first-class polymorphism as found in System F.

Nevertheless, ML has unquestionable strengths as the basis for high-level programming languages. Its implicit polymorphism is extremely convenient for writing concise programs. Functional programming languages such as Haskell and OCaml employ algorithms based on Hindley-Milner type inference and go to great efforts to reduce the need to write type annotations on programs. Whereas the plain Hindley-Milner algorithm supports a limited form of polymorphism in which quantifiers must be top-level and may only be instantiated with monomorphic types, advanced programming techniques often rely on first-class polymorphism, where quantifiers may appear anywhere and may be instantiated with arbitrary polymorphic types, as in System F. However, working directly in System F is painful due to the need for *explicit* type abstraction and application. Alas, type inference, and indeed type checking, is undecidable for System F without type annotations [22, 29].

The primary difficulty in extending ML to support first-class polymorphism is with implicit instantiation of polymorphic types: whenever a variable occurrence is typechecked, any quantified type variables are immediately instantiated with (monomorphic) types. Whereas with plain ML there is

no harm in greedily instantiating type variables, with first-class polymorphism there is sometimes a non-trivial choice to be made over whether to instantiate or not.

The basic Hindley-Milner algorithm [3] restricts the use of polymorphism in types to *type schemes* of the form $\forall \vec{a}. A$ where A does not contain any further polymorphism. This means that, for example, given a function $\text{single} : \forall a. a \rightarrow \text{List } a$, that constructs a list of one element, and a polymorphic function choosing its first argument $\text{choose} : \forall a. a \rightarrow a \rightarrow a$, the expression single choose is assigned the type $\text{List } (a \rightarrow a \rightarrow a)$, for some fixed type a determined by the context in which the expression is used. The type $\text{List } (a \rightarrow a \rightarrow a)$ arises from instantiating the quantifier of single with $a \rightarrow a \rightarrow a$. But what if instead of constructing a list of choice functions at a fixed type, a programmer wishes to construct a list of polymorphic choice functions of type $\text{List } (\forall a. a \rightarrow a \rightarrow a)$? This requires instantiating the quantifier of single with a polymorphic type $\forall a. a \rightarrow a \rightarrow a$, which is forbidden in ML, and indeed the resulting System F type is not even an ML type scheme. However, in a richer language such as System F, the expression single choose could be annotated as appropriate in order to obtain either the type $\text{List } (a \rightarrow a \rightarrow a)$ or the type $\text{List } (\forall a. a \rightarrow a \rightarrow a)$.

All is not lost. By adding a sprinkling of explicit type annotations, in combination with other extensions, it is possible to retain much of the convenience of ML alongside the expressiveness of System F. Indeed, there is a plethora of techniques bridging the expressiveness gap between ML and System F without sacrificing desirable type inference properties of ML [7, 11–14, 24, 25, 27, 28].

However, there is still not widespread consensus on what constitutes a good design for a language combining ML-style type inference with System F-style first-class polymorphism, beyond the typical criteria of decidability, soundness, completeness, and principal typing. As Serrano et al. [25] put it in their PLDI 2018 paper, type inference in the presence of first-class polymorphism is still “a deep, deep swamp” and “no solution (...) with a good benefit-to-weight ratio has been presented to date”. While previous proposals offer considerable expressive power, we nevertheless consider the following combination of design goals to be both compelling and not yet achieved by any prior work:

- **Familiar System F types** Our ideal solution would use exactly the type language of System F. Systems such as HML [13], MLF [11], Poly-ML¹ [7], and QML [24], capture (or exceed) the power of System F, but employ a strict superset of System F’s type language. Whilst in some cases this difference is superficial, we consider that it does increase the burden on the programmer to understand and use these systems effectively, and may also contribute to increasing the syntactic overhead and decreasing the clarity of programs.

¹The name Poly-ML does not appear in the original [7] paper, but was introduced retrospectively [11].

- **Close to ML type inference** Our ideal solution would conservatively extend ML and standard Hindley-Milner type inference, including the (now-standard) *value restriction* [30], without being tied to one particular type inference algorithm. Systems such as MLF and Boxy Types have relied on much more sophisticated type inference techniques than needed in classical Hindley-Milner type inference, and proven difficult to implement or extend further because of their complexity. Other systems, such as GI, are relatively straightforward to implement atop an OutsideIn(X)-style constraint-based type inference algorithm, but would be much more work to add to a standard Hindley-Milner implementation.

- **Low syntactic overhead** Our ideal solution would provide first-class polymorphism without significant departures from ordinary ML-style programming. Early systems [9, 10, 20, 23] showed how to accommodate System F-style polymorphism by associating it with nominal datatype constructors, but this imposes a significant syntactic overhead to make use of these capabilities, which can also affect the readability and maintainability of programs. All previous systems necessarily involve some type annotations as well, which we also desire to minimise as much as possible.

- **Predictable behaviour** Our ideal solution would avoid guessing polymorphism and be specified so that programmers can anticipate where type annotations will be needed. More recent systems, such as HMF [12] and GI [25], use System F types, and are relatively easy to implement, but employ heuristics to guess one of several different polymorphic types, and require programmer annotations if the default heuristic behaviour is not what is needed.

In short, we consider that the problem of extending ML-style type inference with the power of System F is solved as a *technical* problem by several existing systems, but there remains a significant *design* challenge to develop a system that uses *familiar System F types*, is *close to ML type inference*, has *low syntactic overhead*, and has *predictable behaviour*. Of course, these desiderata represent our (considered, but subjective) views as language designers, and others may (and likely will) disagree. We welcome such debate.

Our contribution: FreezeML. In this paper, we introduce FreezeML, a core language extending ML with two System F-like features:

- “frozen” variable occurrences for which polymorphic instantiation is inhibited (written $[x]$ to distinguish them from ordinary variables x whose polymorphic types are implicitly instantiated); and
- type-annotated lambda abstractions $\lambda(x : A). M$.

FreezeML also refines the typing rule for let by:

- restricting let-bindings to have principal types; and
- allowing type annotations on let-bindings.

In FreezeML explicit type annotations are only required on lambda binders used in a polymorphic way, and on let-bindings that assign a non-principal type to a let-bound term; annotations are not required (or allowed) anywhere else. As we shall see in Section 2, the introduction of type-annotated let-bindings and frozen variables allows us to macro-express explicit versions of generalisation and instantiation (the two features that are implicit in plain ML). Thus, unlike ML, although FreezeML still has ML-like variables and let-binding it *also* enjoys explicit encodings of all of the underlying System F features. Correspondingly, frozen variables and type-annotated let-bindings are also central to encoding type abstraction and type application of System F (Section 4.1). Although, as we explain later, our approach is similar in expressiveness to existing proposals such as Poly-ML, we believe its close alignment with System F types and ML type inference are important benefits, and we argue via examples that its syntactic overhead and predictability compare favourably with the state of the art. Nevertheless, further work would need to be done to systematically compare the syntactic overhead and predictability of our approach with existing systems — this criticism, however, also applies to most previous work on new language design ideas.

A secondary technical contribution we make is to repair technical problem faced by FreezeML and some previous systems. In FreezeML, we restrict generalisation to principal types. However, directly incorporating this constraint into the type system results in rules that are syntactically not well-founded. We clarify that the typing relation can still be defined and inductive reasoning about it is still sound. This observation may also apply to other systems, such as HMF [13] and Poly-ML [7], where the same issue arises but was not previously addressed.

Contributions. This paper is a programming language design paper. Though we have an implementation on top of the Links programming language [2]² implementation is not the primary focus. The paper makes the following main contributions:

- A high-level introduction to FreezeML (Section 2).
- A type system for FreezeML as a conservative extension of ML with the expressive power of System F (Section 3).
- Local type-preserving translations back and forth between System F and FreezeML, and a discussion of the equational theory of FreezeML (Section 4).
- A type inference algorithm for FreezeML as an extension of algorithm W [3], which is sound, complete, and yields principal types (Section 5).

Section 6 discusses implementation, Section 7 presents related work and Section 8 concludes.

²<https://github.com/links-lang/links>

2 An Overview of FreezeML

We begin with an informal overview of FreezeML. Recall that the types of FreezeML are exactly those of System F.

Implicit Instantiation. In FreezeML (as in plain ML), when variable occurrences are typechecked, the outer universally quantified type variables in the variable’s type are instantiated implicitly. Suppose a programmer writes `choose id`, where `choose` : $\forall a.a \rightarrow a \rightarrow a$ and `id` : $\forall a.a \rightarrow a$. The quantifier in the type of `id` is implicitly instantiated with an as yet unknown type a , yielding the type $a \rightarrow a$. The type $a \rightarrow a$ is then used to instantiate the quantifier in the type of `choose`, yielding `choose id` : $(a \rightarrow a) \rightarrow (a \rightarrow a)$. The concrete type of a depends on the context in which the expression is used. For instance, if we were to apply `choose id` to an increment function then a would be unified with `Int`. (For the formal treatment of type inference in Section 5 we will be careful to explicitly distinguish between *rigid* type variables, like those bound by the quantifiers in the types of `choose` and `id`, and *flexible* type variables, like the a in the type inferred for the expression `choose id`.)

Explicit Freezing (`[x]`). The programmer may explicitly prevent a variable from having its already existing quantifiers instantiated by using the freeze operator `[−]`. Whereas each ordinary occurrence of `choose` has type $a \rightarrow a \rightarrow a$ for some type a , a frozen occurrence `[choose]` has type $\forall a.a \rightarrow a \rightarrow a$. More interestingly, whereas the term `single choose` has type `List (a → a → a)`, the term `single [choose]` has type `List (∀ a.a → a → a)`. This makes it possible to pass polymorphic arguments to functions that expect them. Consider a function `auto` : $(\forall a.a \rightarrow a) \rightarrow (\forall a.a \rightarrow a)$. Whereas the term `auto id` does not typecheck (because `id` is implicitly instantiated to type $a \rightarrow a$ which does not match the argument type $\forall a.a \rightarrow a$ of `auto`) the term `auto [id]` does.

Explicit Generalisation (`$V`). We can generalise an expression to its principal polymorphic type by binding it to a variable and then freezing it, for instance: `let id = λx.x in poly [id]`, where `poly` : $(\forall a.a \rightarrow a) \rightarrow \text{Int} \times \text{Bool}$. The explicit generalisation operator `$` generalises the type of any value. Whereas the term `λx.x` has type $a \rightarrow a$, the term `$(λx.x)` has type $\forall a.a \rightarrow a$, allowing us to write `poly $(λx.x)`. Explicit generalisation is macro-expressible [6] in FreezeML.

$$\$V \equiv \text{let } x = V \text{ in } [x]$$

We can also define a type-annotated variant:

$$\$^A V \equiv \text{let } (x : A) = V \text{ in } [x]$$

Note that FreezeML adopts the ML value restriction [30]; hence let generalisation only applies to syntactic values.

