

When Subtyping Constraints Liberate^{*†}

A Novel Type Inference Approach for First-Class Polymorphism

LIONEL PARREAUX, HKUST, Hong Kong, China

ALEKSANDER BORUCH-GRUSZECKI, EPFL, Switzerland

ANDONG FAN, HKUST, Hong Kong, China

CHUN YIN CHAU, HKUST, Hong Kong, China

Type inference in the presence of *first-class* or “*impredicative*” *second-order polymorphism* à la System F has been an active research area for several decades, with original works dating back to the end of the 80s. Yet, until now many basic problems remain open, such as how to type check expressions like $(\lambda x. (x\ 123, x\ \text{True}))\ \text{id}$ reliably. We show that a type inference approach based on *multi-bounded polymorphism*, a form of implicit polymorphic subtyping with multiple lower and upper bounds, can help us resolve most of these problems in a uniquely simple and regular way. We define $F_{\{\leq\}}$, a declarative type system derived from the existing theory of implicit coercions by [Cretin and Rémy](#), and we introduce SuperF, a novel algorithm to infer polymorphic multi-bounded $F_{\{\leq\}}$ types while checking user type annotations written in the syntax of System F. We use a recursion-avoiding heuristic to guarantee termination of type inference at the cost of rejecting some valid programs, which thankfully rarely triggers in practice. We show that SuperF is vastly more powerful than all first-class-polymorphic type inference systems proposed so far, significantly advancing the state of the art in type inference for general-purpose programming languages.

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

Additional Key Words and Phrases: type inference, first-class polymorphism, subtyping, constraint solving

1 INTRODUCTION

Consider the function `foo` defined as: `foo f = (f 123, f True)`.

For the past several decades, researchers have been unable to decide which type to infer for this function, should one wish to go beyond the classical Hindley-Milner type inference discipline of ML languages (which yields a unification error here), and should one want to support some form of first-class polymorphism. In particular, no system was proposed where such a term could be given a “satisfactory” type. A satisfactory type for `foo` should let us type check expressions like `foo (fun x → x)` at type $(\text{Int}, \text{Bool})$ and `foo (fun x → Some x)` at type $(\text{Option Int}, \text{Option Bool})$.¹

In this paper, we propose SuperF, a system which can infer the following type for `foo`:

$$\text{foo} : \forall a\ b. ((\text{Int} \rightarrow a) \wedge (\text{Bool} \rightarrow b)) \rightarrow (a, b)$$

which is really just syntax sugar for the following polymorphic type that incorporates several bounds for its extra type parameter `c`:

$$\text{foo} : \forall a\ b\ c\ \{c \leq \text{Int} \rightarrow a, c \leq \text{Bool} \rightarrow b\}. c \rightarrow (a, b)$$

This type is *satisfactory* in that it lets us type check the function call examples above at the desired respective types $(\text{Int}, \text{Bool})$ and $(\text{Option Int}, \text{Option Bool})$.

^{*}The title of this paper took inspiration from that of a talk by Runar Bjarnason: *Constraints Liberate, Liberties Constrain*.

[†]This is a technical report version of the conference paper that appeared in POPL 2024 [[Parreaux et al. 2024](#)].

¹MLsub, which uses Algebraic Subtyping [[Dolan and Mycroft 2017](#)], would infer types $(\text{Int} \vee \text{Bool}, \text{Int} \vee \text{Bool})$ and $(\text{Option}(\text{Int} \vee \text{Bool}), \text{Option}(\text{Int} \vee \text{Bool}))$, which are imprecise and thus unsatisfactory.

Our approach is reminiscent of *bounded polymorphism*, but our type system $F_{\{\leq\}}$ diverges from most research on this topic in the following crucial respects:

- We assume *implicit, erased* polymorphism, meaning that polymorphic types do not need to be introduced and eliminated in terms and do not participate in the term’s runtime semantics. Consequently, we can have subtyping relationships such as $\forall\alpha. \alpha \rightarrow \alpha \leq \text{Int} \rightarrow \text{Int}$, so that any term with the left-hand side type is known to *also* have the right-hand side type.
- We allow lower bounds *as well as* upper bounds, rather than only upper bounds (like in System $F_{<}$) or only lower bounds (like in ML^F), making our system symmetric, but introducing possible inconsistencies between bounds, which need to be dealt with carefully.
- We allow associating an *arbitrary number of bounds* to the same type variable, instead of just one,² reminiscent of the old idea of *constrained types* [Odersky et al. 1999], which was previously studied mostly in the context of ML-style (i.e., not *first-class*) polymorphism.

Fortunately, this combination of type system features was already formally described as part of System F_{cc} by Cretin and Rémy [2014] in their groundbreaking work laying the foundations for general implicit polymorphism with subtyping. This allows us to define $F_{\{\leq\}}$ by translation to F_{cc} and focus on the *algorithmic* aspects of the system in this paper, noting that Cretin and Rémy only presented a *declarative* type system for System F_{cc} and did not investigate type inference.

While SuperF crucially relies on multi-bounded types internally, these types are non-denotable: they are not accessible to users, who are limited to writing type annotations in the syntax of System F, i.e., where polymorphic types cannot have bounds on their quantified type variables. Indeed, to make subtype constraint solving tractable, SuperF internally follows a *polarized* type syntax, whereby bounds are only allowed in *positive* polymorphic types. Annotated types cannot use bounds because these types are both positive *and* negative, since they *provide* a type for the annotated expression but also are *checked* against the expression’s inferred type. This is an atypical property of our system: it means that users cannot always provide explicit type signatures that are as general as the types that are internally inferred by the system. Nevertheless, we argue that in practice users rarely want to use bounds in their type signatures anyway, and that System F types are sufficient for a large number of practical functional programming use cases, including a majority of use cases studied in previous work.

Subtyping for System F types was historically realized in System F_η , which relates polymorphic types by subtyping based on their specificity. F_η has an “ η rule” to subtype function types covariantly in their results and contravariantly in their arguments, which makes System F typing complete with respect to η conversion. Subtyping in F_η is known to be undecidable [Chrzęszcz 1998; Tiuryn and Urzyczyn 1996], and $F_{\{\leq\}}$ is a strict generalization of F_η (i.e., all well-typed System F and F_η terms are well-typed $F_{\{\leq\}}$ terms, but the converse does not hold), so it is natural to conjecture that $F_{\{\leq\}}$ subtyping is also undecidable, making type inference necessarily incomplete.

SuperF is a *terminating* but *incomplete* type inference algorithm for $F_{\{\leq\}}$ that relies on a simple recursion-detection heuristic to raise errors in tricky recursive-looking cases. Without this check, type inference would diverge on terms that exhibit indirectly-recursive structures, such as the Ω combinator ($\text{fun } x \rightarrow x \ x$) ($\text{fun } x \rightarrow x \ x$). Thankfully, we find that this check rarely triggers in practical, real-world examples. On the other hand, even when it succeeds, SuperF does not always infer principal types. This is mainly because it always distributes polymorphic types over function types even though that can sometimes lead to worse outcomes (an example is given as G14 in

²Intuitively, this is equivalent in power to allowing at most one upper bound and at most one lower bound on each type variable together with allowing intersection types in negative positions and union types in positive positions [Parreaux 2020], which is why we could represent the type of `foo` more concisely via an intersection in the example above.

Section 5.5). Whether we always infer principal types in a restricted system without distributivity is currently an open question.

From the programmer’s perspective, the main limitation of SuperF type inference can be understood as: *SuperF never assumes parametrically-polymorphic types for function parameters*. Type annotations must be used when such polymorphic parameters are needed. Nevertheless, SuperF is appreciably more powerful than all other existing systems; we show that it can seamlessly handle most examples previously considered in the literature (a first), and that it does so *without requiring any annotations whatsoever*. Indeed, multi-bounded polymorphic subtyping already covers a lot of additional ground compared to existing unification-based approaches, allowing programs like the one at the beginning of this introduction to type check without having to resort to higher-rank parametric polymorphism where other systems would have to.

Our specific contributions are as follows:

- We describe the main problems of first-class-polymorphic type inference as well as our main ideas to address these problems (Section 2). To the best of our knowledge, although they are quite natural in retrospect, these ideas are novel and have not been developed before.
- We formalize $F_{\{\leq\}}$, a lambda calculus with first-class, implicit, erased, and what we call *multi-bounded* polymorphism (Section 3). We show that the soundness of this system can be obtained by translation into the existing System F_{cc} , whose semantic soundness was mechanically verified in Coq by Cretin [2014].
- We formalize our SuperF type inference system, emphasizing the novel and crucial concept of *subtype extrusion and avoidance* (Section 4). We prove that this type inference system is sound for $F_{\{\leq\}}$ and terminating.
- We present a practical implementation of SuperF (Section 5), which can be tried in a web demo available at <https://hkust-taco.github.io/superf/>. This implementation was evaluated on all examples described in previous work as well as many new and less trivial ones, demonstrating that our system is considerably more powerful than previously proposed approaches. We also show that SuperF infers precise yet concise types for existing ML programs by porting the `List` module from OCaml’s standard library to SuperF.

2 PRESENTATION

We first present our approach from a high-level point of view focused on intuition.

2.1 Motivation

Before we delve into technical details, it is worth considering *why* one might care about first-class polymorphism and why impredicative polymorphic type inference is a useful addition to a functional programming language’s type system.

Early on, Peyton Jones et al. [2007] presented many convincing use cases for *higher-rank polymorphism*, a restricted form of first-class polymorphism. These examples included data structure fusion, encapsulation of state and other effects, existential type encodings, generic programming, folding and mapping functions respecting data structure invariants, internal type class desugaring, etc. Nowadays, higher-rank polymorphism is used pervasively throughout the Haskell ecosystem.

Peyton Jones et al. [2007, §3.4] also presented several use cases for *impredicative* polymorphic type inference (i.e., unrestricted first-class polymorphism, as studied in this paper). Their examples rely on the observation that without impredicative polymorphism, several abstractions that work

well in the monomorphic case start failing in the polymorphic case, requiring tedious manual specialization. For instance, consider the following definitions:

```
revapp : ∀ a. a → (a → b) → b      |      poly : (∀ v. v → v) → (Int, Bool)
revapp x f = f x                       |      poly f = (f 3, f True)
```

While both of these functions fit into classical higher-rank polymorphism, applying one to the other, as in `(revapp (fun x → x) poly)`, requires *impredicatively* instantiating the type variable `a` to type `∀ v. v → v`. A similar thing happens with the fix-point combinator, whose usual type, `fix : (a → a) → a`, cannot be applied to a polymorphic function without impredicative polymorphism.

As another example, recent versions of the Haskell language started supporting abstracting over record fields through a special type class so that, for example, it became possible to write functions working over any input type that contains a field named "size" of type `a`. Allowing this useful abstraction infrastructure to extend to the setting of records with *polymorphic* fields requires impredicative polymorphism, as it requires instantiating `a` in our example with a polymorphic type.

Altogether, supporting first-class polymorphism in its "full glory" has many important real-world benefits, which also explains why the Haskell community has spent over a decade trying to add such generalized support to the Glasgow Haskell Compiler (GHC) [Serrano et al. 2020, 2018; Vytiniotis et al. 2006] and why even languages like Scala 3 are starting to add the feature.³

The solution we propose in this paper is a general one: SuperF supports higher-rank polymorphism as a side effect of its support for unrestricted first-class polymorphism, and all classical examples of predicative higher-rank polymorphism are still supported out of the box in SuperF.

2.2 Parametricity and Subtyping

Recall the definition of function `foo` shown in the introduction: `foo f = (f 123, f True)`

To be polymorphic or not to be. Clearly, function parameter `f` should be polymorphic one way or another in order for `foo` to be well-typed, because `f` is applied to arguments of completely different types `Int` and `Bool`. So an intuitive type for `foo` could have the following form:

```
foo : (∀ a. a → ?) → (?, ?)
```

This leads us to a first problem: what type should the parameter function `f` be expected to *return*? The simplest answer that comes to mind would be to have it return the same type as its input:

```
foo : (∀ a. a → a) → (Int, Bool)
```

While this is possible, it is obviously not general enough — what if the caller wanted to pass to `foo` the function `fun x → Some x`? So far, most previous state-of-the-art approaches to first-class polymorphism already throw in the towel and ask the user to provide an explicit type annotation for `f`. The user could annotate `f` as `∀ a. a → a` or they could annotate it as `∀ a. a → Option a`, or as `∀ a. a → List a`, etc., etc., depending on their intended use of `foo`. Moreover, these parametric polymorphism considerations are not *stable* in the face of small changes. Consider:

```
foo1 f = (f E0, f E1)
```

where `E0` and `E1` are two arbitrary expressions. If `E0` and `E1` happen to be typeable at the same type `τ`, `f` no longer needs to be polymorphic, which allows calling `foo1` with *less powerful* non-polymorphic function arguments of type `τ → R` for some `R`. This may in turn be needed later, depending on the uses of `foo1`. For example, taking `E0 = E1 = fun x → x`, we have:

```
foo2 f = (f (fun x → x), f (fun x → x))
foo2 : ∀ a b. ((a → a) → b) → (b, b)
```

³See *polymorphic function types*: <https://docs.scala-lang.org/scala3/reference/new-types/polymorphic-function-types.html>.

This typing of `foo2` would let us write expressions like `foo2 (fun f → f 1 > 0)`, of type $(\text{Bool}, \text{Bool})$. But the parameter-polymorphic type is still on the table, and is neither more nor less general:

$$\text{foo2} : (\forall a. a \rightarrow a) \rightarrow ((\forall a. a \rightarrow a), (\forall a. a \rightarrow a))$$

which would allow us to type `foo2 id` as $((\forall a. a \rightarrow a), (\forall a. a \rightarrow a))$. The fundamental *instability* of this kind of reasoning gets even worse whenever we pass to `f` argument that *may or may not* be typed identically depending on polymorphism decisions made earlier in type inference, as in:

$$\text{foo3 } f = (f (\text{fun } (x, y) \rightarrow (y, x)), f (\text{fun } (x, y) \rightarrow (x, y)))$$

Above, we could type both arguments to `f` at the same type $\forall a. (a, a) \rightarrow (a, a)$, or we could type them at unrelated types $\forall a b. (a, b) \rightarrow (b, a)$ and $\forall a b. (a, b) \rightarrow (a, b)$ respectively, and again there is no most general choice. Many types seem possible, and there is no obvious way to proceed. This leads us to formulate another major problem: when *should* inferred types be polymorphic and when should they *not* be? In practice, these two problems make the job of a type inference engine for System F unpleasantly cumbersome, and the need for a form of subtyping is felt [Rémy 2005].

Subtype polymorphism and intersection types. Our failure to find a most general type for `foo` was due to a belief that `f` should be assigned a *parametric* polymorphic type. Parametric polymorphism is a form of *infinitary* polymorphism [Aiken and Wimmers 1993; Leivant 1990] – that is, a function `f` of type $\forall a. a \rightarrow a$ can be seen as a function having at the same time all the types $f : \text{Int} \rightarrow \text{Int}; f : \text{Bool} \rightarrow \text{Bool}; f : \text{Option Int} \rightarrow \text{Option Int}$; etc. In type theory, we can express such *overloaded* function types through *intersection types* (\wedge), as in:

$$f : (\text{Int} \rightarrow \text{Int}) \wedge (\text{Bool} \rightarrow \text{Bool}) \wedge (\text{Option Int} \rightarrow \text{Option Int}) \wedge \dots$$

It should now become clear that assigning to the `f` parameter of `foo` an infinitary intersection type is definitely “overkill”. After all, `f` is only applied to `Int` and `Bool` arguments, so a type such as $(\text{Int} \rightarrow \text{S}) \wedge (\text{Bool} \rightarrow \text{T})$ should suffice. Since there are no constraints on `s` and `t`, these can be taken to be type variables of `foo`, yielding the type:

$$\text{foo} : \forall a b. ((\text{Int} \rightarrow a) \wedge (\text{Bool} \rightarrow b)) \rightarrow (a, b)$$

This type (call it τ_0) happens to be a *principal* or *most general* type of `foo`,⁴ meaning that all other types that can be assigned to `foo` can be obtained by widening τ_0 through subtyping.

PROOF. Assume that `foo` has type τ . Since `foo` is a function, τ must be a supertype of $\tau_1 \rightarrow \tau_2$ for some τ_1, τ_2 , i.e., $\tau_1 \rightarrow \tau_2 \leq \tau$ (τ is not necessarily a function type; for instance, it could be \top). Since `foo` returns a pair, we must have $(\sigma_1, \sigma_2) \leq \tau_2$ for some σ_1, σ_2 . Since `f` is applied to an argument of type `Int`, resulting in type σ_1 , we must have $\tau_1 \leq \text{Int} \rightarrow \sigma_1$. Similarly, we must have $\tau_1 \leq \text{Bool} \rightarrow \sigma_2$. We need to show that $\tau_0 = \forall \alpha \beta. ((\text{Int} \rightarrow \alpha) \wedge (\text{Bool} \rightarrow \beta)) \rightarrow (\alpha, \beta) \leq \tau$, a relationship which is implied by $\tau_0 \leq \tau_1 \rightarrow \tau_2$. To show the latter, instantiate the left-hand side by taking $\alpha = \sigma_1$ and $\beta = \sigma_2$, resulting in $((\text{Int} \rightarrow \sigma_1) \wedge (\text{Bool} \rightarrow \sigma_2)) \rightarrow (\sigma_1, \sigma_2) \leq \tau_1 \rightarrow \tau_2$, which works because that is equivalent to: • $\tau_1 \leq (\text{Int} \rightarrow \sigma_1) \wedge (\text{Bool} \rightarrow \sigma_2)$, i.e., $\tau_1 \leq \text{Int} \rightarrow \sigma_1$ and $\tau_1 \leq \text{Bool} \rightarrow \sigma_2$; and • $(\sigma_1, \sigma_2) \leq \tau_2 = (\sigma_1, \sigma_2)$. \square

Notice that τ_0 does not even involve any higher-rank parametric polymorphism! (Instead, it has a higher-rank intersection type.) Yet, a function of this type can be called with parametrically-polymorphic arguments, such as functions of type $\tau_{\text{id}} = \forall c. c \rightarrow c$ which should not come as a surprise, as we saw that τ_{id} stands for $(\text{Int} \rightarrow \text{Int}) \wedge (\text{Bool} \rightarrow \text{Bool}) \wedge \dots$, which is clearly a subtype of $(\text{Int} \rightarrow \text{Int}) \wedge (\text{Bool} \rightarrow \text{Bool})$. Similarly, we can now also call `foo` with argument `fun x → Some x`, of type $\forall c. c \rightarrow \text{Option } c$, which yields a result of type $(\text{Option Int}, \text{Option Bool})$. As a less semantic and more syntactic way to see why the latter works, instantiate `a` to `Option Int` and `b` to `Option Bool`

⁴SuperF does not always infer principal types when they exist, but it does so in many simple cases like this.

in the type of `foo`, and notice that $\forall c. c \rightarrow \text{Option } c$ is indeed a subtype of $(\text{Int} \rightarrow \text{Option Int}) \wedge (\text{Bool} \rightarrow \text{Option Bool})$ because that is the same as saying that $\forall c. c \rightarrow \text{Option } c$ is a subtype of $\text{Int} \rightarrow \text{Option Int}$ *and also* a subtype of $\text{Bool} \rightarrow \text{Option Bool}$, which is true as in each cases we can instantiate the polymorphic type by substituting c accordingly. Furthermore, a function of type τ_0 can *also* be called with monomorphic arguments, of types such as $\top \rightarrow \text{Str}$ (where $\text{Int} \leq \top$ and $\text{Bool} \leq \top$). So it is more general than all System F types we have previously considered for it.

This idea generalizes smoothly to the abstracted example `foo1` mentioned earlier:

$$\text{foo1 } f = (f \ E0, f \ E1) \quad | \quad \text{foo1} : ((T0 \rightarrow a) \wedge (T1 \rightarrow b)) \rightarrow (a, b)$$

where $T0$ and $T1$ are the inferred types of respectively $E0$ and $E1$. Thanks to subtyping, whether the types chosen for $E0$ and $E1$ happen to be compatible or not has become irrelevant. This effectively removes the “discontinuity” in reasoning between these two very similar cases: we no longer have to reason about awkward cases where different, incomparable types suddenly become possible due to small decisions made when inferring the types of subterms.

All in all, picking a *parametrically*-polymorphic type for f makes it more powerful than strictly necessary, and this is a problem because that type occurs *negatively* (i.e., in *input* position), meaning that its excessive power leads to a *weaker* overall type, preventing the discovery of a sufficiently general one. All we had to do was to abandon the infinitary context of parametric polymorphism and assume the finitary context of intersection-based subtype polymorphism instead.

Union types. A dual phenomenon happens with union types in positive positions. Unions give us a natural least upper bound between arbitrary types, even when these types have very different shapes and structures, while other approaches would have to approximate each type to a common and often less precise super type, a process that may not always have good solutions. Consider:

$$\text{bar } x = \text{if } \langle \text{condition} \rangle \text{ then } x \text{ else } \text{id}$$

where `id` has type $\forall a. a \rightarrow a$. Because x and `id` flow together into the result type, most previous approaches would try to *unify* their respective types into a unique result type. But there is no best way of performing such unification: we could make x have the same polymorphic type as `id`, but this would prevent legitimate uses of `bar` such as `bar succ 0` (where `succ`: $\text{Int} \rightarrow \text{Int}$), which could be typed at Int if we instantiated the type of `id` to $\text{Int} \rightarrow \text{Int}$ while typing `bar`. In SuperF, we have:

$$\text{bar} : \forall a. a \rightarrow (a \vee \forall b. b \rightarrow b)$$

which our implementation automatically simplifies by distributing the inner quantifier out:

$$\text{bar} : \forall a b. a \rightarrow (a \vee (b \rightarrow b))$$

and which does not impose any premature decisions on the polymorphism of the x parameter and on the instantiation of `id`’s type. SuperF correctly infers type $\forall b. (\text{Int} \rightarrow \text{Int}) \vee (b \rightarrow b)$ for `bar succ`, which our implementation simplifies to $\text{Int} \rightarrow \text{Int}$ (because $\forall b. b \rightarrow b$ is a subtype of $\text{Int} \rightarrow \text{Int}$).

Going beyond System F. System F only allows expressing the infinitary form of polymorphism (i.e., parametric polymorphism). We prefer inferring types in a slightly more powerful language which smoothes out the rugged edges and discontinuities of System F’s type language. The idea is not new: it goes back at least to ML^F , which used bounded quantification to allow a notion of subsumption (referred to as *generic instance*) to be used in inferred types, so as to abstract over the different possible generalization choices made by users of polymorphic functions. Yet, ML^F still required type annotations in many places, notably for all parameters used polymorphically, and it could not type check the `foo` function variations shown above (except for `foo2`) because it only allowed *one lower bound* per type variable, instead of multiple lower and upper bounds. So we argue that ML^F did not go far enough, and that what we actually need is full-blown polymorphic subtyping and full-blown multi-bounded polymorphism. Because it is a *superset* of System F, we call the type inference system presented in this paper *SuperF*.

On the fluidity of subtyping. Our declarative type system $F_{\{\leq\}}$ allows specifying bounded polymorphic types and unions in arbitrary positions, while the SuperF subset of $F_{\{\leq\}}$, used for type inference, only allows such types in positive positions. Moreover, SuperF types themselves are a superset of System F types. Here is an analogy to justify this design and to explain the ease of our approach compared to the difficulty of inferring pure System F types: one can think of System F types as integer numbers, SuperF types as real numbers, and $F_{\{\leq\}}$ types as imaginary numbers. It is much easier to find *real* solutions to polynomials than it is to find *integer* solutions to them, i.e., diophantine equations. Making an analogy between the problem of finding valid type assignments and the problem of finding polynomial roots, the ruggedness of System F types (integer numbers) forces one into making awkward discrete choices, whereas the fluidity of SuperF subtyping (real numbers) smoothes over these choices, letting one pick most general solutions to sub-problem in a way that composes well. But the best way of *explaining* how to find SuperF type solutions to type inference problems is through the more general framework of $F_{\{\leq\}}$ (imaginary numbers), which generalizes the subtyping relation (numeric domain) further yet.

2.3 SuperF and Implicit Multi-Bounded Polymorphism

We saw that intersections were better at expressing *requirements* on the types of *inputs*, since making these requirements as weak as possible conversely strengthens the overall type. On the other hand, in this paper we are not interested in using intersections in *positive* (output) position, although such intersections can be used to represent various features, such as ad-hoc function overloading. For example, in some programming languages like TypeScript and in semantic subtyping approaches [Frisch et al. 2002], one can write a definition like `dup x = (x, x)` and assign it *several* type signatures, such as $\text{Int} \rightarrow (\text{Int}, \text{Int})$ and $\text{Bool} \rightarrow (\text{Bool}, \text{Bool})$, with the effect of giving `dup` the combined intersection type $(\text{Int} \rightarrow (\text{Int}, \text{Int})) \wedge (\text{Bool} \rightarrow (\text{Bool}, \text{Bool}))$. By contrast, in SuperF, we only assign `dup` the more general type $\forall a. a \rightarrow (a, a)$, which is fine to do in output position as this makes strictly *more* programs type check than picking a specific intersection.⁵

Similarly, SuperF does not support unions in negative positions, though these can be used to represent a structural form of algebraic data types [Castagna et al. 2016; Parreaux and Chau 2022].

Polarized unions & intersection desugaring. When unions are restricted to positive positions and intersections to negative ones, subtype inference remains simple and tractable [Dolan and Mycroft 2017]. In fact, this restricted use of unions and intersections turn out to be equivalent [Parreaux 2020] to a form of bounded polymorphism with both upper and lower bounds on type variables (*multi-bounded polymorphism*). Consider `baz` below and two *equivalent* types for it:

```
bar f x = if f x then f else fun x → x
-- using unions and intersections:
bar : ∀ a b. (a ∧ (b → Bool)) → b → (a ∨ ∀ d. d → d)           -- (A)
-- using multi-bounded polymorphism:
bar : ∀ a b c {a ≤ b → Bool, c ≥ a, c ≥ ∀ d. d → d}. a → b → c  -- (B)
```

PROOF. To show that (A) and (B) are equivalent, we show subtyping in both directions. To show subtyping between two polymorphic types, we: (1) assume the right-hand side type variable bounds; (2) instantiate the type variables on the left-hand side type to types that satisfy the left-hand side type variable bounds; and (3) show subtyping between the underlying bodies of the two polymorphic types. In our case, to show (A) \leq (B), we pick $a = a$ and $b = b$ and thus have to show $(a \wedge (b \rightarrow \text{Bool})) \rightarrow b \rightarrow (a \vee \forall d. d \rightarrow d) \leq a \rightarrow b \rightarrow c$ assuming $a \leq b \rightarrow \text{Bool}$, $c \geq a$, and $c \geq \forall d. d \rightarrow d$, which decomposes by function co- and contravariance to: $\bullet a \leq a \wedge (b \rightarrow \text{Bool})$, i.e., $\circ a \leq a$

⁵Note that not all intersection-overloaded function type signatures can be represented using pure parametric polymorphism. For instance, in TypeScript it is possible to define functions with types such as $(\text{Int} \rightarrow \text{Bool}) \wedge (\text{Bool} \rightarrow \text{Int})$.

(immediate); and $\circ a \leq b \rightarrow \text{Bool}$ (by assumption); $\bullet b \leq b$ (immediate) $\bullet a \vee \forall d. d \rightarrow d \leq c$, i.e., $\circ a \leq c$ (by assumption); and $\circ \forall d. d \rightarrow d \leq c$ (by assumption). To show the other direction $(B) \leq (A)$, we pick $a = a \wedge (b \rightarrow \text{Bool})$, $b = b$, and $c = a \vee \forall d. d \rightarrow d$, which satisfies all the left-hand side type variable bounds and syntactically results in the desired type. \square

Using bounds internally is more convenient than directly using unions and intersections because (1) it facilitates both the implementation and formal specification of type inference; and (2) it allows avoiding some repetition in inferred types, as identical requirements attached to different occurrences of a type variables can be shared as part of a single bound. (For example, $\forall c \{c \leq \forall a. a \rightarrow a\}. (c, c) \rightarrow c$ is more concise than the equivalent $\forall c. (c \wedge (\forall a. a \rightarrow a), c \wedge (\forall a. a \rightarrow a)) \rightarrow c$.) So in this paper we focus on *pure* multi-bounded polymorphism and regard unions and intersections as *syntactic sugar* for the former. Nevertheless, it is often useful to *think of* and *describe* bounded polymorphic types in terms of unions and intersections, which are often more intuitive.

Expressiveness. Extending System F with multi-bounded polymorphism is not insignificant and does bring additional expressiveness to the language. Appendix A.2 presents an example program which is not typeable in System F but is typeable in F_ω , and for which SuperF infers a principal type without the help of any annotations. Note that SuperF is not limited in the ranks of the polymorphic types it infers. For instance, Appendix A.3 exemplifies SuperF type inference for rank 3.

2.4 Type Inference Approach

We now describe type inference in SuperF.

Subtyping over unification. It was historically recognized that subtype inference is notoriously harder than traditional unification-based type inference as in ML. Therefore, it was natural to expect that the difficulty of subtype inference would *compound* with that of first-class-polymorphic type inference, which is probably why the vast majority of previous work stuck with unification — indeed, even the venerable ML^F system, which features *lower-bounded* polymorphism and therefore has a core notion of subtyping (based on generic instances), uses first-order *unification* as the workhorse of its type inference engine. We believe that this is fundamentally a mistake, and that subtype inference in fact *alleviates* most of the intricacies of unification-based approaches to first-class-polymorphic type inference. Indeed, subtyping is better at tracking the *flow of values* through a program, which becomes important in this context.⁶ We argue that ML can get away with approximating data flows in a way that creates undue relationships between unrelated types because its core type system is truly simplistic (essentially simply-typed lambda calculus). But as soon as we move to the first-class-polymorphic setting, such approximations induce all sorts of problems and situations where one needs to make premature polymorphism choices.

Main ideas of SuperF type inference. Generally, our main ideas are as follows:

- Never *infer* parametrically-polymorphic types in negative position, since it is often enough to infer intersection types there.⁷
- Segregate the syntax of types between positive and negative types, similar to the approaches of Dolan and Mycroft [2017]; Pottier [1998], among others.
- Let users provide type annotations that include System F-style polymorphic types, and make sure we can *check* our inferred types *against* these types.

⁶The usefulness of combining subtype-based flow analysis with polymorphism was already noticed by the turn of the century in the context of control-flow analysis, e.g., by Faxén [1997]; Rehof and Fähndrich [2001]; Smith and Wang [2000].

⁷This insight was also leveraged in earlier work by Jim [2000], in a different type inference setting.

- Only infer polymorphic types for lambda abstractions, as we are mostly interested in polymorphic functions, and subtyping with distributivity usually alleviates the need for other forms of polymorphic values in a call-by-value system.⁸ (This point is minor and could easily be relaxed.)
- Offload generated subtyping constraints to a dedicated subtype constraint solver (described in Section 2.5); solve these constraints down to type variable bounds upon generalizing the types of lambda expressions.

The approach described above essentially allows us to infer types for many well-typed System F terms commonly encountered in functional programs, as well as for some terms that are well-typed in our system but ill-typed in System F.

Checking restricted user annotations. While SuperF internally supports multi-bounded polymorphism in positive positions, which allows type inference to use subtyping to delay polymorphism choices, we believe most programmers rarely think in terms of such bounded types. To make checking programmer type annotations against inferred types practical, we *restrict the language of type annotations* to System F types, meaning that users may specify type signatures containing first-class-polymorphic types, but these may *not* specify bounds. Indeed, if we allowed bounds on user-provided type signatures, we would encounter much more difficult constraints to solve – e.g., consider `foo (add : (Int → Int) ∧ (Str → Str)) : τ = ...` where the type annotations are equivalent to giving `foo` type $\forall a \{a \leq \text{Int} \rightarrow \text{Int}, a \leq \text{Str} \rightarrow \text{Str}\}. a \rightarrow \tau$. While type checking the body of `foo` and its uses of the `add` parameter, we would need to deal with constraints of the form $a \leq \sigma$ where a is a rigid type variable with upper bounds; which is equally as hard as dealing with overloading constrains like $(\text{Int} \rightarrow \text{Int}) \wedge (\text{Str} \rightarrow \text{Str}) \leq \sigma$. We are not aware of any constraint-solving approaches able to deal with such types without having to either approximate them⁹ or rely on backtracking, which we argue makes type inference impractical in a real-world context.

Subtype Extrusion. A major difficulty of inferring subtypes in the context of *nested* polymorphism – whether first-class or ML-style – is to avoid mixing up type variables coming from different polymorphism levels. For example, consider typing an expression of the form `fun x → E (fun y → x (y, y))` or equivalently `fun x → let f y = x (y, y) in E f`, where E is some arbitrary subexpression. The variable y introduced by the inner argument function will have type β , for some fresh type variable β , which will be generalized locally as part of the overall type of the inner function, yielding a type of the form $\forall \beta \dots \{ \dots \}. \beta \rightarrow \tau$, assuming $x (y, y)$ returns type τ . Because β is generalized in this nested local function, it is *more polymorphic* than the type of x (call it α), so leaking β into a bound of α as in $\alpha \leq (\beta, \beta) \rightarrow \tau$ would be unsound. To see why, consider that β may later be *instantiated* into several unrelated types as part of its use in the body of E , by which point we would lose the connection between these specific instances and α .

We solve this problem by *extruding* types that are *too polymorphic*, like (β, β) , into some less polymorphic types like (β', β') , where β' is added to the same quantifier as α . Crucially, the extrusion is made to be an *approximation* of the original type by adding the bound $\beta \leq \beta'$ onto the original β .¹⁰ This way, whenever β is later instantiated in E – say, successively to `Int` and `Bool` – then by transitivity we will get constraints $\text{Int} \leq \beta'$ and $\text{Bool} \leq \beta'$, which will ensure that

⁸An example program that is rejected because we only generalize lambdas is given in a Section 2.7 footnote.

⁹MLsub by Dolan and Mycroft [2017] achieves subsumption checking in the presence of such constraints by *under-approximating intersection types*, so for example it treats $(\tau_1 \rightarrow \tau_2) \wedge (\sigma_1 \rightarrow \sigma_2)$ as equivalent to $(\tau_1 \vee \sigma_1) \rightarrow (\tau_2 \wedge \sigma_2)$. This simplifying assumption is *unsound* in the presence of first-class polymorphism, so their technique cannot be used in SuperF. For instance, consider that `id` can be typed as $\forall \alpha. \alpha \rightarrow \alpha$, which is a subtype of $(\text{Int} \rightarrow \text{Int}) \wedge (\text{Str} \rightarrow \text{Str})$, which according to MLsub is a subtype of $(\text{Int} \vee \text{Str}) \rightarrow (\text{Int} \wedge \text{Str})$ – but this type is clearly not a valid type for `id`.

¹⁰If β also occurred negatively in the extruded type, we would also add a *lower* approximant β'' where $\beta'' \leq \beta$ and $\beta'' \leq \beta'$.

the final type inferred for x will be essentially $(\text{Int} \vee \text{Bool}, \text{Int} \vee \text{Bool}) \rightarrow \tau$, the type of a function that accounts for the fact that *both* integers *and* booleans may flow into its parameters.