Explicit Instantiation (`@M`). As in ML, the polymorphic types of variables are implicitly instantiated when type-checking each variable occurrence. Unlike in ML, other terms

can have polymorphic types, which are *not* implicitly instantiated. Nevertheless, we can instantiate a term by binding it to a variable: `let x = head ids in x 42`, where `head : $\forall a. \text{List } (a) \rightarrow a$` returns the first element in a list and `ids : $\text{List } (\forall a. a \rightarrow a)$` is a list of polymorphic identity functions. The explicit instantiation operator `@` supports instantiation of a term without having to explicitly bind it to a variable. For instance, whereas the term `head ids` has type $\forall a. a \rightarrow a$ the term `(head ids)@` in the context of application to `42` has type $\text{Int} \rightarrow \text{Int}$, so `(head ids)@ 42` is well-formed. Explicit instantiation is macro-expressible in FreezeML:

$$M@ \equiv \text{let } x = M \text{ in } x$$

Ordered Quantifiers. Like in System F, but unlike in ML, the order of quantifiers matters. Quantifiers introduced through generalisation are ordered by the sequence in which they first appear in a type. Type annotations allow us to specify a different quantifier order, but variable instantiation followed by generalisation restores the canonical order. For example, if we have functions $f : (\forall a b. a \rightarrow b \rightarrow a \times b) \rightarrow \text{Int}$, $\text{pair} : \forall a b. a \rightarrow b \rightarrow a \times b$, and $\text{pair}' : \forall b a. a \rightarrow b \rightarrow a \times b$, then $f \text{ [pair]}$, $f \text{ \$pair}$, $f \text{ \$pair}'$ have type Int and behave identically, whereas $f \text{ [pair}']$ is ill-typed.

Monomorphic parameter inference. As in ML, function arguments need not have annotations, but their inferred types must be monomorphic, i.e. we cannot typecheck `bad`:

$$\text{bad} = \lambda f. (f \ 42, f \ \text{True})$$

Unlike in ML we can annotate arguments with polymorphic types and use them at different types:

$$\text{poly} = \lambda (f : \forall a. a \rightarrow a). (f \ 42, f \ \text{True})$$

One might hope that it is safe to infer polymorphism by local, compositional reasoning, but that is not the case. Consider the following two functions.

$$\begin{aligned} \text{bad1} &= \lambda f. (\text{poly } [f], (f \ 42) + 1) \\ \text{bad2} &= \lambda f. ((f \ 42) + 1, \text{poly } [f]) \end{aligned}$$

We might reasonably expect both to be typeable by assigning the type $\forall a. a \rightarrow a$ to f . Now, assume type inference is left-to-right. In `bad1` we first infer that f has type $\forall a. a \rightarrow a$ (as $[f]$ is the argument to `poly`); then we may instantiate a to Int when applying f to `42`. In `bad2` we eagerly infer that f has type $\text{Int} \rightarrow \text{Int}$; now when we pass $[f]$ to `poly`, type inference fails. To rule out this kind of sensitivity to the order of type inference, and the resulting incompleteness of our type inference algorithm, we insist that unannotated λ -bound variables be monomorphic. This in turn entails checking monomorphism constraints on type variables and maintaining other invariants (Section 3.2). (One can build more sophisticated systems that defer determining whether a term is polymorphic or not until more information becomes available — both Poly-ML and MLF do, for instance — but we prefer to keep things simple.)

2.1 FreezeML by Example

Figure 1 presents a collection of FreezeML examples that showcase how our system works in practice. We use functions with type signatures shown in Figure 2 (adapted from Serrano et al. [25]). In Figure 1 well-formed expressions are annotated with a type inferred in FreezeML, whilst ill-typed expressions are annotated with \times . Sections A-E of the table are taken from [25]. Section F of the table contains additional examples which further highlight the behaviour of our system. Examples F1-F4 show how to define some of the functions and values in Figure 2 in FreezeML. In FreezeML it is sometimes possible to infer a different type depending on the presence of freeze, generalisation, and instantiation operators. In such cases we provide two copies of an example in Figure 1, the one with extra FreezeML annotations being marked with \bullet . Sometimes explicit instantiation, generalisation, or freezing is mandatory to make an expression well-formed in FreezeML. In such cases there is only one, well-formed copy of an example marked with a \star , e.g. A9 \star . Example F10 \dagger typechecks only in a system without a value restriction due to generalisation of an application.

3 FreezeML via System F and ML

In this section we give a syntax-directed presentation of FreezeML and discuss various design choices that we have made. We wish for FreezeML to be an ML-like call-by-value language with the expressive power of System F. To this end we rely on a standard call-by-value definition of System F, which additionally obeys the value restriction (i.e. only values are allowed under type abstractions). We take mini-ML [1] as a core representation of a call-by-value ML language. Unlike System F, ML separates monotypes from (polymorphic) type schemes and has no explicit type abstraction and application. Polymorphism in ML is introduced by generalising the body of a let-binding, and eliminated implicitly when using a variable. Another crucial difference between System F and ML is that in the former the order of quantifiers in a polymorphic type matters, whereas in the latter it does not. Full definitions of System F and ML, including the syntax, kinding and typing rules, as well as translation from ML to System F, are given in the extended version of this paper [5].

Notations. We write $\text{ftv}(A)$ for the sequence of distinct free type variables of a type in the order in which they first appear in A . For example, $\text{ftv}((a \rightarrow b) \rightarrow (a \rightarrow c)) = a, b, c$. Whenever a kind environment Δ appears as a domain of a substitution or a \forall quantifier, it is allowed to be empty. In such case we identify type $\forall \Delta. H$ with H . We write $\Delta - \Delta'$ for the restriction of Δ to those type variables that do not appear in Δ' . We write $\Delta \# \Delta'$ to mean that the type variables in Δ and Δ' are disjoint. Disjointness is also implicitly required when concatenating Δ and Δ' to Δ, Δ' .

A	POLYMORPHIC INSTANTIATION	B	INFERENCE WITH POLYMORPHIC ARGUMENTS
A1	$\lambda x y. y : a \rightarrow b \rightarrow b$	B1★	$\lambda(f : \forall a. a \rightarrow a).$ $(f \ 1, f \ \text{True}) : (\forall a. a \rightarrow a) \rightarrow \text{Int} \times \text{Bool}$
A1●	$\$(\lambda x y. y) : \forall a. b. a \rightarrow b \rightarrow b$	B2★	$\lambda(xs : \text{List}(\forall a. a \rightarrow a)).$ $\text{poly}(\text{head } xs) : \text{List}(\forall a. a \rightarrow a) \rightarrow \text{Int} \times \text{Bool}$
A2	$\text{choose id} : (a \rightarrow a) \rightarrow (a \rightarrow a)$	D APPLICATION FUNCTIONS	
A2●	$\text{choose [id]} : (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$	D1★	$\text{app poly [id]} : \text{Int} \times \text{Bool}$
A3	$\text{choose [] ids} : \text{List}(\forall a. a \rightarrow a)$	D2★	$\text{revapp [id] poly} : \text{Int} \times \text{Bool}$
A4	$\lambda(x : \forall a. a \rightarrow a). x \ x : (\forall a. a \rightarrow a) \rightarrow (b \rightarrow b)$	D3★	$\text{runST [argST]} : \text{Int}$
A4●	$\lambda(x : \forall a. a \rightarrow a). x \ [x] : (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$	D4★	$\text{app runST [argST]} : \text{Int}$
A5	$\text{id auto} : (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$	D5★	$\text{revapp [argST] runST} : \text{Int}$
A6	$\text{id auto}' : (\forall a. a \rightarrow a) \rightarrow (b \rightarrow b)$	E η -EXPANSION	
A6●	$\text{id [auto]}' : \forall b. (\forall a. a \rightarrow a) \rightarrow (b \rightarrow b)$	E1	$k \ h \ l : \times$
A7	$\text{choose id auto} : (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$	E2★	$k \ $(\lambda x. (h \ x)@) \ l : \forall a. \text{Int} \rightarrow a \rightarrow a$ where $k : \forall a. a \rightarrow \text{List } a \rightarrow a$ $h : \text{Int} \rightarrow \forall a. a \rightarrow a$ $l : \text{List}(\forall a. \text{Int} \rightarrow a \rightarrow a)$
A8	$\text{choose id auto}' : \times$	E3	$r \ (\lambda x \ y. y) : \times$
A9★	$f \ (\text{choose [id]}) \ \text{ids} : \forall a. a \rightarrow a$ where $f : \forall a. (a \rightarrow a) \rightarrow \text{List } a \rightarrow a$	E3●	$r \ $(\lambda x. $(\lambda y. y)) : \text{Int}$ where $r : (\forall a. a \rightarrow \forall b. b \rightarrow b) \rightarrow \text{Int}$
A10★	$\text{poly [id]} : \text{Int} \times \text{Bool}$	F FreezeML PROGRAMS	
A11★	$\text{poly} \ $(\lambda x. x) : \text{Int} \times \text{Bool}$	F1	$\text{id} = \$(\lambda x. x) : \forall a. a \rightarrow a$
A12★	$\text{id poly} \ $(\lambda x. x) : \text{Int} \times \text{Bool}$	F2	$\text{ids} = [[\text{id}]] : \text{List}(\forall a. a \rightarrow a)$
C FUNCTIONS ON POLYMORPHIC LISTS		F3	$\text{auto} = \lambda(x : \forall a. a \rightarrow a). x \ [x] : (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$
C1	$\text{length ids} : \text{Int}$	F4	$\text{auto}' = \lambda(x : \forall a. a \rightarrow a). x \ x : \forall b. (\forall a. a \rightarrow a) \rightarrow b \rightarrow b$
C2	$\text{tail ids} : \text{List}(\forall a. a \rightarrow a)$	F5★	$\text{auto [id]} : \forall a. a \rightarrow a$
C3	$\text{head ids} : \forall a. a \rightarrow a$	F6	$(\text{head ids}) :: \text{ids} : \text{List}(\forall a. a \rightarrow a)$
C4	$\text{single id} : \text{List}(a \rightarrow a)$	F7★	$(\text{head ids})@ \ 3 : \text{Int}$
C4●	$\text{single [id]} : \text{List}(\forall a. a \rightarrow a)$	F8	$\text{choose}(\text{head ids}) : (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$
C5★	$[\text{id}] :: \text{ids} : \text{List}(\forall a. a \rightarrow a)$	F8●	$\text{choose}(\text{head ids})@ : (a \rightarrow a) \rightarrow (a \rightarrow a)$
C6★	$\$(\lambda x. x) :: \text{ids} : \text{List}(\forall a. a \rightarrow a)$	F9	$\text{let } f = \text{revapp [id]} \ \text{in } f \ \text{poly}$ $: \text{Int} \times \text{Bool}$
C7	$(\text{single inc}) \ ++ \ (\text{single id}) : \text{List}(\text{Int} \rightarrow \text{Int})$	F10†	$\text{choose id} \ (\lambda(x : \forall a. a \rightarrow a). \$(\text{auto}' \ x))$ $: (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$
C8★	$g \ (\text{single [id]}) \ \text{ids} : \forall a. a \rightarrow a$ where $g : \forall a. \text{List } a \rightarrow \text{List } a \rightarrow a$		
C9★	$\text{map poly} \ (\text{single [id]}) : \text{List}(\text{Int} \times \text{Bool})$		
C10	$\text{map head} \ (\text{single ids}) : \text{List}(\forall a. a \rightarrow a)$		

Figure 1. Example FreezeML Terms and Types

$\text{head} : \forall a. \text{List } a \rightarrow a$	$\text{id} : \forall a. a \rightarrow a$	$\text{map} : \forall a \ b. (a \rightarrow b) \rightarrow \text{List } a \rightarrow \text{List } b$
$\text{tail} : \forall a. \text{List } a \rightarrow \text{List } a$	$\text{ids} : [\forall a. a \rightarrow a]$	$\text{app} : \forall a \ b. (a \rightarrow b) \rightarrow a \rightarrow b$
$[\] : \forall a. \text{List } a$	$\text{inc} : \text{Int} \rightarrow \text{Int}$	$\text{revapp} : \forall a \ b. a \rightarrow (a \rightarrow b) \rightarrow b$
$(::) : \forall a. a \rightarrow \text{List } a \rightarrow \text{List } a$	$\text{choose} : \forall a. a \rightarrow a \rightarrow a$	$\text{runST} : \forall a. (\forall s. ST \ s \ a) \rightarrow a$
$\text{single} : \forall a. a \rightarrow \text{List } a$	$\text{poly} : (\forall a. a \rightarrow a) \rightarrow \text{Int} \times \text{Bool}$	$\text{argST} : \forall s. ST \ s \ \text{Int}$
$(++) : \forall a. \text{List } a \rightarrow \text{List } a \rightarrow \text{List } a$	$\text{auto} : (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$	$\text{pair} : \forall a \ b. a \rightarrow b \rightarrow a \times b$
$\text{length} : \forall a. \text{List } a \rightarrow \text{Int}$	$\text{auto}' : \forall b. (\forall a. a \rightarrow a) \rightarrow (b \rightarrow b)$	$\text{pair}' : \forall b \ a. a \rightarrow b \rightarrow a \times b$

Figure 2. Type signatures for functions used in the text; adapted from [25].

3.1 FreezeML

FreezeML is an extension of ML with two new features. First, let-bindings and lambda-bindings may be annotated with arbitrary System F types. Second, FreezeML adds a new form $[x]$, called *frozen variables*, for preventing variables from being instantiated.

The syntax of FreezeML is given in Figure 3. (We name the syntactic categories for later use in Section 5.) The types

are the same as in System F. We explicitly distinguish two kinds of type: a monotype (S), is as in ML a type entirely free of polymorphism, and a guarded type (H) is a type with no top-level quantifier (in which any polymorphism is guarded by a type constructor). The terms include all ML terms plus frozen variables ($[x]$) and lambda- and let-bindings with type ascriptions. Values are those terms that may be generalised under the value restriction. They are slightly more general

Type Variables	$\text{TVar} \ni a, b, c$
Type Constructors	$\text{Con} \ni D ::= \text{Int} \mid \text{List} \mid \rightarrow \mid \times \mid \dots$
Types	$\text{Type} \ni A, B ::= a \mid D\bar{A} \mid \forall a.A$
Monotypes	$\text{MType} \ni S, T ::= a \mid D\bar{S}$
Guarded Types	$\text{GType} \ni H ::= a \mid D\bar{A}$
Type Instantiation	$\text{Subst} \ni \delta ::= \emptyset \mid \delta[a \mapsto A]$
Term Variables	$\text{Var} \ni x, y, z$
Terms	$\text{Term} \ni M, N ::= x \mid [x] \mid \lambda x.M$ $\quad \quad \quad \mid \lambda(x : A).M \mid MN$ $\quad \quad \quad \mid \text{let } x = M \text{ in } N$ $\quad \quad \quad \mid \text{let } (x : A) = M \text{ in } N$
Values	$\text{Val} \ni V, W ::= x \mid [x] \mid \lambda x.M$ $\quad \quad \quad \mid \lambda(x : A).M$ $\quad \quad \quad \mid \text{let } x = V \text{ in } W$ $\quad \quad \quad \mid \text{let } (x : A) = V \text{ in } W$
Guarded Values	$\text{GVal} \ni U ::= x \mid \lambda x.M \mid \lambda(x : A).M$ $\quad \quad \quad \mid \text{let } x = V \text{ in } U$ $\quad \quad \quad \mid \text{let } (x : A) = V \text{ in } U$
Kinds	$\text{Kind} \ni K ::= \bullet \mid \star$
Kind Environments	$\text{PEnv} \ni \Delta ::= \cdot \mid \Delta, a$
Type Environments	$\text{TEnv} \ni \Gamma ::= \cdot \mid \Gamma, x : A$

Figure 3. FreezeML Syntax

$\Delta \vdash A : K$	$\text{arity}(D) = n$ $\Delta \vdash A_1 : K$ \dots
$\frac{a \in \Delta}{\Delta \vdash a : \bullet}$	$\frac{\Delta \vdash A_n : K}{\Delta \vdash D\bar{A} : K}$
$\frac{\Delta, a \vdash A : \star}{\Delta \vdash \forall a.A : \star}$	$\frac{\Delta \vdash A : \bullet}{\Delta \vdash A : \star}$

Figure 4. FreezeML Kinding Rules

$\Delta \vdash \delta : \Delta' \Rightarrow_K \Delta''$	$\Delta \vdash \delta : \Delta' \Rightarrow_K \Delta''$	$\Delta, \Delta'' \vdash A : K$
$\frac{\Delta \vdash \emptyset : \cdot \Rightarrow_K \Delta'}{\Delta \vdash \delta[a \mapsto A] : (\Delta', a) \Rightarrow_K \Delta''}$		

Figure 5. FreezeML Instantiation Rules

than the value forms of Standard ML in that they are closed under let binding (as in OCaml). Guarded values are those values that can only have guarded types (that is, all values except those that have a frozen variable in tail position).

The FreezeML kinding judgement $\Delta \vdash A : K$ states that type A has kind K in kind environment Δ . The kinding rules are given in Figure 4. As in ML we distinguish monomorphic types (\bullet) from polymorphic types (\star). Unlike in ML polymorphic types can appear inside data type constructors.

Rules for type instantiation are given in Figure 5. The judgement $\Delta \vdash \delta : \Delta' \Rightarrow \Delta''$ defines a well-formed finite map from type variables in Δ , Δ' into type variables in Δ, Δ'' , such that $\delta(a) = a$ for every $a \in \Delta$. As such, it is only well-defined

$$\begin{aligned} \emptyset(A) &= A & \delta[a \mapsto A](a) &= A \\ \delta(D\bar{A}) &= D(\delta(\bar{A})) & \delta[a \mapsto A](b) &= \delta(b) \\ \delta(\forall a.A) &= \forall c.\delta[a \mapsto c](A), & \text{where } c &\notin \text{ftv}(\delta(b)) \text{ for all } b \neq c \end{aligned}$$

Figure 6. Application of a Type Instantiation in FreezeML

$\Delta; \Gamma \vdash M : A$	VAR $\frac{x : \forall \Delta'. H \in \Gamma}{\Delta \vdash \delta : \Delta' \Rightarrow_{\star} \cdot}$	APP $\frac{\Delta; \Gamma \vdash M : A \rightarrow B \quad \Delta; \Gamma \vdash N : A}{\Delta; \Gamma \vdash MN : B}$
FREEZE $\frac{x : A \in \Gamma}{\Delta; \Gamma \vdash [x] : A}$	$\frac{\Delta \vdash \delta : \Delta' \Rightarrow_{\star} \cdot}{\Delta; \Gamma \vdash x : \delta(H)}$	
LAM $\frac{\Delta; \Gamma, x : S \vdash M : B}{\Delta; \Gamma \vdash \lambda x.M : S \rightarrow B}$	LAM-ASCRIBE $\frac{\Delta; \Gamma, x : A \vdash M : B}{\Delta; \Gamma \vdash \lambda(x : A).M : A \rightarrow B}$	
LET $(\Delta', \Delta'') = \text{gen}(\Delta, A', M) \quad (\Delta, \Delta'', M, A') \Downarrow A$ $\Delta, \Delta''; \Gamma \vdash M : A' \quad \Delta; \Gamma, x : A \vdash N : B$ $\text{principal}(\Delta, \Gamma, M, \Delta'', A')$		$\frac{\Delta; \Gamma \vdash \text{let } x = M \text{ in } N : B$
	LET-ASCRIBE $(\Delta', A') = \text{split}(A, M)$ $\Delta, \Delta'; \Gamma \vdash M : A' \quad \Delta; \Gamma, x : A \vdash N : B$	$\frac{\Delta; \Gamma \vdash \text{let } (x : A) = M \text{ in } N : B$

Figure 7. FreezeML typing rules

if Δ and Δ' are disjoint and Δ and Δ'' are disjoint. Type instantiation accounts for polymorphism by either being restricted to instantiate type variables with monomorphic kinds only (\Rightarrow_{\bullet}) or permitting polymorphic instantiations (\Rightarrow_{\star}). The following rule is admissible

$$\frac{\Delta, \Delta' \vdash A : K \quad \Delta \vdash \delta : \Delta' \Rightarrow_{K'} \Delta''}{\Delta, \Delta'' \vdash \delta(A) : K \sqcup K'}$$

where $\bullet \sqcup \bullet = \bullet$ and $\bullet \sqcup \star = \star \sqcup \bullet = \star \sqcup \star = \star$. We apply type instantiation in a standard way, taking care to account for shadowing of type variables (Figure 6).