The original idea of *subtype extrusion* is due to Parreaux [2020], who informally presented a level-based algorithm for subtype inference with nested polymorphism. This approach was itself inspired by the older level-based unification algorithm implemented for the Caml and later OCaml compilers [Kiselyov 2013; Pottier and Rémy 2005], which finds its origin in the work of Rémy [1992], who initially called these polymorphism levels *ranks*, or “degrés” in French [Rémy 1990]. In this paper, we formalize a simplified version of Parreaux’s extrusion algorithm, which does not use explicit levels but can be implemented with levels as an optimization (our SuperF implementation does so).

2.5 Subtype Constraint Solving

The constraint-solving approach of SuperF is relatively straightforward, which we consider to be one of the main contributions of this paper.

First, we decompose concrete type constraints following the structure of types, for example decomposing $(\tau_1, \tau_2) \leq (\sigma_1, \sigma_2)$ into $\tau_1 \leq \sigma_1$ and $\tau_2 \leq \sigma_2$, and decomposing $\tau_1 \rightarrow \tau_2 \leq \sigma_1 \rightarrow \sigma_2$, by following function parameter contravariance and function result covariance, into $\sigma_1 \leq \tau_1$ and $\tau_2 \leq \sigma_2$. At the same time, we compute the transitive closure of type variable bounds, meaning that when constraining $\alpha \leq \tau$ we also constrain $\sigma \leq \tau$ for each existing lower bound σ of α , and symmetrically for $\tau' \leq \alpha$ and upper bounds. For polymorphic types appearing on the left, of the form $\forall \bar{\alpha} \{ \bar{B} \}. \tau \leq \sigma$, all we need to do is (1) instantiate all quantified type variables $\bar{\alpha}$ to fresh type variables $\bar{\alpha}'$ in τ and B ; (2) solve the bounds $[\alpha \mapsto \alpha'] \bar{B}$ as constraints to make sure they hold; and (3) continue by constraining $[\alpha \mapsto \alpha'] \tau \leq \sigma$.

Type Avoidance. when polymorphic types are found on the *right*, as in $\tau \leq \forall \alpha. \sigma$, we need to turn the α type variable into a so-called *skolem*, i.e., a type variable on which no bounds can be added (because we are not allowed to assume anything about this type). We continue by constraining $\tau \leq \sigma$, during which we may end up trying to add to existing (outer) type variables bounds that refer to this α skolem, which makes no sense outside of the corresponding polymorphic type. Therefore, we need to *approximate* these bounds until they no longer refer to α . A simple solution, which we settle on in this paper, is to *widen* all positive occurrences of α to \top and to *narrow* all its negative occurrences to \perp . As a concrete example, consider passing some argument of type $\alpha \rightarrow \beta$ to a function of type $(\forall \gamma. \gamma \rightarrow \gamma) \rightarrow \tau$, where α and β are some type variables coming from the outside. This leads to $\alpha \rightarrow \beta \leq (\forall \gamma. \gamma \rightarrow \gamma)$ and to constraining $\gamma_{\text{sk.}} \leq \alpha$ and $\beta \leq \gamma_{\text{sk.}}$ where $\gamma_{\text{sk.}}$ is the skolem, resulting in overall inferred bounds $\{\top \leq \alpha, \beta \leq \perp\}$. The argument function, whose type must subtype $\alpha \rightarrow \beta$, would thus essentially be restricted to non-terminating computations (of type \perp) or delayed non-terminating computations like `fun _ → failwith "oops"` (of type $\top \rightarrow \perp$).

Delaying instantiation. Instantiating polymorphic types found on the left *too early* and those found on the right *too late* may lead to unnecessary failures [Zhao et al. 2019]. For example, in $\forall \alpha. \alpha \rightarrow \alpha \leq \forall \beta. \beta \rightarrow \beta$, if we instantiate the left-hand-side first, we proceed to constraining $\alpha \rightarrow \alpha \leq \beta_{\text{sk.}} \rightarrow \beta_{\text{sk.}}$, which leads to bounds $\beta_{\text{sk.}} \leq \alpha \leq \beta_{\text{sk.}}$, and we end up with type-avoided bounds $\top \leq \alpha \leq \perp$ on α , which are inconsistent (thus failing the process). But instantiating the right-hand-side first works out because since α is then a *locally-instantiated variable*, it is *allowed* to refer to $\beta_{\text{sk.}}$, and we can forget about both α and $\beta_{\text{sk.}}$ after the constraining is done – that is, in general, unless α leaked into an outer constraint while constraining, in which case we would have to keep an extruded version of it around, which would fail for the same reason as above.

Distributivity. It is sometimes possible to either *delay* the instantiation or *rush* the skolemization of polymorphic types by using the *distributivity* property, which states that $\forall \alpha. \tau \rightarrow \sigma$ is *equivalent*

by subtyping to $\tau \rightarrow \forall\alpha. \sigma$ when α does not occur in τ (this is called *deep skolemization* by Peyton Jones et al. [2007]). For example, to successfully constrain $\forall\alpha. \text{Int} \rightarrow \alpha \rightarrow \alpha \leq \text{Int} \rightarrow \forall\beta. \beta \rightarrow \beta$, we can either distribute the left-hand side, resulting in the equivalent constraint $\text{Int} \rightarrow \forall\alpha. \alpha \rightarrow \alpha \leq \text{Int} \rightarrow \forall\beta. \beta \rightarrow \beta$, which succeeds, or distribute the right-hand side, resulting in $\forall\alpha. \text{Int} \rightarrow \alpha \rightarrow \alpha \leq \forall\beta. \text{Int} \rightarrow \beta \rightarrow \beta$, which also succeeds.¹¹ This is because in both cases, after distributing, we are able to skolemize the right-hand side before instantiating the left-hand side. Therefore, we try to make use of distributivity whenever possible in SuperF.

Cycle checking. Our implementation does not support recursive types, so it rejects cyclic constraints such as $\alpha \leq \text{Int} \rightarrow \alpha$, using a form of subtyping-aware “occurs-check”. This check prevents typing some programs that could have meaningful types even without recursive types, such as `let rec f a = f`, which can be typed as $\top \rightarrow \top$, as well as $\top \rightarrow \top \rightarrow \top$, as well as $\top \rightarrow \top \rightarrow \top \rightarrow \top$, etc. and in general $\mu\alpha. (\top \rightarrow \alpha)$ where μ is the recursive type binder. However, the loss of expressiveness is not very significant; indeed, ML programmers are already used to “ill-formed” recursive definitions like these being rejected, and their use is of limited interest when more usual forms of type-level recursion can be achieved through standard data type definitions.

Unfortunately, even when we reject cyclic constraints, the subtype constraining process described thus far may still not terminate on inputs involving *indirect* forms of recursion, such as the standard $\Omega = (\lambda x. x x) (\lambda x. x x)$ term and other recursion combinators like $Y f = (\lambda x. f (x x)) (\lambda x. f (x x))$. Therefore, we propose the use of a *heuristic* which we refer to as the “*suspiciously recursive-looking criterion*” (SRLC). Constraints that fail the SRLC are simply rejected conservatively, terminating type inference with a failure.

2.6 Suspiciously Recursive-Looking Criterion (SRLC)

The core challenge for ensuring the termination of type inference is that constraint solving may end up comparing an infinite number of distinct types, which are generated by instantiating and skolemizing universal types, as well as copying types during subtype extrusion.

Going back to the roots. Our idea is to assign what we call a “*root*” to each type and to make sure that all types ever created during constraint solving only have a finite number of roots. We abort constraint solving with an error upon comparing a pair of types with the same roots as another pair of types currently being compared, which effectively bounds the depth of recursive constraining calls. The root of a type variable α is:

- α itself if α already existed at the start of the constraint-solving run;
- β if α was created by instantiating, skolemizing, or extruding a type variable with root β during the current constraint-solving run.

The roots of other types are formed by substituting all type variables with their respective roots.

Termination argument. Because the only way to create new types during constraint solving is to substitute the type variables of existing types with fresh type variables (which by construction yields a type with the same root as the original) or with \top and \perp , it is clear that the number of roots we will ever reach is finite, ensuring the termination of constraint solving.

Practicality. We found that in practice, the SLRC seldom gets in the way of typing correct terms, the only exception being terms that make use of self-application (such as G9 in Table 3, presented in Section 5.5). Crucially, we assume monomorphic recursion as a primitive of the language, meaning that we type recursive definitions by first assigning them a type variable and then constraining

¹¹Note that we do not always have a choice about which side to distribute, unlike here.

that type variable to be a supertype of the inferred result type.¹² In this context, the programs the SRLC rules out are mostly ones that use indirect recursion (in the style of the Y combinator).

2.7 Expressiveness and Limitations

SuperF is uniquely expressive, but it does naturally have limitations.

Expressiveness. HML [Leijen 2009] is a restriction of ML^F where type annotations are restricted to System F types and where all polymorphic parameters must be annotated (even when not used polymorphically within the function body). There are small differences between HML/ ML^F and SuperF in the places where generalization occurs,¹³ which are usually smoothed over by distributivity. Modulo these differences, we conjecture that SuperF subsumes HML, which in turn subsumes ML type inference. So SuperF subsumes ML type inference as well, and type annotations are not needed by SuperF in programs where ML type checking succeeds. While we have not yet proved it, we experimentally verified this claim by porting OCaml’s `List` module to our SuperF implementation, requiring no type annotations (see Section 5.3). ML^F (even in its “shallow” variant [Le Botlan and Rémy 2009]) allows unrestricted type annotations and can therefore accept programs with annotated lower bounds that are not *syntactically* valid in SuperF. Conversely, SuperF accepts programs that do not type check in ML^F . Therefore, SuperF and ML^F are technically incomparable in terms of expressiveness. However, we conjecture that a syntactic restriction of ML^F source programs to using only System F type annotations would be subsumed by SuperF even without the additional polymorphism restriction of HML. Since functional programmers do not normally write type annotations with lower bounds (we are not aware of any functional programs using such annotations in the wild), one could argue that SuperF is more expressive than ML^F in practice.

Limitation: needed annotations. While it is enough to infer non-parametric intersection types in negative positions in many cases, it is not enough when a *parametrically*-polymorphic input type is needed. For instance, the following refactored version of `foo` no longer admits a precise type:

```
fooLet f = let g x = f x in (g 123, g True)
fooLet : ∀ a. ((Int ∨ Bool) → a) → (a, a)
```

This is because `fooLet` makes both arguments to `f` pass through the bottleneck of `g’s x` parameter, thereby merging their originally-separate data flows. It is not possible to recover the precise type of `foo` through type annotations, even in the declarative system.¹⁴ On the other hand, type annotations *can* be used if a different type is desired, for example annotating `(f : ∀ a. a → a)`.

Limitation: restricted type syntax. As explained in the introduction, SuperF is atypical in that it infers types that cannot always be written down (inferred types are generally *non-denotable*) due to the polarity restriction of its internal type syntax. Users may only provide System F type signatures, in which polymorphic types do not have bounds, so that for example the principal type of `foo` cannot be used as an explicit type signature. However, any well-typed SuperF term admits a number of possible System F types. These types are often less precise than the SuperF type, and there is often not a best one to pick, but since it is the user who is providing the type signature, this is not a problem in practice. It appears that users seldom want to use bounds in their annotations anyway (indeed, much of the work following up on ML^F was based on this very assumption). Moreover, in a typical program, a large number of functions, if not most of them, will not be associated with an explicit type signatures. For instance, even though the standard practice in Haskell is to annotate exported top-level definitions with type signatures, typical Haskell

¹²Of course, we can also type check polymorphically-recursive functions as long as explicit type signatures are provided.

¹³For instance, `flip const 42` is generalized (to $\forall 'a . 'a \rightarrow 'a$) in HML but not in SuperF, as it is not syntactically a function.

¹⁴Note that we *cannot* type `g` at $(Int \rightarrow a) \wedge (Bool \rightarrow b)$, the type of `f` in the original `foo`, because $F_{\{\leq\}}$ and a fortiori SuperF deliberately lack an intersection-introduction typing rule (which would make type inference intractable).

code contains many unannotated helper functions, notably in `where` clauses and in non-exported definitions, as well as lambda expressions passed as arguments to other functions or stored in data structures. When working with such local definitions, the full power of SuperF type inference is available, and previously-untypeable functions like `foo` can now be used. Therefore, our restricted type syntax should not be thought of chiefly as a limitation: the SuperF algorithm it enables is a strict improvement over most existing approaches for inferring System F-style first-class polymorphism.

2.8 Practical Considerations: Stability and Error Messages

We now briefly explain why we think SuperF is a useful and practical type inference system.

Robustness and Stability. An important design consideration for first-class-polymorphic type inference is *robustness* against innocuous changes to the program, i.e., type inference *stability* [Le Botlan and Rémy 2009, §4.5]. Fragile type inference leads to suboptimal user experience as it becomes harder for programmers to predict when type annotations will be needed. Thanks to subtyping, which faithfully encodes program data flows without having to approximate types through unification, SuperF achieves a high degree of stability. However, just like for expressiveness, SuperF is neither more nor less stable than ML^F in general. Table 1 summarizes some of the important stability properties of ML^F and SuperF. While inferring precise data flow types generally improves stability, it also makes the system *too expressive* for supporting stability properties that more rigid systems like ML^F possess, as exemplified in M1. A counterexample showing that M1 does not hold in SuperF is the `fooLet` function mentioned in Section 2.7, whose `let`-reduced form has the same (more general) type as `foo`. The reason this is not a problem in ML^F is that ML^F is able to type check neither `foo` nor `fooLet`, nor any similar functions where unannotated parameters are used at distinct types in the function’s body. We conjecture that M1 holds in SuperF if we make it a side condition that `x` be used *linearly* in `a2`, though proving this is quite challenging and we have not done so yet. Note that the side condition of MS5 is crucial, as neither ML^F nor SuperF are stable under general η expansion. Here, “`f` has function type” should be understood in a restrictive, syntactic sense: for instance, neither a type variable with function type bounds nor a polymorphic type with a function type body is a *function type* as required by the condition. The side condition that `x` should not appear under lambdas in S3 can be intuitively understood as follows: lambdas introduce nested polymorphic scopes, so for `x` to interact with the corresponding nested polymorphic type variables without leaking them and causing imprecision, `x` may need to be parametrically polymorphic, but we do not infer parametrically polymorphic types for function parameters.

Table 1. Example *typeability equivalences* demonstrating some stability properties of ML^F and SuperF. In each case, the left-hand side is typeable *if and only if* the right-hand side is, with equivalent types. Rows labeled with “MS*” show properties shared by both systems while those labeled with “M*” show properties enjoyed by ML^F but not SuperF and those labeled with “S*” show properties enjoyed by SuperF but not ML^F. Metavariables `s`, `t`, `f`, `g` range over terms while `x`, `y` range over variables.

M1	Let conversion	<code>let x = s in t</code>	$[x \mapsto s]t$	<code>x</code> occurs in <code>t</code>
MS2	Redefine application	<code>s t</code>	<code>(fun f x → f x) s t</code>	
MS3	Reorder arguments	<code>s t₁ t₂</code>	<code>(fun x y → s (y x)) t₂ t₁</code>	
MS4	Curryfication	<code>s (t₁, t₂)</code>	<code>(fun x y → s (x, y)) t₁ t₂</code>	
MS5	η conversion	<code>f</code>	<code>fun x → f x</code>	<code>f</code> has function type
S1	Tail conversion	<code>(if t then f else g) x</code>	<code>if t then f x else g x</code>	
S2	Head conversion	<code>x (if s then t₁ else t₂)</code>	<code>if s then x t₁ else x t₂</code>	
S3	Rename argument	<code>(fun x → t) y</code>	$[x \mapsto y]t$	<code>x</code> does not occur under lambdas in <code>t</code>

<u>Core syntax</u>		<u>Contexts</u>	
<i>Type</i>	$\tau, \sigma ::= \alpha \mid \tau \rightarrow \tau \mid \forall V\{\Sigma\}. \tau \mid \langle \tau \rangle \mid \boxed{\tau}_x \mid \top$	<i>Typing ctx.</i>	$\Gamma ::= \epsilon \mid \Gamma \cdot (x : \tau) \mid \Gamma \cdot B$
<i>Bound</i>	$B ::= \alpha \leq \tau \mid \tau \leq \alpha$	<i>Bounds ctx.</i>	$\Xi, \Sigma ::= \epsilon \mid \Xi \cdot B$
<i>Term</i>	$s, t ::= x \mid \lambda x. t \mid t t \mid () \mid t : \tau \mid \mathbf{let} \ x = t \ \mathbf{in} \ t$	<i>Variables</i>	$V, W, X ::= \epsilon \mid V \cdot \alpha$

Fig. 1. Syntax of types, terms, and contexts.

User-friendly type error messages. A non-obvious advantage of type inference systems like SuperF based on subtyping (in the school of *algebraic subtyping* [Dolan 2017]) rather than unification is that they are easily adapted to track and propagate *type provenance information* relating all inferred types to the relevant source code locations from which these types arose. This information can in turn be used to produce precise and informative type error messages [Bhanuka et al. 2023]. In the first-class polymorphism setting, one has to additionally take care of the interaction between rigid variables and skolems, which can leak and cause problems down the line. SuperF is well-equipped to deal with these problems thanks to its fine-grained tracking of data flows.

3 DECLARATIVE TYPE SYSTEM

We now present the declarative $F_{\{\leq\}}$ type system formally.

The syntax of $F_{\{\leq\}}$ is presented in Figure 1. The only non-standard syntax is that of polymorphic types $\forall V\{\Sigma\}. \tau$, which include a set of bounds Σ on the quantified variables V , and the *box* $\boxed{\tau}_x$ and *unbox* $\langle \tau \rangle$ forms, which are explained next. For simplicity of presentation, we do not have special forms for concrete features such as product and sum types, as their addition is straightforward; the bare-bone features presented here are sufficient to demonstrate our main ideas.

Notations and Shorthands. Notation \overline{E}_i^i denotes a repetition of $i = 0$ to n occurrences of E_i , and we omit i when it is unambiguous. We also make use of the following notations:

$$V \notin X \triangleq \overline{\alpha \notin X}^{\alpha \in V}; \quad \mathcal{B}(\Gamma) \triangleq \overline{B}^{B \in \Gamma}; \quad \perp \triangleq \forall \alpha. \alpha; \quad \tau_1 \vee \tau_2 \triangleq \bigvee \{ \tau_1, \tau_2 \}; \quad \bigvee \{ \overline{\tau}_i^i \} \triangleq \forall \alpha \{ \overline{\tau}_i \leq \alpha^i \}. \alpha \quad (\alpha \text{ fresh})$$

3.1 Declarative Typing Rules

System $F_{\{\leq\}}$ is a mostly standard polymorphic type system whose main distinctive features are *multi-bounded polymorphic types* and *boxed types*. The latter prevents unannotated function parameters from being used polymorphically, which reflects the limitations of the SuperF algorithm.

The declarative typing rules of $F_{\{\leq\}}$ are presented in Figure 2. Rules **T-UNIT**, **T-VAR**, **T-APP**, **T-SUBS**, and **T-LET** are all standard. **T-ABS** extends the typing context with a *boxed type* $\boxed{\tau}_x$ for its parameter x , which prevents the parameter from being used polymorphically without an annotation, as we shall see next. **T-UNBOX** allows erasing leftover boxes whose variables are no longer in scope, using the box erasure syntax $\tau \setminus x$ formally defined and exemplified in Appendix B.1. **T-ASC** expects the annotated expression to have type $\langle \tau \rangle$, which is a type used to *unbox* polymorphic types, making them available for use. Rule **T-FORALL** allows *generalizing* the type τ of an arbitrary term t .¹⁵ There are two requirements for generalization: first, the free type variables of τ must not occur in Γ , which is standard; second, we require that the quantification $\forall V\{\Sigma\}$ be *consistent* in Γ , written $\mathcal{B}(\Gamma) \vdash \forall V\{\Sigma\}$ **cons.** and read “in context Γ , the quantification of variables V with bounds Σ is consistent”. We define consistency as the existence of a solution to Σ , i.e. a substitution of V such that the variable bounds Σ hold in $\mathcal{B}(\Gamma)$. This makes sure no parts of a typing derivation

¹⁵While SuperF only generalizes lambda expressions, it is simpler to present generalization as an orthogonal rule in $F_{\{\leq\}}$.

$$\begin{array}{c}
\boxed{\Gamma \vdash t : \tau} \\
\text{T-UNIT} \\
\frac{}{\Gamma \vdash () : \top} \\
\text{T-VAR} \\
\frac{\Gamma(x) = \tau}{\Gamma \vdash x : \tau} \\
\text{T-ABS} \\
\frac{\Gamma \cdot (x : \boxed{\tau_1}_x) \vdash t : \tau_2}{\Gamma \vdash \lambda x. t : \tau_1 \rightarrow \tau_2} \\
\text{T-UNBOX} \\
\frac{\Gamma \vdash t : \tau \quad x \notin \text{dom}(\Gamma)}{\Gamma \vdash t : (\tau \setminus x)} \\
\text{T-APP} \\
\frac{\Gamma \vdash t_0 : \tau_1 \rightarrow \tau_2 \quad \Gamma \vdash t_1 : \tau_1}{\Gamma \vdash t_0 t_1 : \tau_2} \\
\text{T-ASC} \\
\frac{\Gamma \vdash t : \langle \tau \rangle}{\Gamma \vdash (t : \tau) : \tau} \\
\text{T-SUBS} \\
\frac{\Gamma \vdash t : \tau_1 \quad \mathcal{B}(\Gamma) \vdash \tau_1 \leq \tau_2}{\Gamma \vdash t : \tau_2} \\
\text{T-FORALL} \\
\frac{\Gamma \cdot \Sigma \vdash t : \tau \quad V \notin \text{FV}(\Gamma) \quad \mathcal{B}(\Gamma) \vdash \forall V \{ \Sigma \} \text{ cons.}}{\Gamma \vdash t : \forall V \{ \Sigma \}. \tau} \\
\text{T-LET} \\
\frac{\Gamma \vdash t_1 : \tau_1 \quad \Gamma \cdot (x : \tau_1) \vdash t_2 : \tau_2}{\Gamma \vdash \text{let } x = t_1 \text{ in } t_2 : \tau_2}
\end{array}$$

Fig. 2. Declarative typing rules.

$$\begin{array}{c}
\boxed{\Xi \vdash \tau \leq \tau} \\
\text{S-TOP} \\
\frac{}{\Xi \vdash \tau \leq \top} \\
\text{S-VARREFL} \\
\frac{}{\Xi \vdash \alpha \leq \alpha} \\
\text{S-HYP} \\
\frac{B \in \Xi}{\Xi \vdash B} \\
\text{S-UNBOX1} \\
\frac{\Xi \vdash \tau \leq \sigma}{\Xi \vdash \boxed{\tau}_x \leq \langle \sigma \rangle} \\
\text{S-UNBOX2} \\
\frac{}{\Xi \vdash \tau \leq \langle \tau \rangle} \\
\text{S-CONGBOUND} \\
\frac{}{\Xi \vdash \boxed{\alpha}_x \leq \boxed{\tau}_x} \\
\text{S-TRANS} \\
\frac{\Xi \vdash \tau_0 \leq \tau_1 \quad \Xi \vdash \tau_1 \leq \tau_2}{\Xi \vdash \tau_0 \leq \tau_2} \\
\text{S-FUN} \\
\frac{\Xi \vdash \tau_0 \leq \tau_1 \quad \Xi \vdash \tau_2 \leq \tau_3}{\Xi \vdash \tau_1 \rightarrow \tau_2 \leq \tau_0 \rightarrow \tau_3} \\
\text{S-CONGFUN} \\
\frac{}{\Xi \vdash \boxed{\tau_1 \rightarrow \tau_2}_x \leq \tau_1 \rightarrow \boxed{\tau_2}_x} \\
\text{S-FORALL-R} \\
\frac{V \notin \text{FV}(\tau) \quad \Xi \vdash \forall V \{ \Sigma \} \text{ cons.}}{\Xi \vdash \tau \leq \forall V \{ \Sigma \}. \tau} \\
\text{S-FORALL-COV} \\
\frac{\Xi \vdash \forall V \{ \Sigma \} \text{ cons.} \quad \Xi \cdot \Sigma \vdash \tau \leq \sigma}{\Xi \vdash \forall V \{ \Sigma \}. \tau \leq \forall V \{ \Sigma \}. \sigma} \\
\text{S-FORALL-L} \\
\frac{}{\Xi \vdash \forall \bar{\alpha}_i^i \{ \Sigma \}. \tau \leq [\bar{\alpha}_i \mapsto \bar{\tau}_i^i] \tau} \\
\text{S-FORALL-DISTR} \\
\frac{\Xi \vdash \forall V \{ \Sigma \} \text{ cons.} \quad V \notin \text{FV}(\tau_1)}{\Xi \vdash \forall V \{ \Sigma \}. \tau_1 \rightarrow \tau_2 \leq \tau_1 \rightarrow \forall V \{ \Sigma \}. \tau_2}
\end{array}$$

Fig. 3. Declarative subtyping rules.

can rely on unsound assumptions. In particular, subtyping assumptions cannot be used to carry subtyping proofs in the same way as dependent type systems and GADTs can be used to carry type equality or subtyping proofs [Boruch-Gruszecki et al. 2022; Parreaux et al. 2019; Scherer and Rémy 2015]. We disallow inconsistent constraints because these could be used to make the type system accept obviously ill-typed terms. For example, $\forall \alpha \{ \text{Nat} \leq \alpha \leq \text{Int} \}. \alpha \rightarrow \alpha$ (a shorthand for $\forall \alpha \{ \text{Nat} \leq \alpha, \alpha \leq \text{Int} \}. \alpha \rightarrow \alpha$) is consistent because $\text{Nat} \leq \text{Int}$. It defines an identity type that works on any type between Nat and Int . On the other hand, $\forall \alpha \beta \{ \top \leq \alpha \leq \beta, \beta \leq \text{Int} \}. \alpha \rightarrow \beta$ is *inconsistent* because the (transitive) bounding on α and β , $\top \leq \text{Int}$, cannot be derived.

3.2 Declarative Subtyping Rules

The subtyping rules of $F_{\{\leq\}}$ are presented in Figure 3. Rules **S-TOP**, **S-VARREFL**, **S-TRANS**, and **S-FUN** are standard. **S-HYP** is used to leverage any type variable bounds present in the typing context. **S-FORALL-R** introduces polymorphic types in the middle of subtyping. Its premise makes sure that the added quantified type variables are *not* in the type being quantified – the idea is that these type variables will appear later through the combined use of **S-FORALL-COV**¹⁶ to widen a type under quantifiers, and **S-FORALL-L**, to instantiate the type variables of a polymorphic type. This way, we can use the subtyping rules to rearrange the quantifiers of a type – for example, we can

¹⁶**S-FORALL-COV** should not be confused with “kernel” rules in the context of explicit polymorphism *à la* System F_{\leq} . Unlike these systems, $F_{\{\leq\}}$ allows *polymorphic subtyping*, so that relationships like $\forall \alpha \{ \alpha \leq \text{Int} \}. \tau \leq \forall \alpha \{ \alpha \leq \text{Nat} \}. \tau$ do hold.

derive subtyping relationships such as $\forall\alpha\gamma.\forall\beta\{\beta \leq \alpha\}.\alpha \rightarrow \beta \rightarrow \gamma \leq \forall\beta\alpha\{\beta \leq \alpha\}.\alpha \rightarrow \beta \rightarrow \top$. **S-FORALL-DISTR** describes the distributivity of polymorphic types over arrow types; note that the other direction is already admissible by covariance of function results and uses of **S-FORALL-R/S-FORALL-COV/S-FORALL-L**. All these rules require that the corresponding quantifications be consistent, to avoid introducing bad types during subtyping. Finally, **S-UNBOX1** and **S-UNBOX2** are the rules that make type ascription work. The latter is used when the ascribed type needs not be unboxed, and the former is used to open a boxed function parameter type. Note that boxes only block polymorphic types and they congrue with type variables and functions:¹⁷ **S-CONGBOUND** allows *upcasting* a boxed type variable to one of its bounds and **S-CONGFUN** allows pushing a box into the result type.

Derived lattice types. As hinted previously, we can *encode* bottom and union types.

THEOREM 3.1. $\perp \triangleq \forall\alpha.\alpha$ is the bottom of (\leq) , i.e., for all τ **wf**, we have $\perp \leq \tau$.

THEOREM 3.2. For all $\overline{\tau_i}^i$ **wf** where $\alpha \notin FV(\overline{\tau_i}^i)$, the type defined as $\bigvee\{\overline{\tau_i}^i\} \triangleq \forall\alpha\{\overline{\tau_i} \leq \alpha\}.\alpha$ is the least upper bound of all $\overline{\tau_i}^i$, i.e.,

- (A) $\bigvee\{\overline{\tau_i}^i\}$ is an upper bound of $\overline{\tau_i}^i$, meaning that for all j we have $\tau_j \leq \bigvee\{\overline{\tau_i}^i\}$; and
- (B) for all σ **wf** such that σ is an upper bound of $\overline{\tau_i}^i$, we have $\bigvee\{\overline{\tau_i}^i\} \leq \sigma$.

Remark 1. By contrast, we cannot encode intersections in the same way as unions, because that would require the use of existential quantification, which is not supported by $F_{\{\leq\}}$ nor by the underlying System F_{cc} . On the other hand, using universal quantification, we can still encode intersections that occur *negatively* in some outer type, i.e., if α occurs negatively in $\tau[\alpha]$, then $\tau[\sigma_0 \wedge \sigma_1]$ can be encoded as $\forall\alpha\{\alpha \leq \sigma_0, \alpha \leq \sigma_1\}.\tau[\alpha]$, where $\alpha \notin FV(\tau) \cup BV(\tau[\cdot])$.

Example 3.3. Consider term $t = \lambda x. \mathbf{let} \ y = (\mathit{choose} \ x \ (\lambda z. z) : \forall\alpha.\alpha \rightarrow \alpha) \ \mathbf{in} \ (y \ 0, y \ \mathbf{True})$ in context $\Xi = (\mathit{choose} : \forall\beta.\beta \rightarrow \beta \rightarrow \beta)$. As expected intuitively, t can be typed in $F_{\{\leq\}}$ at type $(\forall\alpha.\alpha \rightarrow \alpha) \rightarrow (\text{Int}, \text{Bool})$. First note that $\lambda z. z$ can be typed at $\forall\alpha.\alpha \rightarrow \alpha$. We can also assign to parameter x the type $\forall\alpha.\alpha \rightarrow \alpha$, so x in the lambda body is typed at $\boxed{\forall\alpha.\alpha \rightarrow \alpha}_x$. We need to pass arguments of mismatched boxed and unboxed types to *choose*. This is no problem: instantiate *choose*'s type variable β to the union $\boxed{\forall\alpha.\alpha \rightarrow \alpha}_x \vee (\forall\alpha.\alpha \rightarrow \alpha)$, letting us type the call *choose* $x \ (\lambda z. z)$, which returns the same union type. We can then upcast the result by *instantiating* the union type to $\langle \forall\alpha.\alpha \rightarrow \alpha \rangle$ (since each side is a subtype of it) and use that type in **T-Asc**.

Example 3.4. Consider again the `fooLet` function from Section 2.7. Observe that we *cannot* type `g` polymorphically enough because it lacks a type annotation on `f` to that effect.

3.3 Metatheory

To ensure the soundness of $F_{\{\leq\}}$, we translate our terms into valid System F_{cc} terms. System F_{cc} is a very general declarative type system of so-called *implicit coercions* (which are essentially polymorphic subtyping assumptions parameterized with a context), designed by **Cretin [2014]**; **Cretin and Rémy [2012, 2014]**. The soundness of our translation is stated in the following theorem:

THEOREM 3.5. If $\vdash t : \tau$, then $\vdash^{F_{cc}} \tilde{t} : \llbracket \tau \rrbracket$ for some System F_{cc} term \tilde{t} .

The definition of type translation $\llbracket \tau \rrbracket$ and the proof of the above theorem are enclosed in Appendix C. While the theorem only allows translating terms well-typed in the empty context, a term well-typed in an arbitrary Γ can be translated after wrapping it in abstractions to capture its free variables.

¹⁷In an extended system with sum and product types, we would define similar rules for pushing boxes inside them.

<p><i>Positive type</i> $\tau^+ ::= a \mid \tau^- \rightarrow \tau^+ \mid \forall \bar{\alpha}\{\Sigma\}. \tau^+ \mid \top$</p> <p><i>Negative type</i> $\tau^- ::= a \mid \tau^+ \rightarrow \tau^- \mid \forall \alpha. \tau^- \mid \top$</p> <p><i>Polarity</i> $\pm ::= + \mid -$</p> <p><i>Contexts</i> $\Gamma ::= \epsilon \mid \Gamma \cdot (x : \tau^+)$ $\phi ::= \Xi \mid \Xi \cdot \triangleright \phi$</p>	<p><i>Bound</i> $B ::= a \leq \tau^- \mid \tau^+ \leq a$</p> <p><i>Type variable</i> $a, b ::= \alpha \mid \alpha_\alpha$</p> <p><i>Constraint</i> $C ::= \tau^+ \leq^\phi \tau^-$</p> <p>$\Xi^? ::= \Xi \mid \mathbf{err}$ $\Delta ::= \epsilon \mid \Delta \cdot C$</p>
---	---

Fig. 4. Polarized syntax of SuperF. Notation: we write $\tau^+ \leq \tau^-$ as a shorthand for $\tau^+ \leq^\epsilon \tau^-$.

4 TYPE INFERENCE

We now formally describe the SuperF type inference algorithm.

4.1 Restricted Polar Syntax

The syntax of types used by SuperF is presented in Figure 4. We now separate *positive* types τ^+ from *negative* types τ^- . This is reminiscent of the way that Le Botlan and Rémy disallow polymorphic types with non-trivial bounds “below arrows” in Shallow ML^F [Le Botlan and Rémy 2009].¹⁸ We require that all type variables quantified in negative positions be boundless.

Type variables a, b include *rooted variables* α_β . Each α_β is a distinct, standalone variable with two components: a name α and a root β . We have $FV(\alpha_\beta) = \{\alpha_\beta\}$. We use the root as metadata to ensure type inference terminates (see Section 4.3). As a shorthand, we define $\alpha_{\beta_\gamma} \triangleq \alpha_\gamma$.¹⁹

4.2 Type Inference Rules

The type inference rules of SuperF are presented in Figure 5. Judgment $\Gamma \vdash t : \tau^+ \Rightarrow \Delta$ is read “under context Γ , term t can be typed at τ^+ by generating constraints Δ ”. The constraints in the output context Δ are not resolved yet and may thus be inconsistent. The *subtype constraining* judgment $V, W \vdash \Xi \gg \Delta \gg \Sigma^?$, read “under rigid variables V , skolems W , and assumptions Ξ , the constraints Δ produce a set of bounds or an error $\Sigma^?$ ” and presented in the next section, is used to ensure that they are not. Type inference derivations make use of *freshness* premises. We only consider *well-formed* derivations, which are those where ‘ α *fresh*’ occurs at most once for each α . We write $S_1 \# S_2$ to denote that the sets S_1 and S_2 are disjoint.