The FreezeML judgement $\Delta; \Gamma \vdash M : A$ states that term M has type A in kind environment Δ and type environment Γ ; its rules are shown in Figure 7. These rules are adjusted with respect to ML to allow full System F types everywhere except in the types of variables bound by unannotated lambdas, where only monotypes are permitted.

As in ML, the VAR rule implicitly instantiates variables. The \star in the judgement $\Delta \vdash \delta : \Delta' \Rightarrow_{\star} \cdot$ indicates that the type variables in Δ' may be instantiated with polymorphic types. The FREEZE rule differs from the VAR rule only in that it suppresses instantiation. In the LAM rule, the restriction to a syntactically monomorphic argument type ensures that an argument cannot be used at different types inside the body of

$(\Delta, \Delta', M, A') \Downarrow A$	
$M \in \text{GVal}$	$\Delta \vdash \delta : \Delta' \Rightarrow \bullet \cdot \quad M \notin \text{GVal}$
$(\Delta, \Delta', M, A') \Downarrow \forall \Delta'. A'$	$(\Delta, \Delta', M, A') \Downarrow \delta(A')$
$\text{gen}(\Delta, A, M) = \begin{cases} (\Delta', \Delta') & \text{if } M \in \text{GVal} \\ (\cdot, \Delta') & \text{otherwise} \end{cases}$ <p style="text-align: center;">where $\Delta' = \text{ftv}(A) - \Delta$</p> $\text{split}(\forall \Delta'. H, M) = \begin{cases} (\Delta, H) & \text{if } M \in \text{GVal} \\ (\cdot, \forall \Delta'. H) & \text{otherwise} \end{cases}$ <p>$\text{principal}(\Delta, \Gamma, M, \Delta', A') =$ $\Delta' = \text{ftv}(A') - \Delta$ and $\Delta, \Delta'; \Gamma \vdash M : A'$ and (for all $\Delta'', A'' \mid$ if $\Delta'' = \text{ftv}(A'') - \Delta$ and $\Delta, \Delta''; \Gamma \vdash M : A''$ then there exists δ such that $\Delta \vdash \delta : \Delta' \Rightarrow_{\star} \Delta''$ and $\delta(A') = A''$)</p>	

Figure 8. FreezeML auxiliary definitions

a lambda abstraction. However, the type of an unannotated lambda abstraction may subsequently be generalised. For example, consider the expression $\text{poly } \$(\lambda x. x)$. The parameter x cannot be typed with a polymorphic type; giving the syntactic monotype a to x yields type $a \rightarrow a$ for the lambda-abstraction. The $\$$ operator then generalises this to $\forall a. a \rightarrow a$ as the type of argument passed to poly . The LAM-ASCRIBE rule allows an argument to be used polymorphically inside the body of a lambda abstraction. The APP rule is standard.

Let Bindings. Because we adopt the value restriction, the LET rule behaves differently depending on whether or not M is a guarded value (cf. GVal syntactic category in Figure 3). The choice of whether to generalise the type of M is delegated to the judgement $(\Delta, \Delta', M, A') \Downarrow A$, where A' is the type of M and Δ'' are the generalisable type variables of M , i.e. $\Delta'' = \text{ftv}(A') - \Delta$. The \Downarrow judgement determines A , the type given to x while type-checking N . If M is a guarded value, we generalise and have $A = \forall \Delta'. A'$. If M is not a guarded value, we have $A = \delta(A')$, where δ is an instantiation with $\Delta \vdash \delta : \Delta' \Rightarrow \bullet \cdot$. This means that instead of abstracting over the unbound type variables Δ'' of A' , we instantiate them *monomorphically*. We further discuss the need for this behaviour in Section 3.2.

The gen judgement used in the LET rule may seem surprising – its first component is unused whilst the second component is identical in both cases and corresponds to the generalisable type variables of A' . Indeed, the first component of gen is irrelevant for typing but it is convenient for writing the translation from FreezeML to System F (Figure 11 in Section 4.2), where it is used to form a type abstraction, and in the type inference algorithm (Figure 16 in Section 5.4), where it allows us to collapse two cases into one.

The LET rule requires that A' is the principal type for M . This constraint is necessary to ensure completeness of our type inference algorithm; we discuss it further in Section 3.2. The relation principal is defined in Figure 8.

The LET-ASCRIBE rule is similar to the LET rule, but instead of generalising the type of M , it uses the type A supplied via an annotation. As in LET, A' denotes the type of M . However, the annotated case admits non-principal types for M . The split operator enforces the value restriction. If M is a guarded value, A' must be a guarded type, i.e. we have $A' = H$ for some H . We then have $A = \forall \Delta'. H$. If M is not a guarded value split requires $A' = A$ and $\Delta' = \cdot$. This means that *all* toplevel quantifiers in A must originate from M itself, rather than from generalising it.

Every valid typing judgement in ML is also a valid typing judgement in FreezeML.

Theorem 1. *If $\Delta; \Gamma \vdash M : S$ in ML then $\Delta; \Gamma \vdash M : S$ in FreezeML.*

(The exact derivation can differ due to differences in the kinding rules and the principality constraint on the LET rule.)

3.2 Design Considerations

Monomorphic instantiation in the LET rule. Recall that the LET rule enforces the value restriction by instantiating those type variables that would otherwise be quantified over. Requiring these type variables to be instantiated with monotypes allows us to avoid problems similar to the ones outlined in Section 2. Consider the following two functions.

$$\text{bad3} = \lambda(\text{bot} : \forall a. a). \text{let } f = \text{bot bot in } (\text{poly } [f], (f \ 42) + 1)$$

$$\text{bad4} = \lambda(\text{bot} : \forall a. a). \text{let } f = \text{bot bot in } ((f \ 42) + 1, \text{poly } [f])$$

Since we do not generalise non-values in let-bindings due to the value restriction, in both of these examples f is initially assigned the type a rather than the most general type $\forall a. a$ (because bot bot is a non-value). Assuming type inference proceeds from left to right then type inference will succeed on bad3 and fail on bad4 for the same reasons as in Section 2. In order to rule out this class of examples, we insist that non-values are first generalised and then instantiated with monomorphic types. Thus we constrain a to only unify with monomorphic types, which leads to type inference failing on both bad3 and bad4.

Our guiding principle is “never guess polymorphism”. While our system permits instantiation of quantifiers with polymorphic types – per VAR rule – it does not permit polymorphic instantiations of type variables inside the type environment. The high-level invariant that FreezeML uses to ensure that this principle is not violated is that any (as yet) unknown types appearing in the type environment (which maps term variables to their currently inferred types) during type inference must be explicitly marked as monomorphic. The only means by which inference can introduce unknown types into the type environment are through unannotated

lambda-binders or through not generalising let-bound variables. By restricting these cases to be monomorphic we ensure in turn that any unknown type appearing in the type environment must be explicitly marked as monomorphic.

Principal Type Restriction. The LET rule requires that when typing $\text{let } x = M \text{ in } N$, the type A' given to M must be principal. Consider the program

$$\text{bad5} = \text{let } f = \lambda x.x \text{ in } [f] \quad 42$$

On the one hand, if we infer the type $\forall a.a \rightarrow a$ for f , then bad5 will fail to type check as we cannot apply a term of polymorphic type (instantiation is only automatic for variables). However, given a traditional declarative type system one might reasonably propose $\text{Int} \rightarrow \text{Int}$ as a type for f , in which case bad5 would be typeable – albeit a conventional type inference algorithm would have difficulty inferring a type for it. In order to ensure completeness of our type inference algorithm in the presence of generalisation and freeze, we bake principality into the typing rule for let, similarly to [7, 13, 15, 27]. This means that the only legitimate type that f may be assigned is the most general one, that is $\forall a.a \rightarrow a$.

One may think of side-stepping the problem with bad5 by always instantiating terms that appear in application position (after all, it is always a type error for an uninstantiated term of polymorphic type to appear in application position). But then we can exhibit the same problem with a slightly more intricate example.

$$\text{bad6} = \text{let } f = \lambda x.x \text{ in id } [f] \quad 42$$

The principality condition is also applied in the non-generalising case of the LET rule, meaning that we must instantiate the principal type for M rather than an arbitrary one. Otherwise, we could still type bad4 by assigning bot bot type $\forall a.a \rightarrow a$. In the LET rule Δ' would be empty, making instantiation a no-op.

Well-foundedness. The alert reader may already have noticed a complication resulting from the principal type restriction: $\text{principal}(\Delta, \Gamma, M, \Delta', A')$ contains a negative occurrence of the typing relation, in order to express that Δ', A' is a “most general” solution for Δ'', A'' among all possible derivations of $\Delta, \Delta''; \Gamma \vdash M : A''$. This negative occurrence means that *a priori*, the rules in Figures 7 and 8 do not form a proper inductive definition.

This is a potentially serious problem, but it can be resolved easily by observing that the rules, while not syntactically well-founded, can be stratified. Instead of considering the rules in Figures 7 and 8 as a *single* inductive definition, we consider them to determine a *function* $\mathcal{J}[\![\!-\!] \!]$ from terms M to triples (Δ, Γ, A) . The typing relation is then defined as $\Delta; \Gamma \vdash M : A \iff (\Delta, \Gamma, A) \in \mathcal{J}[\![M]\!]$. We can easily prove by induction on M that $\mathcal{J}[\![M]\!]$ is well-defined. Furthermore, we can show that the inference rules in Figure 7 hold and are invertible. When reasoning about typing judgements, we

can proceed by induction on M and use inversion. It is also sound to perform recursion over typing derivations provided the principal assumption is not needed; we indicate this by greying out this assumption (for example in Figure 11). We give full details and explain how this reasoning is performed in the extended version of this paper [5].

Type Variable Scoping. A type annotation in FreezeML may contain type variables that is not bound by the annotation. In contrast to many other systems, we do not interpret such variables existentially, but allow binding type variables across different annotations. In an expression $\text{let } (x : A) = M \text{ in } N$, we therefore consider the toplevel quantifiers of A bound in M , meaning that they can be used freely in annotations inside M , rather like GHC’s scoped type variables [21]. However, this is only true for the generalising case, when M is a guarded value. In the absence of generalisation, any polymorphism in the type A originates from M directly (e.g., because M is a frozen variable). Hence, if M is not a guarded value no bound variables of A are bound in M .