Rules **I-TOP**, **I-VAR**, **I-ASC** and **I-LET** are straightforward. Importantly, note that the $(t : \tau^-)$ term in rule **I-ASC** only permits τ^- which are syntactically both positive and negative types (i.e., System F types), since τ^- also occurs where the judgment syntactically requires a positive type. Rule **I-APP** infers the types of both parts t_1 and t_2 of an application $t_1 t_2$, *assumes* some function type $\tau_2^+ \rightarrow \alpha$ for the former (where α is a fresh type variable), and finally returns α as a result along with the constraint that the type inferred for t_1 indeed is a subtype of $\tau_2^+ \rightarrow \alpha$. Rule **I-ABS** is the most interesting. It type checks the body of a lambda in a context Γ extended by $(x : \alpha)$ where α is fresh, which so far is standard. Then, it *resolves* the inferred constraints Δ , reducing them down to some bounds Ξ . These comprise bounds on the type variable α and on type variables created during lambda body typing or while resolving the constraints, as well as bounds on *outer* type variables (bound by other lambdas). The goal is to generalize the resulting polymorphic type “as much as possible” without leaking references to locally-quantified type variables into the outer polymorphic context, as we explained in Section 2.4. The *rigid variable* set V contains all variables which should

¹⁸Since Shallow ML^F only supports type variable lower bounds (no upper bounds) and since lower bounds are positive positions, this syntactic restriction is similar to having a polarized type syntax similar to ours.

¹⁹The core syntax only supports type variables a of the forms α and α_β . This shorthand defines the meaning of syntax α_a . It reflects the fact that the *root* of the variable stays the same even as a fresh type variable is duplicated into a new type variable. Without the shorthand, we would need to explicitly destructure a , which quickly becomes too verbose.

$$\boxed{\Gamma \vdash t : \tau^+ \Rightarrow \Delta}$$

$$\begin{array}{c}
\text{I-TOP} \\
\hline
\Gamma \vdash () : \top \Rightarrow \epsilon
\end{array}
\quad
\begin{array}{c}
\text{I-VAR} \\
\hline
\Gamma(x) = \tau^+ \\
\Gamma \vdash x : \tau^+ \Rightarrow \epsilon
\end{array}
\quad
\begin{array}{c}
\text{I-ASC} \\
\hline
\Gamma \vdash t : \tau^+ \Rightarrow \Delta \\
\Gamma \vdash (t : \tau^-) : \tau^- \Rightarrow \Delta \cdot (\tau^+ \leq \tau^-)
\end{array}$$

$$\begin{array}{c}
\text{I-LET} \\
\hline
\Gamma \vdash t_1 : \tau_1^+ \Rightarrow \Delta_0 \quad \Gamma \cdot (x : \tau_1^+) \vdash t_2 : \tau_2^+ \Rightarrow \Delta_1 \\
\Gamma \vdash \text{let } x = t_1 \text{ in } t_2 : \tau_2^+ \Rightarrow \Delta_0 \cdot \Delta_1
\end{array}
\quad
\begin{array}{c}
\text{I-APP} \\
\hline
\alpha \text{ fresh} \quad \Gamma \vdash t_1 : \tau_1^+ \Rightarrow \Delta_1 \quad \Gamma \vdash t_2 : \tau_2^+ \Rightarrow \Delta_2 \\
\Gamma \vdash t_1 t_2 : \alpha \Rightarrow \Delta_1 \cdot \Delta_2 \cdot (\tau_1^+ \leq \tau_2^+ \rightarrow \alpha)
\end{array}$$

$$\begin{array}{c}
\text{I-ABS} \\
\hline
\alpha \text{ fresh} \quad V \supseteq FV(\Gamma) \quad V \# FV(\Delta) \setminus FV(\Gamma) \\
\Gamma \cdot (x : \alpha) \vdash t : \tau^+ \Rightarrow \Delta \quad V, \epsilon \vdash \epsilon \gg \Delta \gg \Xi \quad \text{split}_V(\text{uproot}(\Xi)) = (\Xi_0, \bar{y}, \Xi_1) \\
\Gamma \vdash \lambda x. t : \forall \alpha \bar{y} \{ \Xi_1 \}. \alpha \rightarrow \tau^+ \Rightarrow \Xi_0
\end{array}$$

Fig. 5. Type inference rules of SuperF.

not be bound in the abstraction's type, and which will be bound in some outer context. This includes all type variables already in Γ , as well as additional fresh variables which are used by extrusion as approximants (see Section 4.3). To ensure constraining terminates, some type variables in Ξ may have been ascribed with *roots*; these ascriptions are unnecessary after subtype constraining is done, so they can be removed by *uproot*.

Definition 4.1 (Uprooting). Let $\text{uproot}(\alpha_\beta) \triangleq \alpha$, which naturally extends to τ and Ξ by congruence.

Definition 4.2 (Context splitting). Let us generalize the set difference operator ' \setminus ' to contexts. We *split* a context between outer bounds (those that mention *only* outer variables) and the rest using: $\text{split}_V(\Xi) \triangleq (\text{outer}_V(\Xi), FV(\Xi) \setminus V, \text{inner}_V(\Xi))$

$$\text{outer}_V(\Xi) \triangleq \frac{}{\tau^+ \leq \sigma^-} \in \Xi, FV(\tau^+ \leq \sigma^-) \subseteq V \quad \text{inner}_V(\Xi) \triangleq \Xi \setminus \text{outer}_V(\Xi)$$

4.3 Subtype Constraining

The subtype constraining rules are defined in Figure 6.

Constraining algorithm. We first describe how to interpret the subtype constraining rules as an algorithm. Our algorithm uses the first variable context both as an input and an output: we input rigid variables from the context (V) and get back approximants from extrusion (V'), combining both in the judgment as a single context ($V \cdot V'$). The former is an immutable set of rigid variables while the latter is a supply of rigid variables to be used as approximants in extrusion, formally represented by the subset premises in **I-ABS** and **C-FORALL-R**. In an implementation of constraining, V' would be a mutable set. In a constraining derivation ($V \cdot V'$), $W \vdash \Xi \gg \Delta \gg \Xi'$, V' and Ξ' are the only output; everything else is an input. The order in which the constraining rules are applied does not affect the correctness of the corresponding constraining judgments, but our algorithm always applies the rules in the order given in the figure. (Crucially, **C-FORALL-R** should be applied *before* **C-FORALL-L** when possible, as explained in Section 2.5.) In order to algorithmically construct a derivation of constraining given the inputs (V, W, Ξ, Δ), we begin by attempting to use the first rule whose conclusion shape matches the inputs. If the matching rule has constraining premises, we recurse on them. If deriving any of them results in **err**, the algorithm backtracks and uses **C-FAIL** to construct the derivation. Otherwise, we construct a successful derivation. In particular, if the SRLC (premise of **C-FLEX-L/C-FLEX-R**) fails, the output is **err**.

Acyclicity. For type inference to be sound, the constructed bounds graphs must remain acyclic, because F_{cc} does not support “full” recursive types (i.e., ones where recursive occurrences may appear in bounds). We define the semantic *acyclic* check for this purpose in Appendix B.2. That definition does not count two kinds of “harmless” cycles:

$$\boxed{V, W \vdash \Xi \gg \Delta \gg \Xi'}$$

Given rigid variables V , skolems W , and assuming bounds Ξ ,
the constraints in Δ can be solved by introducing new bounds Ξ' (if $\Xi' \neq \text{err}$)

$$\begin{array}{c}
\text{C-EMPTY} \\
\hline
V, W \vdash \Xi \gg \epsilon \gg \epsilon \\
\\
\text{C-FUN} \\
\hline
V, W \vdash \Xi_1 \gg \Delta \cdot (\tau_3^+ \leq^{\triangleright\phi} \tau_1^-) \cdot (\tau_2^+ \leq^{\triangleright\phi} \tau_4^-) \gg \Xi_2 \\
\hline
V, W \vdash \Xi_1 \gg \Delta \cdot (\tau_1^- \rightarrow \tau_2^+ \leq^{\phi} \tau_3^- \rightarrow \tau_4^-) \gg \Xi_2 \\
\\
\text{C-FLEX-L} \\
B = (a \leq \tau^-) \quad \text{root}(B) \notin \text{roots}(\phi) \quad \text{acyclic}(\Xi; B) \\
\hline
a \notin V \cdot W \quad V, W \vdash \Xi; B \gg \Delta \cdot (\tau^+ \leq^{B \cdot \phi} \tau^-) \stackrel{(\tau^+ \leq a) \in \Xi}{\gg} \Xi' \\
\hline
V, W \vdash \Xi \gg \Delta \cdot (a \leq^{\phi} \tau^-) \gg \Xi' \cdot B \\
\\
\text{C-RIGID-L} \\
a \in V \quad V, W \vdash \tau^- \rightsquigarrow (\Sigma, \sigma^-) \quad V, W \vdash \Xi \gg \Delta \cdot \Sigma \gg \Xi' \\
\hline
V, W \vdash \Xi \gg \Delta \cdot (a \leq^{\phi} \tau^-) \gg \Xi' \cdot (a \leq \sigma^-) \\
\\
\text{C-EMPTY} \\
\hline
V, W \vdash \Xi \gg \epsilon \gg \epsilon \\
\\
\text{C-TOP} \\
\hline
V, W \vdash \Xi \gg \Delta \cdot (\tau \leq^{\phi} \top) \gg \Xi' \\
\\
\text{C-VARREFL} \\
\hline
V, W \vdash \Xi \gg \Delta \gg \Xi' \\
\hline
V, W \vdash \Xi \gg \Delta \cdot (a \leq^{\phi} a) \gg \Xi' \\
\\
\text{C-FLEX-R} \\
B = (\tau^+ \leq a) \quad \text{root}(B) \notin \text{roots}(\phi) \quad \text{acyclic}(\Xi; B) \\
\hline
a \notin V \cdot W \quad V, W \vdash \Xi; B \gg \Delta \cdot (\tau^+ \leq^{B \cdot \phi} \tau^-) \stackrel{(a \leq \tau^-) \in \Xi}{\gg} \Xi' \\
\hline
V, W \vdash \Xi \gg \Delta \cdot (\tau^+ \leq^{\phi} a) \gg \Xi' \cdot B \\
\\
\text{C-RIGID-R} \\
a \in V \quad V, W \vdash \tau^+ \rightsquigarrow (\Sigma, \sigma^+) \quad V, W \vdash \Xi \gg \Delta \cdot \Sigma \gg \Xi' \\
\hline
V, W \vdash \Xi \gg \Delta \cdot (\tau^+ \leq^{\phi} a) \gg \Xi' \cdot (\sigma^+ \leq a) \\
\\
\text{C-FORALL-L} \\
\beta \text{ fresh} \quad \sigma^- = [\alpha \mapsto \beta_\alpha] \tau^- \quad V' \supseteq FV(\tau^+ \leq \tau^-) \setminus W \quad V' \# W \cdot \beta_\alpha \\
V', W \cdot \beta_\alpha \vdash \epsilon \gg \tau^+ \leq^{\triangleright\phi} \sigma^- \gg \Xi \quad V, W \vdash \Xi_0 \gg \Delta \cdot \text{outer}_{V'}(\Xi)^\phi \gg \Xi_1 \\
\hline
V, W \vdash \Xi_0 \gg \Delta \cdot (\tau^+ \leq^{\phi} \forall \alpha. \tau^-) \gg \Xi_1 \\
\\
\text{C-FORALL-R} \\
\beta_\alpha \text{ fresh} \quad V, W \vdash \Xi_0 \gg \Delta \cdot ([\alpha \mapsto \beta_\alpha] \bar{B})^\phi \cdot ([\alpha \mapsto \beta_\alpha] \tau^+ \leq^{\phi} \tau^-) \gg \Xi_1 \\
\hline
V, W \vdash \Xi_0 \gg \Delta \cdot (\forall \bar{\alpha} \{ \bar{B} \}. \tau^+ \leq^{\phi} \tau^-) \gg \Xi_1 \\
\\
\text{C-FAIL} \\
\hline
V, W \vdash \Xi \gg \Delta \gg \text{err}
\end{array}$$

$$\boxed{V, W \vdash \tau^\pm \rightsquigarrow (\Sigma, \sigma^\pm)}$$

Given rigid variables V and skolems W , the extrusion of τ^\pm is σ^\pm in Σ .

$$\begin{array}{c}
\text{X-1} \\
\hline
FV(\tau^\pm) \subseteq V \\
\hline
V, W \vdash \tau^\pm \rightsquigarrow (\epsilon, \tau^\pm) \\
\\
\text{X-2} \\
F = FV(\tau^\pm) \setminus (V \cdot W) \quad (\beta_a, a, \gamma_a) \text{ X-fresh} \quad \{\beta_a, \gamma_a\} \subseteq V^{a \in F} \\
\rho = [a \mapsto \beta_a, a \mapsto \gamma_a]^{a \in F}, [b \mapsto \perp, b \mapsto \top]^{b \in W} \quad W \vdash \tau^\pm \text{ X-ok} \\
\hline
V, W \vdash \tau^\pm \rightsquigarrow ((\beta_a \leq a) \cdot (a \leq \gamma_a))^{a \in F}, \rho^\pm \tau^\pm
\end{array}$$

Fig. 6. Subtype constraining and extrusion rules of SuperF.

- *Immediate* type variable bound cycles. Indeed, the variables involved in an immediate cycle are simply equivalent (and could be unified). For instance, the judgment $\text{acyclic}((\alpha \leq \beta) \cdot (\beta \leq \alpha))$ holds. A type like $\forall \alpha \beta \{ \alpha \leq \beta, \beta \leq \alpha \}. (\alpha \rightarrow \beta) \rightarrow (\alpha \rightarrow \beta)$ is equivalent to $\forall \alpha. (\alpha \rightarrow \alpha) \rightarrow (\alpha \rightarrow \alpha)$ and does not contain a ‘real’ or *semantic* cycle.
- *Contravariant* cycles like in type $\forall \alpha \beta \{ \alpha \leq \alpha \rightarrow \beta \}. \alpha \rightarrow \beta$, which would be inferred for term $\lambda x. x x$, or equivalently $\forall \alpha \beta. ((\alpha \rightarrow \beta) \wedge \alpha) \rightarrow \beta$, which is equivalent and also syntactically cyclic. Such syntactically cyclic bounds do not actually lead to cyclic reasoning because they do not induce cycles in the polar traversal of types. To understand this, the reader will want to refer to the definition of polar traversal in Appendix B.2 as the *reach*[±] dual functions.

Constraint annotations. Constraints have the form $\tau^+ \leq^{\phi} \tau^-$, where ϕ is a set of subtyping relationships, some possibly guarded by \triangleright , which are currently in the process of being constrained.

This annotation has two uses: First, it allows catching, in **C-SKIP**, constraints that go through *immediate* type variable cycles (consider what happens with, e.g., $\Delta = (\alpha \leq \beta) \cdot (\beta \leq \alpha) \cdot (\alpha \leq \perp)$). Note that \in does not look past \triangleright in ϕ , so that **C-SKIP** is prevented from triggering on previously seen constraints after going through a type constructor (see how **C-FUN** adds \triangleright in its sub-constraints). Indeed, we do not wish the constraining algorithm to admit recursive types because our declarative system $F_{\{\leq\}}$ does not support them, but this would happen if we did not guard the constraints in \triangleright while going into type constructors. The effect of \triangleright is to prevent **C-SKIP** from picking up constraints that are being solved at an outer constructor level. Second, it allows the SRLC to catch constraints applying on roots currently being constrained (as an example, see Appendix A.5). The reason why ϕ cannot be a simple *context* in the constraining judgment is that constraining is specified as a “tail-recursive” worklist algorithm, where subderivations do not necessarily correspond to implied constraints, since they also concern “sibling” constraints in Δ . We also write $(\tau^+ \leq \sigma^-)^\phi$ as a shorthand for $\tau^+ \leq^\phi \sigma^-$.

Constraining rules. Rules **C-EMPTY**, **C-TOP**, **C-VARREFL**, and **C-FUN** are straightforward. There are four nontrivial rules for type variables. **C-FLEX-L** and **C-FLEX-R** register new bounds on locally-flexible type variables and traverse the transitive consequences of these bounds. They perform a cyclicity check to ensure that cyclic bounds such as $\alpha \leq \tau \rightarrow \alpha$ are rejected and a roots-check to guarantee termination. Moreover, they augment the sub-constraint annotations with the bound B currently being propagated. Finally, they output B as part of the result. **C-RIGID-L** and **C-RIGID-R** register *extruded* bounds on *locally-rigid* variables. Locally-rigid variables are those bound by outer lambdas as well as those currently rigidified by an outer application of **C-FORALL-R**; extrusion ensures that inner variables do not leak through the bounds of outer ones. **C-FORALL-L** instantiates the type variables and bounds of a polymorphic type found on the left. **C-FORALL-R** skolemizes a quantified variable found on the right-hand side, as described in Section 2.5. We first compute the set of *free variables* V' involved in the subsequent $\tau^+ \leq \tau^-$ comparison in order to locally rigidify them. We recurse on that comparison, which gives us the intermediate result Ξ . Finally, we continue constraining, adding the delayed rigid variable bounds $outer_{V'}(\Xi)$ as new constraints to solve. We can discard $inner_{V'}(\Xi)$ because these bounds only concern purely local variables that were created while constraining the skolemized right-hand-side type. For instance, when constraining $\forall \alpha. \alpha \rightarrow \alpha \leq \forall \beta. \beta \rightarrow \beta$, these inner bounds will contain $\alpha'_\alpha \leq \beta'_\beta$ and $\beta'_\beta \leq \alpha'_\alpha$, which are irrelevant outside the subderivation. Note that there is no need to relate V' and V : it is enough to only have all free type variables in $\tau^+ \leq \tau^-$. We suspend the constraining of any type variables *other than* those that are locally introduced while instantiating nested polymorphic types. In the second subderivation, we resume constraining the delayed bounds on V' induced by $\tau^+ \leq \sigma^-$.

Polar substitutions. Notation $\rho = [a^- \mapsto \tau^-, a^+ \mapsto \tau^+]$ denotes a *polar substitution* ρ , which is a substitution that maps the *positive* occurrences of a type variable a to some negative type τ^- and the *negative* occurrences of a to a possibly different positive type τ^+ . A polar substitution ρ is applied as either $\rho^+(\sigma^+)$ or $\rho^-(\sigma^-)$ depending on the assumed polarity of σ^\pm . (Since the syntaxes of positive and negative types overlap, it is not always possible to tell their polarity without context.)

Extrusion. Extrusion is a critical part of polymorphic type inference. Its task is to avoid adding to a type variable bounds that are “more polymorphic” than the type variable itself. We distinguish between *locally-rigid* type variables V and *skolems* W : the former are *less* polymorphic than the current polymorphism level, while the latter are *more* polymorphic. Skolems are introduced by **C-FORALL-R** and represent arbitrary future instantiations; less polymorphic types should thus not be allowed to mention them. Rules **C-RIGID-L** and **C-RIGID-R** use extrusion to ensure that the bounds added to rigid type variables mention neither flexible variables nor skolems. Extrusion

itself is defined with two rules: **X-1** and **X-2**. The first one simply filters out types that already only mention rigid variables. In **X-2**, set F is the list of flexible variables to extrude. Flexible variables are “more polymorphic” than rigid variables, thus flexible variables should not appear in the bounds of rigid variables. Thankfully, we are allowed to add new bounds to flexible variables, so that we can create less polymorphic *approximants* for them to replace their uses in the bounds of rigid variables. An approximant is related to the original polymorphic type variable through the latter’s bounds. Concretely, each $a \in F$ is extruded into a lower approximant β_a (for the negative occurrences of a) and an upper approximant γ_a (for the positive occurrences of a), and we register the bounds $\beta_a \leq a$, $a \leq \gamma_a$ as constraints — these should be viewed as bounds registered on a , as the approximants cannot themselves refer to the more polymorphic a in their own bounds. $[a- \mapsto \beta_a, a+ \mapsto \gamma_a]$ denotes the corresponding polar substitution. The **X-fresh** premise denotes that each a is always approximated into the *same* approximants in the whole constraining derivation and that each approximant is an otherwise fresh rigid variable. We also use an **X-ok** premise to ensure that the bounds of universal types remain consistent after extrusion.²⁰ In contrast to flexible variables, skolems are “too polymorphic” and their occurrences can only be approximated to either \top or \perp , depending on each occurrence’s polarity. We cannot register approximants for skolems because we are not allowed to add bounds to skolems, which denote *unknown, arbitrary* types. Both skolems and flexible variables are approximated by extrusion such that the rigid variable bound resulting from **C-RIGID-L** and **C-RIGID-R** is at least as constraining as the corresponding type from the original constraint. For instance, if α is rigid and β is flexible, then the constraint $\alpha \leq \top \rightarrow \beta$ results in a bound $\alpha \leq \top \rightarrow \gamma_\beta$ such that $\gamma_\beta \leq \beta$, where γ_β is the *lower approximant* of β , a special locally-rigid variable. The approximant is bound at the same polymorphism level as α , making it equally polymorphic as α and generalized alongside it, for example in **I-Abs** (see the description of **I-Abs** in Section 4.2). Approximants always end up being bounded by the opposite bounds of their root variable: the lower approximant of β inherits β ’s upper bounds, and the upper approximant of β inherits β ’s lower bounds, which results from adding the approximant bounds as *constraints* to be solved later (rather than as solved bounds). This ensures that bounding a flexible variable with its approximants yields consistent bounds, correctly making constraints like $\alpha \rightarrow \alpha \leq \forall \beta. \beta \rightarrow \beta$ fail.

Example 4.3. Assuming $f : (\forall \alpha. \alpha \rightarrow \alpha) \rightarrow \beta$, notice that **let** $x = \lambda y. y$ **in** $f (\lambda z. x z)$ can be typed in SuperF because $\forall \gamma. \gamma \rightarrow \gamma$ is inferred as the type of x and subsequently used parametrically in $x z$. On the other hand, $(\lambda x. f (\lambda z. x z)) (\lambda y. y)$ cannot be typed by SuperF because we never assume parametrically-polymorphic types for parameters. In this case, the type variable α assigned to z and the type variable γ assumed for the result of the $x z$ application are leaked outside the scope of the λz abstraction because they are involved in constraint $\beta \leq \alpha \rightarrow \gamma$, where β is the type variable assigned to x . Because of this constraint, both α and γ need to be extruded to the lower polymorphic level of β . Therefore, $\lambda z. x z$, whose inferred type is, $\alpha \rightarrow \gamma$, cannot be assigned the polymorphic type required by f . Note also that after η -contracting $\lambda z. x z$ to x , $(\lambda x. f x) (\lambda y. y)$ become typeable in SuperF, which shows that SuperF, like $F_{\{\leq\}}$, is not stable under η expansion.

Termination and the SRLC. The Suspiciously Recursive-Looking Criterion (SRLC) is formally represented as the two *roots* premises in **C-FLEX-L** and **C-FLEX-R**. As explained in Section 2.5, the goal of the SRLC is to ensure that for all inputs, deriving subtype constraining always terminates. We provide an intuitive proof sketch for the termination theorem in Appendix D.1.

²⁰This is currently implemented as a rather ad-hoc and incomplete check, given by Definition B.8 in Appendix B.3, which crudely ensures that polymorphic types will not get inconsistent bounds after extrusion by checking that the skolems being extruded are not reachable from both the upper bounds and lower bounds of the same type variable. If this check fails, type inference conservatively fails.

Definition 4.4 (Roots). Define $root(\alpha_\beta) \triangleq \beta$, $root(\alpha) \triangleq \alpha$ and, when τ is not a type variable, $root(\tau) \triangleq \overline{[a \mapsto root(a)]^{a \in FV(\tau)}} \tau$. In addition, define $root(\tau \leq \sigma) = root(\tau) \leq root(\sigma)$, $roots(\Xi) \triangleq \{root(\tau \leq \sigma) \mid (\tau \leq \sigma) \in \Xi\}$, and $roots(\Xi \triangleright \phi) \triangleq roots(\Xi) \cup roots(\phi)$.

Order of constraining type variables. The fact that C-FLEX-L and C-FLEX-R are tried before C-RIGID-L and C-RIGID-R is important to avoid non-termination of the constraining algorithm. In a previous version of this paper [Parreaux et al. 2024], we mistakenly had the rules in the opposite order (contrary to what our implementation actually does), which is fixed in this version.

4.4 Distributivity

Distributivity, which is supported by $F_{\{\leq\}}$, adds flexibility to the type inference system. The basic idea of using distributivity in type inference is described in the two following rules:

$$\begin{array}{c}
 \text{C-FORALL-DISTR-R} \\
 \frac{\alpha \notin FV(\tau_1^+) \quad \Xi_1 \gg \Delta \cdot (\tau_1^+ \leq^\phi \forall \alpha. \tau_1^+ \rightarrow \tau_2^-) \gg \Xi_2}{\Xi_1 \gg \Delta \cdot (\tau_1^+ \leq^\phi \tau_1^+ \rightarrow \forall \alpha. \tau_2^-) \gg \Xi_2}
 \end{array}
 \qquad
 \begin{array}{c}
 \text{C-FORALL-DISTR-L} \\
 \frac{\text{split}_{FV(\tau_1^-)}(\Sigma) = (\Sigma_0, _, \Sigma_1) \quad V_0 = V \cap FV(\Sigma_0) \quad V_1 = V \setminus V_0 \quad \Xi_1 \gg \Delta \cdot (\forall V_0\{\Sigma_0\}. \tau_1^- \rightarrow \forall V_1\{\Sigma_1\}. \tau_2^+ \leq^\phi \tau_3^-) \gg \Xi_2}{\Xi_1 \gg \Delta \cdot (\forall V\{\Sigma\}. \tau_1^- \rightarrow \tau_2^+ \leq^\phi \tau_3^-) \gg \Xi_2}
 \end{array}$$

C-FORALL-DISTR-R is straightforward: it simply pushes quantifiers out of arrow-type right-hand sides so as to make skolemization happen as early as possible. **C-FORALL-DISTR-L** does the reverse, pushing quantifiers into arrow-type right-hand sides to delay their instantiation. The latter is more complicated to specify because it has to deal with splitting multi-bound quantifiers, filtering out those type variables whose bounds transitively refer to the arrow type’s left-hand side and which can therefore not be distributed. These rules are straightforwardly generalized to look deeply inside arrow types, so that they can push quantifiers across multiple nested arrow types; while this is implemented in our prototype, we omit the formal details here for the sake of brevity.

Limitations. Distributivity reasoning in SuperF is unlike the rest of the type inference system, which is why we describe it in its own subsection. The application of all other type inference rules is “uncontroversial”, in the sense that there is no obvious better thing to do in the context where each rule applies. On the other hand, whether to apply distributivity is a choice, and the choice is non-obvious: while applying it is *usually* beneficial (so in practice, we always apply it eagerly whenever possible), it can occasionally lead to *worse* type inference outcomes, yielding errors when a type inference run without distributivity would have succeeded.

For example, consider the following program, which contains a type ascription on subterm a :

```
test f = let a () = f () in a :  $\tau \rightarrow (\forall b. b \rightarrow b)$ 
```

When typing $f()$, SuperF creates a fresh type variable α and constrains the type of f to be a subtype of $\tau \rightarrow \alpha$. But because the type variable of f is *less polymorphic*, that type is *extruded* into $\tau \rightarrow \beta_\alpha$ where $\alpha \geq \beta_\alpha$ and β_α is a fresh type variable at the same (outer) polymorphism level as f . Binding a is then assigned the generalized function type $\tau = \forall \alpha\{\beta_\alpha \leq \alpha\}. \tau \rightarrow \alpha$. Then, when typing the ascription in the body of the let binding, a constraint of the form $\tau \leq \tau \rightarrow \sigma_{id}$ is generated. SuperF finds that the left-hand side is a polymorphic type and the right-hand side is an arrow whose right-hand side is polymorphic, so it has two choices:

- Distribute the latter out (i.e., what the implementation does), resulting in $\tau \leq \forall \gamma. \tau \rightarrow \gamma \rightarrow \gamma$, which can be skolemized to $\tau \leq \tau \rightarrow \gamma^{sk} \rightarrow \gamma^{sk}$, after which the left-hand side is finally instantiated to $\tau \rightarrow \alpha'$ where $\alpha' \geq \beta_\alpha$. So the constraint results in adding upper bound $\tau \rightarrow \gamma^{sk} \rightarrow \gamma^{sk}$ to the local flexible variable α' and propagating that bound to β_α . But β_α is a less polymorphic (rigid) type variable, so the bound is extruded to $\tau \rightarrow \tau \rightarrow \perp$, an approximation that may later lead to type errors.

- Refrain from distributing and instead instantiate the LHS, resulting in $\top \rightarrow \alpha' \leq \top \rightarrow \sigma_{\text{id}}$ where $\alpha' \geq \beta_\alpha$, which decomposes to $\alpha' \leq \sigma_{\text{id}}$. This results in adding bound σ_{id} to α' and propagating that bound to β_α . Notice that this time, no extrusion/avoidance is performed, because we kept the right-hand side nested polymorphic type intact.

So while distributive reasoning is often useful, there are cases where it is in fact counter-productive. The “dummy” LHS polymorphic type τ , which just wraps a plain type variable with a single lower bound, leads SuperF to distribute the RHS although that is not necessary to solve the constraint successfully, which in turn leads to type extrusion. To remain predictable, SuperF does not perform any non-local reasoning on whether performing distributivity is likely to be beneficial or not, and so it fails in this example. It is an open question whether we can determine when distributivity is unnecessary, or at least whether we can design a heuristic to avoid with high accuracy the cases when it is likely harmful.

4.5 Metatheory

The following theorem captures the main correctness criterion for type inference:

THEOREM 4.5 (SOUNDNESS OF TYPE INFERENCE). *If $\vdash t : \tau^+ \Rightarrow \Delta$ and $\vdash \epsilon \gg \Delta \gg \Xi'$, then we have $\vdash t : \forall V\{\Xi\}. \tau^+$, where $\Xi = \text{uproot}(\Xi')$, and $V = \text{FV}(\Xi)$.*

If type inference infers a type for a closed term, we obtain a typing that can be translated to System F_{cc} . We state the accompanying correctness theorem of constraining as follows:

Definition 4.6 (Constraint entailment). With $\Xi \vdash \Delta$ we denote that $\overline{\Xi \vdash \tau \leq \sigma}^{(\tau \leq \sigma) \in \Delta}$.

THEOREM 4.7 (SOUNDNESS OF CONSTRAINING). *Let $V, \bar{\alpha} \vdash \epsilon \gg \Delta \gg \Sigma'$ and $\Sigma = \text{uproot}(\Sigma')$. Then for all $\bar{\tau}$ where $\text{FV}(\bar{\tau}) \# \bar{\beta}$ and $\bar{\beta} = \text{FV}(\Sigma) \setminus \{V \cdot \bar{\alpha}\}$ we have*

- $\text{outer}_V(\Sigma) \cdot \text{inner}_V([\bar{\alpha} \mapsto \bar{\tau}]\Sigma) \vdash [\bar{\alpha} \mapsto \bar{\tau}]\Delta$ and
- $\text{outer}_V(\Sigma) \vdash \forall \bar{\beta}\{\text{inner}_V([\bar{\alpha} \mapsto \bar{\tau}]\Sigma)\} \text{ cons.}$

This theorem is only stated for derivations whose assumptions are empty ϵ ; this matches how constraining is used in **I-ABS** and **C-FORALL-R**. Intuitively, the results Σ of such derivations satisfy two properties: they entail (or resolve) the input constraints Δ , and their inner flexible variable bounds $\text{inner}_V(\Sigma)$ are consistent if we assume the outer rigid variable bounds $\text{outer}_V(\Sigma)$. These properties hold for all skolem replacement $[\bar{\alpha} \mapsto \bar{\tau}]$, which is crucially important for **C-FORALL-R**. We prove Theorems 4.5 and 4.7 in Appendix D.

5 PRACTICAL IMPLEMENTATION OF SUPERF

In this section, we describe MLscriptF, our implementation of SuperF in the existing MLscript language, and empirically evaluate it on existing test suites and examples from previous work literature. There are a few inessential differences between MLscriptF and SuperF as formalized in Section 4, the main ones being that MLscriptF:

- supports type checking recursive definitions, which is not shown in Section 4 but is straightforward: we bind the definition’s name to a fresh type variable while typing the definition’s body and then constrain that type variable to be a supertype of the inferred body type. We also refrain from generalizing the functions constituting the bodies of these recursive definitions, as that would lead to polymorphism extrusion and unnecessary failures.²¹

²¹Another way of understanding this is that we do not support type inference for polymorphically-recursive functions.

- does not generalize lambdas nested under other lambdas (as in curried functions); doing so is redundant since distributivity can always be used to push polymorphic types back inside arrow types when needed. More subtly, we also do not want to incur *too much* polymorphism in the presence of recursive terms, as that would degrade the type checker’s performance.
- uses explicit *polymorphism levels* to track extrusion and type avoidance, as described in Section 2.4. This can be seen as an optimization (to minimize the number of type traversals) that does not affect the functional properties of the system.
- has top-level `def` bindings with call-by-name semantics which are always generalized.

5.1 Implementation and Validation on MLscript Test Suite

We have implemented our approach in MLscript, a new ML-family language currently in active development. The MLscript codebase already had a significant amount of regression tests written for it, amounting to about 8000 lines of test *code* (not counting comments nor blank lines), which runs (in parallel) in about 7 seconds on a recent x86 MacBook laptop.²² Many of these tests produce statistics, such as the number of calls to the subtype constraining method, to check that the amount of work remains reasonable. About 70% of these tests were written before implementing first-class polymorphism (FCP); they test various aspects of the language through small functions constructed in it. This makes us confident that the addition of our FCP technique to an existing language is practical and does not introduce performance problems in the form of, for example, pathological cases that would blow up the time spent type checking a given file.

5.2 Validation on ML^F Test Suite

The test suite of ML^F was graciously provided to us by its authors [Le Botlan and Rémy \[2009\]](#). It provided us with a lot of interesting and tricky examples, which SuperF all managed to type check. Because ML^F is the previously-unbeaten champion in FCP type inference expressiveness, we are quite satisfied that we could handle its test suite better than ML^F itself, which required manual type annotations in more places.

5.3 Validation on Existing OCaml Code

We have experimentally evaluated our claim that SuperF subsumes ML type inference by porting the `List` module implementation from OCaml’s standard library (`list.ml`), which is about 500 lines of non-empty, non-comment lines of ML code. Our implementation

- did not require any type annotations;
- inferred relatively concise/compact and readable types;
- inferred types that were often more precise than the OCaml ones, thanks to subtyping; and
- inferred types which were successfully *checked against* the explicit type signatures provided in the OCaml standard library’s `list.mli` interface file.

These claims are substantiated in supplementary material (see the file `OCamlList.md`).

5.4 Evaluation on Examples from The Literature

Table 2 presents all the examples that we could gather from previous work on FCP. (Most of them were already summarized by [Emrich et al. \[2020\]](#) for their paper on FreezeML.) We observe that

²²All running times in this section include the time to actually *run* the tests (which is done via JavaScript code generation and executed through *nodejs*), as MLscript is a real programming language and not just a type checker.

Table 2. Typeability comparisons in different systems. ‘●’ means the term type checks; ‘●*’ means its inferred type is not as general as it could be (assuming SuperF subtyping); ‘⦿’ means the term type checks only after adding some type-free *polymorphism annotations* (e.g., FreezeML’s freezing annotations); ‘○’ means a type error is raised. SuperF–D is SuperF *without distributivity*, while SuperF has distributivity.