Note that given the let binding above, where A has the shape $\forall \Delta.H$, there is no ambiguity regarding which of the type variables in Δ result from generalisation and which originate from M itself. If M is a guarded value, its type is guarded, too, and hence all variables in Δ result from generalisation. Conversely, if $M \notin \text{GVal}$, then there is no generalisation at all.

Due to the unambiguity of the binding behaviour in our system with the value restriction, we can define a purely syntax-directed well-formedness judgement for verifying that types in annotations are well-kinded and respect the intended scoping of type-variables. We call this property well-scopedness, and it is a prerequisite for type inference. The corresponding judgement is $\Delta \Vdash M$, checking that in M , the type annotations are well-formed with respect to kind environment Δ (Figure 9). The main subtlety in this judgement is in how Δ grows when we encounter *annotated* let-bindings. For annotated lambdas, we just check that the type annotation is well-formed in Δ but do not add any type variables in Δ . For plain let, we just check well-scopedness recursively. However, for annotated let-bindings, we check that the type annotation A is well-formed, and we check that M is well-scoped *after extending Δ with the top-level type variables of A* . This is sensible because in the LET-ASCRIBE rule, these type variables (present in the type annotation) are introduced into the kind environment when type checking M . In an unannotated let, in contrast, the generalisable type variables are not mentioned in M , so it does not make sense to allow them to be used in other type annotations inside M .

As a concrete example of how this works, consider an explicitly annotated let-binding of the identity function: $\text{let } (f : \forall a.a \rightarrow a) = \lambda(x : a).x \text{ in } N$, where the a type annotation on x is bound by $\forall a$ in the type annotation on f . However, if we left off the $\forall a.a \rightarrow a$ annotation on f ,

$\frac{}{\Delta \Vdash [x]}$	$\frac{}{\Delta \Vdash x}$	$\frac{\Delta \Vdash M}{\Delta \Vdash \lambda x.M}$
$\frac{\Delta \Vdash A : \star}{\Delta \Vdash M}$	$\frac{\Delta \Vdash M}{\Delta \Vdash N}$	$\frac{\Delta \Vdash M \quad \Delta \Vdash N}{\Delta \Vdash \mathbf{let} \ x = M \ \mathbf{in} \ N}$
$\frac{}{\Delta \Vdash \lambda(x : A).M}$	$\frac{}{\Delta \Vdash MN}$	$\frac{}{\Delta \Vdash \mathbf{let} \ x = M \ \mathbf{in} \ N}$
$\frac{\Delta \Vdash A : \star \quad (\Delta', A') = \mathbf{split}(A, M) \quad \Delta, \Delta' \Vdash M \quad \Delta \Vdash N}{\Delta \Vdash \mathbf{let} \ (x : A) = M \ \mathbf{in} \ N}$		

Figure 9. Well-Scopedness of FreezeML Terms

then the a annotation on x would be unbound. This also means that in expressions, we cannot let type annotations α -vary freely; that is, the previous expression is α -equivalent to $\mathbf{let} \ (f : \forall b.b \rightarrow b) = \lambda(x : b).x \ \mathbf{in} \ N$ but not to $\mathbf{let} \ (f : \forall b.b \rightarrow b) = \lambda(x : a).x \ \mathbf{in} \ N$. This behaviour is similar to other proposals for scoped type variables [21].

“Pure” FreezeML. In a hypothetical version of FreezeML without the value restriction, a purely syntactic check on $\mathbf{let} \ (x : A) = M \ \mathbf{in} \ N$ is not sufficient to determine which top-level quantifiers of A are bound in M . In the expression

$$\begin{aligned} &\mathbf{let} \ (f : \forall a b.a \rightarrow b \rightarrow b) = \\ &\quad \mathbf{let} \ (g : \forall b.a \rightarrow b \rightarrow b) = \lambda y z.z \ \mathbf{in} \ \mathbf{id} \ [g] \\ &\mathbf{in} \ N \end{aligned}$$

the outer \mathbf{let} generalises a , unlike the subsequent variable b , which arises from the inner \mathbf{let} binding. The well-scopedness judgement would require typing information. Moreover, the LET-ASC rule would have to nondeterministically split the type annotation A into $\forall \Delta', \Delta''. H$, such that Δ' contains those variables to generalise (a in the example), and Δ'' contains those type variables originating from M directly (b in the example). Similarly, type inference would have to take this splitting into account.

Instantiation strategies. In FreezeML (and indeed ML) the only terms that are implicitly instantiated are variables. Thus (head ids) 42 is ill-typed and we must insert the instantiation operator $@$ to yield a type-correct expression: (head ids) $@$ 42. It is possible to extend our approach to perform *eliminator instantiation*, whereby we implicitly instantiate terms appearing in monomorphic elimination position (in particular application position), and thus, for instance, infer a type for bad5 without compromising completeness.

Another possibility is to instantiate all terms, except those that are explicitly frozen or generalised. Here, it also makes sense to extend the $[-]$ operator to act on arbitrary terms, rather than just variables. We call this strategy *pervasive instantiation*. Like eliminator instantiation, pervasive instantiation infers a type for (head ids) 42. However, pervasive

$$\begin{aligned} \mathcal{E}[[x]] &= [x] \\ \mathcal{E}[[\lambda x^A.M]] &= \lambda(x : A).\mathcal{E}[[M]] \\ \mathcal{E}[[MN]] &= \mathcal{E}[[M]] \ \mathcal{E}[[N]] \\ \mathcal{E}[[\Lambda a.V^B]] &= \mathbf{let} \ (x : \forall a.B) = (\mathcal{E}[[V]])@ \ \mathbf{in} \ [x] \\ \mathcal{E}[[M^{\forall a.B} A]] &= \mathbf{let} \ (x : B[A/a]) = (\mathcal{E}[[M]])@ \ \mathbf{in} \ [x] \end{aligned}$$

Figure 10. Translation from System F to FreezeML

instantiation requires inserting explicit generalisation where it was previously unnecessary. Moreover, pervasive instantiation complicates the meta-theory, requiring two mutually recursive typing judgements instead of just one.

The formalism developed in this paper uses variable instantiation alone, but our implementation also supports eliminator instantiation. We defer further theoretical investigation of alternative strategies to future work.

4 Relating System F and FreezeML

In this section we present type-preserving translations mapping System F terms to FreezeML terms and vice versa. We also briefly discuss the equational theory induced on FreezeML by these translations.

4.1 From System F to FreezeML

Figure 10 defines a translation $\mathcal{E}[[-]]$ of System F terms into FreezeML. The translation depends on types of subterms and is thus formally defined on derivations, but we use a shorthand notation in which subterms are annotated with their type (e.g., in $\Lambda a.V^B$, B indicates the type of V).

Variables are frozen to suppress instantiation. Term abstraction and application are translated homomorphically.

Type abstraction $\Lambda a.V$ is translated using an annotated let-binding to perform the necessary generalisation. However, we cannot bind x to the translation of V directly as only *guarded* values may be generalised but $\mathcal{E}[[V]]$ may be an unguarded value (concretely, a frozen variable). Hence, we bind x to $(\mathcal{E}[[V]])@$, which is syntactic sugar for $\mathbf{let} \ y = \mathcal{E}[[V]] \ \mathbf{in} \ y$. This expression is indeed a guarded value. We then freeze x to prevent immediate instantiation. Type application MA , where M has type $\forall a.B$, is translated similarly to type abstraction. We bind x to the result of translating M , but only after instantiating it. The variable x is annotated with the intended return type $B[A/a]$ and returned frozen.

Explicit instantiation is strictly necessary and the following, seemingly easier translation is incorrect.

$$\mathcal{E}[[M^{\forall a.B} A]] \neq \mathbf{let} \ (x : B[A/a]) = \mathcal{E}[[M]] \ \mathbf{in} \ [x]$$

The term $\mathcal{E}[[M]]$ may be a frozen variable or an application, whose type cannot be implicitly instantiated to type $B[A/a]$.

For any System F value V (i.e., any term other than an application), $\mathcal{E}[[V]]$ yields a FreezeML value (Figure 3).

Each translated term has the same type as the original.

$$\begin{array}{c}
C \left[\frac{x : A \in \Gamma}{\Delta; \Gamma \vdash [x] : A} \right] = x \quad C \left[\frac{\Delta; \Gamma, x : S \vdash M : B}{\Delta; \Gamma \vdash \lambda x.M : S \rightarrow B} \right] = \lambda x^S.C[[M]] \quad C \left[\frac{\Delta; \Gamma, x : A \vdash M : B}{\Delta; \Gamma \vdash \lambda(x : A).M : A \rightarrow B} \right] = \lambda x^A.C[[M]] \\
C \left[\frac{x : \forall \Delta'. H \in \Gamma \quad \Delta \vdash \delta : \Delta' \Rightarrow_{\star} \cdot}{\Delta; \Gamma \vdash x : \delta(H)} \right] = x \delta(\Delta') \quad C \left[\frac{\Delta; \Gamma \vdash M : A \rightarrow B \quad \Delta; \Gamma \vdash N : A}{\Delta; \Gamma \vdash MN : B} \right] = C[[M]] C[[N]] \\
C \left[\frac{\begin{array}{c} (\Delta', \Delta'') = \text{gen}(\Delta, A', M) \quad (\Delta, \Delta'', M, A') \Downarrow A \\ \Delta, \Delta''; \Gamma \vdash M : A' \quad \Delta; \Gamma, x : A \vdash N : B \\ \text{principal}(\Delta, \Gamma, M, \Delta'', A') \end{array}}{\Delta; \Gamma \vdash \text{let } x = M \text{ in } N : B} \right] = \text{let } x^A = \Lambda \Delta'. C[[M]] \text{ in } C[[N]] = C \left[\frac{\begin{array}{c} (\Delta', A') = \text{split}(A, M) \\ \Delta, \Delta'; \Gamma \vdash M : A' \\ \Delta; \Gamma, x : A \vdash N : B \end{array}}{\Delta; \Gamma \vdash \text{let } (x : A) = M \text{ in } N : B} \right]
\end{array}$$

Figure 11. Translation from FreezeML to System F

Theorem 2 (Type preservation). *If $\Delta; \Gamma \vdash M : A$ in System F then $\Delta; \Gamma \vdash \mathcal{E}[[M]] : A$ in FreezeML.*

4.2 From FreezeML to System F

Figure 11 gives the translation of FreezeML to System F. The translation depends on types of subterms and is thus formally defined on derivations. Frozen variables in FreezeML are simply variables in System F. A plain (i.e., not frozen) variable x is translated to a type application $x \delta(\Delta')$, where $\delta(\Delta')$ stands for the pointwise application of δ to Δ' . Here, δ and Δ' are obtained from x 's type derivation in FreezeML; Δ' contains all top-level quantifiers of x 's type. This makes FreezeML's implicit instantiation of non-frozen variables explicit. Lambda abstractions and applications translate directly. Let-bindings in FreezeML are translated as generalised let-bindings in System F where $\text{let } x^A = M \text{ in } N$ is syntactic sugar for $(\lambda x^A.N) M$. Here, generalisation is repeated type abstraction.

Each translated term has the same type as the original.

Theorem 3 (Type preservation). *If $\Delta; \Gamma \vdash M : A$ holds in FreezeML then $\Delta; \Gamma \vdash C[[M]] : A$ holds in System F.*

4.3 Equational reasoning

We can derive and verify equational reasoning principles for FreezeML by lifting from System F via the translations. We write $M \simeq N$ to mean M is observationally equivalent to N whenever $\Delta; \Gamma \vdash M : A$ and $\Delta; \Gamma \vdash N : A$. At a minimum we expect β -rules to hold, and indeed they do; the twist is that they involve substituting a different value depending on whether the variable being substituted for is frozen or not.

$$\begin{array}{l}
\text{let } x = V \text{ in } N \quad \simeq N[\$V / [x], (\$V)@ / x] \\
\text{let } (x : A) = V \text{ in } N \quad \simeq N[\$^A V / [x], (\$^A V)@ / x] \\
(\lambda x.M) V \quad \simeq M[V / [x], V@ / x] \\
(\lambda(x : A).M) V \quad \simeq M[V / [x], V@ / x]
\end{array}$$

If we perform type-erasure then these rules degenerate to the standard ones. We can also verify that η -rules hold.

$$\begin{array}{l}
\text{let } x = U \text{ in } x \quad \simeq U \quad \text{let } x = [y] \text{ in } x \quad \simeq y \\
\text{let } (x : A) = U \text{ in } x \quad \simeq U \quad \text{let } (x : A) = [y] \text{ in } x \quad \simeq y \\
\lambda x.M x \quad \simeq M \quad \lambda(x : A).M [x] \quad \simeq M
\end{array}$$

5 Type Inference

In this section we present a sound and complete type inference algorithm for FreezeML. The style of presentation is modelled on that of Leijen [13].

5.1 Type Variables and Kinds

When expressing type inference algorithms involving first-class polymorphism, it is crucial to distinguish between object language type variables, and meta language type variables that stand for unknown types required to solve the type inference problem. This distinction is the same as that between *eigenvariables* and *logic variables* in higher-order logic programming [18]. We refer to the former as *rigid* type variables and the latter as *flexible* type variables. For the purposes of the algorithm we will explicitly separate the two by placing them in different kind environments.

As in the rest of the paper, we let Δ range over fixed kind environments in which every type variable is monomorphic (kind \bullet). In order to support, for instance, applying a function to a polymorphic argument, we require flexible variables that may be unified with polymorphic types. For this purpose we introduce refined kind environments ranged over by Θ . Type variables in a refined kind environment may be polymorphic (kind \star) or monomorphic (kind \bullet). In our algorithms we place rigid type variables in a fixed environment Δ and flexible type variables in a refined environment Θ . Refined kind environments (Θ) are given by the following grammar.

$$\text{KEnv } \Theta ::= \cdot \mid \Theta, a : K$$

We often implicitly treat fixed kind environments \bar{a} as refined kind environments $\bar{a} : \bullet$. The refined kinding rules are given in Figure 12.

$\Theta \vdash A : K$	CONS $\text{arity}(D) = n$ $\Theta \vdash A_1 : K$
TYVAR $a : K \in \Theta$... FORALL $\Theta, a : \bullet \vdash A : \star$
$\Theta \vdash a : K$	UPCAST $\Theta \vdash A : \bullet$
$\Theta \vdash \Gamma$	EXTEND $\Theta \vdash \Gamma \quad \Theta \vdash A : \star$ (for all $a \in \text{ftv}(A) \mid a : \bullet \in \Theta$)
EMPTY $\Theta \vdash \cdot$	$\Theta \vdash \Gamma, x : A$

Figure 12. Refined Kinding Rules

$\Delta \vdash \theta : \Theta \Rightarrow \Theta'$	$\Delta \vdash \theta : \Theta' \Rightarrow \Theta$	$\Delta, \Theta \vdash A : K$
$\Delta \vdash \emptyset : \cdot \Rightarrow \Theta$	$\Delta \vdash \theta[a \mapsto A] : (\Theta', a : K) \Rightarrow \Theta$	

Figure 13. Type Substitutions

The key difference with respect to the object language kinding rules is that type variables can now be polymorphic. Rather than simply defining kinding of type environments point-wise the EXTEND rule additionally ensures that all type variables appearing in a type environment are monomorphic. This restriction is crucial for avoiding guessing of polymorphism. More importantly, it is also key to ensuring that typing judgements are stable under substitution. Without it it would be possible to substitute monomorphic type variables with types containing nested polymorphic variables, thus introducing polymorphism into a monomorphic type.

We generalise typing judgements $\Delta; \Gamma \vdash M : A$ to $\Theta; \Gamma \vdash M : A$, adopting the convention that $\Theta \vdash \Gamma$ and $\Theta \vdash A$ must hold as preconditions.

5.2 Type Substitutions

In order to define the type inference algorithm we will find it useful to define a judgement for type substitutions θ , which operate on flexible type variables, unlike type instantiations δ , which operate on rigid type variables. The type substitution rules are given in Figure 13. The rules are as in Figure 7, except that the kind environments on the right of the turnstile are refined kind environments and rather than the substitution having a fixed kind, the kind of each type variable must match up with the kind of the type it binds.

We write ι_Θ for the identity type substitution on Θ , omitting the subscript when clear from context.

$$\iota = \emptyset \quad \iota_{\Theta, a:K} = \iota_\Theta[a \mapsto a]$$

Composition of type substitutions is standard.

$$\theta \circ \emptyset = \theta \quad \theta \circ \theta'[a \mapsto A] = (\theta \circ \theta')[a \mapsto \theta(A)]$$

The rules shown in Figure 14 are admissible and we make use of them freely in our algorithms and proofs.

S-IDENTITY $\Delta \vdash \iota_\Theta : \Theta \Rightarrow \Theta$	S-WEAKEN $\Delta \vdash \theta : \Theta \Rightarrow \Theta'$ $\Delta, \Delta' \vdash \theta : \Theta \Rightarrow \Theta', \Theta''$
S-COMPOSE $\Delta \vdash \theta : \Theta' \Rightarrow \Theta''$ $\Delta \vdash \theta' : \Theta \Rightarrow \Theta'$ $\Delta \vdash \theta \circ \theta' : \Theta \Rightarrow \Theta''$	S-STRENGTHEN $\Delta \vdash \theta : \Theta \Rightarrow \Theta'$ $\text{ftv}(\theta) \# \Delta', \Theta''$ $\Delta - \Delta' \vdash \theta : \Theta \Rightarrow \Theta' - \Theta''$

Figure 14. Properties of Substitution

$\text{unify} : (\text{PEnv} \times \text{KEnv} \times \text{Type} \times \text{Type}) \rightarrow (\text{KEnv} \times \text{Subst})$

$\text{unify}(\Delta, \Theta, a, a) =$
return (Θ, ι)

$\text{unify}(\Delta, (\Theta, a : K), a, A) =$
let $\Theta_1 = \text{demote}(K, \Theta, \text{ftv}(A) - \Delta)$
assert $\Delta, \Theta_1 \vdash A : K$
return $(\Theta_1, \iota[a \mapsto A])$

$\text{unify}(\Delta, (\Theta, a : K), A, a) =$
let $\Theta_1 = \text{demote}(K, \Theta, \text{ftv}(A) - \Delta)$
assert $\Delta, \Theta_1 \vdash A : K$
return $(\Theta_1, \iota[a \mapsto A])$

$\text{unify}(\Delta, \Theta, D \bar{A}, D \bar{B}) =$
let $(\Theta_1, \theta_1) = (\Theta, \iota)$
let $n = \text{arity}(D)$
for $i \in 1..n$
let $(\Theta_{i+1}, \theta_{i+1}) =$
let $(\Theta', \theta') = \text{unify}(\Delta, \Theta_i, \theta_i(A_i), \theta_i(B_i))$
return $(\Theta', \theta' \circ \theta_i)$
return $(\Theta_{n+1}, \theta_{n+1})$

$\text{unify}(\Delta, \Theta, \forall a.A, \forall b.B) =$
assume fresh c
let $(\Theta_1, \theta') = \text{unify}(\Delta, \Theta, A[c/a], B[c/b])$
assert $c \notin \text{ftv}(\theta')$
return (Θ_1, θ')

$\text{demote}(\star, \Theta, \Delta) = \Theta$

$\text{demote}(\bullet, \cdot, \Delta) = \cdot$

$\text{demote}(\bullet, (\Theta, a : K), \Delta) = \text{demote}(\bullet, \Theta, \Delta), a : \bullet \quad (a \in \Delta)$

$\text{demote}(\bullet, (\Theta, a : K), \Delta) = \text{demote}(\bullet, \Theta, \Delta), a : K \quad (a \notin \Delta)$

Figure 15. Unification Algorithm

5.3 Unification

A crucial ingredient for type inference is unification. The unification algorithm is defined in Figure 15. It is partial in that it either returns a result or fails. Following Leijen [13] we explicitly indicate the successful return of a result X by writing return X . Failure may be either explicit or implicit (in the case that an auxiliary function is undefined).

The algorithm takes a quadruple (Δ, Θ, A, B) of a fixed kind environment Δ , a refined kind environment Θ , and types A and B , such that $\Delta, \Theta \vdash A, B$. It returns a unifier, that is, a pair (Θ', θ) of a new refined kind environment Θ' and a type substitution θ , such that $\Delta \vdash \theta : \Theta \Rightarrow \Theta'$.