	MLF	FPH	HMF	HML	GI	FreezeML	SuperF–D	SuperF
A. Polymorphic instantiation								
A1 $\lambda x. \lambda y. y$	●	●*	●	●	●	●*	●	●
A2 choose id	●	●*	●*	●	●	●*	●	●
A3 choose nil ids	●	●	●	●	●	●	●	●
A4 $\lambda x. x x$	○	○	○	○	○	○	●	●
A5 id auto	●	●	●	●	●	●	●	●
A6 id auto'	●	●	●	●	●	●*	●	●
A7 choose id auto	●	○	●	●	●	●	●	●
A8 choose id auto'	○	○	○	○	○	○	●	●
A9 f (choose id) ids	●	●	○	●	○	⦿	●	●
where $f : \forall \alpha. (\alpha \rightarrow \alpha) \rightarrow \text{List } \alpha \rightarrow \alpha$								
A10 poly id	●	●	●	●	●	⦿	●	●
A11 poly $(\lambda x. x)$	●	●	●	●	●	⦿	●	●
A12 id poly $(\lambda x. x)$	●	●	●	●	●	⦿	●	●
B. Inference with polymorphic arguments								
B1 $\lambda f. (f \top, f \text{ True})$	○	○	○	○	○	○	●	●
B2 $\lambda xs. \text{poly}(\text{head } xs)$	●	○	○	○	○	○	●	●
C. Functions on polymorphic lists								
C1 length ids	●	●	●	●	●	●	●	●
C2 tail ids	●	●	●	●	●	●	●	●
C3 head ids	●	●	●	●	●	●	●	●
C4 single id	●	●	●*	●	●	●*	●	●
C5 cons id ids	●	●	●	●	●	⦿	●	●
C6 cons $(\lambda x. x)$ ids	●	●	●	●	●	⦿	●	●
C7 append (single inc) (single id)	●	●	●	●	●	●	●	●
C8 g (single id) ids	●	●	○	●	○	⦿	●	●
where $g : \forall \alpha. \text{List } \alpha \rightarrow \text{List } \alpha \rightarrow \alpha$								
C9 map poly (single id)	●	●	○	●	○	⦿	●	●
C10 map head (single ids)	●	●	●	●	●	●	●	●
D. Application functions								
D1 app poly id	●	●	●	●	●	⦿	●	●
D2 revapp id poly	●	●	●	●	●	⦿	●	●
D3 runST argST	●	●	●	●	●	⦿	●	●
D4 app runST argST	●	●	●	●	●	⦿	●	●
D5 revapp argST runST	●	●	●	●	●	⦿	●	●
E. η-expansion								
E1 $k \ h \ l$	●	○	○	○	○	○	●	●
E2 $k (\lambda x. h \ x) \ l$	●	●	●	●	●	⦿	●	●
where $k : \forall \alpha. \alpha \rightarrow \text{List } \alpha \rightarrow \alpha$ $h : \text{Int} \rightarrow \forall \alpha. \alpha \rightarrow \alpha$ $l : \text{List } (\forall \alpha. \text{Int} \rightarrow \alpha \rightarrow \alpha)$								
E3 $r (\lambda x. \lambda y. y)$	●	○	○	●	○	⦿	●	●
where $r : (\forall \alpha. \alpha \rightarrow \forall \beta. \beta \rightarrow \beta) \rightarrow \text{Int}$								
F. FreezeML paper additions								
F5 auto id	●	●	●	●	●	⦿	●	●
F6 cons (head ids) ids	●	○	●	●	●	●	●	●
F7 head ids 3	●	●	●	●	●	⦿	●	●
F8 choose (head ids)	●	●	●	●	●	●	●	●
F9 let $f = \text{revapp id in } f \text{ poly}$	●	○	○	●	○	⦿	●	●
F10 choose id $(\lambda x. \text{auto}' \ x)$	●	○	○	○	○	○	●	●

SuperF seamlessly handles *all* these previous examples without the need for any type or other polymorphism/freezing annotations, with or without distributivity.

head : $\forall \alpha. \text{List } \alpha \rightarrow \alpha$	id : $\forall \alpha. \alpha \rightarrow \alpha$	map : $\forall \alpha \beta. (\alpha \rightarrow \beta) \rightarrow \text{List } \alpha \rightarrow \text{List } \beta$
tail : $\forall \alpha. \text{List } \alpha \rightarrow \text{List } \alpha$	ids : $\text{List } (\forall \alpha. \alpha \rightarrow \alpha)$	app : $\forall \alpha \beta. (\alpha \rightarrow \beta) \rightarrow \alpha \rightarrow \beta$
nil : $\forall \alpha. \text{List } \alpha$	inc : $\text{Int} \rightarrow \text{Int}$	revapp : $\forall \alpha \beta. \alpha \rightarrow (\alpha \rightarrow \beta) \rightarrow \beta$
cons : $\forall \alpha. \alpha \rightarrow \text{List } \alpha \rightarrow \text{List } \alpha$	choose : $\forall \alpha. \alpha \rightarrow \alpha \rightarrow \alpha$	runST : $\forall \alpha. (\forall s. \text{ST } s \alpha) \rightarrow \alpha$
single : $\forall \alpha. \alpha \rightarrow \text{List } \alpha$	poly : $(\forall \alpha. \alpha \rightarrow \alpha) \rightarrow (\text{Int}, \text{Bool})$	argST : $\forall s. \text{ST } s \text{Int}$
append : $\forall \alpha. \text{List } \alpha \rightarrow \text{List } \alpha \rightarrow \text{List } \alpha$	fst : $\forall \alpha \beta. (\alpha, \beta) \rightarrow \alpha$	zero : ChurchInt
length : $\forall \alpha. \text{List } \alpha \rightarrow \text{Int}$	auto : $(\forall \alpha. \alpha \rightarrow \alpha) \rightarrow (\forall \alpha. \alpha \rightarrow \alpha)$	succ : $\text{ChurchInt} \rightarrow \text{ChurchInt}$
const : $\forall \alpha \beta. \alpha \rightarrow \beta \rightarrow \alpha$	auto' : $\forall \beta. (\forall \alpha. \alpha \rightarrow \alpha) \rightarrow (\beta \rightarrow \beta)$	$\text{ChurchInt} = \forall \alpha. (\alpha \rightarrow \alpha) \rightarrow (\alpha \rightarrow \alpha)$
$z = \lambda f. \lambda x. x$	$s = \lambda n. \lambda f. \lambda x. f (n f x)$	$n3 = s (s (s z))$

Fig. 7. Type signatures and definitions for terms used in Table 2 and Table 3.

Table 3. Typeability comparisons (continued). Legends are in Table 2; ‘x’ means non-termination or crash.

	MLF	FPH	HMF	HML	GI	FreezeML	SuperF-D	SuperF
G. SuperF additions (this paper)								
G1A $z : \text{ChurchInt}$	●	●	●	●	●	●	●	●
G2 s	●*	●*	●*	●*	●*	●	●*	●
G2A $s : \text{ChurchInt} \rightarrow \text{ChurchInt}$	○	○	○	○	○	●	○	●
G3A $n3 : \text{ChurchInt}$	●	●	●	●	○	●	○	●
G4A $(\text{fun } () \rightarrow n3 \ n3) : () \rightarrow \text{ChurchInt}$	●	●	●	●	○	○	○	●
G5 $\text{fst} (\text{fst} (\text{fst} (n3 (\lambda x. (x, \emptyset)) 1)))$	○	○	○	○	○	○	●	●
G6 $(\text{succ} (\text{succ} \text{zero})) (\text{succ} (\text{succ} \text{zero}))$	●	●	●	●	○	●	●	●
G7 $(s (s z)) (s (s z))$	●	●	●	●	○	●	●	●
G8 let rec to_church $n = \text{if } n == 0 \text{ then } z \text{ else else } s (\text{to_church } (n - 1))$	●*	●*	●*	●*	○	●	●*	●
G8A $\text{to_church} : \text{Int} \rightarrow \text{ChurchInt}$	○	○	○	○	○	●	○	●
G9 let rec $i \ x = \text{if True then } x \text{ else } i \ i \ x$	○	○	○	○	x	○	○	○
G10 $(\lambda x. x \ x) (\lambda x. x \ x)$	○	○	○	○	○	○	○	○
G11 $\text{auto} (\text{auto}' \ \text{id})$	●	●	●	●	●	○	○	●
G12 $(\lambda y. \text{let tmp} = y \ \text{id in } y \ \text{const}) (\lambda x. x \ x)$	○	○	○	x	○	○	●	●
G13 $(\lambda k. (k (\lambda x. x), k (\lambda x. \text{single } x))) (\lambda f. (f \ 1, f \ \text{True}))$	○	○	○	○	○	○	○	○
G14 $(\lambda f. \text{let } a() = f \ \text{id in } (a()) : \text{Int} \rightarrow (\forall 'b. 'b \rightarrow 'b)) (\text{const } (\text{const } \ \text{id}))$	●	○	○	○	○	○	●	○

5.5 New Examples

Because the examples presented in previous work were too simple to truly test the limits of SuperF, we created a number of more intricate expressions in Table 3, several of which come directly from the test suite of ML^F. Note that SuperF does not use bidirectional type checking, so for an annotated term $t : \tau$ to type check, t itself must type check. Therefore, the annotation in, e.g., $s : \text{ChurchInt} \rightarrow \text{ChurchInt}$ does *not* help type inference — it only checks that the type inferred for s can be instantiated to the less precise type $\text{ChurchInt} \rightarrow \text{ChurchInt}$.

SuperF cannot type check examples G2A, G4A, G8A, and G11 without distributivity. For example, in G11, distributivity pushes a quantifier into the result type of auto' , making it polymorphic, so it can be passed to auto . G4A is interesting because it is type checked successfully by MLF, FPH, and HML but not by SuperF without distributivity. This is due to SuperF only generalizing lambda abstractions, whereas these other systems generalize even applications (see also Section 2.7).

G14 is a contrived example manually constructed to demonstrate that distributivity in conjunction with subtype extrusion and type annotations sometimes leads to worse outcomes: SuperF distributes $\text{Int} \rightarrow (\forall 'b. 'b \rightarrow 'b)$ into $\forall 'b. \text{Int} \rightarrow ('b \rightarrow 'b)$, leading to the premature skolemization of $'b$ and to more extrusion than without distributing.

The definition in G9 uses itself recursively at different types and therefore needs polymorphic recursion, for which SuperF would require a type annotation — adding the explicit type signature

‘ $i : \forall a. a \rightarrow a$ ’ makes it type check. Interestingly, it is alternatively possible to type check this example using recursive types, which our implementation also supports (behind a feature flag).

6 RELATED WORK

We now discuss some related work.

6.1 First-Class Polymorphism in ML-style languages

First-class polymorphism has been studied extensively in the setting of ML-style languages, where polymorphism is normally implicit.

Boxed polymorphism and related systems. *Poly-ML* [Garrigue and Rémy 1997] and its predecessors [Läufer and Odersky 1994; O’Toole and Gifford 1989; Rémy 1994] augment ML type schemes with boxed polytypes. In Poly-ML, boxed polytypes wrap type schemes into simple ML types using the syntax $[\sigma]$, preventing implicit instantiation. The introduction form of boxed polytypes requires explicit annotation of the type scheme ($[t : \sigma]$). Boxed polytypes must be eliminated explicitly at the use site with the syntax $\langle t \rangle$ to instantiate the type scheme. The approach never infers first-class-polymorphic types, but only propagate the polytypes provided by the user. However, this juxtaposition of ML-style type inference and System F-style type checking means that the order of the terms being processed can affect the result of type inference, since extra type information may be obtained from type inference. When a let-binding or an application is typed, it is possible that one of the two subterms must be typed before the other for the other to be typeable. Since there is no way to tell in advance if a type inference path would succeed, the algorithm must backtrack and explore the other paths when it fails. The number of such paths is usually exponential in the size of the program, leading to combinatorial explosion. For example, while typing $\lambda f. \langle f \rangle (g f)$, where $g = \lambda f : [\forall \alpha. \alpha \rightarrow \alpha]. \langle f \rangle f$, typing the former $\langle f \rangle$ would fail since the type of f is a type variable, not a boxed polytype. However, typing the argument $(g f)$ first would force the type of f to be $[\forall \alpha. \alpha \rightarrow \alpha]$, which then allows $\langle f \rangle$ to be typed. To address this, the authors chose to reject programs where any of the paths may fail. This rejects a wide range of otherwise typeable programs and requires explicit type annotations on lambda parameters. The authors claim to have achieved complete type inference by incorporating the notion of “labels” into the declarative system, where each boxed polytype is associated with a label. A user-provided type annotation implicitly quantifies over all labels. A boxed polytype can only be opened when its label does not occur anywhere else, which implies that it must originate from an annotation. In SuperF, the unannotated version of g can be typed as $\lambda f. f f : \forall \alpha \beta \{ \alpha \leq \alpha \rightarrow \beta \}. \alpha \rightarrow \beta$, and the unannotated version of the problematic term above as $\lambda f. f (g f) : \forall \alpha \beta \gamma \{ \alpha \leq \alpha \rightarrow \beta, \alpha \leq \beta \rightarrow \gamma \}. \alpha \rightarrow \gamma$, which is a subtype of the only annotated version typeable in Poly-ML, $\lambda f : [\forall \alpha. \alpha \rightarrow \alpha]. \langle f \rangle (g f)$, with type $[\forall \alpha. \alpha \rightarrow \alpha] \rightarrow [\forall \alpha. \alpha \rightarrow \alpha]$.

ML^F. Le Botlan and Rémy [2003] introduced ML^F, which has been the “undefeated champion” of type inference for first-class polymorphism *until now*. ML^F uses a more flexible type language than System F in that it includes bounded quantification, whereby types are ordered based on their polymorphism “power”, akin to our subtyping relation. ML^F takes inspiration from boxed polymorphism and attempts to alleviate the verbose introduction and elimination of boxes. This is achieved by introducing two forms of type variable bindings, namely rigid bindings ($\alpha \Rightarrow \sigma$) and flexible bindings ($\alpha \geq \sigma$), which can appear in a context called the “prefix” and in binders for bounded polymorphism. The usual subsumption relation is split into an abstraction relation Ξ and an instance relation \sqsubseteq . The Ξ relation is used to implicitly abstract polytypes as type variables with a rigid binding in the prefix, which allows the polymorphism to propagate by preventing implicit instantiation (cf. boxed polymorphism in Poly-ML). For example, consider the function

$a = \lambda z. \omega^\dagger z$, where $\omega^\dagger = \lambda(x : \sigma_{id}) x x$. It is clear that z must have a polymorphic type, but it is not used polymorphically in a , its polymorphism is merely passed to ω^\dagger . Abstracting σ_{id} as α in the prefix allows the body of a to be typed by $(\alpha \Rightarrow \sigma_{id}) z : \alpha \vdash \omega^\dagger z : \alpha$ without ever guessing polymorphism. Then a can be typed as $\forall(\alpha \Rightarrow \sigma_{id}) \alpha \rightarrow \alpha$. Moreover, as seen in the type of a , rigid bindings keep track of sharing, which is crucial to type inference as it allows instantiation to be delayed. To reveal and instantiate an abstraction, the programmer must explicitly provide a type annotation that corresponds to the rigid binding. Contrary to the reversible Ξ abstraction, the instance relation \sqsubseteq is used to irreversibly instantiate a type scheme, subject to the flexible bindings. Intuitively, rigid bindings represent the types of terms that must be polymorphic, while flexible bindings represent the types that may be instantiated. For example, the principal type of ω^\dagger is $\forall(\alpha \Rightarrow \sigma_{id}, \alpha' \geq \sigma_{id}) \alpha \rightarrow \alpha'$, where the parameter type α must be polymorphic since it is used polymorphically in the body, but the result type α' may be instantiated to any instance of σ_{id} . In ML^F , lambda parameters that are used polymorphically require annotations, since type inference does not guess second-order types. SuperF does not guess second-order polymorphic types either, but uses multiple bounds (which can be represented as intersections) to accommodate for finitarily-polymorphic uses of unannotated parameters.

Le Botlan and Rémy note that ML^F infers principal types for types which are typeable in their system, but that typing some terms may require type annotations – and in fact, different annotations may lead to incomparable principal types. For example, they remark that $\lambda x. x x$ is not typeable in ML^F while both $\lambda(x : \forall\alpha.\alpha) x x$ and $\lambda(x : \forall\alpha.\alpha \rightarrow \alpha) x x$ are, but neither has a more general type (their types are unrelated). On the other hand, $\lambda x. x x$ is typeable in SuperF, inferring type $\forall\alpha\beta. ((\alpha \rightarrow \beta) \wedge \alpha) \rightarrow \beta$ (i.e., $\forall\alpha\beta\gamma\{\gamma \leq \alpha \rightarrow \beta, \gamma \leq \alpha\}. \gamma \rightarrow \beta$ or equivalently $\forall\alpha\beta\{\alpha \leq \alpha \rightarrow \beta\}. \alpha \rightarrow \beta$), which is a subtype of *both* ML^F types mentioned above.

The bindings in prefixes and quantifiers in ML^F resemble constrained types in SuperF. However, a few distinctive features in SuperF contribute to its greater expressiveness over ML^F . Firstly, function types are contravariant in their parameter types in SuperF, whereas they are invariant in their parameter types in ML^F . Contravariance in the subtyping relation is crucial to type inference. For example, consider the two terms $f = \lambda(x : \sigma_{id} \rightarrow \sigma_{auto}). x$ and $g = \lambda(x : \sigma_{auto}). id$, where $id = \lambda x. x$, $\sigma_{id} = \forall\alpha. \alpha \rightarrow \alpha$ and $\sigma_{auto} = \sigma_{id} \rightarrow \sigma_{id}$. The application $f g$ is typeable in a straightforward manner in SuperF thanks to subtyping. The same term is, however, not typeable in ML^F . Secondly, as a result of the restriction to monotypes under type constructors, ML^F requires the introduction of “administrative” type variable bindings even when these bindings are used once and do not represent any shared information. Let us consider the following example. The inferred type of id is the expected $\sigma_{id} = \forall\alpha. \alpha \rightarrow \alpha$ in both SuperF and ML^F . One may wish to instantiate α to some concrete type, say σ_{id} . In SuperF, id can be typed as $\sigma_{id} \rightarrow \sigma_{id}$, which is $\alpha \rightarrow \alpha$ with σ_{id} substituted for α . However, in ML^F , since only monotypes are allowed under an arrow, the type scheme σ_{id} must be abstracted in a binding as $\forall(\alpha \Rightarrow \sigma_{id}) \alpha \rightarrow \alpha$. One may expect id to also be typeable as $\forall(\alpha \Rightarrow \sigma_{id}) \forall(\beta \Rightarrow \sigma_{id}) \alpha \rightarrow \beta$, or even the more general $\sigma_{id} \rightarrow \sigma_{id}$, which is a shorthand for $\forall(\alpha \Rightarrow \sigma_{id}) \forall(\alpha' \geq \sigma_{id}) \alpha \rightarrow \alpha'$. This is, however, not the case as ML^F keeps track of sharing even after the $\alpha \geq \perp$ binding has been instantiated and rigidified to $\alpha \Rightarrow \sigma_{id}$. The arguably counterintuitive pseudo-subtyping relation in ML^F imposes some burden on users, and may catch them by surprise. Finally, SuperF allows multiple bounds on the same variable, including multiple upper bounds and multiple lower bounds, whereas ML^F only allows one lower or rigid bound on a variable. Supporting multiple bounds for each type variable allows us to effectively infer intersection types in negative positions (such as function parameters) and union types in positive positions (such as function results).

Variations on ML^F . Rémy and Jakobowski [2007] proposed a graphical representation of ML^F types that allows for an efficient unification algorithm. Le Botlan and Rémy further investigated a “shallow” version of ML^F where type annotations and rigid bounds are restricted to System F types [Le Botlan and Rémy 2009] which, while less expressive, still has interesting properties, and more importantly admits simpler semantics and metatheory. Finally, Rémy and Jakobowski [2012] designed xML^F , a “Church-style” intermediate language for ML^F which would be amenable to intrinsically-typed compilation as in compilers like the Glasgow Haskell Compiler (GHC).

Restrictions on ML^F . *HML* [Leijen 2009] is a simplification of ML^F restricting type variable bindings to flexible bounds (discarding rigid bounds). *HML* retains the expressiveness of by requiring annotations on polymorphic function arguments, which is marginally higher than the requirement in ML^F . The author believes that this requirement to programmers is justified by the simpler types they work with. Like ML^F , *HML* is robust against most program transformations, including abstracting and inlining let-bindings, and abstraction with higher-order functions (most notably with `app` and `revapp`). An exception to this is η -expansion, which may require a type annotation on a polymorphic parameter type. The author believes that robustness is an important property as it forms the basis of equational reasoning. *SuperF* enjoys even more robustness thanks to there being no requirement to annotate polymorphic parameters, thanks to constrained types, which makes η -expansion robust as well as the other transformations.

HMF [Leijen 2008] is another, more drastic restriction of ML^F discarding flexible bindings, i.e., a simple extension of the Hindley-Milner type system with System F types, supporting type inference through a straightforward extension of algorithm W. The author believes that a simple type system is beneficial in that it eases the burden on programmers in understanding the types, and in simplifying the metatheory and the implementation of type inference. To achieve this, *HMF* always predicatively instantiates ambiguous applications eagerly. Annotations are thus required when impredicative instantiation is needed, in addition to the usual annotation requirement on function parameters that are used polymorphically. This has the adverse effect of introducing a dependency of typability on the order of parameter, which is alleviated in some cases by considering application chains all at once, instead of treating each of their constituent applications individually. This system has the major drawbacks that type inference is less stable than in ML^F -style systems.

FPH [Vytiniotis et al. 2008] is yet another proposed type inference system which uses System F types but can still be understood as a restriction of ML^F . *QML* [Russo and Vytiniotis 2009] extends ML types with universal and existential quantified types, which co-exist with type schemes. While type schemes are *implicitly* introduced by let-bindings and eliminated when applications happen, quantified types are *explicitly* introduced and eliminated by type annotations. For example, let $id = \lambda x. x$ in id will be typed as a type scheme $\Pi(\alpha)(\alpha \rightarrow \alpha)$. To have a first-class-polymorphic type, an explicit type annotation is necessary as let $pid = \{\forall\alpha. \alpha \rightarrow \alpha\}$ in pid , here pid has the type $\Pi(\epsilon)(\forall\alpha. \alpha \rightarrow \alpha)$. On the other hand, it is also necessary to provide an explicit type annotation to eliminate this quantified type as $(pid \{\forall\alpha. \alpha \rightarrow \alpha\})$ 3. Like many systems, *QML* also does not guess polymorphic types for function parameters. Parameters can only be polymorphic when they are annotated with quantified types. As an example, let $poly = \lambda f. \text{let } y = f \{\forall\alpha. \alpha \rightarrow \alpha\} \text{ in } (y \ 1, y \ \text{true})$ in $poly$ has the type $(\forall\alpha. \alpha \rightarrow \alpha) \rightarrow (Int, Bool)$. The treatment of existentially-quantified types is similar, requiring explicit type annotations to introduce and eliminate these types.

FreezeML. *FreezeML* [Emrich et al. 2020, 2022] introduces a “freeze” operator to explicitly disable instantiation. Otherwise, like in ML, polymorphic types are implicitly instantiated. For example, *single id* where $single : \forall a. a \rightarrow \text{List } a$ and $id : \forall a. a \rightarrow a$ is typed as $\text{List } (a \rightarrow a)$ since id ’s type is implicitly instantiated. In contrast, with id frozen using syntax $[id]$, then $single [id]$ is typed

as $\text{List } (\forall a. a \rightarrow a)$. With the frozen operator and some macro expansions, users may explicitly control whether the types of terms are generalized or instantiated. FreezeML’s type inference algorithm extends the traditional algorithm W and achieves principal type inference of first-class polymorphism. However, it requires a significant number of annotations for polymorphic code: users need to manually control the explicit freezing, instantiation, and generalization of polymorphic types, *in addition to* the need of annotating the full type of polymorphic function parameters, which altogether amounts to a significant burden for programmers. It seems SuperF is more general than FreezeML – indeed, we conjecture that *most* (if not *all*) FreezeML programs can be typed by SuperF even after stripping all their freezing and instantiation annotations. For example, in FreezeML, the example E2 can only be typed with two explicit operators (as $k \$(\lambda x. (h x))@ l$) while SuperF requires no annotations at all. Table 2 in appendix shows that such annotations are required for the majority of examples from the existing literature, although these examples are all quite simple. FreezeML annotations, while more concise than having to specify full types, are arguably not very intuitive. The E2 annotations above are needed to explicitly adapt the quantification of h from $\text{Int} \rightarrow \forall a. a \rightarrow a$ to $\forall a. \text{Int} \rightarrow a \rightarrow a$, so it becomes compatible with a type that only happens to *flow into* the same result variable. Users are effectively reduced to retyping by η expansion and adding freezing/instantiation in a way that may seem quite obscure to whoever reads the code afterward. Furthermore, FreezeML does not support reordering quantifiers implicitly. For example, $\forall \alpha \beta. \alpha \rightarrow \beta$ and $\forall \alpha \beta \gamma. \alpha \rightarrow \beta$ and $\forall \beta \alpha. \alpha \rightarrow \beta$ are all considered distinct, incompatible types [Emrich et al. 2022]. One may use explicit generalization to restore a canonical order of quantifiers, but this still results in a quite inflexible system where many valid programs fail due to syntactic details, rather than semantic ones. While FreezeML can type the `foo` examples from the introduction using various sets of annotations, it cannot assign a principal type to it, so there is no best set of annotations to use for that example. In particular, no set of annotations can let one type check both `foo (fun x → x)` and `foo (fun x → Some x)` with the expected result types, whereas our system type checks both “out of the box”, without any annotations.

First-class polymorphism for Haskell. *Boxy types* [Vytiniotis et al. 2006] were used to implement GHC’s original `ImpredicativeTypes` extension. The system incorporates bidirectional type inference by distinguishing types to be inferred and types to be checked using boxes. Boxy-types retains the idea of not guessing polymorphic types, but refines it by inferring polytypes for function arguments with locally known type information. The system as presented was very fragile against local transformations and thus difficult to use in practice, which is why it was eventually replaced by the more predictable *Quick Look*. For example, typability is not preserved under η -expansion.

Guarded instantiation (GI) [Serrano et al. 2018] focuses on simplicity and tries to balance between complexity, expressiveness, and annotation burden. In GI, polymorphic instantiation only happens on guarded type variables. The guardedness is decided in function applications by examining whether the type variable is guarded by type constructors in any parameter types, for example both type variables in $\forall a b. (a \rightarrow b) \rightarrow \text{List } a \rightarrow \text{List } b$ are guarded but neither in $\forall a b. a \rightarrow b$ is. Unguarded type variables may only be instantiated by monomorphic types or polymorphic types without top-level quantifiers, such as $\text{Int} \rightarrow \forall a. a$. Note that only the types of those parameters that are given corresponding arguments are considered. As an example, in `append id` where `append` : $\forall a. a \rightarrow \text{List } a \rightarrow \text{List } a$, though a is guarded in the second parameter type, in this partial application only the first parameter type is considered for the guardedness, so a cannot be polymorphically instantiated. In `append id ids`, the second argument is now given so a is considered guarded and therefore may be instantiated polymorphically. Type instantiation in GI takes no account of the context of the call and each argument in a call is treated independently between different function

applications. As usual, for lambdas, GI will not guess polymorphic parameters. Parameters may only be polymorphic when annotated.

Quick Look [Serrano et al. 2020] is a type inference system for Haskell that further simplifies the first-class-polymorphic type system. The authors claim that it is eminently practical for Haskell programming. However, it is quite inexpressive when compared to most systems mentioned previously, and in particular it removes all forms of subtyping (so-called “deep subsumption”) from Haskell, needing manual η expansions instead even to just reorder or distribute quantifiers. Since the official support of Quick Look in GHC around 2021, users of GHC have already complained about its inflexibility, as it makes them manually eta-expand code in many places where they did not need to before and where it appears “clear” that the compiler should be able to figure things out.²³

Predicative higher-rank polymorphic type inference. Type inference can become much easier if one is to abandon the *impredicative* setting of first-class polymorphic type inference, i.e., if one separates the syntax of types between polymorphic and monomorphic types, only allowing the latter to be substituted for type variables. Odersky and Läufer [1996] demonstrated an example of this approach in their seminal work on type inference for higher-rank polymorphism. Many works have followed up on the idea, bringing various improvements [Dunfield and Krishnaswami 2013, 2019; Peyton Jones et al. 2007; Vytiniotis et al. 2006; Xie 2021; Xie and Oliveira 2018; Xue and Oliveira 2021] many of which relying on the idea of bidirectional typing [Dunfield and Krishnaswami 2021; Dunfield and Pfenning 2004].

Stability considerations. Le Botlan and Rémy [2009] originally investigated the stability of first-class polymorphic type inference in the context of ML^F . As Table 1 shows, we support most of their stability guarantees and some more. Bottu and Eisenberg [2021] recently tried to characterize such a notion of stability for Haskell. They developed 6 “similarity” properties and showed that in the context of GHC, all combinations of *shallow or deep* and *eager or lazy* instantiation break some of these similarities, though the *lazy and shallow* combination breaks the least. Of these similarities, one does not hold in SuperF due to the restricted language of annotations (similarity 2, about inferred types being usable as type signatures). The rest of these similarities trivially pass the stability criterion in SuperF, although SuperF is *deep* (and *not shallow*) as well as being *lazy*.²⁴

6.2 Type Inference for System F

Boehm [1985] and Pfenning [1988] proposed partial type inference procedures based on second-order unification which is undecidable but has a practical semi-algorithm due to Huet [1975]. Their systems required placeholder annotations (without type information) for all type abstractions and type applications, which they called a *partial type inference* approach. Mitchell [1988] introduced the now-usual notion of *type containment* for System F, a notion of subtyping involving polymorphic types (what was known as *subtyping* at the time did not include polymorphic subtyping). He showed how to make System F complete with respect to η expansion in System F_η , meaning all terms that System F can type after η expanding some of their subterms are typeable in System F_η through subsumption. Type inference for System F_η was later shown undecidable, due to an undecidable type containment relation Wells [1996]. Rémy [2005] proposed a simpler partial type inference system for System F based on a similar notion of type containment, but the approach was not very satisfactory, pushing Rémy to investigate the more promising approach of ML^F .

²³For example, see the following highly-upvoted reddit comment thread entitled “Was simplified subsumption worth it for industry Haskell programmers?” (accessed June 2022): <https://redd.it/ujpzx3>.

²⁴The reason Bottu and Eisenberg feel the need to be *shallow* in their approach rather than *deep* can be attributed to the presence of Haskell-specific features in their system, such as explicit type applications, and to their lack of subtyping.

In parallel, there has been much work on systems like F_{\leq} , which combined first-class polymorphism, type variable bounds, and subtyping, and even also sometimes included intersections [Pierce 1997]. However, these systems usually assume *explicit* polymorphism, whereby polymorphic types have to be introduced and eliminated in terms, and where subtyping between them is “structural,” i.e., there are no relationships between types such as $\forall a. a \rightarrow a$ and $\text{Int} \rightarrow \text{Int}$. This is by contrast to the traditional ML approach to polymorphism, which is implicit.

In this paper, we follow the ML tradition of implicit polymorphism, which is more appropriate as a foundation for generalizing the ML type system to first-class polymorphism. More specifically, we follow an approach closely inspired by *algebraic subtyping* [Dolan 2017], and its reformulation as *Simple-sub* [Parreaux 2020], including its novel approach to nested polymorphism with extrusion in the presence of subtyping.

6.3 Polyvariant Flow Analysis

Palsberg and Pavlopoulou [1998] showed that *polyvariant* analysis (also known as *context-sensitive* analysis) can be related formally to a subtyping system with union, intersection, and recursive types. Unions model sets of abstract values and intersections model each usage of an abstract value. Their system conspicuously does not feature polymorphism.

Faxén [1997]; Smith and Wang [2000] propose inferring polymorphic types rather than intersections for function definitions, which is more flexible and composable as it can process unrelated definitions separately, whereas the approach based solely on intersections is a global process. Smith and Wang develop a polymorphic system with subtype constraint solving designed for flow analysis. Their constraint closure rules are evocative of our **C-FLEX-L/C-FLEX-R**, **C-FUN**, and **C-FORALL-L**. However, they have no rule analogous to **C-FORALL-R**, nor do they need any notions of type avoidance and extrusion, which are proper to a *type inference* view of constraint solving (as opposed to the restrictive view of flow analysis). They give functions polymorphic types that contain all locally-inferred *unsolved subtyping constraints*.²⁵ The fact they do not solve constraints locally means they do not check for the consistency of inferred types (so they may fail to report type errors at definition sites) and it means their constraints are resolved over and over again every time the polymorphic types that capture them are instantiated. Similarly to us, their basic system is non-terminating and they introduce a termination condition which “detects a certain kind of self-referential flow in the constraints”. They show that termination can be guaranteed by merging some instantiations in this case. Wang and Smith [2001] adapt this approach to object-oriented programming but abandon the first-class polymorphism part.

Rehof and Fähndrich [2001] also rely on polymorphic subtyping; they design a control-flow analysis that can be reduced to the problem of context-free language reachability. They represent constraint sets as substitutions to avoid copying constraints around and obtain better scalability. Whether this insight can be applied to SuperF remains to be determined.

6.4 Intersection and mixed type systems

Jim’s polarized type system called “P” [Jim 2000] mixes universal quantification with intersection types in a way that strongly resembles ours, restricting the former to positive type positions and the latter to negative type positions. However, to retain decidability, the types of P are much more restricted syntactically. For example, they do not admit universal quantifiers nor intersections on the right of arrow types. P has special typing rules for variable and lambda applications. The former uses in its premise a *rewritten* lambda term (*not* a subterm of the term being typed), which makes

²⁵Faxén proposes simplifying these constraints incrementally, but this does not change the problem fundamentally.

is difficult to explain and characterize well-typed terms on an intuitive level.²⁶ Moreover, since the syntax of negative types does not include polymorphic types, users cannot check their programs against provided System F type signatures.

6.5 Implicit Coercion Constraints

Motivated by the goal of designing a type system that could generalize all previous bounded quantification approaches (mainly ML-style constrained types [Odersky et al. 1999], System $F_{<}$, [Pierce 1991], and ML^F [Le Botlan and Rémy 2003]), Cretin and Rémy [2014] developed a calculus of *implicit coercions* called System F_{cc} . Coercion constraints are very expressive and designed in a way that is more general than bounded quantification — notably, they do not require the use of a dedicated typing rule for generalization, lifting generalization to the coercion relation (which is a form of generalized subtyping). Thankfully, it is straightforward to encode our *multi-bounded polymorphism* type system into F_{cc} , as we demonstrate in Appendix C. Scherer and Rémy [2015] later investigated versions of coercion constraints allowing abstracting over possibly-inconsistent coercions, to support applications like GADTs [Xi et al. 2003].

7 CONCLUSION AND FUTURE WORK

We presented SuperF, a novel type inference algorithm for first-class polymorphism based on *multi-bounded polymorphism*. SuperF is uniquely expressive, stable under small program changes, and allows for understandable error messages to be reported. We also presented $F_{\{\leq\}}$, a simple type system that serves as the declarative basis of SuperF and that is founded on the theory of implicit coercions from the existing F_{cc} . SuperF was implemented as part of a real-world programming language and evaluated on previous test suites as well as examples from the literature.

SuperF could serve as a starting point for more advanced type systems features, such as higher-order subtyping, existentials and generalized algebraic data types (GADTs), as well as dependent types. For example, SuperF could provide an adequate algorithmic basis to the cDOT calculus of Boruch-Gruszecki et al. [2022] and to the λ_I^V calculus of Xue and Oliveira [2021], who left algorithmic formulations of their systems for future work. Unlike SuperF, λ_I^V uses *predicative* polymorphism. Restricting SuperF to predicative type instantiation should be straightforward and may allow dropping the SRLC. We would also like to support *unrestricted* type annotations in the future, which may be helped by integrating some form of bidirectional typing.

ACKNOWLEDGMENTS

We would like to thank Didier Rémy, Didier Le Botlan, Bruno Oliveira, Stephen Dolan, and the anonymous reviewers for their helpful feedback. This research was partially funded by Hong Kong Research Grant Council project number 26208821.

DATA AVAILABILITY STATEMENT

The artifact of this paper, which was granted the “Functional” and “Reusable” badges, is available on Zenodo at <https://zenodo.org/records/8424750> [Parreaux et al. 2023]. The corresponding Github repository is available at <https://github.com/hkust-taco/superf> and the system can be tried in a web demo at <https://hkust-taco.github.io/superf/>.

²⁶Generally, type inference can be greatly simplified if it is allowed to perform some steps of reduction in the program being considered. For example, ML type inference can be expressed as a simple type unification pass after reducing all let bindings. However, such *rewriting-based* approaches are not realistic and are not used in practice [Le Botlan and Rémy 2009].

REFERENCES

- Alexander Aiken and Edward L. Wimmers. 1993. Type Inclusion Constraints and Type Inference. In *Proceedings of the Conference on Functional Programming Languages and Computer Architecture (Copenhagen, Denmark) (FPCA '93)*. Association for Computing Machinery, New York, NY, USA, 31–41. <https://doi.org/10.1145/165180.165188> \hookrightarrow page 5
- Ishan Bhanuka, Lionel Parreaux, David Binder, and Jonathan Immanuel Brachthäuser. 2023. Getting into the Flow: Towards Better Type Error Messages for Constraint-Based Type Inference. *Proc. ACM Program. Lang.* 7, OOPSLA2, Article 237 (oct 2023), 29 pages. <https://doi.org/10.1145/3622812> \hookrightarrow page 14
- Richard S. Bird and Ross Paterson. 1999. De Bruijn Notation as a Nested Datatype. *J. Funct. Program.* 9, 1 (jan 1999), 77–91. <https://doi.org/10.1017/S0956796899003366> \hookrightarrow page 38
- Hans-J. Boehm. 1985. Partial polymorphic type inference is undecidable. In *26th Annual Symposium on Foundations of Computer Science (sfcs 1985)*. 339–345. <https://doi.org/10.1109/SFCS.1985.44> \hookrightarrow page 31
- Aleksander Boruch-Gruszecki, Radosław Waundefinedko, Yichen Xu, and Lionel Parreaux. 2022. A Case for DOT: Theoretical Foundations for Objects with Pattern Matching and GADT-Style Reasoning. *Proc. ACM Program. Lang.* 6, OOPSLA2, Article 179 (oct 2022), 30 pages. <https://doi.org/10.1145/3563342> \hookrightarrow pages 15 and 33
- Gert-Jan Bottu and Richard A. Eisenberg. 2021. Seeking Stability by Being Lazy and Shallow: Lazy and Shallow Instantiation is User Friendly. In *Proceedings of the 14th ACM SIGPLAN International Symposium on Haskell (Virtual, Republic of Korea) (Haskell 2021)*. Association for Computing Machinery, New York, NY, USA, 85–97. <https://doi.org/10.1145/3471874.3472985> \hookrightarrow page 31
- Giuseppe Castagna, Tommaso Petrucciani, and Kim Nguyen. 2016. Set-theoretic types for polymorphic variants. In *Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming (ICFP 2016)*. Association for Computing Machinery, Nara, Japan, 378–391. <https://doi.org/10.1145/2951913.2951928> \hookrightarrow page 7
- Jacek Chrząszcz. 1998. Polymorphic subtyping without distributivity. In *Mathematical Foundations of Computer Science 1998*, Luboš Brim, Jozef Gruska, and Jiří Zlatuška (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 346–355. \hookrightarrow page 2
- Julien Cretin. 2014. *Erasable coercions: a unified approach to type systems*. Theses. Université Paris-Diderot - Paris VII. <https://tel.archives-ouvertes.fr/tel-00940511> \hookrightarrow pages 3 and 16
- Julien Cretin and Didier Rémy. 2012. On the Power of Coercion Abstraction. In *Proceedings of the 39th ACM Symposium on Principles of Programming Languages (POPL 2012)*. Philadelphia, PA, USA. To appear. \hookrightarrow page 16
- Julien Cretin and Didier Rémy. 2014. System F with Coercion Constraints. In *Logics In Computer Science (LICS)*. ACM. \hookrightarrow pages 1, 2, 16, and 33
- Stephen Dolan. 2017. *Algebraic subtyping*. Ph.D. Dissertation. \hookrightarrow pages 14 and 32
- Stephen Dolan and Alan Mycroft. 2017. Polymorphism, subtyping, and type inference in MLsub. *ACM SIGPLAN Notices* 52, 1 (Jan. 2017), 60–72. <https://doi.org/10.1145/3093333.3009882> \hookrightarrow pages 1, 7, 8, and 9
- Jana Dunfield and Neel Krishnaswami. 2021. Bidirectional Typing. *ACM Comput. Surv.* 54, 5, Article 98 (may 2021), 38 pages. <https://doi.org/10.1145/3450952> \hookrightarrow page 31
- Jana Dunfield and Neelakantan R. Krishnaswami. 2013. Complete and Easy Bidirectional Typechecking for Higher-Rank Polymorphism. *SIGPLAN Not.* 48, 9 (Sept. 2013), 429–442. <https://doi.org/10.1145/2544174.2500582> \hookrightarrow page 31
- Jana Dunfield and Neelakantan R. Krishnaswami. 2019. Sound and Complete Bidirectional Typechecking for Higher-Rank Polymorphism with Existentials and Indexed Types. *Proc. ACM Program. Lang.* 3, POPL, Article 9 (Jan. 2019), 28 pages. <https://doi.org/10.1145/3290322> \hookrightarrow page 31
- Jana Dunfield and Frank Pfenning. 2004. Tridirectional Typechecking. In *Proceedings of the 31st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (Venice, Italy) (POPL '04)*. Association for Computing Machinery, New York, NY, USA, 281–292. <https://doi.org/10.1145/964001.964025> \hookrightarrow page 31
- 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 Conference on Programming Language Design and Implementation (London, UK) (PLDI 2020)*. Association for Computing Machinery, New York, NY, USA, 423–437. <https://doi.org/10.1145/3385412.3386003> \hookrightarrow pages 24 and 29
- Frank Emrich, Jan Stolarek, James Cheney, and Sam Lindley. 2022. Constraint-Based Type Inference for FreezeML. *Proc. ACM Program. Lang.* 6, ICFP, Article 111 (aug 2022), 26 pages. <https://doi.org/10.1145/3547642> \hookrightarrow pages 29 and 30
- Karl-Filip Faxén. 1997. Polyvariance, polymorphism and flow analysis. In *Analysis and Verification of Multiple-Agent Languages*, Mads Dam (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 260–278. \hookrightarrow pages 8 and 32
- A. Frisch, G. Castagna, and V. Benzaken. 2002. Semantic subtyping. In *Proceedings 17th Annual IEEE Symposium on Logic in Computer Science*. 137–146. <https://doi.org/10.1109/LICS.2002.1029823> \hookrightarrow page 7
- Jacques Garrigue and Didier Rémy. 1997. Extending ML with semi-explicit higher-order polymorphism. In *Theoretical Aspects of Computer Software*, Martín Abadi and Takayasu Ito (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 20–46. \hookrightarrow page 27
- P. Giannini and S.R. Della Rocca. 1988. Characterization of typings in polymorphic type discipline. In *[1988] Proceedings. Third Annual Symposium on Logic in Computer Science*. 61–70. <https://doi.org/10.1109/LICS.1988.5101> \hookrightarrow page 38

- G.P. Huet. 1975. A unification algorithm for typed λ -calculus. *Theoretical Computer Science* 1, 1 (1975), 27–57. [https://doi.org/10.1016/0304-3975\(75\)90011-0](https://doi.org/10.1016/0304-3975(75)90011-0) \hookrightarrow page 31
- Trevor Jim. 2000. A Polar Type System.. In *ICALP Satellite Workshops*. Citeseer, 323–338. \hookrightarrow pages 8 and 32
- Oleg Kiselyov. 2013. Efficient generalization with levels (Okmij Blog). <http://okmij.org/ftp/ML/generalization.html#levels>. Accessed: 2020-06-30. \hookrightarrow page 10
- Konstantin Läufer and Martin Odersky. 1994. Polymorphic Type Inference and Abstract Data Types. *ACM Trans. Program. Lang. Syst.* 16, 5 (sep 1994), 1411–1430. <https://doi.org/10.1145/186025.186031> \hookrightarrow page 27
- Didier Le Botlan and Didier Rémy. 2003. MLF: Raising ML to the power of System F. In *Proceedings of the Eighth ACM SIGPLAN International Conference on Functional Programming*. 27–38. \hookrightarrow pages 27, 28, 33, and 38
- Didier Le Botlan and Didier Rémy. 2009. Recasting MLF. *Information and Computation* 207, 6 (2009), 726–785. <https://doi.org/10.1016/j.ic.2008.12.006> \hookrightarrow pages 12, 13, 17, 24, 29, 31, and 33
- Daan Leijen. 2008. HMF: Simple Type Inference for First-Class Polymorphism. In *Proceedings of the 13th ACM SIGPLAN International Conference on Functional Programming (Victoria, BC, Canada) (ICFP '08)*. Association for Computing Machinery, New York, NY, USA, 283–294. <https://doi.org/10.1145/1411204.1411245> \hookrightarrow page 29
- Daan Leijen. 2009. Flexible Types: Robust Type Inference for First-Class Polymorphism. In *Proceedings of the 36th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (Savannah, GA, USA) (POPL '09)*. Association for Computing Machinery, New York, NY, USA, 66–77. <https://doi.org/10.1145/1480881.1480891> \hookrightarrow pages 12 and 29
- Daniel Leivant. 1990. Discrete Polymorphism. In *Proceedings of the 1990 ACM Conference on LISP and Functional Programming (Nice, France) (LFP '90)*. Association for Computing Machinery, New York, NY, USA, 288–297. <https://doi.org/10.1145/91556.91675> \hookrightarrow page 5
- John C. Mitchell. 1988. Polymorphic type inference and containment. *Information and Computation* 76, 2 (1988), 211–249. [https://doi.org/10.1016/0890-5401\(88\)90009-0](https://doi.org/10.1016/0890-5401(88)90009-0) \hookrightarrow page 31
- Martin Odersky and Konstantin Läufer. 1996. Putting Type Annotations to Work. In *Proceedings of the 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (St. Petersburg Beach, Florida, USA) (POPL '96)*. Association for Computing Machinery, New York, NY, USA, 54–67. <https://doi.org/10.1145/237721.237729> \hookrightarrow page 31
- Martin Odersky, Martin Sulzmann, and Martin Wehr. 1999. Type inference with constrained types. *Theory and Practice of Object Systems* 5, 1 (1999), 35–55. \hookrightarrow pages 2 and 33
- J. W. O’Toole and D. K. Gifford. 1989. Type Reconstruction with First-Class Polymorphic Values. In *Proceedings of the ACM SIGPLAN 1989 Conference on Programming Language Design and Implementation (Portland, Oregon, USA) (PLDI '89)*. Association for Computing Machinery, New York, NY, USA, 207–217. <https://doi.org/10.1145/73141.74836> \hookrightarrow page 27
- Jens Palsberg and Christina Pavlopoulou. 1998. From Polyvariant Flow Information to Intersection and Union Types. In *Proceedings of the 25th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (San Diego, California, USA) (POPL '98)*. Association for Computing Machinery, New York, NY, USA, 197–208. <https://doi.org/10.1145/268946.268963> \hookrightarrow page 32
- Lionel Parreaux. 2020. The Simple Essence of Algebraic Subtyping: Principal Type Inference with Subtyping Made Easy (Functional Pearl). *Proc. ACM Program. Lang.* 4, ICFP, Article 124 (Aug. 2020), 28 pages. <https://doi.org/10.1145/3409006> \hookrightarrow pages 2, 7, 10, and 32
- Lionel Parreaux, Aleksander Boruch-Gruszecki, Andong Fan, and Chun Yin Chau. 2023. *When Subtyping Constraints Liberate: A Novel Type Inference Approach for First-Class Polymorphism (Artifact)*. <https://doi.org/10.5281/zenodo.8424750> \hookrightarrow page 33
- Lionel Parreaux, Aleksander Boruch-Gruszecki, Andong Fan, and Chun Yin Chau. 2024. When Subtyping Constraints Liberate: A Novel Type Inference Approach for First-Class Polymorphism. *Proc. ACM Program. Lang.* 8, POPL, Article 48 (jan 2024), 33 pages. <https://doi.org/10.1145/3632890> \hookrightarrow pages 1 and 22
- Lionel Parreaux, Aleksander Boruch-Gruszecki, and Paolo G. Giarrusso. 2019. Towards Improved GADT Reasoning in Scala. In *Proceedings of the Tenth ACM SIGPLAN Symposium on Scala (London, United Kingdom) (Scala '19)*. Association for Computing Machinery, New York, NY, USA, 12–16. <https://doi.org/10.1145/3337932.3338813> \hookrightarrow page 15
- Lionel Parreaux and Chun Yin Chau. 2022. MLstruct: Principal Type Inference in a Boolean Algebra of Structural Types. *Proc. ACM Program. Lang.* 6, OOPSLA2, Article 141 (Oct 2022), 30 pages. <https://doi.org/10.1145/3563304> \hookrightarrow page 7
- Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Mark Shields. 2007. Practical Type Inference for Arbitrary-Rank Types. *J. Funct. Program.* 17, 1 (Jan. 2007), 1–82. <https://doi.org/10.1017/S0956796806006034> \hookrightarrow pages 3, 11, 31, 38, and 39
- Frank Pfenning. 1988. Partial Polymorphic Type Inference and Higher-Order Unification. In *Proceedings of the 1988 ACM Conference on LISP and Functional Programming (Snowbird, Utah, USA) (LFP '88)*. Association for Computing Machinery, New York, NY, USA, 153–163. <https://doi.org/10.1145/62678.62697> \hookrightarrow page 31
- Benjamin C Pierce. 1991. *Programming with intersection types and bounded polymorphism*. Ph. D. Dissertation. Citeseer. \hookrightarrow page 33

- Benjamin C. Pierce. 1997. Intersection types and bounded polymorphism. *Mathematical Structures in Computer Science* 7, 2 (1997), 129–193. <https://doi.org/10.1017/S096012959600223X> ↪ page 32
- François Pottier. 1998. *Type Inference in the Presence of Subtyping: from Theory to Practice*. Research Report RR-3483. INRIA. <https://hal.inria.fr/inria-00073205> ↪ page 8
- François Pottier and Didier Rémy. 2005. *The Essence of ML Type Inference*. 389–489. ↪ page 10
- Jakob Rehof and Manuel Fähndrich. 2001. Type-Based Flow Analysis: From Polymorphic Subtyping to CFL-Reachability. In *Proceedings of the 28th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (London, United Kingdom) (POPL '01). Association for Computing Machinery, New York, NY, USA, 54–66. <https://doi.org/10.1145/360204.360208> ↪ pages 8 and 32
- Didier Rémy. 1990. *Algèbres Touffues. Application au Typage Polymorphe des Objets Enregistrements dans les Langages Fonctionnels*. Thèse de doctorat. Université de Paris 7. ↪ page 10
- Didier Rémy. 1992. *Extension of ML type system with a sorted equation theory on types*. Research Report RR-1766. INRIA. <https://hal.inria.fr/inria-00077006> Projet FORMEL. ↪ page 10
- Didier Rémy. 1994. Programming Objects with ML-ART, an Extension to ML with Abstract and Record Types. In *Proceedings of the International Conference on Theoretical Aspects of Computer Software (TACS '94)*. Springer-Verlag, Berlin, Heidelberg, 321–346. ↪ page 27
- Didier Rémy. 2005. Simple, Partial Type-Inference for System F Based on Type-Containment. *SIGPLAN Not.* 40, 9 (Sept. 2005), 130–143. <https://doi.org/10.1145/1090189.1086383> ↪ pages 5 and 31
- Didier Rémy and Boris Yakobowski. 2007. A Graphical Presentation of MLF Types with a Linear-Time Unification Algorithm. In *Proceedings of the 2007 ACM SIGPLAN International Workshop on Types in Languages Design and Implementation* (Nice, Nice, France) (TLDI '07). Association for Computing Machinery, New York, NY, USA, 27–38. <https://doi.org/10.1145/1190315.1190321> ↪ page 29
- Claudio Russo and Dimitrios Vytiniotis. 2009. QML: Explicit First-Class Polymorphism for ML. In *ML '09 Proceedings of the 2009 ACM SIGPLAN workshop on ML* (ml '09 proceedings of the 2009 acm sigplan workshop on ml ed.). ACM New York, NY, USA, 3–14. <https://www.microsoft.com/en-us/research/publication/qml-explicit-first-class-polymorphism-for-ml/> ↪ page 29
- Didier Rémy and Boris Yakobowski. 2012. A church-style intermediate language for MLF. *Theoretical Computer Science* 435 (2012), 77–105. <https://doi.org/10.1016/j.tcs.2012.02.026> Functional and Logic Programming. ↪ page 29
- Gabriel Scherer and Didier Rémy. 2015. Full Reduction in the Face of Absurdity. In *Programming Languages and Systems*, Jan Vitek (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 685–709. ↪ pages 15 and 33
- Alejandro Serrano, Jurriaan Hage, Simon Peyton Jones, and Dimitrios Vytiniotis. 2020. A Quick Look at Impredicativity. *Proc. ACM Program. Lang.* 4, ICFP, Article 89 (aug 2020), 29 pages. <https://doi.org/10.1145/3408971> ↪ pages 4 and 31
- Alejandro Serrano, Jurriaan Hage, Dimitrios Vytiniotis, and Simon Peyton Jones. 2018. Guarded Impredicative Polymorphism. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Philadelphia, PA, USA) (PLDI 2018). Association for Computing Machinery, New York, NY, USA, 783–796. <https://doi.org/10.1145/3192366.3192389> ↪ pages 4 and 30
- Scott F. Smith and Tiejun Wang. 2000. Polyvariant Flow Analysis with Constrained Types. In *Programming Languages and Systems*, Gert Smolka (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 382–396. ↪ pages 8 and 32
- Jerzy Tiuryn and Pawel Urzyczyn. 1996. The Subtyping Problem for Second-Order Types is Undecidable. In *Proceedings of the Eleventh Annual IEEE Symposium on Logic in Computer Science (LICS 1996)* (New Brunswick, NJ, USA). IEEE Computer Society Press, 74–85. ↪ page 2
- Dimitrios Vytiniotis, Stephanie Weirich, and Simon Peyton Jones. 2006. Boxy Types: Inference for Higher-Rank Types and Impredicativity. In *Proceedings of the Eleventh ACM SIGPLAN International Conference on Functional Programming* (Portland, Oregon, USA) (ICFP '06). Association for Computing Machinery, New York, NY, USA, 251–262. <https://doi.org/10.1145/1159803.1159838> ↪ pages 4, 30, and 31
- Dimitrios Vytiniotis, Stephanie Weirich, and Simon Peyton Jones. 2008. FPH: First-Class Polymorphism for Haskell. In *Proceedings of the 13th ACM SIGPLAN International Conference on Functional Programming* (Victoria, BC, Canada) (ICFP '08). Association for Computing Machinery, New York, NY, USA, 295–306. <https://doi.org/10.1145/1411204.1411246> ↪ page 29
- Tiejun Wang and Scott F. Smith. 2001. Precise Constraint-Based Type Inference for Java. In *ECOOP 2001 — Object-Oriented Programming*, Jørgen Lindskov Knudsen (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 99–117. ↪ page 32
- Joseph B. Wells. 1996. *Type inference for System F with and without the eta rule*. Ph.D. Dissertation. <http://ezproxy.ust.hk/login?url=https://www.proquest.com/dissertations-theses/type-inference-system-f-without-eta-rule/docview/304323713/se-2> Copyright - Database copyright ProQuest LLC; ProQuest does not claim copyright in the individual underlying works; Last updated - 2023-02-24. ↪ page 31
- Hongwei Xi, Chiyang Chen, and Gang Chen. 2003. Guarded Recursive Datatype Constructors. In *Proceedings of the 30th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (New Orleans, Louisiana, USA) (POPL '03).

- Association for Computing Machinery, New York, NY, USA, 224–235. <https://doi.org/10.1145/604131.604150> ↔ pages 33 and 38
- Ningning Xie. 2021. Higher-rank polymorphism : type inference and extensions. <http://hdl.handle.net/10722/307011> ↔ page 31
- Ningning Xie and Bruno C. d. S. Oliveira. 2018. Let Arguments Go First. In *Programming Languages and Systems*, Amal Ahmed (Ed.), Springer International Publishing, Cham, 272–299. ↔ page 31
- Mingqi Xue and Bruno C.d.S. Oliveira. 2021. A dependently typed calculus with polymorphic subtyping. *Science of Computer Programming* 208 (2021), 102655. <https://doi.org/10.1016/j.scico.2021.102655> ↔ pages 31 and 33
- Jinxu Zhao, Bruno C. d. S. Oliveira, and Tom Schrijvers. 2019. A Mechanical Formalization of Higher-Ranked Polymorphic Type Inference. *Proc. ACM Program. Lang.* 3, ICFP, Article 112 (jul 2019), 29 pages. <https://doi.org/10.1145/3341716> ↔ page 10

A ADDITIONAL SUPERF EXAMPLES

In this appendix, we provide some additional examples related to SuperF type inference.

A.1 Concrete Example of Need for Annotations

Consider a function `mapExprs`: $(\forall a. \text{Expr } a \rightarrow \text{Expr } a) \rightarrow \text{Program} \rightarrow \text{Program}$, which presumably transforms the typed subexpressions of a program by applying an argument function at an arbitrary number of different expression types that are not known statically (`Expr` could be, for instance, the usual typed expression GADT [Xi et al. 2003] or an expression data type with type-level de Bruijn indices [Bird and Paterson 1999; Peyton Jones et al. 2007]). The following SuperF definition fails to infer a sufficiently polymorphic type in this case:

```
mapAndPrintSize f = mapExprs (fun e → print_int (size e) ; f e)
```

Indeed, in this case, parameter `f` is assigned a type variable and the local type variable assigned to `e` is extruded through it. So calling `mapAndPrintSize id` yields an error, as the local `a` skolem in `Expr a` (from the definition of `mapExprs`) flows out to the extruded type variable where it undergoes type avoidance, resulting in `Expr ⊥` and `Expr ⊤` before flowing back into `Expr a`, which fails because `⊥` and `⊤` do not match `a`. In cases like this, a simple type annotation resolves the problem:

```
mapAndPrintSize (f: ∀ a. Expr a → Expr a) =
  mapExprs (fun e → print_int (size e); f e)

mapAndPrintSize id -- ok
```

Stronger forms of this restriction are present in all previous first-class and higher-rank polymorphism type inference approaches we are aware of, including ML^F . The latter requires that *all parameters used polymorphically* must be annotated, whereas we only require annotations for parameters used *parametrically*-polymorphically (i.e., *non-parametric* – or *finitary* – polymorphism can be left unannotated thanks to intersection type inference). Other approaches generally require annotations for parameters that are to be assigned a polymorphic type at all, which is even more restrictive. For instance, consider that in both SuperF and ML^F , the following definition is accepted:

```
mapAndPrintHello f = print_string "Hello!" ; mapExprs f

mapAndPrintHello : (∀ a. Expr a → Expr a) → Program → Program
```

A.2 Expressiveness of SuperF

Consider the following `monster` term along with some helper definitions:

```
I x = x
K x y = x
auto x = x x

monster = (fun y → (let tmp = y I in y K)) auto
```

which is not typeable in System F (though it is typeable in F_ω) [Giannini and Della Rocca 1988], leading Le Botlan and Rémy [2003] to conjecture that it is not typeable in their ML^F type system either.

But there is nothing *fundamentally* hard to type about this example. System F can in fact type `monster` after it is made to take just one step of reduction (leading to `let tmp = self I in self K`). The reason System F cannot type the original term is that System F lacks the expressiveness needed to accurately describe a set of constraints inferred from the program's data flow.

SuperF infers the following types for this program, without the help of any type annotations:

```
I : ∀ a. a → a
K : ∀ a. a → ⊤ → a
```

```

auto : ∀ a b. ((a → b) ∧ a) → b
monster : τ → ∀ a. a → τ → a

```

To understand how this works, note that the `fun y → (let tmp = y I in y K)` subterm of `monster` is given inferred type $\forall r. ((\forall a. a \rightarrow \tau \rightarrow a) \rightarrow r) \wedge ((\forall a. a \rightarrow a) \rightarrow \tau) \rightarrow r$, which clearly shows that parameter `y` is inferred to be a function expected to yield some type `r` when applied to `K`, of type $\text{KT} = \forall a. a \rightarrow \tau \rightarrow a$, and to yield any type τ when applied to `I` (since the result of that application is discarded), with the overall result type being `r`. Then, applying that subterm to `auto` is simply a matter of checking that `auto` can receive both `K` and `I` as individual arguments while yielding `r` as a result of the former, which is expressed as the subtyping constraint $\forall a. ((a \rightarrow b) \wedge a) \rightarrow b \leq \text{KT} \rightarrow r$, which decomposes into $b \leq r$ and $\text{KT} \leq (a\emptyset \rightarrow b) \wedge a\emptyset$, into $a1 \rightarrow \tau \rightarrow a1 \leq a\emptyset \rightarrow b$ (i.e., $a\emptyset \leq a1$ and $\tau \rightarrow a1 \leq b$) and $\text{KT} \leq a\emptyset$, which after simplification yields $r = \tau \rightarrow \text{KT}$, the `monster` type shown above. We describe the general idea of how such subtyping constraints are solved in Section 2.5 and rigorously formalize the process in Section 4.3.

It is easy to show that these types are principal using the same approach as in Section 2.2 — intuitively, all the constraints we solved above are indeed *necessary* consequences of the term being well-typed, and thus the constrained types we infer are essentially representations of every constraint any valid type must satisfy.

A.3 Unlimited-rank Type Inference

SuperF is not limited in the rank of the polymorphic types it infers. In practice, it always infers parametrically-polymorphic types for odd ranks and intersection types for even ranks.

Consider the following elaboration on our introduction example:

```

foo f = (f 123, f True)

bar f = (f (fun x → x), f (fun x → Some x))

test = bar foo

```

SuperF successfully infers the following types:

```

foo : ∀ a b. ((Int → a) ∧ (Bool → b)) → (a, b)

bar : ∀ a b. (((∀ c. c → c) → a) ∧ (∀ d. d → Some[d]) → b) → (a, b)

test : (Int, Bool), (Option Int, Option Bool)

```

Notice that the type inferred for `bar` has rank 3, as it contains \forall types behind two arrow type left-hand sides [Peyton Jones et al. 2007, §3.1].

A.4 Type Error Messages

Assuming `k : (∀ a. a → a) → int`, the expression `k (fun x → x + 1)` results in the following error message in SuperF, where the flow of rigid type variable `a` into parameter `x` is clear:

```

[ERROR] Type mismatch in application:
  1.8: k (fun x -> x + 1)
— type `a` is not an instance of type `int`
  1.5: let k: (forall a. a -> a) -> int
— Note: constraint arises from reference:
  1.8: k (fun x -> x + 1)

```

By contrast, here is the error generated by GHC, an industry-strength compiler that has supported bidirectionally-typed higher-rank polymorphism (a strictly easier feature to implement, which is sufficient for this example) for many years:

- Couldn't match expected type 'Int' with actual type 'a'
'a' is a rigid type variable bound by
a type expected by the context:
forall a. a -> a
at <interactive>:4:3-20
- In the first argument of '(+)', namely '(x :: Int)'
In the expression: (x :: Int) + 1
In the first argument of 'k', namely '(\ x -> (x :: Int) + 1)'
- Relevant bindings include x :: a (bound at <interactive>:4:5)

The more tricky case to handle is when a type variable leaks out of its scope. SuperF keeps track of where leaks happen (which correspond to *type extrusions*) and reports them in the resulting error messages as possible locations where type annotations may be needed. For instance, assuming the same type for k as above, the term `(fun f -> k (fun x -> f x)) id` yields the following error:

```
[ERROR] Type error in application
  1.233:      (fun f -> k (fun x -> f x)) id
— type variable `a` leaks out of its scope
  1.229:      def k: (forall a. a -> a) -> int
— back into type variable `a`
  1.229:      def k: (forall a. a -> a) -> int
— Adding a type annotation to any of the following terms
  may help resolve the problem
— • this application:
  1.233:      (fun f -> k (fun x -> f x)) id
```

While the suggested annotation location does not refer specifically to parameter `f`, which is the most natural place to add a type annotation, it does refer to the relevant application where the leak takes place. Here is the rather verbose and unhelpful error produced by GHC:

- Couldn't match type 'a0' with 'a'
Expected: a -> a
Actual: a0 -> a0
because type variable 'a' would escape its scope
This (rigid, skolem) type variable is bound by
a type expected by the context:
forall a. a -> a
at <interactive>:5:11-22
- In the expression: f x
In the first argument of 'k', namely '(\ x -> f x)'
In the expression: k (\ x -> f x)
- Relevant bindings include
x :: a (bound at <interactive>:5:14)
f :: a0 -> a0 (bound at <interactive>:5:4)

A.5 Ω and the SRLC

Here we give an example constraint-solving run which involves the suspiciously recursive-looking criterion (SRLC).

Consider the problem of type checking term Ω , i.e., `(fun x -> x x) (fun x -> x x)`.

At some point in the main constraint-solving run, we get the constraint $\tau \leq \beta$, where $\tau = \forall \alpha \{ \alpha \leq \alpha \rightarrow \alpha \}. \alpha \rightarrow \alpha$ and $\beta \leq \beta \rightarrow \beta$

This leads to:

- add lower bound τ to β
- constrain $\tau \leq \beta \rightarrow \beta$ (UB of β)
- instantiate τ to $\gamma_\alpha \rightarrow \gamma_\alpha$ where $\gamma_\alpha \leq \gamma_\alpha \rightarrow \gamma_\alpha$ (we write γ_α to show that γ 's shadow is α)
- constrain $\gamma_\alpha \rightarrow \gamma_\alpha \leq \beta \rightarrow \beta$
- constrain $\beta \leq \gamma_\alpha$ (LHS of function types)
- constrain $\beta \leq \gamma_\alpha \rightarrow \gamma_\alpha$ (UB of γ_α)
- constrain $\tau \leq \gamma_\alpha \rightarrow \gamma_\alpha$ (LB of β)

In turn, this leads to another round of exactly the same constraints until we reach:

- constrain $\tau \leq \delta_\alpha \rightarrow \delta_\alpha$

But both pairs of types $(\tau, \delta_\alpha \rightarrow \delta_\alpha)$ and $(\tau, \gamma_\alpha \rightarrow \gamma_\alpha)$ have the same root $(\tau, \alpha \rightarrow \alpha)$, so the SRLC kicks in and aborts the constraint-solving run.

B EXTRA FORMAL DEFINITIONS AND EXAMPLES

We report on extra formal definitions as well as examples of the formal definitions here.

B.1 Box Erasure

Definition B.1. Box erasure, written $\tau \setminus x$, is defined inductively on the syntax of types and bounds contexts as follows:

$$\begin{aligned}
 \top \setminus x &\triangleq \top \\
 \alpha \setminus x &\triangleq \alpha \\
 (\tau_1 \rightarrow \tau_2) \setminus x &\triangleq (\tau_1 \setminus x) \rightarrow (\tau_2 \setminus x) \\
 \forall V\{\Xi\}. \tau \setminus x &\triangleq \forall V\{\{\Xi \setminus x\}\}. (\tau \setminus x) \\
 \langle \tau \rangle \setminus x &\triangleq \langle \tau \setminus x \rangle \\
 \boxed{\tau}_x \setminus x &\triangleq \tau \setminus x \\
 \boxed{\tau}_y \setminus x &\triangleq \boxed{\tau \setminus x}_y \quad (x \neq y) \\
 \epsilon \setminus x &\triangleq \epsilon \\
 \Xi \cdot (\tau \leq \sigma) \setminus x &\triangleq (\Xi \setminus x) \cdot ((\tau \setminus x) \leq (\sigma \setminus x))
 \end{aligned}$$

Example B.2. Consider the following program:

```

test = fun f → fun x → f x

test (fun res → res) id

```

Note that to type the body of `test`, we need to wrap the type τ its parameter `x` into a box $\boxed{\tau}_x$, which means that parameter `f` must be a function that accepts such boxed arguments. Therefore, taking (arbitrarily, for the example) $\tau = \forall \gamma. \gamma \rightarrow \gamma$, the type of `test` could be, for instance, $\sigma = \forall \alpha. (\boxed{\tau}_x \rightarrow \alpha) \rightarrow \alpha$, or it could be $\sigma' = \forall \alpha \beta \{\beta \leq \boxed{\tau}_x\}. (\beta \rightarrow \alpha) \rightarrow \alpha$. In both cases, all the x -branded boxes of this type can be removed at the *use site* of `test`, as in $\sigma \setminus x = \forall \alpha. (\tau \rightarrow \alpha) \rightarrow \alpha$ and $\sigma' \setminus x = \forall \alpha \beta \{\beta \leq \tau\}. (\beta \rightarrow \alpha) \rightarrow \alpha$, which allows passing in the polymorphic identity function `id` for `x` and using the box-free result polymorphically.

Remark 2. Naturally, in all developments, we use Barendregt's convention and assume all variables in the program are distinct.

B.2 Syntactic and Semantic Acyclicity

A context is *syntactically acyclic* when its bounds graph is acyclic. By contrast, a context is *semantically acyclic* (or just *acyclic*) when all its bounds graph cycles (if it has any) are “harmless”, meaning that traversing the bounds graph while following polarities will not lead to a cyclic path. The simplest example of that would be $\forall \alpha \{ \alpha \leq \top \rightarrow \alpha \}$. α – here, the upper bound of α , which is $\top \rightarrow \alpha$, is not reachable from traversing the bounds graph starting from α at a positive polarity (which is the body of the positive polymorphic type, always considered positive), so it is irrelevant and the bounds graph of this type is semantically acyclic despite being syntactically cyclic.

In general, we require all inferred types to be acyclic because supporting cyclic contexts would require support for full-fledged recursive types, which System F_{cc} does not yet possess.²⁷

Definition B.3 (Acyclicity). Acyclicity is defined through the *reach* judgment, which is the smallest relation satisfying the following rules (where equality is taken to mean mutual set inclusion):

$$\boxed{\text{acyclic}(\Xi) \triangleq \text{for all } a, \{a^+, a^-\} \# \text{reach}^+(a, \Xi, 0) \text{ and } a^- \notin \text{reach}^-(a, \Xi, 0)} \quad \text{reach}^\pm(\top, \Xi, n) = \emptyset$$

$$\text{reach}^\pm(\tau \rightarrow \sigma, \Xi, n) = \text{reach}^\mp(\tau, \Xi, 1) \cup \text{reach}^\pm(\sigma, \Xi, 1) \quad \text{reach}^\pm(\forall V\{\Sigma\}. \tau, \Xi, n) = \text{reach}^\pm(\tau, \Xi \cdot \Sigma, n)$$

$$\frac{}{a^\pm \in \text{reach}^\pm(a, \Xi, 1)} \quad \frac{(\tau \leq a) \in \Xi \quad b \in \text{reach}^+(\tau, \Xi, n)}{b \in \text{reach}^+(a, \Xi, n)} \quad \frac{(a \leq \tau) \in \Xi \quad b \in \text{reach}^-(\tau, \Xi, n)}{b \in \text{reach}^-(a, \Xi, n)}$$

Where we use \mp as a shorthand for $+$ when \pm is $-$ and $-$ when \pm is $+$.