A type variable unifies with itself, yielding the identity substitution. Due to the use of explicit kind environments, there is no need for an explicit occurs check to avoid unification of a type variable a with a type A including recursive occurrences of a . Unification of a flexible variable a with a type A implicitly performs an occurs check by checking that the type substituted for a is well-formed in an environment (Δ, Θ_1) that does not contain a . A polymorphic flexible variable unifies with any other type, as is standard. A monomorphic flexible variable only unifies with a type A if A may be *demoted* to a monomorphic type. The auxiliary demote function converts any polymorphic flexible variables in A to monomorphic flexible variables in the refined kind environment. This demotion is sufficient to ensure that further unification cannot subsequently make A polymorphic. Unification of data types is standard, checking that the data type constructors match, and recursing on the substructures. Following Leijen [13], unification of quantified types ensures that forall-bound type variables do not escape their scope by introducing a fresh rigid (skolem) variable and ensuring it does not appear in the free type variables of the substitution.

Theorem 4 (Unification is sound). *If $\Delta, \Theta \vdash A, B : K$ and $\text{unify}(\Delta, \Theta, A, B) = (\Theta', \theta)$ then $\theta(A) = \theta(B)$ and $\Delta \vdash \theta : \Theta \Rightarrow \Theta'$.*

Theorem 5 (Unification is complete and most general). *If $\Delta \vdash \theta : \Theta \Rightarrow \Theta'$ and $\Delta, \Theta \vdash A : K$ and $\Delta, \Theta \vdash B : K$ and $\theta(A) = \theta(B)$, then $\text{unify}(\Delta, \Theta, A, B) = (\Theta'', \theta')$ where there exists θ'' satisfying $\Delta \vdash \theta'' : \Theta'' \Rightarrow \Theta'$ such that $\theta = \theta'' \circ \theta'$.*

5.4 The Inference Algorithm

The type inference algorithm is defined in Figure 16. It is partial in that it either returns a result or fails. The algorithm takes a quadruple $(\Delta, \Theta, \Gamma, M)$ of a fixed kind environment Δ , a refined kind environment Θ , a type environment Γ , and a term M , such that $\Delta; \Theta \vdash \Gamma$. If successful, it returns a triple (Θ', θ, A) of a new refined kind environment Θ' , a type substitution θ , such that $\Delta \vdash \theta : \Theta \Rightarrow \Theta'$, and a type A such that $\Delta, \Theta' \vdash A : \star$.

The algorithm is an extension of algorithm W [3] adapted to use explicit kind environments Δ, Θ . Inferring the type of a frozen variable is just a matter of looking up its type in the type environment. As usual, the type of a plain (unfrozen) variable is inferred by instantiating any polymorphism with fresh type variables. The returned identity type substitution is weakened accordingly. Crucially, the argument type inferred for an unannotated lambda abstraction is monomorphic. If on the other hand the argument type is

```

infer : (PEnv × KEnv × TEnv × Term) → (KEnv × Subst × Type)
infer(Δ, Θ, Γ, [x]) =
  return (Θ, ι, Γ(x))
infer(Δ, Θ, Γ, x) =
  let ∀ā.H = Γ(x)
  assume fresh  $\bar{b}$ 
  return ((Θ,  $\bar{b} : \star$ ), ι, H[ $\bar{b}/\bar{a}$ ])
infer(Δ, Θ, Γ, λx.M) =
  assume fresh a
  let (Θ1, θ[a ↦ S], B) = infer(Δ, (Θ, a : ●), (Γ, x : a), M)
  return (Θ1, θ, S → B)
infer(Δ, Θ, Γ, λ(x : A).M) =
  let (Θ1, θ, B) = infer(Δ, Θ, (Γ, x : A), M)
  return (Θ1, θ, A → B)
infer(Δ, Θ, Γ, M N) =
  let (Θ1, θ1, A') = infer(Δ, Θ, Γ, M)
  let (Θ2, θ2, A) = infer(Δ, Θ1, θ1(Γ), N)
  assume fresh b
  let (Θ3, θ3[b ↦ B]) = unify(Δ, (Θ2, b : ★), θ2(A'), A → b)
  return (Θ3, θ3 ∘ θ2 ∘ θ1, B)
infer(Δ, Θ, Γ, let x = M in N) =
  let (Θ1, θ1, A) = infer(Δ, Θ, Γ, M)
  let Δ' = ftv(θ1) - Δ
  let (Δ'', Δ''') = gen((Δ, Δ'), A, M)
  let Θ'1 = demote(●, Θ1, Δ''')
  let (Θ2, θ2, B) = infer(Δ, Θ'1 - Δ'', θ1(Γ), x : ∀Δ'' . A, N)
  return (Θ2, θ2 ∘ θ1, B)
infer(Δ, Θ, Γ, let (x : A) = M in N) =
  let (Δ', A') = split(A, M)
  let (Θ1, θ1, A1) = infer((Δ, Δ'), Θ, Γ, M)
  let (Θ2, θ'2) = unify((Δ, Δ'), Θ1, A', A1)
  let θ2 = (θ'2 ∘ θ1)
  assert ftv(θ2) # Δ'
  let (Θ3, θ3, B) = infer(Δ, Θ2, (θ2(Γ), x : A), N)
  return (Θ3, θ3 ∘ θ2, B)

```

Figure 16. Type Inference Algorithm

annotated with a type, then we just use that type directly. For applications we use the unification algorithm to check that the function and argument match up. Generalisation is performed for unannotated let-bindings in which the let-binding is a guarded value. For unannotated let-bindings in which the let-binding is not a guarded value, generalisation is suppressed and any ungeneralised flexible type variables are demoted to be monomorphic. When a let-binding is annotated with a type then rather than performing generalisation we use the annotation, taking care to account for any polymorphism that is already present in the inferred type for M

using split, and checking that none of the quantifiers escape by inspecting the codomain of θ_2 .

Theorem 6 (Type inference is sound). *If $\Delta, \Theta \vdash \Gamma$ and $\Delta \Vdash M$ and $\text{infer}(\Delta, \Theta, \Gamma, M) = (\Theta', \theta, A)$ then $\Delta, \Theta'; \theta(\Gamma) \vdash M : A$ and $\Delta \vdash \theta : \Theta \Rightarrow \Theta'$.*

Theorem 7 (Type inference is complete and principal). *Let $\Delta \Vdash M$ and $\Delta, \Theta \vdash \Gamma$. If $\Delta \vdash \theta : \Theta \Rightarrow \Theta'$ and $\Delta, \Theta'; \theta(\Gamma) \vdash M : A$, then $\text{infer}(\Delta, \Theta, \Gamma, M) = (\Theta'', \theta', A')$ where there exists θ'' satisfying $\Delta \vdash \theta'' : \Theta'' \Rightarrow \Theta'$ such that $\theta = \theta'' \circ \theta'$ and $\theta''(A') = A$.*

6 Implementation

We have implemented FreezeML as an extension of Links. This exercise was mostly routine. In the process we addressed several practical concerns and encountered some non-trivial interactions with other features of Links. In order to keep this paper self-contained we avoid concrete Links syntax, but instead illustrate the ideas of the implementation in terms of extensions to the core syntax used in the paper.

In ASCII we render $[x]$ as $\sim x$. For convenience, Links builds in the generalisation $\$$ and instantiation operators $@$.

In practice (in Links and other functional languages), it is often convenient to include a type signature for a function definition rather than annotations on arguments. Thus

$$\begin{aligned} f &: \forall a. A \rightarrow B \rightarrow C \\ f \ x \ y &= M \\ N \end{aligned}$$

is treated as:

$$\text{let } (f : \forall a. A \rightarrow B \rightarrow C) = \lambda(x : A). \lambda(y : B). M \text{ in } N$$

Though x and y are not themselves annotated, A and B may be polymorphic, and may mention a .

Given that FreezeML is explicit about the order of quantifiers, adding support for explicit type application [4] is straightforward. We have implemented this feature in Links.

Links has an implicit subkinding system used for various purposes including classifying base types in order to support language-integrated query [16] and distinguishing between linear and non-linear types in order to support session typing [17]. In plain FreezeML, if we have $\text{poly} : (\forall a. a \rightarrow a) \rightarrow \text{Int} \times \text{Bool}$ and $\text{id} : \forall a. a \rightarrow a$, then we may write $\text{poly} [\text{id}]$. The equivalent in Links also works. However, the type inferred for the identity function in Links is not $\forall a. a \rightarrow a$, but rather $\forall (a : \circ). a \rightarrow a$, where the subkinding constraint \circ captures the property that the argument is used linearly. Given this more refined type for id the term $\text{poly} [\text{id}]$ no longer type-checks. In this particular case one might imagine generating an implicit coercion (a function that promises to use its argument linearly may be soundly treated as a function that may or may not use its argument linearly). In general one has to be careful to be explicit about the kinds of type variables when working with first-class polymorphism.

Similar issues arise from the interaction between first-class polymorphism and Links's effect type system [16].

Existing infrastructure for subkinding in the implementation of Links was helpful for adding support for FreezeML as we exploit it for tracking the monomorphism / polymorphism distinction. However, there is a further subtlety: in FreezeML type variables of monomorphic kind may be instantiated with (though not unified with) polymorphic types; this behaviour differs from that of other kinds in Links.

The Links source language allows the programmer to explicitly distinguish between rigid and flexible type variables. Flexible type variables can be convenient to use as wild-cards during type inference. As a result, type annotations in Links are slightly richer than those admitted by the well-scopedness judgement of Figure 9. It remains to verify the formal properties of the richer system.

7 Related Work

There are many previous attempts to bridge the gap between ML and System F. Some systems employ more expressive types than those of System F; others implement heuristics in the type system to achieve a balance between increased complexity of the system and reducing the number of necessary type annotations; finally, there are systems like ours that eschew such heuristics for the sake of simplifying the type system further. Users then have to state their intentions explicitly, potentially resulting in more verbose programs.

Expressive Types. MLF [11] (sometimes stylised as ML^F) is considered to be the most expressive of the conservative ML extensions so far. MLF achieves its expressiveness by going beyond regular System F types and introducing polymorphically bounded types, though translation from MLF to System F and vice versa remains possible [11, 12]. MLF also extends ML with type annotations on lambda binders. Annotations on binders that are *used* polymorphically are mandatory, since type inference will not guess second-order types. This is required to maintain principal types.

HML [14] is a simplification of MLF. In HML all polymorphic function arguments require annotations. It significantly simplifies the type inference algorithm compared to MLF, though polymorphically bounded types are still used.

Heuristics. HMF [13] contrasts with the above systems in that it only uses regular System F types (disregarding order of quantifiers). Like FreezeML, it only allows principal types for let-bound variables, and type annotations are needed on all polymorphic function parameters. HMF allows both instantiation and generalisation in argument positions, taking n-ary applications into account. The system uses weights to select between less and more polymorphic types. Whole lambda abstractions require an annotation to have a polymorphic return type. Such term annotations are *rigid*, meaning they suppress instantiation and generalisation. As instantiation

is implicit in HMF, rigid annotations can be seen as a means to freeze arbitrary expressions.