The “flag” parameter n can be 0 or 1 and indicates whether we have traversed at least one type constructor (if not, any cycle is a spurious “immediate” cycle).

Example B.4. Context $\Xi = (\alpha \rightarrow \beta) \rightarrow \gamma \leq \alpha$ is not acyclic because we have:

$$\begin{aligned}
 \text{reach}^+(\alpha, \Xi, 0) &= \text{reach}^+((\alpha \rightarrow \beta) \rightarrow \gamma, \Xi, 0) \\
 &= \text{reach}^-(\alpha \rightarrow \beta, \Xi, 1) \cup \text{reach}^+(\gamma, \Xi, 1) \\
 &= \text{reach}^+(\alpha, \Xi, 1) \cup \{\gamma^+\} \\
 &= \{\alpha^+, \gamma^+\} \\
 &\ni \alpha^+
 \end{aligned}$$

Example B.5. Context $\Xi = \alpha \leq \top \rightarrow \alpha$ is not acyclic because we have:

$$\begin{aligned}
 \text{reach}^-(\alpha, \Xi, 0) &= \text{reach}^-(\top \rightarrow \alpha, \Xi, 0) \\
 &= \text{reach}^+(\top, \Xi, 1) \cup \text{reach}^-(\alpha, \Xi, 1) \\
 &= \{\alpha^-\} \\
 &\ni \alpha^-
 \end{aligned}$$

Definition B.6 (Syntactic acyclicity).

$$\text{synacyc}(\Xi) \triangleq \text{for all } a, \{a^+, a^-\} \# \text{reach}^+(a, \Xi, 0) \cup \text{reach}^-(a, \Xi, 0)$$

Example B.7. Context $\Xi = \alpha \leq \alpha \rightarrow \top$ is semantically acyclic because $\text{reach}^+(\alpha, \Xi, 0) = \emptyset$ and $\text{reach}^-(\alpha, \Xi, 0) = \alpha^+$. It can be rewritten into the equivalent $\Xi' = \alpha \leq \alpha' \wedge (\alpha' \rightarrow \top)$ (α' fresh), syntax sugar for $\Xi' = \alpha \leq \beta', \beta' \leq \alpha', \beta' \leq \alpha' \rightarrow \top$ (α', β' fresh), which is more obviously semantically acyclic – but is still *syntactically cyclic*, as $\text{reach}^-(\alpha, \Xi', 0) = \alpha^+$.

²⁷We conjecture that it would not be hard to extend F_{cc} with full recursive types, but we have not yet tried doing so.

B.3 Extrusion Safety Check

We first define *collect* that returns all the bounds nesting in the input type:

$$\begin{aligned}
\text{collect}(\top) &= \epsilon \\
\text{collect}(\alpha) &= \epsilon \\
\text{collect}(\forall\alpha. \tau^-) &= \text{collect}(\tau^-) \\
\text{collect}(\tau^\mp \rightarrow \tau^\pm) &= \text{collect}(\tau^\mp) \cdot \text{collect}(\tau^\pm) \\
\text{collect}(\forall V\{\Sigma\}. \tau^+) &= V \cdot \text{collect}(\tau^+) \cdot \overline{\text{collect}(\sigma^+)^{(\sigma^+ \leq \alpha) \in \Sigma}} \cdot \overline{\text{collect}(\sigma^-)^{(\alpha \leq \sigma^-) \in \Sigma}} \cdot \Sigma
\end{aligned}$$

We say τ^\pm is safe to extrude with the skolem set W if the type satisfies the following **X-OK** check:

Definition B.8. $W \vdash \tau^\pm$ **X-OK** \triangleq for all $\alpha \in S$, $W \# \text{reach}^+(\alpha, \mathcal{B}(S), 1)$ or $W \# \text{reach}^-(\alpha, \mathcal{B}(S), 1)$ where $S = \text{collect}(\tau^\pm)$. We ignore polarities on type variables returned by *reach* $^\pm$.

C TRANSLATION TO F_{cc}

In this appendix we present the translation from our system to System F_{cc} . We will make two small adjustments to our system. First: we define the translation for an “erased” version of the system which doesn’t allow either the $\boxed{\tau}_x$ or the $\langle \tau \rangle$ type form, doesn’t allow the $(x : \tau)$ term form, doesn’t have the following rules: **T-UNBOX**, **T-ASC**, **S-UNBOX1**, **S-UNBOX2**, **S-CONGBOUND**, **S-CONGFUN**, and has the following version of **T-ABS**:

$$\begin{array}{c}
\text{T-ABS} \\
\frac{\Gamma \cdot (x : \tau_1) \vdash t : \tau_2}{\Gamma \vdash \lambda x. t : \tau_1 \rightarrow \tau_2}
\end{array}$$

Clearly a valid derivation in the normal system is also valid in the erased system, after we erase $\boxed{\tau}_x$ and $\langle \tau \rangle$ to τ and $(x : \tau)$ to x .

Second, we will need to add type well-formedness to our system, since System F_{cc} requires type variables to be bound in context.

We extend contexts Γ and Ξ with type bindings: Γ may have the form $\Gamma \cdot \alpha$, likewise for Ξ .

Definition C.1. Quantifier consistency $\Xi \vdash \forall \alpha_i^i \{ \overline{B}_j^j \}$ **cons.** is defined as both $\Xi \vdash \forall \alpha_i^i \{ \overline{B}_j^j \}$, \top **wf** and $\Xi \vdash \theta \overline{B}_j^j$, where $\theta = [\alpha_i \mapsto \tau_i]^i$ for some $\overline{\tau}_i^i$.

Figure 8 shows rules of type well-formedness $\Gamma \vdash \tau$ **wf**. As every bounds context Ξ is also grammatically a typing context Γ , we will also write $\Xi \vdash \tau$ **wf**.

$\boxed{\Gamma \vdash \tau \text{ wf}}$

$$\begin{array}{c}
\begin{array}{ccc}
\text{W-TOP} & \text{W-VAR} & \text{W-FUN} \\
\frac{}{\Gamma \vdash \top \text{ wf}} & \frac{\alpha \in \Gamma}{\Gamma \vdash \alpha \text{ wf}} & \frac{\Gamma \vdash \tau_1 \text{ wf} \quad \Gamma \vdash \tau_2 \text{ wf}}{\Gamma \vdash \tau_1 \rightarrow \tau_2 \text{ wf}}
\end{array} \\
\\
\text{W-FORALL} \\
\frac{\Gamma \cdot \overline{\alpha}_i^i \vdash \sigma_1 \text{ wf} \quad \Gamma \cdot \overline{\alpha}_i^i \vdash \sigma_2 \text{ wf} \quad \overline{(\sigma_1 \leq \sigma_2) \in \overline{B}}}{\Gamma \vdash \forall \overline{\alpha}_i^i \{ \overline{B}_j^j \}. \tau \text{ wf}} \quad \Gamma \cdot \overline{\alpha}_i^i \cdot \overline{B}_j^j \vdash \tau \text{ wf}
\end{array}$$

Fig. 8. Type well-formedness.

Since we now desire types to be well-formed, we need to make small changes to the typing and subtyping rules: variables need to be bound in contexts as appropriate and well-formedness needs to be required at certain points. Figures 9 and 10 show these changes.

$$\boxed{\Gamma \vdash t : \tau}$$

$$\begin{array}{c} \text{T-ABS} \\ \Gamma \cdot (x : \tau_1) \vdash t : \tau_2 \quad \Gamma \vdash \tau_1 \mathbf{wf} \\ \hline \Gamma \vdash \lambda x. t : \tau_1 \rightarrow \tau_2 \end{array} \qquad \begin{array}{c} \text{T-FORALL} \\ \Gamma \cdot V \cdot \Sigma \vdash t : \tau \quad V \notin FV(\Gamma) \quad \mathcal{B}(\Gamma) \vdash \forall V \{ \Sigma \} \mathbf{cons.} \\ \hline \Gamma \vdash t : \forall V \{ \Sigma \}. \tau \end{array}$$

Fig. 9. Adjustments to typing rules.

$$\boxed{\Xi \vdash \tau \leq \sigma}$$

$$\begin{array}{c} \text{S-TOP} \\ \Xi \vdash \tau \mathbf{wf} \\ \hline \Xi \vdash \tau \leq \top \end{array} \qquad \begin{array}{c} \text{S-VARREFL} \\ \alpha \in \Xi \\ \hline \Xi \vdash \alpha \leq \alpha \end{array} \qquad \begin{array}{c} \text{S-FORALL-COV} \\ \Xi \vdash \forall V \{ \Sigma \} \mathbf{cons.} \quad \Xi \cdot V \cdot \Sigma \vdash \tau \leq \sigma \\ \hline \Xi \vdash \forall V \{ \Sigma \}. \tau \leq \forall V \{ \Sigma \}. \sigma \end{array}$$

$$\begin{array}{c} \text{S-FORALL-L} \\ \Xi \vdash \forall V \{ \Sigma \} \mathbf{cons.} \quad \overline{\Xi \vdash \tau_i \mathbf{wf}^i} \quad \Xi \vdash [\overline{\alpha_i \mapsto \tau_i^i}] \Sigma \\ \hline \Xi \vdash \forall \overline{\alpha_i^i} \{ \Sigma \}. \tau \leq [\overline{\alpha_i \mapsto \tau_i^i}] \tau \end{array}$$

Fig. 10. Adjustments to subtyping rules.

The intent behind these changes is to allow the following *regularity* lemmas to hold.

LEMMA C.2. *Let $\Gamma \mathbf{wf}$. Then $\Gamma \vdash t : \tau$ implies that $\Gamma \vdash \tau \mathbf{wf}$.*

LEMMA C.3. *Let $\Gamma \mathbf{wf}$. Then $\Gamma \vdash \tau \leq \sigma$ implies that $\Gamma \vdash \tau \mathbf{wf}$ and $\Gamma \vdash \sigma \mathbf{wf}$.*

In the following sections, we implicitly assume that all contexts are well-formed.

C.1 Translation

In this subsection, we show the exact translation to F_{cc} . We write variables standing for F_{cc} objects with a tilde: $\tilde{\alpha}, \tilde{\Gamma}$. We try keep the metavariables uniform, e.g. the metavariable for F_{cc} terms is \tilde{t} , not \tilde{a} ; the only exception is when we reference variables used to instantiate F_{cc} rules, where we use the F_{cc} metavariable with a tilde (e.g. “Use COERWEAK with $\tilde{\Sigma} = \dots$ ”).

Where useful, we write F_{cc} judgments with a clarification on the turnstile: $\tilde{\Gamma} \vdash^{F_{cc}} \tilde{t} : \tilde{\tau}$.

F_{cc} includes a notion of type (kind, proposition) *equality* $\tilde{\tau} =^{F_{cc}} \tilde{\sigma}$, which should not be confused with the two types being the same. For instance, we have $\pi_1 \langle \tilde{\tau}, \tilde{\sigma} \rangle =^{F_{cc}} \tilde{\tau}$, but the two types are distinct objects. (This equality must be explicitly used by appropriate F_{cc} rules.) To avoid confusion, we add a clarification on the equality sign when using this form of equality.

We use Θ as a metavariable for mappings from our type variables to F_{cc} types.

We use $\llbracket \cdot \rrbracket$ to denote the translation function. Where useful, we add a metavariable in bold font in a superscript to clarify what sort of object is produced by the translation. We define the translation function as follows:

$$\llbracket \forall V \{ \Sigma \} \rrbracket_{\Theta} \triangleq \{ \tilde{\alpha} : \llbracket V \rrbracket^{\mathbf{K}} \mid \llbracket \Sigma \rrbracket_{\Theta \circ \llbracket V, \tilde{\alpha} \rrbracket \Theta} \}$$

$$\begin{aligned} \llbracket \alpha \cdot V; \tilde{\tau} \rrbracket^\Theta &\triangleq \llbracket \alpha; \tilde{\tau} \rrbracket^\Theta \circ \llbracket V; \pi_2 \tilde{\tau} \rrbracket^\Theta \\ \llbracket \alpha; \tilde{\tau} \rrbracket^\Theta &\triangleq [\alpha \mapsto \pi_1 \tilde{\tau}] \end{aligned}$$

$$\begin{aligned} \llbracket \alpha \cdot V \rrbracket^\kappa &\triangleq \star \times \llbracket V \rrbracket^\kappa \\ \llbracket \alpha \rrbracket^\kappa &\triangleq \star \times 1 \end{aligned}$$

$$\begin{aligned} \llbracket (\tau \leq \sigma) \cdot \Sigma \rrbracket_\Theta &\triangleq \llbracket \tau \leq \sigma \rrbracket_\Theta \wedge \llbracket \Sigma \rrbracket_\Theta \\ \llbracket \epsilon \rrbracket_\Theta &\triangleq \top \end{aligned}$$

$$\llbracket \tau \leq \sigma \rrbracket_\Theta \triangleq [\llbracket \tau \rrbracket_\Theta \triangleright \llbracket \sigma \rrbracket_\Theta]$$

$$\llbracket \tau \rightarrow \sigma \rrbracket_\Theta \triangleq \llbracket \tau \rrbracket_\Theta \rightarrow \llbracket \sigma \rrbracket_\Theta$$

$$\llbracket \forall V \{ \Sigma \}. \tau \rrbracket_\Theta \triangleq \forall (\tilde{\alpha} : \llbracket \forall V \{ \Sigma \} \rrbracket_\Theta) \llbracket \tau \rrbracket_{\Theta \circ \llbracket V; \tilde{\alpha} \rrbracket^\Theta}$$

$$\llbracket \alpha \rrbracket_\Theta \triangleq \Theta(\alpha)$$

$$\llbracket \top \rrbracket_\Theta \triangleq \top$$

$$(\epsilon; \Theta) \sim (\emptyset; \Theta)$$

$$\frac{(\Gamma; \Theta'_1) \sim (\tilde{\Gamma}; \Theta_2)}{\llbracket (x : \tau) \cdot \Gamma; \Theta_1 \rrbracket \sim \llbracket (x : \llbracket \tau \rrbracket_\Theta, \tilde{\Gamma}; \Theta_2) \rrbracket} \quad \frac{(\Gamma; \Theta'_1) \sim (\tilde{\Gamma}; \Theta_2)}{\llbracket (V \cdot \Sigma \cdot \Gamma; \Theta_1) \rrbracket \sim \llbracket (\tilde{\alpha} : \llbracket \forall V \{ \Sigma \} \rrbracket_{\Theta_1}, \tilde{\Gamma}; \Theta_2) \rrbracket}$$

$$\Gamma \sim (\tilde{\Gamma}; \Theta) \triangleq (\Gamma; \text{id}) \sim (\tilde{\Gamma}; \Theta)$$

C.2 Core theorem

In this subsection we show the core translation theorem. While we restrict our derivations to only use type variables bound in the context, this is not a significant loss of expressiveness: any derivation using the original rules can easily be transformed into a derivation obeying the well-formedness restriction by simply binding the type variables in the outermost context.

THEOREM C.4. *If $\vdash t : \tau$, then $\vdash^{F_{cc}} \tilde{t} : \llbracket \tau \rrbracket$ for some F_{cc} term \tilde{t} .*

PROOF. A special case of Lemma C.5. □

LEMMA C.5. *Let $\Gamma \sim (\tilde{\Gamma}; \Theta)$. Then $\Gamma \vdash t : \tau$ implies that $\tilde{\Gamma} \vdash^{F_{cc}} \tilde{t} : \llbracket \tau \rrbracket_\Theta$ for some F_{cc} term \tilde{t} .*

PROOF. By induction on the typing derivation.

Case T-UNIT. Holds if $()$ is translated to an arbitrary well-typed closed F_{cc} term, for instance the identity function, or a pair of identity functions.

Case T-VAR, T-ABS, T-APP. Trivial.

Case T-SUBS. By TERMCOER with empty $\tilde{\Sigma}$. Coercion can be shown with Lemma C.11.

Case T-FORALL. By TERMCOER with $\tilde{\Sigma} = \tilde{\alpha} : \llbracket \forall V \{ \Sigma \} \rrbracket_\Theta^\kappa$. Typing can be shown by the IH. Coercion can be shown by COERGEN and by Lemma C.12. □

C.3 Coercion lemmas

The propositions below follow from a simple inspection of the appropriate definitions.

PROPOSITION C.6. *If $\Gamma \vdash \tau$ wf, then $FV(\tau) \subseteq FV(\Gamma)$.*

PROPOSITION C.7. *Let $\Gamma \sim (\tilde{\Gamma}; \Theta)$. Then $\alpha \in FV(\Gamma)$ implies that $\tilde{\Gamma} \vdash \Theta(\alpha) : \star$.*

PROPOSITION C.8. *Let $\Gamma \sim (\tilde{\Gamma}; \Theta)$. Then $\Gamma \vdash \tau$ wf implies that $\tilde{\Gamma} \vdash \llbracket \tau \rrbracket_{\Theta} : \star$.*

We define an auxiliary translation which maps sequences of types to tuples of F_{cc} types.

$$\begin{aligned} \llbracket \tau \cdot \bar{\sigma} \rrbracket_{\Theta}^{(\tau \times \bar{\sigma})} &\triangleq \langle \llbracket \tau \rrbracket_{\Theta}, \llbracket \bar{\sigma} \rrbracket_{\Theta}^{(\tau \times \bar{\sigma})} \rangle \\ \llbracket \tau \rrbracket_{\Theta}^{(\tau \times \tau)} &\triangleq \langle \llbracket \tau \rrbracket_{\Theta}, \langle \rangle \rangle \end{aligned}$$

PROPOSITION C.9. *Let $\Gamma \sim (\tilde{\Gamma}; \Theta)$. Then $\Gamma \vdash \forall V\{B_i^i\}. \tau$ wf and $\overline{\Gamma} \vdash \theta B_i^i$ (where $\text{dom}(\theta) = \{V\}$) implies that $\tilde{\Gamma} \vdash^{F_{cc}} \llbracket \theta V \rrbracket_{\Theta}^{(\tau \times \tau)} : \llbracket V \rrbracket^{\mathbf{K}}$.*

LEMMA C.10. *Let $\Gamma \sim (\tilde{\Gamma}; \Theta)$. Then $(\tau \leq \sigma) \in \Gamma$ implies that $\tilde{\Gamma} \vdash^{F_{cc}} \llbracket \tau \rrbracket_{\Theta} \triangleright \llbracket \sigma \rrbracket_{\Theta}$.*

PROOF. By structural induction on Γ .

Case $\Gamma = \Gamma' \cdot (x : \tau')$. Then $(\tau \leq \sigma) \in \Gamma'$ and we conclude by the IH and F_{cc} weakening.

Case $\Gamma = \Gamma' \cdot V \cdot \Sigma$ s.t. $(\tau \leq \sigma) \in \Gamma'$. Likewise.

Case $\Gamma = \Gamma' \cdot V \cdot \Sigma$ s.t. $(\tau \leq \sigma) \in \Sigma$. Then $\tilde{\Gamma} = \tilde{\Gamma}'$, $\tilde{\alpha} : \llbracket \forall V\{\Sigma\} \rrbracket_{\Theta_0}$ and $\Theta = \Theta_0 \circ \llbracket V; \tilde{\alpha} \rrbracket^{\Theta}$. Since $(\tau \leq \sigma) \in \Sigma$, we have

$$\llbracket \forall V\{\Sigma\} \rrbracket_{\Theta_0} = \{\tilde{\beta} : \llbracket V \rrbracket^{\mathbf{K}} \mid \dots \wedge \llbracket \tau \rrbracket_{\Theta'_0} \triangleright \llbracket \sigma \rrbracket_{\Theta'_0} \wedge \dots\}$$

such that $\tilde{\beta}$ is not free in any type in the image of Θ_0 and $\Theta'_0 = \Theta_0 \circ \llbracket V; \tilde{\beta} \rrbracket^{\Theta}$. Note that by TYPEUNPACK and TYPEVAR, we have $\tilde{\Gamma}' \vdash \tilde{\alpha} : \llbracket V \rrbracket^{\mathbf{K}}$. By PROPRES and PROPANDPROJ, we then have $\tilde{\Gamma}' \vdash^{F_{cc}} \llbracket \tilde{\beta} \mapsto \tilde{\alpha} \rrbracket (\llbracket \tau \rrbracket_{\Theta'_0} \triangleright \llbracket \sigma \rrbracket_{\Theta'_0})$. We observe that $\llbracket \tilde{\beta} \mapsto \tilde{\alpha} \rrbracket (\llbracket \tau \rrbracket_{\Theta'_0} \triangleright \llbracket \sigma \rrbracket_{\Theta'_0})$ is the same as $\llbracket \tau \rrbracket_{\Theta} \triangleright \llbracket \sigma \rrbracket_{\Theta}$ and we conclude with COERPROP. \square

LEMMA C.11. *Let $\Gamma \sim (\tilde{\Gamma}; \Theta)$. Then $\Gamma \vdash \tau \leq \sigma$ implies that $\tilde{\Gamma} \vdash^{F_{cc}} \llbracket \tau \rrbracket_{\Theta} \triangleright \llbracket \sigma \rrbracket_{\Theta}$.*

PROOF. By induction on the subtyping derivation.

Case S-TOP, S-VARREFL. By COERTOP and COERREFL respectively.

Case S-TRANS, S-FUN. By the IH and respectively either COERTRANS or COERARR.

Case S-HYP. By Lemma C.10.

Case S-FORALL-R. By COERWEAK with $\tilde{\Sigma} = \tilde{\alpha} : \llbracket \forall V\{\Sigma\} \rrbracket_{\Theta}^{\mathbf{K}}$ and COERGEN with $\tilde{\sigma} = \tilde{\alpha}$.

Case S-FORALL-COV. Then $\tau = \forall V\{\Sigma\}. \tau'$ and $\sigma = \forall V\{\Sigma\}. \sigma'$. By COERWEAK with $\tilde{\Sigma} = \tilde{\alpha} : \llbracket \forall V\{\Sigma\} \rrbracket_{\Theta}^{\mathbf{K}}$ and COERTRANS, it suffices to show that both $\tilde{\Gamma}, \tilde{\alpha} : \llbracket \forall V\{\Sigma\} \rrbracket_{\Theta}^{\mathbf{K}} \vdash^{F_{cc}} \llbracket \tau \rrbracket_{\Theta} \triangleright \llbracket \sigma' \rrbracket_{\Theta}$ and $\tilde{\Gamma} \vdash^{F_{cc}} (\tilde{\alpha} : \llbracket \forall V\{\Sigma\} \rrbracket_{\Theta}^{\mathbf{K}} \vdash \llbracket \sigma' \rrbracket_{\Theta}) \triangleright \llbracket \sigma \rrbracket_{\Theta}$. The former can be shown by the IH. The latter can be shown by COERGEN as long as we have $\tilde{\Gamma} \vdash \llbracket \forall V\{\Sigma\} \rrbracket_{\Theta}^{\mathbf{K}}$. As a premise of S-FORALL-COV, we have $\Gamma \vdash \forall V\{\Sigma\}$, which means we can conclude by the same argument as in Lemma C.12. (Note that in this argument we use the IH rather than the entire Lemma C.11).

Case S-FORALL-L. Similarly to the previous case: by COERGEN, using Lemma C.12 and the consistency premise of S-FORALL-L.

Case S-FORALL-DISTR. Then $\tau = \forall V\{\Sigma\}. \tau_1 \rightarrow \tau_2$ and $\sigma = \tau_1 \rightarrow \forall V\{\Sigma\}. \tau_2$.

By COERWEAK with $\tilde{\Sigma} = \tilde{\alpha} : \llbracket \forall V\{\Sigma\} \rrbracket_{\Theta}^{\mathbf{K}}$ and COERTRANS, it suffices to show both of:

$$\begin{aligned} \tilde{\Gamma}, \tilde{\alpha} : \llbracket \forall V\{\Sigma\} \rrbracket_{\Theta}^{\mathbf{K}} \vdash^{F_{cc}} \llbracket \forall V\{\Sigma\}. \tau_1 \rightarrow \tau_2 \rrbracket_{\Theta} \triangleright \llbracket \tau_1 \rightarrow \tau_2 \rrbracket_{\Theta} \\ \tilde{\Gamma} \vdash^{F_{cc}} (\tilde{\alpha} : \llbracket \forall V\{\Sigma\} \rrbracket_{\Theta}^{\mathbf{K}} \vdash \llbracket \tau_1 \rightarrow \tau_2 \rrbracket_{\Theta}) \triangleright \llbracket \tau_1 \rightarrow \forall V\{\Sigma\}. \tau_2 \rrbracket_{\Theta} \end{aligned}$$

The former can be derived via COERINST with $\tilde{\sigma} = \tilde{\alpha}$. The latter simplifies by COERARR to:

$$\begin{aligned} \tilde{\Gamma}, \tilde{\alpha} : \llbracket \forall V \{ \Sigma \} \rrbracket_{\Theta}^{\kappa} \vdash^{F_{cc}} \tau_1 \triangleright \tau_1 \\ \tilde{\Gamma} \vdash^{F_{cc}} (\tilde{\alpha} : \llbracket \forall V \{ \Sigma \} \rrbracket_{\Theta}^{\kappa} \vdash \tau_2) \triangleright \llbracket \forall V \{ \Sigma \} \rrbracket_{\Theta}^{\kappa} \tau_2 \end{aligned}$$

Which hold respectively by reflexivity (COERREFL) and COERGEN. Similarly to previous cases, the premise of COERGEN can be shown using Lemma C.12 and the consistency premise of S-FORALL-DISTR. \square

LEMMA C.12. *Let $\Gamma \sim (\tilde{\Gamma}; \Theta)$. Then $\Gamma \vdash \forall V \{ \Sigma \}$ **cons.** implies that $\tilde{\Gamma} \vdash^{F_{cc}} \llbracket \forall V \{ \Sigma \} \rrbracket_{\Theta}^{\kappa}$.*

PROOF. Let $\llbracket \forall V \{ \Sigma \} \rrbracket_{\Theta}^{\kappa} = \{ \tilde{\alpha} : \llbracket V \rrbracket^{\kappa} \mid \llbracket \Sigma \rrbracket_{\Theta'} \}$. Our goal is, by definition, the same as $\tilde{\Gamma} \vdash^{F_{cc}} \exists \llbracket \forall V \{ \Sigma \} \rrbracket_{\Theta}^{\kappa}$. By PROPEXI and TYPEPACK, it suffices to show that we have $\tilde{\tau}$ such that $\tilde{\Gamma} \vdash^{F_{cc}} \tilde{\tau} : \llbracket V \rrbracket$ and $\tilde{\Gamma} \vdash^{F_{cc}} [\tilde{\alpha} \mapsto \tilde{\tau}] (\llbracket \Sigma \rrbracket_{\Theta'})$. Let θ be the substitution from $\Gamma \vdash \forall V \{ \Sigma \}$ and let $\bar{\tau}$ be its image. We pick $\tilde{\tau} = \llbracket \bar{\tau} \rrbracket_{\Theta}^{(\tau \times \tau)}$. Note that $\tilde{\tau}$ is a tuple of F_{cc} types.

We evidently have $\tilde{\Gamma} \vdash^{F_{cc}} \tilde{\tau} : \llbracket V \rrbracket$, since every type in $\bar{\tau}$ is encoded into a type of F_{cc} kind \star and $\tilde{\tau} = \llbracket \bar{\tau} \rrbracket_{\Theta}^{(\tau \times \tau)}$ is a type tuple of same size as V , meaning its F_{cc} kind is $\llbracket V \rrbracket$.

The remaining goal is $\tilde{\Gamma} \vdash^{F_{cc}} [\tilde{\alpha} \mapsto \tilde{\tau}] (\llbracket \Sigma \rrbracket_{\Theta'})$. We have $[\tilde{\alpha} \mapsto \tilde{\tau}] (\llbracket \Sigma \rrbracket_{\Theta'}) =^{F_{cc}} \llbracket \theta \Sigma \rrbracket_{\Theta}$ by F_{cc} type equality and by the definition of our translation: projections of $\tilde{\alpha}$ occur in $\llbracket \Sigma \rrbracket_{\Theta'}$ where variables from V occur in Σ . Since $\tilde{\tau} = \llbracket \bar{\tau} \rrbracket_{\Theta}^{(\tau \times \tau)} = \llbracket \theta V \rrbracket_{\Theta}^{(\tau \times \tau)}$, we have the desired F_{cc} type equality.

Therefore by PROPEQ and (repeated) PROPANPAIR, it suffices to show that $\tilde{\Gamma} \vdash^{F_{cc}} \llbracket \theta \sigma_1 \leq \theta \sigma_2 \rrbracket_{\Theta}$ for every $(\sigma_1 \leq \sigma_2) \in \Sigma$. As we have $\Gamma \vdash \forall V \{ \Sigma \}$, for each such $(\sigma_1 \leq \sigma_2)$ we have $\Gamma \vdash \theta \sigma_1 \leq \theta \sigma_2$, and we can conclude by (repeated) Lemma C.11. \square

D TYPE INFERENCE CORRECTNESS PROOFS

In this section, we sketch the proof of the core correctness theorem for type inference. First, we need to slightly modify I-ABS to account for the well-formedness requirement that our System F_{cc} translation introduced. Concretely, we add the binding for α in the premise:

$$\frac{\text{I-ABS} \quad \alpha \text{ fresh} \quad \Gamma \cdot \alpha \cdot (x : \alpha) \vdash t : \tau^+ \Rightarrow \Delta \quad V \supseteq FV(\Gamma) \quad V \# FV(\Delta) \quad V, \epsilon \vdash \epsilon \gg \Delta \gg \Xi \quad \text{split}_V(\text{uproot}(\Xi)) = (\Xi_0, \bar{y}, \Xi_1)}{\Gamma \vdash \lambda x. t : \forall \bar{y} \{ \Xi_1 \}. \alpha \rightarrow \tau^+ \Rightarrow \Xi_0}$$

We also need standard weakening and permutation lemmas.

LEMMA D.1 (WEAKENING). *Let Γ, Γ' be well-formed contexts such that $\Gamma = \Gamma_1 \cdot \Gamma_2$ and $\Gamma' = \Gamma_1 \cdot \Gamma_0 \cdot \Gamma_2$. Then:*

- $\Gamma \vdash t : \tau$ implies that $\Gamma' \vdash t : \tau$
- $\Gamma \vdash \tau \leq \sigma$ implies that $\Gamma' \vdash \tau \leq \sigma$
- $\Gamma \vdash \forall V \{ \Sigma \}$ **cons.** implies that $\Gamma' \vdash \forall V \{ \Sigma \}$ **cons.**
- $\Gamma \vdash \tau$ **wf** implies that $\Gamma' \vdash \tau$ **wf**

PROOF. By mechanical inspection of the definition of typing, subtyping, quantifier consistency, type well-formedness. No judgment inspects the number of bindings in the context. \square

LEMMA D.2 (PERMUTATION). *Let Γ, Γ' be well-formed contexts such that Γ' is a permutation of Γ . Then:*

- $\Gamma \vdash t : \tau$ implies that $\Gamma' \vdash t : \tau$
- $\Gamma \vdash \tau \leq \sigma$ implies that $\Gamma' \vdash \tau \leq \sigma$

- $\Gamma \vdash \forall V\{\Sigma\}$ **cons.** implies that $\Gamma' \vdash \forall V\{\Sigma\}$ **cons.**
- $\Gamma \vdash \tau$ **wf** implies that $\Gamma' \vdash \tau$ **wf**

PROOF. By mechanical inspection of the definition of typing, subtyping, quantifier consistency, type well-formedness. No judgment inspects the order of bindings in the context. \square

We can now restate and prove Theorem 4.5:

THEOREM (SOUNDNESS OF TYPE INFERENCE). *If $t : \tau^+ \Rightarrow \Delta$ and $\vdash \epsilon \gg \Delta \gg \Xi'$, then we have $\vdash t : \forall V\{\Xi\}. \tau^+$, where $\Xi = \text{uproot}(\Xi')$, and $V = FV(\Xi)$.*

PROOF. By Lemma D.3 we have $V \cdot \Xi \vdash t : \tau^+$. By the second conclusion of Theorem 4.7 we have $\vdash \forall V\{\Xi\}$ **cons.**. (The constraining premise we have has empty skolems, so the skolem replacement is the identity.) We can then conclude by **T-FORALL**. \square

LEMMA D.3. *Let $\Gamma \vdash t : \tau^+ \Rightarrow \Delta$. Then $V, \epsilon \vdash \epsilon \gg \Delta \gg \Xi'$ implies that $\Gamma \cdot \bar{\gamma} \cdot \Xi \vdash t : \tau^+$, where $\Xi = \text{uproot}(\Xi')$ and $\bar{\gamma} = FV(\Xi) \setminus V$.*

PROOF. By induction on the type inference derivation. All cases but **I-ABS** are trivial, as $\Xi \vdash \Delta$ by the first conclusion of Theorem 4.7. (Again, the constraining premise on which we use the theorem has empty skolems.)

In **I-ABS**, we have:

$$\begin{aligned} t &= \lambda x. t_0 \quad \tau^+ = \forall \bar{\gamma}\{\Xi_1\}. \alpha \rightarrow \tau_0^+ \quad \Delta = \Xi_0 \\ \Gamma \cdot \alpha \cdot (x : \alpha) &\vdash t_0 : \tau_0^+ \Rightarrow \Delta' \quad V', \epsilon \vdash \epsilon \gg \Delta' \gg \Xi' \\ \text{where } V' &\supseteq FV(\Gamma) \quad V' \# FV(\Delta) \quad \text{split}_{V'}(\Xi') = (\Xi_0, \bar{\gamma}, \Xi_1) \end{aligned}$$

We desire $\Gamma \cdot \Xi_0 \vdash \lambda x. t_0 : \forall \bar{\gamma}\{\Xi_1\}. \alpha \rightarrow \tau_0^+$, which by **T-FORALL** and **T-ABS** reduces to

$$\begin{aligned} \Xi_\Gamma \cdot \Xi_0 &\vdash \forall \bar{\gamma}\{\Xi_1\} \text{ **cons.**} \\ \Gamma \cdot \Xi_0 \cdot \bar{\gamma} \cdot \Xi_1 &\cdot (x : \alpha) \vdash t_0 : \tau_0^+ \end{aligned}$$

The former holds by weakening the second conclusion of Theorem 4.7. The latter holds by permutating the conclusion of the IH. \square

We now turn to prove the correctness of the subtype constraint solving, which is by far the most interesting and tricky part of the type inference algorithm.

D.1 Termination of constraining

We now justify the termination of subtype constraint solving. This property trivially implies the termination of type inference as a whole since the type inference rules are defined by structural induction on program terms. We only provide a sketch of the proof.

We need auxiliary definitions. The following definition states an invariant maintained by normal constraining derivations.

Definition D.4. V, W, Ξ is a coherent constraining context triple if V is disjoint from W and for all $B \in \Xi$, at least one hand of B is a type variable a such that $a \notin V \cdot W$.

We say such a triple is rooted in $\bar{\alpha}$ if $\bar{\alpha} \subseteq \{\text{roots}(a) \mid a \in FV(V \cdot W) \cup FV(\Xi_0)\}$.

Definition D.5. Constraints Δ are *simplifiable* (denoted $\gg \Delta \gg$) in roots $\bar{\alpha}$ if for all coherent constraining context triples V, W, Ξ rooted in $\bar{\alpha}$ there is a finite derivation of $V, W \vdash \Xi \gg \Delta \gg \Sigma^?$ for some $\Sigma^?$.

LEMMA D.6 (WORKLIST EXTENSION). *Let $V, W \vdash \Xi_0 \gg \Delta_2 \gg \Sigma_2^?$. Then $\gg \Delta_1 \gg$ in roots $\{\text{roots}(a) \mid a \in FV(V \cdot W) \cup FV(\Xi_0)\}$ implies that $V, W \vdash \Xi_0 \gg \Delta_1 \cdot \Delta_2 \gg \Sigma_3^?$ for some $\Sigma_3^?$.*

PROOF. Follows from the definitions. \square

As a terminology note, we say a constraining derivation “solves” a constraint C if it has a subderivation such that its goals Δ end in C .

THEOREM (TERMINATION OF CONSTRAINING). *For all Δ and all coherent constraining context triples V, W, ϵ , there is a finite derivation of $V, W \vdash \epsilon \gg \Delta \gg \Xi^?$ for some $\Xi^?$.*

PROOF SKETCH. The intuition behind the sketch is that constraining, if viewed as an algorithm, always terminates (i.e., there is a finite constraining derivation), since it either reduces the constraints it still has to traverse or it reduces the number of roots in Δ which still can be visited. Concretely, we sketch how to show that the premises of each rule are smaller according to some measure.

The measure that we use is an ordered pair: its first element is the amount of roots in Δ which still can be visited according to their corresponding ϕ , and its second element is the size of Δ .

The proof is by induction on the size of Δ , calculated using our measure. The proof’s cases are organized according to which constraining rule applies first. Note that in all cases, if the derivation of any premise outputs **err**, we use **C-FAIL** instead of the rule from the current case.

Rule **C-EMPTY** trivially terminates, while rules **C-TOP**, **C-SKIP**, **C-FUN** and **C-FORALL-L** all have a premise with a strictly smaller Δ .

Rules **C-RIGID-L** and **C-RIGID-R** add the result of extrusion Σ to the constraints. By the IH, we have $\gg \Delta \gg$ in current roots. We sketch how to show that we also have $\gg B \gg$ in current roots for all $B \in \Sigma$.

For all $B \in \Sigma$, one hand of B is a rigid variable, and another is a flexible variable. This means that only **C-SKIP**, **C-FLEX-L** and **C-FLEX-R** can apply. If B is not skipped and **C-RIGID-L** or **C-RIGID-R** applies, then B is added to ϕ and all the constraints added to Δ have a rigid variable on at least one hand. The latter means that either a basic rule (**C-TOP**, **C-VARREFL** or **C-SKIP**) or one of the variable rules (**C-FLEX-L**, **C-FLEX-R**, **C-RIGID-L**, **C-RIGID-R**) applies to each such new constraint. The variable rules cannot apply endlessly, since every constraint they add to Δ is also added to ϕ (indirectly in the case of **C-RIGID-L** and **C-RIGID-R**). Then, we have $\gg B \gg$ in current roots as desired.

Hence, by Lemma D.6 we have $\gg \Sigma \gg$ and $\gg \Delta \cdot \Sigma \gg$, which lets us conclude the case.

Rules **C-FLEX-L** and **C-FLEX-R** include the SRLC premise. By the IH we have $\gg \Delta \gg$. The SRLC ensures that each of the constraints added in the premise of **C-FLEX-L** and **C-FLEX-R** is smaller according to our measure. Hence, for each such constraint C we also have $\gg C \gg$, letting us conclude the case by Lemma D.6.

Finally, we consider rule **C-FORALL-R**. Constraint $\tau^+ \leq^\phi \sigma^-$ is smaller than $\Delta \cdot (\tau^+ \leq^\phi \forall \alpha. \tau^-)$, so by the IH there is a finite derivation of $V', W \cdot \beta_\alpha \vdash \epsilon \gg \tau^+ \leq^\phi \sigma^- \gg \Xi^?$ for some $\Xi^?$. If $\Xi^? = \mathbf{err}$ then we use **C-FAIL**, otherwise $\Xi^? = \Xi$.

Now we focus on the second constraining premise. By the IH, there is a finite derivation of $V, W \cdot \beta_\alpha \vdash \Xi_0 \gg \tau^+ \leq^\phi \sigma^- \gg \dots$ (the result is not relevant except that we may need to propagate **err** with **C-FAIL**). We will refer to this derivation as \mathfrak{J}_2 and to first premise’s derivation as \mathfrak{J}_1 . We will sketch how to show there is a finite derivation of $V, W \cdot \beta_\alpha \vdash \Xi_0 \gg \text{outer}_{V'}(\Xi) \gg \dots$, which can be constructed from \mathfrak{J}_1 and \mathfrak{J}_2 . The intuition is that temporarily freezing some type variables, like in the first premise’s derivation \mathfrak{J}_1 , delays solving constraints which can be solved on the spot, like in \mathfrak{J}_2 . In a sense, \mathfrak{J}_2 merely solves more constraints than \mathfrak{J}_1 . In particular, all flexible variables appearing in \mathfrak{J}_1 can be treated as also appearing in \mathfrak{J}_2 (we can act as though both derivations always pick the same variable names). Additionally, note that the only constraints solved in \mathfrak{J}_1

which are not solved in \mathfrak{J}_2 are constraints on rigid variables in $V' \setminus V$ solved by **C-RIGID-L** or **C-RIGID-R**.

We sketch how to construct the desired derivation. Every $B \in \text{outer}_{V'}(\Xi)$ was added by a subderivation of \mathfrak{J}_1 ending in **C-RIGID-L** or **C-RIGID-R**. WLOG assume the former for a particular $B = a \leq \tau_a$. Then that subderivation of \mathfrak{J}_1 solves some $B_1 = a \leq \tau_{a,1}$.

If $a \in V$, then \mathfrak{J}_2 also solves B_1 via **C-RIGID-L**. Otherwise, if a is an approximant for b in \mathfrak{J}_1 , then (and only then) B_1 was added to the constraints in \mathfrak{J}_1 by **C-RIGID-L** or **C-RIGID-R** and no subderivation in \mathfrak{J}_2 solves the exact same goal. We proceed by distinguishing two cases based on how the constraints added by **C-RIGID-L** and **C-RIGID-R** must be solved.

If $\tau_{a,1}$ is also an approximant a' in \mathfrak{J}_1 , then a and a' are the lower and upper approximants for b . The desired subderivation of $a \leq a'$ can be constructed because all constraints on a and a' solved in \mathfrak{J}_1 are instead constraints on b , a flexible variable, in \mathfrak{J}_2 . Hence, \mathfrak{J}_2 solves the constraints arising from the desired derivation needing to solve $a \leq a'$.

Otherwise, if $\tau_{a,1}$ is not an approximant in \mathfrak{J}_1 , then it is an upper bound of b and \mathfrak{J}_2 solves $b \leq \tau_{a,1}$.

Finally, if a is not an approximant in \mathfrak{J}_1 , then \mathfrak{J}_2 has a subderivation which solves $a \leq \tau_{a,1}$ via **C-FLEX-L**.

In all these cases, in the desired derivation we need to construct a (sub)derivation solving $a \leq \tau_a$. Compared to $\tau_{a,1}$ in \mathfrak{J}_2 , some type variables a_a may have been replaced by approximants when B_1 was extruded in \mathfrak{J}_1 . If such type variables are flexible in the desired derivation, i.e., absent from V , then they approximate temporarily frozen variables (ones in $V' \setminus V$). Hence, all constraints on these approximants extruded by \mathfrak{J}_1 corresponds to constraints on their underlying type variable solved by \mathfrak{J}_2 . While the constraints may not necessarily be solved in the same order in \mathfrak{J}_2 and in the desired derivation, the order of constraints does not truly matter. As Appendix ?? explains, if we reorder some goals then the derivation as a whole still traverses and solves the same constraints.

Hence, there is a finite derivation of $V, W \cdot \beta_\alpha \vdash \Xi_0 \gg \text{outer}_{V'}(\Xi) \gg \dots$, which can also be derived in W because $\beta_\alpha \notin FV(\Xi)$. Therefore, by Lemma D.6 there is also a finite derivation of $V, W \vdash \Xi_0 \gg \Delta \cdot \text{outer}_{V'}(\Xi) \gg \Sigma^?$ for some $\Sigma^?$. If $\Sigma^? = \mathbf{err}$ then we use **C-FAIL**, otherwise we can conclude with **C-FORALL-R**. \square

D.2 Soundness of constraining

We now justify that subtype constraining, when it succeeds, leads to valid and consistent subtyping, which allows us to use it in the proof of type inference soundness.

D.2.1 Full propagation. To prove type inference soundness, we define a new inductive relation on bounds context, called *full propagation*, which is a restricted form of subtyping without a transitivity rule and which corresponds closely with what a constraint closure and resolution algorithm should compute as its output context. We show that our subtype constraining algorithm implies full propagation in its output context and then we show that full propagation itself implies that the context's subtyping relationships can be derived *and* that they are consistent. We show consistency by exhibiting a solution to the bounds in the form of the *lower-bounds union* (LBU).

Definition D.7 (Full propagation). We say V, Σ is *fully-propagated* if $V, W, \epsilon \vdash \Sigma \dashv \Sigma$, where the judgment is defined in Figure 11. We only consider well-formed full propagation derivations where $V \# W$.

D.2.2 Concrete types.

$$\boxed{V, W, \Sigma \vdash \tau^+ \leq \tau^- \vdash \Sigma'}$$

We define $V, W, \Sigma \vdash \bar{B} \vdash \Sigma' \triangleq \overline{V, W, \Sigma \vdash B \vdash \Sigma'}$

$$\begin{array}{c}
\text{P-TOP} \\
\hline
V, W, \Sigma \vdash \tau^+ \leq \top \vdash \Sigma' \\
\\
\text{P-RIGID-L} \\
\frac{\alpha \in V \quad \rho = [\beta^- \mapsto \sigma^-, \beta^+ \mapsto \sigma^{+\beta \in FV(\tau^-)} \setminus V] \quad (\alpha \leq \rho \tau^-) \in \Sigma'}{V \vdash \tau^- \text{ X-ok} \quad FV(\rho \tau^-) \subseteq V \quad \overline{V, W, \Sigma \vdash (\sigma^- \leq \beta) \cdot (\beta \leq \sigma^+) \vdash \Sigma'}^{\beta \in FV(\tau^-)} \setminus V}{V, W, \Sigma \vdash \alpha \leq \tau^- \vdash \Sigma'} \\
\\
\text{P-RIGID-R} \\
\frac{\alpha \in V \quad \rho = [\beta^- \mapsto \sigma^-, \beta^+ \mapsto \sigma^{+\beta \in FV(\tau^+)} \setminus V] \quad (\rho \tau^+ \leq \alpha) \in \Sigma'}{V \vdash \tau^+ \text{ X-ok} \quad FV(\rho \tau^+) \subseteq V \quad \overline{V, W, \Sigma \vdash (\sigma^- \leq \beta) \cdot (\beta \leq \sigma^+) \vdash \Sigma'}^{\beta \in FV(\tau^+)} \setminus V}{V, W, \Sigma \vdash \tau^+ \leq \alpha \vdash \Sigma'} \\
\\
\text{P-FLEX} \\
\frac{\alpha \notin V \cdot W \quad (\tau^+ \leq \sigma^-) \in \Sigma' \quad (\tau^+ \leq \sigma^-) \in \{(\alpha \leq \sigma^-), (\tau^+ \leq \alpha)\}}{\alpha \notin V \cdot W \quad (\tau^+ \leq \sigma^-) \notin \Sigma' \quad (\tau^+ \leq \sigma^-) \in \Sigma} \quad \frac{\alpha \notin V \cdot W \quad (\tau^+ \leq \sigma^-) \in \Sigma' \quad (\tau^+ \leq \sigma^-) \in \{(\alpha \leq \sigma^-), (\tau^+ \leq \alpha)\}}{V, W, \Sigma \vdash CLB_V(\tau^+, \Sigma') \leq \sigma^{-'} \vdash \Sigma'}^{\sigma^{-'} \in CUB_V(\sigma^-, \Sigma')}{V, W, \Sigma \vdash \tau^+ \leq \sigma^- \vdash \Sigma'} \\
\\
\text{P-FORALL-L} \\
\frac{V, W, \Sigma \vdash [\bar{\alpha} \mapsto \bar{\pi}] \Xi \vdash \Sigma' \quad V, W, \Sigma \vdash [\bar{\alpha} \mapsto \bar{\pi}] \tau^+ \leq \sigma^- \vdash \Sigma'}{V, W, \Sigma \vdash \forall \bar{\alpha} \{ \Xi \}. \tau^+ \leq \sigma^- \vdash \Sigma'} \\
\\
\text{P-FORALL-R} \\
\frac{\beta \text{ fresh} \quad V' \supseteq FV(\tau^+ \leq \sigma^-) \setminus W \quad V' \# W \cdot \beta \text{ acyclic}(\Xi) \quad V', W \cdot \beta, \epsilon \vdash \tau^+ \leq [\bar{\alpha} \mapsto \beta] \sigma^- \vdash \Xi \quad V', W \cdot \beta, \epsilon \vdash \Xi \vdash \Xi \quad V, W, \Sigma \vdash \text{outer}_{V'}(\Xi) \vdash \Sigma'}{V, W, \Sigma \vdash \tau^+ \leq \forall \alpha. \sigma^- \vdash \Sigma'}
\end{array}$$

Fig. 11. The polarized full-propagation rules. Note that in **P-RIGID-L** and **P-RIGID-R**, σ^- and σ^+ are both positive and negative, i.e. System F types. In **P-FORALL-L**, π can be either positive or negative as long as τ^+ remains a valid positive type after substitutions.

Definition D.8 (Left-Concrete Types). A left-concrete type is *not* a flexible type variable (that may be under universal quantifiers).

$$\top \text{ concrete}_V^+ \quad \tau_1 \rightarrow \tau_2 \text{ concrete}_V^+ \quad \frac{a \in V}{a \text{ concrete}_V^+} \quad \frac{\tau^+ \text{ concrete}_V^+}{\forall \bar{\alpha} \{ \Sigma \}. \tau^+ \text{ concrete}_V^+}$$

Definition D.9 (Right-Concrete Types). A right-concrete type is *not* a flexible type variable (that may be under spurious universal quantifiers such as $\forall \alpha. \beta$ where β is flexible).

$$\top \text{ concrete}_V^- \quad \tau_1 \rightarrow \tau_2 \text{ concrete}_V^- \quad \frac{a \in V}{a \text{ concrete}_V^-} \quad \frac{\alpha \in FV(\tau^-)}{\forall \alpha. \tau^- \text{ concrete}_V^-} \quad \frac{\tau^- \text{ concrete}_V^-}{\forall \alpha. \tau^- \text{ concrete}_V^-}$$

Definition D.10 (Concrete Lower Bounds (CLB)). Assume universal types are well-formed and no name collision of universal quantifiers. $CLB_V(\tau, \Xi) = go_V^+(\tau, \Xi, \epsilon)$. go^+ returns τ itself if τ is

left-concrete or the union of *left-concrete* lower bounds of τ in Ξ that can be transitively reached.

$$\begin{aligned} go_V^+(\tau^+, \Xi, W) &= \tau^+ && \text{if } \tau^+ \text{ \textbf{concrete}}_V^+ \\ go_V^+(a, \Xi, W) &= \sqrt{go_V^+(\tau^+, \Xi, W \cdot a)}^{(\tau^+ \leq a) \in \Xi, \tau^+ \notin (W \cdot a)} && \text{otherwise} \\ go_V^+(\forall \bar{\alpha}\{\Sigma\}. \tau^+, \Xi, W) &= \forall \bar{\alpha}\{\Sigma\}. go_V^+(\tau^+, \Sigma \cdot \Xi, W) && \text{otherwise} \end{aligned}$$

Example D.11. $CLB_\epsilon(a, (\forall \alpha \beta \gamma \{\alpha \rightarrow \beta \leq \gamma, \beta \rightarrow \alpha \leq \gamma, \alpha \leq \beta\}. \gamma \leq a)) = \forall \alpha \beta \gamma \{\alpha \rightarrow \beta \leq \gamma, \beta \rightarrow \alpha \leq \gamma, \alpha \leq \beta\}. (\alpha \rightarrow \beta) \vee (\beta \rightarrow \alpha)$.

Example D.12. Let $\Xi = (\forall \beta \{\Xi'\}. \beta \leq \alpha) \cdot (\tau \leq \alpha) \cdot (\sigma' \leq \gamma)$ and $\Xi' = (\alpha \leq \beta) \cdot (\sigma \leq \beta) \cdot (\gamma \leq \beta)$ where τ, σ, σ' are left-concrete. Then $CLB_\epsilon(\alpha, \Xi) = \tau \vee \forall \beta \{\Xi'\}. (\sigma \vee \sigma')$.

Definition D.13 (Concrete Upper Bounds (CUB)). $CUB_V(a, \Xi) = go_V^-(a, \Xi, \epsilon)$. go^- calculates a set of all *right-concrete* upper bounds of a in Ξ that can be transitively reached.

$$\begin{aligned} go_V^-(\tau^-, \Xi, W) &= \tau^- && \text{if } \tau^- \text{ \textbf{concrete}}_V^- \\ go_V^-(a, \Xi, W) &= \overline{go_V^-(\tau^-, \Xi, W \cdot a)}^{(a \leq \tau^-) \in \Xi, \tau^- \notin (W \cdot a)} && \text{otherwise} \\ go_V^-(\forall \alpha. \tau^-, \Xi, W) &= go_V^-(\tau^-, \Xi, W) && \text{otherwise} \end{aligned}$$

The following properties of CLB and CUB are obvious:

PROPOSITION D.14. $CLB_V(\tau, \Xi)$ is always defined.

PROPOSITION D.15. $CUB_V(\tau, \Xi)$ is always defined.

PROPOSITION D.16. $CLB_V(\tau, \Xi) = CLB_V(\tau, \Xi \cdot (\tau \leq \sigma))$

PROPOSITION D.17. $CUB_V(\tau, \Xi) = CUB_V(\tau, \Xi \cdot (\sigma \leq \tau))$

D.2.3 Lower Bound Union (LBU).

Definition D.18 (Declarative Lower Bound Union). For all Ξ and V and set of proper substitutions ρ , if every substitution $(\alpha \mapsto \tau) \in \rho$ satisfies $\tau = \rho CLB_V(\alpha, \Xi)$, then $\rho \in isLBU_V^\Xi$.

Definition D.19 (Algorithmic Lower Bound Union). We obtain the algorithmic lower bound union by $mkLBU(\Xi, V, \epsilon)$ which is a function call that returns a list of substitutions:

$$mkLBU(\Xi, V, W) = \overline{\alpha \mapsto [mkLBU(\Xi, V, W \cdot \alpha)] CLB_V(\alpha, \Xi)}^{\alpha \in FV(\Xi) \setminus (V \cdot W)}$$

PROPOSITION D.20. *If acyclic*(Ξ), then:

- (1) $mkLBU(\Xi, V, \epsilon) \in isLBU_V^\Xi$
- (2) $mkLBU(\Xi, V, \epsilon)$ is always defined.

PROOF SKETCH. For each type variable $a \in FV(\Xi) \setminus V$, its CLB can always be calculated because each type variable is only accessed at most once in the bounds graph. Due to the acyclicity condition of Ξ , there are no cyclic occurrences of α transitively reachable in the concrete lower bounds of α , so no infinite substitution of LBU would happen when recursively substituting occurrences of other type variables in the concrete lower bounds, and we can disregard α in the recursive computation of lower bound union of other type variables. \square

PROPOSITION D.21. *If acyclic*(Ξ), then there always exists a proper substitution ρ in $isLBU_V^\Xi$.

PROOF. We provide $mkLBU(\Xi, V, \epsilon)$ as the witness. \square

Notation: we write LBU_V^Ξ as an arbitrary instance in $isLBU_V^\Xi$.

D.2.4 Transitive reachability of bounds.

Definition D.22 (Transitively reachable bounds). Given a rigid type variable set V , we denote a flexible type variable α transitively reaches a lower bound τ^+ in context Σ through a set of flexible type variables $\bar{\beta}$ as $\tau^+ \leq_V^* \alpha \in \Sigma \sim \bar{\beta}$.

$$\frac{\tau^+ \leq \alpha \in \Sigma \quad \alpha \notin V}{\tau^+ \leq_V^* \alpha \in \Sigma \sim \alpha} \quad \frac{\tau^+ \leq_V^* \beta \in \Sigma \sim \bar{\gamma} \quad \beta \leq \alpha \in \Sigma \quad \alpha, \beta \notin V}{\tau^+ \leq_V^* \alpha \in \Sigma \sim \bar{\gamma} \cdot \beta}$$

Similarly, we denote an upper bound that is transitively reachable as $\alpha \leq_V^* \tau^- \in \Sigma \sim \bar{\beta}$.

D.2.5 Propagation context adjustment.

LEMMA D.23. *If for all $a \notin V$, $\tau_i^- \in \text{CUB}_V(\tau^-, \Xi')$, we have $V, W, \Xi \cdot (a \leq \tau^-) \vdash \text{CLB}_V(a, \Xi') \leq \tau_i^- \dashv \Xi'$ and acyclic($\Xi \cdot \Xi' \cdot (a \leq \tau^-)$), then $V, W, \Xi \vdash \text{CLB}_V(a, \Xi') \leq \tau_i^- \dashv \Xi' \cdot (a \leq \tau^-)$.*

LEMMA D.24. *If for all $a \notin V$, $\tau_i^- \in \text{CUB}_V(a, \Xi')$, we have $V, W, \Xi \cdot (\tau^+ \leq a) \vdash \text{CLB}_V(\tau^+, \Xi') \leq \tau_i^- \dashv \Xi'$ and acyclic($\Xi \cdot \Xi' \cdot (\tau^+ \leq a)$), then $V, W, \Xi \vdash \text{CLB}_V(\tau^+, \Xi') \leq \tau_i^- \dashv \Xi' \cdot (\tau^+ \leq a)$.*

PROOF SKETCH OF D.23. Essentially, the types on both sides of the propagation relationship to prove are all concrete. The typical case is when two function types, say $l_1 \rightarrow l_2$ and $r_1 \rightarrow r_2$, are on both sides. Due to the acyclicity condition, a would never appear (or be transitively reached) in the positive/negative positions of l_2/r_2 and negative/positive positions of l_1/r_1 . For the positions otherwise, it is possible that a appears due to syntactically cyclic upper bounds, e.g. $l_1 \rightarrow l_2 \leq a \leq \tau^- = a \rightarrow \top$, we need to propagate $a \leq l_1$ in $\Xi' \cdot (a \leq \tau^-)$, while $a \leq \tau^-$ does not create any new concrete lower bound for a so it is equivalent to propagating it in Ξ' , which is already known. Note that it is impossible that $a \leq l_1$ is propagated by **P-SKIP** when $l_1 = \tau^- = a \rightarrow \top$ because this violates the acyclicity condition. So generally, we can inductively reconstruct the whole derivation of the assumed propagation relations with the contexts adjusted as that of the goals. For **D.24**, the proof idea is the same, while there is no syntactic cyclicity in lower bounds. \square

LEMMA D.25. *If for all τ^+ and $a \notin V$:*

- (1) $V, W, \Xi \cdot (\tau^+ \leq a) \vdash \Xi' \dashv \Xi'$
- (2) $V, W, \Xi \cdot (\tau^+ \leq a) \vdash \sigma^+ \leq \sigma^- \dashv \Xi'$
- (3) acyclic($\Xi \cdot \Xi' \cdot (\tau^+ \leq a)$)
- (4) For all $\tau_i^- \in \text{CUB}_V(a, \Xi')$, we have $V, W, \Xi \cdot (\tau^+ \leq a) \vdash \text{CLB}_V(\tau^+, \Xi') \leq \tau_i^- \dashv \Xi'$

then $V, W, \Xi \vdash \sigma^+ \leq \sigma^- \dashv \Xi' \cdot (\tau^+ \leq a)$.

LEMMA D.26. *If for all τ^- and $a \notin V$:*

- (1) $V, W, \Xi \cdot (a \leq \tau^-) \vdash \Xi' \dashv \Xi'$
- (2) $V, W, \Xi \cdot (a \leq \tau^-) \vdash \sigma^+ \leq \sigma^- \dashv \Xi'$
- (3) acyclic($\Xi \cdot \Xi' \cdot (a \leq \tau^-)$)
- (4) For all $\tau_i^- \in \text{CUB}_V(\tau^-, \Xi')$, we have $V, W, \Xi \cdot (a \leq \tau^-) \vdash \text{CLB}_V(a, \Xi') \leq \tau_i^- \dashv \Xi'$

then $V, W, \Xi \vdash \sigma^+ \leq \sigma^- \dashv \Xi' \cdot (a \leq \tau^-)$.

PROOF SKETCH OF D.26 (D.25 IS SYMMETRIC). By induction on the full propagation derivation from (2). Intuitively, this lemma means bounds are still consistent even with a new upper bound to propagate after checking that the new bound is consistent and propagated to the original bounds.

Cases P-TOP, P-VARREFL, P-FUN Immediate.

Cases P-RIGID-L/P-RIGID-R WLOG we consider **P-RIGID-L**. We are to prove $V, W, \Xi \vdash b \leq \sigma^- \dashv \Xi' \cdot (a \leq \tau^-)$. Since $a \notin V$ and $b \in V$, $a \neq b$. We conclude with IH.

Case P-SKIP By assumption, we have $(\sigma^+ \leq \sigma^-) \in (\Xi \cdot (a \leq \tau^-))$ and $(\sigma^+ \leq \sigma^-) \notin \Xi'$. Then the only case we need to consider is when $(\sigma^+ \leq \sigma^-) = (a \leq \tau^-)$, otherwise we conclude with **P-SKIP** since $(\sigma^+ \leq \sigma^-) \in \Xi$.

By **P-FLEX**, we need to prove for all $\tau_i^- \in CUB_V(\tau^-, \Xi' \cdot (a \leq \tau^-))$, we have $V, W, \Xi \vdash CLB_V(a, \Xi' \cdot (a \leq \tau^-)) \leq \tau_i^- \dashv \Xi' \cdot (a \leq \tau^-)$.

It is immediate to see that $CLB_V(a, \Xi' \cdot (a \leq \tau^-)) = CLB_V(a, \Xi')$ since $a \leq \tau^-$ does not introduce new lower bounds for a .

Consider τ^- . If τ^- is right-concrete, then $CUB_V(\tau^-, \Xi' \cdot (a \leq \tau^-)) = CUB_V(\tau^-, \Xi') = \tau^-$. Otherwise, we still have $CUB_V(\tau^-, \Xi' \cdot (a \leq \tau^-)) = CUB_V(\tau^-, \Xi')$ because τ^- , as a flexible type variable (possibly under spurious universal quantifiers), could only be visited once.

Now, the goal is to prove for all $\tau_i^- \in CUB_V(\tau^-, \Xi')$, $V, W, \Xi \vdash CLB_V(a, \Xi') \leq \tau_i^- \dashv \Xi' \cdot (a \leq \tau^-)$.

By assumption, we have for all $\tau_i^- \in CUB_V(\tau^-, \Xi')$, $V, W, \Xi \cdot (a \leq \tau^-) \vdash CLB_V(a, \Xi') \leq \tau_i^- \dashv \Xi'$. We conclude with **D.23**.

Case P-FLEX

For this case we have:

$$\frac{b \notin V \quad (\sigma^+ \leq \sigma^-) \in \{(b \leq \sigma^-), (\sigma^+ \leq b)\} \quad (\sigma^+ \leq \sigma^-) \in \Xi' \quad \frac{V, W, \Xi \cdot (a \leq \tau^-) \vdash CLB_V(\sigma^+, \Xi') \leq \sigma_i^- \dashv \Xi' \sigma_i^- \in CUB_V(\sigma^-, \Xi')}{V, W, \Xi \cdot (a \leq \tau^-) \vdash \sigma^+ \leq \sigma^- \dashv \Xi'}}{V, W, \Xi \cdot (a \leq \tau^-) \vdash \sigma^+ \leq \sigma^- \dashv \Xi'}$$

We denote CLB types $CLB_V(\sigma^+, \Xi')$ as C_l , $CUB_V(\sigma^-, \Xi')$ as C_r , $CLB_V(\sigma^+, \Xi' \cdot (a \leq \tau^-))$ as C'_l , and $CUB_V(\sigma^-, \Xi' \cdot (a \leq \tau^-))$ as C'_r .

By **P-FLEX**, the goal is to prove for all $\sigma_i^- \in C'_r$, we have $V, W, \Xi \vdash C'_l \leq \sigma_i^- \dashv \Xi' \cdot (a \leq \tau^-)$. We discuss the concreteness of τ^- and σ^+/σ^- to discover how the new bound $a \leq \tau^-$ influences the relationships between C_r/C'_r and C_l/C'_l . We use the notation \equiv to describe the equivalency between CLB types, i.e. a CLB type can be replaced by an equivalent one in propagation judgments.

Consider the following cases:

- $(\sigma^+ \leq \sigma^-) = (b \leq \sigma^-)$.

By IH, for all $\sigma_i^- \in C_r$, we have $V, W, \Xi \vdash C_l \leq \sigma_i^- \dashv \Xi' \cdot (a \leq \tau^-)$.

By **D.23**, for all $\tau_i^- \in CUB_V(\tau^-, \Xi')$, we have $V, W, \Xi \vdash CLB_V(a, \Xi') \leq \tau_i^- \dashv \Xi' \cdot (a \leq \tau^-)$ (*).

Consider the concreteness of τ^- and σ^- :

- τ^- **concrete** $_{\bar{V}}$ and σ^- **concrete** $_{\bar{V}}$. We conclude with the IH since $C_l = C'_l, C_r = C'_r$.
- τ^- **concrete** $_{\bar{V}}$ and $\neg(\sigma^-$ **concrete** $_{\bar{V}})$. We have $C_l = C'_l$. If σ^- cannot transitively reach a as an upper bound in $\Xi' \cdot (a \leq \tau^-)$, then we conclude with IH since σ^- does not introduce new upper bounds to propagate, i.e. $C_r = C'_r$. Otherwise, σ^- transitively reaches a , i.e. in the bounds graph (with the new bound $(a \leq \tau^-)$ marked red):

$$b \leq \sigma^- \leq^* a \leq \tau^-$$

then $C'_r = C_r \cup \{\tau^-\}$ since $CUB_V(\tau^-, \Xi' \cdot (a \leq \tau^-)) = \{\tau^-\}$. For all type in C_r , we conclude with IH. We want $V, W, \Xi \vdash C_l \leq \tau^- \dashv \Xi' \cdot (a \leq \tau^-)$, which is implied by (*) since a transitively has the lower bound b .

- $\neg(\tau^-$ **concrete** $_{\bar{V}})$ and σ^- **concrete** $_{\bar{V}}$. We have $C_r = C'_r = \{\sigma^-\}$. If b cannot transitively reach τ^- , then we conclude with IH since $(a \leq \tau^-)$ does not introduce new lower bounds to propagate, i.e. $C_l = C'_l$. Otherwise, b transitively reaches τ^- , i.e.:

$$a \leq \tau^- \leq^* b \leq \sigma^-$$

then $C'_l \equiv C_l \vee CLB_V(a, \Xi')$. The bounds C_l being less than each type in C_r are already propagated by IH. We want $V, W, \Xi \vdash CLB_V(a, \Xi') \leq \sigma^- \dashv \Xi' \cdot (a \leq \tau^-)$, which is implied by (*) since τ^- transitively has the upper bound σ^- .

- $\neg(\tau^- \text{ concrete}_V^-)$ and $\neg(\sigma^- \text{ concrete}_V^-)$. For the cases when (1) both τ^-/σ^- could not transitively reach b/a , (2) only one of τ^-/σ^- transitively reaches b/a , we apply the same reasoning in the former cases. When both τ^-/σ^- transitively reaches b/a ,

$$\begin{array}{ccc} a & \leq & \tau^- \\ \overset{*}{\forall} \downarrow & & \downarrow \overset{*}{\wedge} \\ \sigma^- & \geq & b \end{array}$$

$C'_r = C_r \cdot CUB_V(\tau^-, \Xi')$ and $C'_l \equiv C_l \vee CLB_V(a, \Xi')$. By (*), those relationships are propagated since adding $a \leq \tau^-$ would unify a, τ^-, b , and σ^- .

- $(\sigma^+ \leq \sigma^-) = (\sigma^+ \leq b)$. We apply analogous reasoning to the case above. We discuss the concreteness of τ^- and σ^+/σ^- and whether b/σ^- can transitively reach a/τ^- to discover how the new bound $a \leq \tau^-$ influences the relationships between C_r/C'_r and C_l/C'_l .

Case P-FORALL-L By IH.

Case P-FORALL-R By the IH on the *third* constraining premise; other premises are unchanged. \square

LEMMA D.27. *If for all τ^- and $FV(\tau^- \leq \tau^+) \subseteq V$ and $V, W, \Xi \vdash \sigma^+ \leq \sigma^- \dashv \Xi'$ then $V, W, \Xi \vdash \sigma^+ \leq \sigma^- \dashv \Xi' \cdot (\tau^+ \leq \tau^-)$.*

D.2.6 Correctness of CLB and CUB.

LEMMA D.28. *For all τ^- and $a \notin V$ and Σ and Σ' , if*

- (1) *for all $\sigma^+ \leq b \in \Sigma$ and $b \leq \sigma^- \in \Sigma'$, we have $V, W, \Sigma \vdash \sigma^+ \leq \sigma^- \dashv \Sigma'$; and*
- (2) *for all $\tau^+ \leq a \in \Sigma'$, we have $V, W, \Sigma \vdash \tau^+ \leq \tau^- \dashv \Sigma'$*

then for all $\tau_i^- \in CUB_V(\tau^-, \Sigma')$, we have $V, W, \Sigma \vdash CLB_V(a, \Sigma') \leq \tau_i^- \dashv \Sigma'$.

LEMMA D.29. *For all τ^+ and $a \notin V$ and Σ and Σ' , if*

- (1) *for all $b \leq \sigma^- \in \Sigma$ and $\sigma^+ \leq b \in \Sigma'$, we have $V, W, \Sigma \vdash \sigma^+ \leq \sigma^- \dashv \Sigma'$; and*
- (2) *for all $a \leq \tau^- \in \Sigma'$, we have $V, W, \Sigma \vdash \tau^+ \leq \tau^- \dashv \Sigma'$*

then for all $\tau_i^- \in CUB_V(a, \Sigma')$, we have $V, W, \Sigma \vdash CLB_V(\tau^+, \Sigma') \leq \tau_i^- \dashv \Sigma'$.