Several systems for first-class polymorphism were proposed in the context of the Haskell programming language. These systems include boxy types [27], FPH [28], and GI [25]. The Boxy Types system, used to implement GHC’s `ImpredicativeTypes` extension, was very fragile and thus difficult to use in practice. Similarly, the FPH system – based on MLF – was simpler but still difficult to implement in practice. GI is the latest development in this line of research. Its key ingredient is a heuristic that restricts polymorphic instantiation, based on whether a variable occurs under a type constructor and argument types in an application. Like HMF, it uses System F types, considers n -ary applications for typing, and requires annotations both for polymorphic parameter and return types. However, only top-level type variables may be re-ordered. The authors show how to combine their system with the `OutsideIn(X)` [26] constraint-solving type inference algorithm used by the Glasgow Haskell Compiler. They also report a prototype implementation of GI as an extension to GHC with encouraging experience porting existing Hackage packages that use rank- n polymorphism.

Explicitness. Some early work on first-class polymorphism was based on the observation that polymorphism can be encapsulated inside nominal types [9, 10, 20, 23].

The QML [24] system explicitly distinguishes between polymorphic schemes and quantified types and hence does not use plain System F types. Type schemes are used for ML let-polymorphism and introduced and eliminated implicitly. Quantified types are used for first-class polymorphism, in particular for polymorphic function arguments. Such types must always be introduced and eliminated explicitly, which requires stating the full type and not just instantiating the type variables. All polymorphic instantiations must therefore be made explicitly by annotating terms at call sites. Neither let- nor λ -bound variables can be annotated with a type.

Poly-ML [7] is similar to QML in that it distinguishes two incompatible sorts of polymorphic types. Type schemes arise from standard ML generalisation; (boxed) polymorphic types are introduced using a dedicated syntactic form which requires a type annotation. Boxed polymorphic types are considered to be simple types, meaning that a type variable can be instantiated with a boxed polymorphic type, but not with a type scheme. Terms of a boxed type are not instantiated implicitly, but must be opened explicitly, resulting in instantiation. Unlike QML, the instantiated type is deduced from the context, rather than requiring an annotation.

Unlike FreezeML, Poly-ML supports inferring polymorphic parameter types for unannotated lambdas, but this is limited to situations where the type is unambiguously determined by the context. This is achieved by using *labels*, which track whether polymorphism was guessed or confirmed by a type annotation. Whereas FreezeML has type

annotations on binders, Poly-ML has type annotations on terms and propagates them using the label system.

In Poly-ML, the example $\lambda x.\text{auto } x$ typechecks, guessing a polymorphic type for x ; FreezeML requires a type annotation on x . In FreezeML the program `let id = $\lambda x.x$ in let c = id 3 in auto [id]` typechecks, whereas in Poly-ML a type annotation is required (in order to convert between $\forall a.a \rightarrow a$ and $[\forall a.a \rightarrow a]$). However, Poly-ML could be extended with a new construct for introducing boxed polymorphism without a type annotation, using the principal type instead. With such a change it is possible to translate from FreezeML into this modified version of Poly-ML without inserting any new type annotations (see the extended version of this paper [5]).

The extended version of this paper [5] contains an example-based comparison of FreezeML, GI, MLF, HMF, FPH, and HML.

Instantiation as subsumption. In FreezeML instantiation induces a natural subtyping relation such that $A \leq B$ iff B is an instance of A . In other systems (e.g. [20]) such a subtyping relation applies implicitly to all terms via a subsumption rule. This form of subsumption is fundamentally incompatible with frozen variables, which explicitly suppress instantiation, enabling fine-grained control over exactly where instantiation occurs. Nonetheless, subsumption comes for free on unfrozen variables and potentially elsewhere if one adopts more sophisticated instantiation strategies.

8 Conclusions

In this paper, we have introduced FreezeML as an exercise in language design for reconciling ML type inference with System F-style first-class polymorphism. We have also implemented FreezeML as part of the Links programming language [2], which uses a variant of ML type inference extended with row types, and has a kind system readily adapted to check that inferred function arguments are monotypes.

Directions for future work include extending FreezeML to accommodate features such as higher-kinds, GADTs, and dependent types, as well as exploring different implicit instantiation strategies. It would also be instructive to rework our formal account using the methodology of Gundry et al. [8] and use that as the basis for mechanised soundness and completeness proofs.

Acknowledgments

This work was supported by EPSRC grant EP/K034413/1 ‘From Data Types to Session Types—A Basis for Concurrency and Distribution’, ERC Consolidator Grant Skye (grant number 682315), by an LFCS internship, and by an ISCF Metrology Fellowship grant provided by the UK government’s Department for Business, Energy and Industrial Strategy (BEIS). We are grateful to James McKinna, Didier Rémy, Andreas Rossberg, and Leo White for feedback and to anonymous reviewers for constructive comments.

References

- [1] Dominique Clément, Joëlle Despeyroux, Th. Despeyroux, and Gilles Kahn. 1986. A Simple Applicative Language: Mini-ML. In *LISP and Functional Programming*. 13–27.
- [2] Ezra Cooper, Sam Lindley, Philip Wadler, and Jeremy Yallop. 2006. Links: Web Programming Without Tiers. In *FMCO (Lecture Notes in Computer Science)*, Vol. 4709. Springer, 266–296. <http://links-lang.org/>
- [3] Luís Damas and Robin Milner. 1982. Principal Type-Schemes for Functional Programs. In *POPL*. ACM Press, 207–212.
- [4] Richard A. Eisenberg, Stephanie Weirich, and Hamidhasan G. Ahmed. 2016. Visible Type Application. In *ESOP (Lecture Notes in Computer Science)*, Vol. 9632. Springer, 229–254.
- [5] Frank Emrich, Sam Lindley, Jan Stolarek, James Cheney, and Jonathan Coates. 2019. *FreezeML: Complete and Easy Type Inference for First-Class Polymorphism (extended version)*. Technical Report. arXiv:2004.00396.
- [6] Matthias Felleisen. 1991. On the Expressive Power of Programming Languages. *Sci. Comput. Program.* 17, 1-3 (1991), 35–75.
- [7] Jacques Garrigue and Didier Rémy. 1999. Semi-Explicit First-Class Polymorphism for ML. *Inf. Comput.* 155, 1-2 (1999), 134–169.
- [8] Adam Gundry, Conor McBride, and James McKinna. 2010. Type Inference in Context. In *MSFP@ICFP*. ACM, 43–54.
- [9] Mark P. Jones. 1997. First-class Polymorphism with Type Inference. In *POPL*. ACM Press, 483–496.
- [10] Konstantin Läufer and Martin Odersky. 1994. Polymorphic Type Inference and Abstract Data Types. *ACM Trans. Program. Lang. Syst.* 16, 5 (1994), 1411–1430.
- [11] Didier Le Botlan and Didier Rémy. 2003. ML^F: raising ML to the power of System F. In *ICFP*. ACM, 27–38.
- [12] Daan Leijen. 2007. A type directed translation of MLF to system F. In *ICFP*. ACM, 111–122.
- [13] Daan Leijen. 2008. HMF: simple type inference for first-class polymorphism. In *ICFP*. ACM, 283–294.
- [14] Daan Leijen. 2009. Flexible types: robust type inference for first-class polymorphism. In *POPL*. ACM, 66–77.
- [15] Xavier Leroy and Michel Mauny. 1991. Dynamics in ML. In *FPCA*. Springer, 406–426.
- [16] Sam Lindley and James Cheney. 2012. Row-based effect types for database integration. In *TLDI*. ACM, 91–102.
- [17] Sam Lindley and J Garrett Morris. 2017. Lightweight Functional Session Types. In *Behavioural Types: from Theory to Tools*. River Publishers, 265–286.
- [18] Dale Miller. 1992. Unification Under a Mixed Prefix. *J. Symb. Comput.* 14, 4 (1992), 321–358. [https://doi.org/10.1016/0747-7171\(92\)90011-R](https://doi.org/10.1016/0747-7171(92)90011-R)
- [19] Robin Milner, Mads Tofte, and David MacQueen. 1997. *The Definition of Standard ML (Revised)*. MIT Press.
- [20] Martin Odersky and Konstantin Läufer. 1996. Putting Type Annotations to Work. In *POPL*. ACM Press, 54–67.
- [21] Simon Peyton Jones and Mark Shields. 2002. Lexically scoped type variables. Unpublished. <https://www.microsoft.com/en-us/research/publication/lexically-scoped-type-variables/>.
- [22] Frank Pfenning. 1993. On the Undecidability of Partial Polymorphic Type Reconstruction. *Fundam. Inform.* 19, 1/2 (1993), 185–199.
- [23] Didier Rémy. 1994. Programming Objects with ML-ART, an Extension to ML with Abstract and Record Types. In *TACS (Lecture Notes in Computer Science)*, Vol. 789. Springer, 321–346.
- [24] Claudio V. Russo and Dimitrios Vytiniotis. 2009. QML: Explicit First-class Polymorphism for ML. In *ML*. ACM, 3–14.
- [25] Alejandro Serrano, Jurriaan Hage, Dimitrios Vytiniotis, and Simon Peyton Jones. 2018. Guarded impredicative polymorphism. In *PLDI*. ACM, 783–796.
- [26] Dimitrios Vytiniotis, Simon L. Peyton Jones, Tom Schrijvers, and Martin Sulzmann. 2011. OutsideIn(X) Modular type inference with local assumptions. *J. Funct. Program.* 21, 4-5 (2011), 333–412.
- [27] Dimitrios Vytiniotis, Stephanie Weirich, and Simon L. Peyton Jones. 2006. Boxy types: inference for higher-rank types and impredicativity. In *ICFP*. ACM, 251–262.
- [28] Dimitrios Vytiniotis, Stephanie Weirich, and Simon L. Peyton Jones. 2008. FPH: first-class polymorphism for Haskell. In *ICFP*. ACM, 295–306.
- [29] J. B. Wells. 1994. Typability and Type-Checking in the Second-Order lambda-Calculus are Equivalent and Undecidable. In *LICS*. IEEE Computer Society, 176–185.
- [30] Andrew K. Wright. 1995. Simple Imperative Polymorphism. *Lisp and Symbolic Computation* 8, 4 (1995), 343–355.