PROOF SKETCH OF D.28 (D.29 IS SIMILAR). By the definition of CLB,

$$\begin{aligned} CLB_V(a, \Sigma') &= go_V^+(a, \Sigma', \epsilon) \\ &= \bigvee \overline{go_V^+(\tau^+, \Sigma', \{a\})}^{(\tau^+ \leq a) \in \Sigma', \tau^+ \notin \{a\}} \end{aligned}$$

By P-FORALL-L, for each $(\tau^+ \leq a) \in \Sigma'$ and $\tau^+ \notin \{a\}$, we need to prove for all $\tau_i^- \in CUB_V(\tau^-, \Sigma')$, we have $V, W, \Sigma \vdash go_V^+(\tau^+, \Sigma', \{a\}) \leq \tau_i^- \dashv \Sigma'$. We conclude by invoking Lemma D.30 with the assumption that CLB and CUB always terminate and the following properties: for all $b \leq a \in \Xi'$ where $a, b \notin V$, if $\tau^- \leq_V^* b \in \Xi' \sim \bar{\beta}$ then $V, W, \Xi \vdash \tau^+ \leq a \dashv \Xi'$. \square

LEMMA D.30. *For all $\tau^-, \tau_0^+, \bar{\alpha}, a \notin V, \Sigma$, and Σ' , if*

- (1) $\tau_0^+ \leq_V^* a \in \Sigma' \sim \bar{\alpha}$; and
- (2) *for all $\sigma^+ \leq b \in \Sigma$ and $b \leq \sigma^- \in \Sigma'$, we have $V, W, \Sigma \vdash \sigma^+ \leq \sigma^- \dashv \Sigma'$; and*
- (3) *for all τ^+ such that $\tau^+ \leq_V^* a \in \Sigma' \sim \bar{\beta}$, we have $V, W, \Sigma \vdash \tau^+ \leq \tau^- \dashv \Sigma'$*

then for all $\tau_i^- \in \text{CUB}_V(\tau^-, \Sigma')$, if there exists k and l such that the recursion depth of the CUB call is k and the recursion depth of the following go^+ call is l , then we have $V, W, \Sigma \vdash go_V^+(\tau_0^+, \Sigma', \bar{\alpha}) \leq \tau_i^- \dashv \Sigma'$.

PROOF SKETCH OF D.30. We prove by induction on l and case analysis on the concreteness of each τ^+ (IH1):

- $l = 0$ and τ^+ **concrete** $_V^+$.

Then $go_V^+(\tau^+, \Sigma', \{a\}) = \tau^+$. By assumption, we have $V, W, \Sigma \vdash \tau^+ \leq \tau^- \dashv \Sigma'$ (*). We now perform induction on k (IH2).

- $k = 0$, i.e. $\tau^- = \tau_i^-$ and τ^- **concrete** $_V^-$, we conclude immediately.
- $k > 0$, i.e. $\neg(\tau^- \text{concrete}_V^-)$ and τ^- transitively has the concrete upper bound τ_i^- . We then perform induction on the propagation assumption (*) (IH3).

Case P-TOP Impossible as \top is concrete.

Cases P-VARREFL, P-FUN, P-RIGID-R Impossible since τ^- is not concrete.

Case P-RIGID-L $\tau^+ = b$, $b \in V$, $(b \leq \tau') \in \Sigma'$, $FV(\tau') \in V$, and $V, W, \Sigma \vdash \tau' \leq \tau^- \dashv \Sigma'$. By $FV(\tau') \in V$, τ' **concrete** $_V^+$. By IH3, for all $\tau_i^- \in \text{CUB}_V(\tau^-, \Sigma')$, we have $V, W, \Sigma \vdash \tau' \leq \tau_i^- \dashv \Sigma'$. We conclude with **P-RIGID-L**.

Case P-SKIP Since τ^+ **concrete** $_V^+$, we can only have $\tau^- = b$, $b \notin V$, $(\tau^+ \leq b) \notin \Sigma'$, and $(\tau^+ \leq b) \in \Sigma$. There is a τ' that is non-right-concrete and transitively has the concrete upper bound τ_i^- and $(b \leq \tau') \in \Sigma'$. By assumption (1), we have $V, W, \Sigma \vdash \tau^+ \leq \tau' \dashv \Sigma'$. We conclude with IH2 since τ' goes through $k - 1$ non-right-concrete types to reach τ_i^- .

Case P-FLEX Immediate.

Case P-FORALL-L By IH3.

Case P-FORALL-R $\tau^- = \forall b. \sigma^-$ and $b \notin FV(\sigma^-)$ and $\forall b. \sigma^-$ is not concrete. Notice that $\forall b. \sigma^-$ is essentially a flexible type variable (b') under spurious universal quantifiers, we have $(\tau^+ \leq b') \in \Xi$ and also $(\tau^+ \leq b') \in \text{outer}_{V'}(\Xi)$ since $V' \supseteq FV(\tau^+ \leq b')$. We conclude with IH3.

- $l > 0$ and $\tau^+ = \forall \bar{\alpha}\{\Xi\}$. τ' and $\neg(\tau' \text{concrete}_V^+)$.

Then $go_V^+(\forall \bar{\alpha}\{\Xi\}, \tau', \Sigma', \{a\}) = \forall \bar{\alpha}\{\Xi\}$. $go_V^+(\tau', \Sigma' \cdot \Xi, \{a\})$.

By assumption, $V, W, \Sigma \vdash \forall \bar{\alpha}\{\Xi\}. \tau' \leq \tau^- \dashv \Sigma'$.

We apply similar reasoning as the former case ($l = 0$) by induction on the number of non-right-concrete types that CUB traverses and the propagation derivation. In case **P-FORALL-L**, we have Ξ propagated after substitution. We use this and **P-FORALL-L** to conclude.

- $l > 0$ and $\tau^+ = b$ and $\neg(b \text{concrete}_V^+)$.

Then $go_V^+(b, \Sigma', \{a\}) = \bigvee \overline{go_V^+(\tau', \Sigma', \{a, b\})}^{(\tau' \leq b) \in \Sigma', \tau' \notin \{a, b\}}$. We conclude with **P-FORALL-L** and IH1 since each τ' traverses $l - 1$ non-left-concrete types to reach a left-concrete type. \square

Definition D.31. $\phi \subseteq \Xi \triangleq \overline{B \in \Xi}^{B \in \phi}$. Note that \in does not look past \triangleright in ϕ .

D.2.7 Constraining implies full propagation.

LEMMA D.32. If $V, W \vdash \Xi_0 \gg \Delta \gg \Xi'_0$ and $\overline{\phi \subseteq \Xi_0}^{(\tau^+ \leq \phi \tau^-) \in \Delta}$ and *acyclic*($\Xi \cdot \Xi'$) and $V \# W$, then:

- (1) for all $\tau^+ \leq \phi \tau^- \in \Delta$, we have $V, W, \Xi \vdash \text{uproot}(\tau^+ \leq \tau^-) \dashv \Xi'$.

- (2) (a) for all $\tau^+ \leq a \in \Xi$ and $a \leq \tau^- \in \Xi'$ where $a \notin V$, we have $V, W, \Xi \vdash \tau^+ \leq \tau^- \dashv \Xi'$; and
 (b) for all $a \leq \tau^- \in \Xi$ and $\tau^+ \leq a \in \Xi'$ where $a \notin V$, we have $V, W, \Xi \vdash \tau^+ \leq \tau^- \dashv \Xi'$.
 (3) $V, W, \Xi \vdash \Xi' \dashv \Xi'$

where $\text{uproot}(\Xi_0) = \Xi$ and $\text{uproot}(\Xi'_0) = \Xi'$.

PROOF. By induction on the constraining derivation. The cases where the proof is sketched are clearly marked. We implicitly *uproot* types with roots when they are used in propagation judgments.

Cases C-EMPTY, C-TOP, C-VARREFL, C-FUN, C-FAIL Immediate.

Case C-SKIP (2) and (3) holds by the IH. For (1), WLOG we consider $\Delta = \Delta' \cdot (a \leq \tau^-)$. If $(a \leq \tau^-) \in \Xi'$ then we conclude with the IH's (3). Otherwise, we conclude with **P-SKIP** since $(a \leq \tau^-) \in \phi$ and $\phi \subseteq \Xi$.

Case C-FLEX-L/C-FLEX-R WLOG we consider **C-FLEX-L**.

$$\frac{B = (a \leq \tau^-) \quad \text{root}(B) \notin \text{roots}(\phi) \quad \text{acyclic}(\Xi \cdot B)}{a \notin V \cdot W \quad V, W \vdash \Xi \cdot B \gg \Delta \cdot (\tau^+ \leq^{B \cdot \phi} \tau^-) \xrightarrow{(\tau^+ \leq a) \in \Xi} \gg \Xi'} \gg \Xi'$$

$$\frac{}{V, W \vdash \Xi \gg \Delta \cdot (a \leq^\phi \tau^-) \gg \Xi' \cdot B}$$

- (1) The goal is $V, W, \Xi \vdash \Delta \cdot (a \leq \tau^-) \dashv \Xi' \cdot (a \leq \tau^-)$.
 To propagate $(a \leq \tau^-)$, since $(a \leq \tau^-) \in (\Xi' \cdot (a \leq \tau^-))$ and $a \notin V$, by **P-FLEX**, we want:

for all $\tau_i^- \in \text{CUB}_V(\tau^-, \Xi' \cdot (a \leq \tau^-))$, we have

$$V, W, \Xi \vdash \text{CLB}_V(a, \Xi' \cdot (a \leq \tau^-)) \leq \tau_i^- \dashv \Xi' \cdot (a \leq \tau^-).$$

By the IH's (2) (b), we have $\overline{V, W, \Xi \cdot (a \leq \tau^-) \vdash \tau^+ \leq \tau^- \dashv \Xi'}^{(\tau^+ \leq a) \in \Xi'}$.

By **D.28**, using the above as the second premise and the IH's (2) (a) as the first premise, we have for all $\tau_i^- \in \text{CUB}_V(\tau^-, \Xi')$, we have

$$V, W, \Xi \cdot (a \leq \tau^-) \vdash \text{CLB}_V(a, \Xi') \leq \tau_i^- \dashv \Xi' (**).$$

By **D.26** and the IH's (3), we adjust the context Ξ' such that for all $\tau_i^- \in \text{CUB}_V(\tau^-, \Xi')$, we have $V, W, \Xi \vdash \text{CLB}_V(a, \Xi') \leq \tau_i^- \dashv \Xi' \cdot (a \leq \tau^-) (*)$.

Since τ^- would only be accessed once by CUB no matter its concreteness, $\text{CUB}_V(\tau^-, \Xi' \cdot (a \leq \tau^-)) = \text{CUB}_V(\tau^-, \Xi')$. Since a would only be accessed once by CLB, $\text{CLB}_V(a, \Xi' \cdot (a \leq \tau^-)) = \text{CLB}_V(a, \Xi')$. We conclude that $(a \leq \tau^-)$ is propagated with $(*)$ by the two equations above.

To propagate Δ , we conclude with the IH and $(**)$ and **D.26** that adjusts the context.

- (2) By the IH and **D.26**, we have the property for $\Xi, (a \leq \tau^-)$ and Ξ' , and we want the same for Ξ and Ξ' , $(a \leq \tau^-)$. This requires $\overline{V, W, \Xi \vdash \tau^+ \leq \tau^- \dashv \Xi' \cdot (a \leq \tau^-)}^{(\tau^+ \leq a) \in \Xi}$.

By the IH's (1), we have $\overline{V, W, \Xi \cdot (a \leq \tau^-) \vdash \tau^+ \leq \tau^- \dashv \Xi'}^{(\tau^+ \leq a) \in \Xi}$.

Then by **D.26** (using $(**)$ and the IH for the premises),

we have $\overline{V, W, \Xi \vdash \tau^+ \leq \tau^- \dashv \Xi' \cdot (a \leq \tau^-)}^{(\tau^+ \leq a) \in \Xi}$, which concludes.

- (3) The goal is $V, W, \Xi \vdash \Xi' \cdot (a \leq \tau^-) \dashv \Xi' \cdot (a \leq \tau^-)$, which requires propagating Ξ' and $(a \leq \tau^-)$ under adjusted contexts. We start from the IH's (3) and reason like in (1).

Case C-RIGID-L/C-RIGID-R WLOG we consider **C-RIGID-L**.

$$\frac{\text{C-RIGID-L} \quad a \in V \quad V, W \vdash \tau^- \rightsquigarrow (\Sigma, \sigma^-) \quad V, W \vdash \Xi \gg \Delta \cdot \Sigma \gg \Xi'}{V, W \vdash \Xi \gg \Delta \cdot (a \leq^\phi \tau^-) \gg \Xi' \cdot (a \leq \sigma^-)}$$

- (1) The goal is to show $V, W, \Xi \vdash \Delta \cdot (a \leq \tau^-) \dashv \Xi' \cdot (a \leq \sigma^-)$.

The first subgoal is to show that Δ is propagated.

By IH's (1), we have $V, W, \Xi \vdash \Delta \cdot \Sigma \dashv \Xi'$ where Σ is the output bounds of extrusion. Since σ^- is the output of extrusion, $FV(\sigma^-) \subseteq V$. Thus we can adjust the output context and obtain $V, W, \Xi \vdash \Delta \cdot \Sigma \dashv \Xi' \cdot (a \leq \sigma^-)$ since $(a \leq \sigma^-)$ creates no new type to propagate by D.27, which concludes this subgoal.

The second subgoal is to show that $V, W, \Xi \vdash a \leq \tau^- \dashv \Xi' \cdot (a \leq \sigma^-)$.

We invert the extrusion derivation. For case X-1, the goal is trivial since $\sigma^- = \tau^-$ and $FV(\tau^-) \subseteq V$ which allows us to conclude with P-RIGID-L by picking ρ as empty. For case X-2, $\sigma^- = \rho^- \tau^-$ where ρ is the polarized substitution that extrude all free types variables in $FV(\tau^-) \setminus V$. We conclude with P-RIGID-L by using ρ as the substitution and $V, W, \Xi \vdash \Sigma \dashv \Xi' \cdot (a \leq \sigma^-)$. We can reuse the X-ok premise of X-2. Note that since $V \# W$, we have $W \cap FV(\tau^-) = W \cap (FV(\tau^-) \setminus V)$. For all $b \in W \cap FV(\tau^-)$ we trivially have $V, W, \Xi \vdash (\perp \leq b) \cdot (b \leq \top) \dashv \Xi' \cdot (a \leq \sigma^-)$ for the positive and negative approximants of the skolem b .

- (2) We conclude with the IH's (2) since $a \in V$.

- (3) The goal is to show $V, W, \Xi \vdash \Xi' \cdot (a \leq \sigma^-) \dashv \Xi' \cdot (a \leq \sigma^-)$.

By IH's (3), we have $V, W, \Xi \vdash \Xi' \dashv \Xi'$. As in (1), $FV(\sigma^-) \subseteq V$ and we can adjust the output context to obtain $V, W, \Xi \vdash \Xi' \dashv \Xi' \cdot (a \leq \sigma^-)$.

By P-RIGID-L, we have $V, W, \Xi \vdash a \leq \sigma^- \dashv \Xi' \cdot (a \leq \sigma^-)$ by picking empty ρ .

Case C-FORALL-L/C-FORALL-R All properties hold by the IH. Note that in the first premise of C-FORALL-R, the cache is locked so the free part of the cache is empty, which is a subset of the empty input context, i.e. $\triangleright \phi \subseteq \epsilon$. This allows us to invoke the IH on that premise. \square

COROLLARY D.33. *If $V, W \vdash \epsilon \gg \Delta \gg \Xi'_0$, then $V, W, \epsilon \vdash \Xi' \dashv \Xi'$ where $\uparrow \text{proot}(\Xi'_0) = \Xi'$.*

PROOF. Corollary of D.32. \square

D.2.8 Full propagation implies subtyping and consistent bounds.

Notation: we write $V, W \vdash \sigma \leq \tau \dashv \Sigma$ as a shorthand for $V, W, \epsilon \vdash \sigma \leq \tau \dashv \Sigma$.

LEMMA D.34. *If $V, \bar{\alpha} \vdash \tau \leq \sigma \dashv \Sigma$ and $V, \bar{\alpha} \vdash \Sigma \dashv \Sigma$ and acyclic(Σ), then we have*

$$\text{outer}_V(\Sigma) \vdash LBU_{V \cdot \bar{\alpha}}^{\text{inner}_V(\Sigma)}(\tau \leq \sigma).$$

PROOF SKETCH. By induction on the full propagation derivation. Let $\rho = LBU_{V \cdot \bar{\alpha}}^{\text{inner}_V(\Sigma)}$.

Case P-TOP. By S-TOP.

Case P-VARREFL. By reflexivity of subtyping.

Case P-FUN. By IH and S-FUN.

Case P-RIGID-L.

$$\frac{a \in V \quad \theta = \overline{[b^- \mapsto \sigma^-, b^+ \mapsto \sigma^+ \cdot b \in FV(\tau^-) \setminus V]} \quad (a \leq \theta \tau^-) \in \Sigma}{\frac{W \vdash \tau^- \text{ X-ok} \quad FV(\theta \tau^-) \subseteq V \quad \overline{V, \bar{\alpha} \vdash (\sigma^- \leq b) \cdot (b \leq \sigma^+) \dashv \Sigma} \cdot b \in FV(\tau^-) \setminus V}{V, \bar{\alpha} \vdash a \leq \tau^- \dashv \Sigma}}$$

Since we have both $FV(\theta \tau^-) \subseteq V$ and $a \in V$, it is the case that $(a \leq \theta \tau^-) \in \text{outer}_V(\Sigma)$ and also that $\rho(a \leq \theta \tau^-) = a \leq \theta \tau^-$.

By S-HYP, we have $\text{outer}_V(\Sigma) \vdash \rho(a \leq \theta \tau^-)$. Then by the IH, we have the following.

$$\overline{\text{outer}_V(\Sigma) \vdash \rho(\sigma^- \leq b) \cdot \rho(b \leq \sigma^+)} \cdot b \in FV(\tau^-) \setminus V$$

With the polarized substitutions θ and the subtyping derivations above, we can have $outer_V(\Sigma) \vdash \rho(\theta\tau^- \leq \tau^-)$ by D.38, which allows us to conclude with S-TRANS.

Case P-RIGID-R. Symmetric to the former case.

Case P-SKIP Impossible.

Case P-FLEX.

$$\frac{a \notin V \quad (\tau \leq \sigma) \in \{(a \leq \sigma), (\tau \leq a)\} \quad (\tau \leq \sigma) \in \Sigma \quad \frac{}{V, \bar{\alpha} \vdash CLB_V(\tau, \Sigma) \leq \sigma' \dashv \Sigma}^{\sigma' \in CUB_V(\sigma, \Sigma)}}{V, \bar{\alpha} \vdash \tau \leq \sigma \dashv \Sigma}$$

- $(\tau \leq \sigma) = (a \leq \sigma)$. We want $outer_V(\Sigma) \vdash \rho(a \leq \sigma)$.

After expanding ρa , we want:

$$outer_V(\Sigma) \vdash \rho CLB_V(a, inner_V(\Sigma)) \leq \rho \sigma$$

By the IH, we have (*):

$$\frac{}{outer_V(\Sigma) \vdash \rho CLB_V(a, \Sigma) \leq \rho \sigma'}^{\sigma' \in CUB_V(\sigma, \Sigma')}$$

At this point we can conclude with the following auxilliary lemma.

LEMMA D.35. $\frac{}{outer_V(\Sigma) \vdash \rho CLB_V(a, \Sigma) \leq \rho \sigma'}^{\sigma' \in CUB_V(\sigma, \Sigma')}$, where $\rho = LBU_{V \cdot \bar{\alpha}}^{inner_V(\Sigma)}$. Then $outer_V(\Sigma) \vdash \rho(a \leq \sigma)$.

PROOF. By induction on σ . We inspect its concreteness at the same time.

- σ **concrete**⁻. Then $CUB_V(\sigma, \Sigma') = \sigma$. The lemma's premise is sufficient to conclude, since $FV(outer_V(\Sigma)) \subseteq V$ and CLB never traverses rigid type variables, i.e., $CLB_V(a, \Sigma) = CLB_V(a, inner_V(\Sigma))$.
- $\neg(\sigma$ **concrete**⁻) and $\sigma = b$ and $a \notin V$. Then $\rho \sigma = \rho CLB_V(b, inner_V(\Sigma))$. It is intuitive to see that $outer_V(\Sigma) \vdash \rho CLB_V(a, inner_V(\Sigma)) \leq \rho CLB_V(b, inner_V(\Sigma))$ when $(a \leq b) \in \Sigma'$ because a is b 's immediate lower bound, any concrete lower bound of the former is included in the union of concrete lower bounds of the latter.
- $\neg(\sigma$ **concrete**⁻) and $\sigma = \forall \alpha. \sigma_0$ and $\alpha \notin FV(\sigma_0)$. We conclude by the IH on σ and $\alpha \notin FV(\sigma')$, as $CUB_V(\sigma, \Sigma') = CUB_V(\sigma_0, \Sigma')$. \square
- $(\tau \leq \sigma) = (\tau \leq a)$. The intuition is similar to the second bullet point above: as a direct subtype of a , the concrete lower bounds of τ must be included by that of a .

Case P-FORALL-L.

$$\frac{V, \bar{\alpha} \vdash [\bar{\beta} \mapsto \pi'] \Xi \dashv \Sigma \quad V, \bar{\alpha} \vdash [\bar{\beta} \mapsto \pi'] \tau \leq \sigma \dashv \Sigma}{V, \bar{\alpha} \vdash \forall \bar{\beta} \{ \Xi \}. \tau \leq \sigma \dashv \Sigma}$$

We want $outer_V(\Sigma) \vdash \rho(\forall \bar{\beta} \{ \Xi \}. \tau \leq \sigma)$, or equivalently $outer_V(\Sigma) \vdash (\forall \bar{\beta} \{ \rho \Xi \}. \rho \tau \leq \rho \sigma)$. By the IH, we have $outer_V(\Sigma) \vdash \rho[\bar{\beta} \mapsto \pi'] \Xi$, or equivalently $outer_V(\Sigma) \vdash \rho[\bar{\beta} \mapsto \rho \pi'] \Xi$. Then by S-FORALL-L, we have $outer_V(\Sigma) \vdash \forall \bar{\beta} \{ \rho \Xi \}. \rho \tau \leq [\bar{\beta} \mapsto \rho \pi'] \rho \tau$ (**). Further, by the IH we have $outer_V(\Sigma) \vdash \rho([\bar{\beta} \mapsto \pi'] \tau \leq \sigma)$, or equivalently $outer_V(\Sigma) \vdash \rho([\bar{\beta} \mapsto \rho \pi'] \tau \leq \sigma)$. Since $\bar{\beta}$ is locally bound, we have $\rho[\bar{\beta} \mapsto \rho \pi'] \tau = [\bar{\beta} \mapsto \rho \pi'] \rho \tau$. Then we can conclude by S-TRANS.

Case P-FORALL-R.

$$\frac{\begin{array}{l} b \text{ fresh} \quad V' \supseteq FV(\tau \leq \sigma) \quad V' \# \bar{\alpha} \cdot b \\ \text{acyclic}(\Xi) \quad V', \bar{\alpha} \cdot b \vdash \tau \leq [a \mapsto b]\sigma \vdash \Xi \\ V', \bar{\alpha} \cdot b \vdash \Xi \vdash \Xi \quad V, \bar{\alpha} \vdash \text{outer}_{V'}(\Xi) \vdash \Sigma \end{array}}{V, \bar{\alpha} \vdash \tau \leq \forall a. \sigma \vdash \Sigma}$$

The goal is to show $\text{outer}_V(\Sigma) \vdash \rho\tau \leq \forall a. \rho\sigma$.

By the IH, we have

$$\text{outer}_{V'}(\Xi) \vdash \rho'\tau \leq \rho'[a \mapsto b]\sigma,$$

where $\rho' = \text{LBU}_{V', \bar{\alpha} \cdot b}^{\text{inner}_{V'}(\Xi)}$.

Since the domain of ρ' is $FV(\text{inner}_{V'}(\Xi)) \setminus (V' \cdot \bar{\alpha} \cdot b)$, which is disjoint with $FV(\tau \leq [a \mapsto \pi']\sigma)$, we then equivalently have that:

$$\text{outer}_{V'}(\Xi) \vdash \tau \leq [a \mapsto b]\sigma \quad (\text{I}).$$

By the IH, we have

$$\text{outer}_V(\Sigma) \vdash \rho(\text{outer}_{V'}(\Xi)) \quad (\text{II}).$$

Given (I) and (II), by subtyping entailment under substitutions (D.39) we have:

$$\text{outer}_V(\Sigma) \vdash \rho\tau \leq \rho[a \mapsto b]\sigma \quad (\text{III}).$$

We now show a chain of subtyping relationships and connect them by **S-TRANS** to conclude:

$$\begin{aligned} \text{outer}_V(\Sigma) \vdash \rho\tau \leq \forall b. \rho\tau & \quad (\text{by S-FORALL-R and } b \text{ is fresh}) \\ & \leq \forall b. \rho[a \mapsto b]\sigma \quad (\text{by S-FORALL-COV and III}) \\ & = \forall a. \rho\sigma \quad (\text{by } \alpha\text{-equivalency}) \quad \square \end{aligned}$$

LEMMA D.36. *If $V, \bar{\alpha} \vdash \tau \leq \sigma \vdash \Sigma$ and $\text{acyclic}(\Sigma)$, then we have*

$$\text{outer}_V(\Sigma) \cdot \text{inner}_V(\Sigma) \vdash \tau \leq \sigma.$$

PROOF. By induction on the full propagation derivation.

Case P-VARREFL By **S-VARREFL**.

Case P-FUN By the IH and **S-FUN**.

Case P-RIGID-L, P-RIGID-R By **S-HYP** & **S-TRANS** & D.38.

Case P-FLEX By **S-HYP**. Note that $(\tau \leq \sigma) \notin \text{outer}_V(\Sigma)$.

Case P-FORALL-L By the IH, **S-FORALL-L** & **S-TRANS**.

Case P-FORALL-R

$$\frac{\begin{array}{l} b \text{ fresh} \quad V' \supseteq FV(\tau \leq \sigma) \quad V' \# \bar{\alpha} \cdot b \\ \text{acyclic}(\Xi) \quad V', \bar{\alpha} \cdot b \vdash \tau \leq [a \mapsto b]\sigma \vdash \Xi \\ V', \bar{\alpha} \cdot b \vdash \Xi \vdash \Xi \quad (\text{I}) \quad V, \bar{\alpha} \vdash \text{outer}_{V'}(\Xi) \vdash \Sigma \end{array}}{V, \bar{\alpha} \vdash \tau \leq \forall a. \sigma \vdash \Sigma}$$

The goal is to show $\text{outer}_V(\Sigma) \cdot \text{inner}_V(\Sigma) \vdash \tau \leq \forall a. \sigma$. Note that Ξ may contain new flexible variables “created” while deriving (I); we define those variables as $X \equiv FV(\Xi) \setminus \{V' \cdot \bar{\alpha} \cdot b\}$. Note that $FV(\Xi) \setminus \{V' \cdot \bar{\alpha} \cdot b\} = FV(\text{inner}_{V'}(\Xi)) \setminus \{V' \cdot \bar{\alpha} \cdot b\}$.

By the IH, we have:

$$\text{outer}_{V'}(\Xi) \cdot \text{inner}_{V'}(\Xi) \vdash \tau \leq [a \mapsto b]\sigma \quad (\text{IV}).$$

By the IH, we have:

$$\text{outer}_V(\Sigma) \cdot \text{inner}_V(\Sigma) \vdash \text{outer}_{V'}(\Xi) \quad (\text{II}).$$

By [D.34](#) (reusing premise (I)), we have $outer_{V'}(\Xi) \vdash LBU_{V \cdot \bar{\alpha} \cdot b}^{inner_{V'}(\Xi)}(inner_{V'}(\Xi))$, which, if we use $LBU_{V \cdot \bar{\alpha} \cdot b}^{inner_{V'}(\Xi)}$ as the substitution, gives us

$$outer_{V'}(\Xi) \vdash \forall X\{inner_{V'}(\Xi)\} \mathbf{cons.} \quad (\text{III}).$$

We next show several subtyping relationships. By [S-FORALL-R](#) and (III), we have $outer_{V'}(\Xi) \vdash \tau \leq \forall X\{inner_{V'}(\Xi)\}$. $\forall b. \tau$ where b is a fresh type variable. Note that $FV(\tau) \subseteq V' \# X \cdot b$, as X are fresh flexible variables from (I) and b itself is fresh. Then by [S-FORALL-COV](#), we have

$$outer_{V'}(\Xi) \vdash \forall X\{inner_{V'}(\Xi)\}. \forall b. \tau \leq \forall X\{inner_{V'}(\Xi)\}. \forall b. [a \mapsto b]\sigma.$$

By [S-FORALL-L](#), we have

$$outer_{V'}(\Xi) \vdash \forall X\{inner_{V'}(\Xi)\}. \forall b. [a \mapsto b]\sigma \leq \forall b. [a \mapsto b]\sigma.$$

Note that $FV(\sigma) \subseteq V' \# X$ and the range of ω avoids X and b .

Finally, we chain the above relationships with [S-TRANS](#) and rename b to a by α -equivalency to have $outer_{V'}(\Xi) \vdash \tau \leq \forall a. \sigma$. We conclude with (II) and subtyping entailment ([D.37](#)). \square

D.2.9 Auxiliary subtyping lemmas.

LEMMA D.37 (SUBTYPING ENTAILMENT). *If $\Xi \vdash \tau \leq \sigma$ and $\Sigma \vdash \Xi$, then $\Sigma \vdash \tau \leq \sigma$.*

PROOF. By $\Sigma \vdash \Xi$, we have $\Sigma \vdash \forall \epsilon\{\Xi\} \mathbf{cons.}$ (I). By weakening, we have $\Sigma \cdot \Xi \vdash \tau \leq \sigma$ (II). We show a chain of subtyping relationships and connect them by [S-TRANS](#) to conclude:

$$\begin{aligned} \Sigma \vdash \tau \leq \forall \epsilon\{\Xi\}. \tau & & (\text{by S-FORALL-R and (I)}) \\ & \leq \forall \epsilon\{\Xi\}. \sigma & (\text{by S-FORALL-COV and (II)}) \\ & \leq \sigma & (\text{by S-FORALL-L}) \quad \square \end{aligned}$$

LEMMA D.38 (POLARIZED TYPE SUBSTITUTION). *If $\theta = \overline{[\alpha^- \mapsto \sigma^-, \alpha^+ \mapsto \sigma^+]}$ where for each $\beta \in \bar{\beta}$, β^- is substituted as \perp and β^+ is substituted as \top in θ , and $\Sigma \vdash \rho(\sigma^- \leq \alpha) \cdot \rho(\alpha \leq \sigma^+)$, then:*

- (1) *if $\bar{\beta} \vdash \tau^+ \mathbf{X-ok}$ and $\Xi \vdash \rho\tau^+ \leq \rho\tau^+$, then $\Sigma \vdash \rho(\tau^+ \leq \theta^+\tau^+)$.*
- (2) *if $\bar{\beta} \vdash \tau^- \mathbf{X-ok}$ and $\Xi \vdash \rho\tau^- \leq \rho\tau^-$, then $\Sigma \vdash \rho(\theta^-\tau^- \leq \tau^-)$.*

PROOF SKETCH. By induction on the polarized type τ^\pm for (1) (τ^- for (2)). Polarized substitutions replace a type variable at the negative position by its subtype (i.e. σ^-), and symmetrically for the positive position. $\Xi \vdash \rho\tau^\pm \leq \rho\tau^\pm$ makes sure τ^\pm itself has consistent bounds under some substitution ρ , and $\tau^\pm \mathbf{X-ok}$ guarantees that polarized substitutions do not introduce inconsistent polymorphic types. \square

LEMMA D.39 (SUBTYPING ENTAILMENT UNDER SUBSTITUTIONS). *If $\Xi \vdash \tau \leq \sigma$ and $\Sigma \vdash \rho\Xi$ and $dom(\rho) \subseteq FV(\tau \leq \sigma)$ and ρ avoids capturing and ρ is idempotent, then $\Sigma \vdash \rho(\tau \leq \sigma)$.*

PROOF. By induction on the subtyping judgment. For cases [S-FORALL-R](#), [S-FORALL-COV](#), and [S-FORALL-DISTR](#), we expand the bound consistency premise $\Xi \vdash \forall V\{\Sigma'\} \mathbf{cons.}$ to subtyping judgments $\Xi \vdash \rho'\Sigma'$ where $dom(\rho') = V$ according to the definition of consistency. We invoke the IH on those subtyping judgments and have $\Sigma \vdash \rho\rho'\Sigma'$. To prove the goal, we need $\Sigma \vdash \forall V\{\rho\Sigma'\} \mathbf{cons.}$ We use $\rho\rho'$ as the substitutions of V , and $\Sigma \vdash \rho\rho'\Sigma'$ is implied by $\Sigma \vdash \rho\rho'\Sigma'$ because ρ is capture-avoiding and idempotent. \square

D.2.10 Soundness of constraining.

THEOREM D.40 (SOUNDNESS OF CONSTRAINING). *Let $V, \bar{\alpha} \vdash \epsilon \gg \Delta \gg \Sigma'$ and $\Sigma = \text{uproot}(\Sigma')$. Then for all $\bar{\tau}$ where $FV(\bar{\tau}) \# \bar{\beta}$ and $\bar{\beta} = FV(\Sigma) \setminus \{V \cdot \bar{\alpha}\}$ we have*

- $\text{outer}_V(\Sigma) \cdot \text{inner}_V([\bar{\alpha} \mapsto \bar{\tau}]\Sigma) \vdash [\bar{\alpha} \mapsto \bar{\tau}]\Delta$ and
- $\text{outer}_V(\Sigma) \vdash \forall \bar{\beta} \{ \text{inner}_V([\bar{\alpha} \mapsto \bar{\tau}]\Sigma) \}$ **cons.**

PROOF. By Lemmas D.32 and D.36, we have $\text{outer}_V(\Sigma) \cdot \text{inner}_V(\Sigma) \vdash \Delta$. By Lemma D.39, we have $\text{outer}_V(\Sigma) \cdot \text{inner}_V([\bar{\alpha} \mapsto \bar{\tau}]\Sigma) \vdash [\bar{\alpha} \mapsto \bar{\tau}]\Delta$. Notice that $[\bar{\alpha} \mapsto \bar{\tau}]\text{outer}_V(\Sigma) = \text{outer}_V(\Sigma)$. By Lemmas D.32 and D.34, we have $\text{outer}_V(\Sigma) \vdash \text{LBU}_{V \cdot \bar{\alpha}}^{\text{inner}_V(\Sigma)} \text{inner}_V(\Sigma)$. By Lemma D.39, we have $\text{outer}_V(\Sigma) \vdash [\bar{\alpha} \mapsto \bar{\tau}]\text{LBU}_{V \cdot \bar{\alpha}}^{\text{inner}_V(\Sigma)} \text{inner}_V(\Sigma)$. Given the fact that Σ contains no bounds on skolems and skolems in the bounds of rigid type variables are extruded, we equivalently have:

$$\text{outer}_V(\Sigma) \vdash [\bar{\alpha} \mapsto \bar{\tau}]\text{LBU}_{V \cdot \bar{\alpha}}^{\text{inner}_V(\Sigma)} \text{inner}_V([\bar{\alpha} \mapsto \bar{\tau}]\Sigma).$$

This allows us to conclude that $\text{outer}_V(\Sigma) \vdash \forall \bar{\beta} \{ \text{inner}_V([\bar{\alpha} \mapsto \bar{\tau}]\Sigma) \}$ **cons.** □

Received 2023-07-11; accepted 2023-11-07