

Local Type Inference for Polarised System F with Existentials

ANONYMOUS AUTHOR(S)

Change! Type inference for polymorphic types has been extensively studied, with the Damas-Hindley-Milner algorithm being the most well-known approach for handling parametric polymorphism. However, supporting impredicative existential types has remained an open challenge due to undecidability barriers. This paper presents F_{\exists}^{\pm} , a type system extending the Call-By-Push-Value calculus with first-class existential and universal quantifiers that can be instantiated impredicatively. We give a local type inference algorithm for F_{\exists}^{\pm} that avoids the undecidability issues through a novel combination of unification and anti-unification techniques. A key contribution is providing a declarative specification of the inference algorithm based on the F_{\exists}^{\pm} type system and proving the algorithm's soundness and completeness. We also explore extending the inference framework to bidirectional typechecking, different subtyping regimes, and bounded polymorphism. This work opens up new avenues for bringing the benefits of impredicative existential types to programming languages through principled type inference.

CCS Concepts: • **Do Not Use This Code** → **Generate the Correct Terms for Your Paper**; *Generate the Correct Terms for Your Paper*; Generate the Correct Terms for Your Paper; Generate the Correct Terms for Your Paper.

Additional Key Words and Phrases: Type Inference, System F, Call-by-Push-Value, Polarized Typing, Focalisation, Subtyping

ACM Reference Format:

Anonymous Author(s). 2018. Local Type Inference for Polarised System F with Existentials. *J. ACM* 37, 4, Article 111 (August 2018), 43 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Over the last half-century, there has been considerable work on developing type inference algorithms for programming languages, mostly aimed at solving the problem of *type polymorphism*.

That is, in pure polymorphic lambda calculus—System F [Girard 1971; Reynolds 1974]—the polymorphic type $\forall\alpha. A$ has a big lambda $\Lambda\alpha. e$ as an introduction form, and an explicit type application $e \{A\}$ as an elimination form. This is an extremely simple and compact type system, whose rules fit on a single page, but whose semantics are advanced enough to accurately represent concepts such as parametricity and representation independence. However, System F by itself is unwieldy as a programming language. The fact that the universal type $\forall\alpha. A$ has explicit eliminations means that programs written using polymorphic types will need to be stuffed to the gills with type annotations explaining how and when to instantiate the quantifiers.

Therefore, most work on type inference has been aimed at handling type instantiations implicitly—we want to be able to use a polymorphic function like $\text{len} : \forall\alpha. [\alpha] \rightarrow \text{Int}$ at any concrete list type without explicitly instantiating the quantifier in len 's type. That is, we want to write $\text{len} [1, 2, 3]$ instead of writing $\text{len} \{\text{Int}\} [1, 2, 3]$.

The way this is typically done is by using a *subtyping* relation, induced by the polymorphic instantiation: the type $\forall\alpha. A$ is a *subtype* of all of its instantiations. So we wish to be free to use

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Association for Computing Machinery.

0004-5411/2018/8-ART111 \$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

the same polymorphic function `len` at many different types such as $\text{[Int]} \rightarrow \text{Int}$, $\text{[Bool]} \rightarrow \text{Int}$, $\text{[Int]} \times \text{Bool} \rightarrow \text{Int}$, and so on. However, the subtyping can be used nondeterministically: whenever we see a polymorphic type $\forall \alpha. A$, we know it is a subtype of *any* of its instantiations. To turn this into an algorithm, we have to actually make some choices.

The most famous of the algorithms for handling this is the Damas-Hindley-Milner (DHM) algorithm [Damas et al. 1982; Hindley 1969; Milner 1978]. DHM deals with this choice using *unification*. Whenever we would have had to introduce a particular concrete type in the specification, the DHM algorithm introduces a *unification variable*, and incrementally instantiates this variable as more and more type constraints are observed.

However, the universal quantifier is not the only quantifier! Dual to the universal quantifier \forall is the existential quantifier \exists . Even though existential types have an equal logical status to universal types, they have been studied much less frequently in the literature. The most widely-used algorithm for handling existentials is the algorithm of Odersky and Laufer [Läufer et al. 1994]. This algorithm treats existentials in a second-class fashion: they are not explicit connectives of the logic but rather are tied to datatype declarations. As a result, packing and unpacking existential types is tied to constructor introductions and constructor eliminations in pattern matches. This allows Damas-Milner inference to continue to work almost unchanged, but it does come at the cost of losing first-class existentials and also of restricting the witness types to monomorphic types.

There has been a limited amount of work on support for existential types as a first-class concept. In an unpublished draft, Leijen [2006] studied an extension of Damas-Milner inference in which type schemes contain alternating universal and existential quantifiers. Quantifiers still range over monotypes (i.e., types with no quantifiers inside of them), and higher-rank polymorphism is not permitted: one cannot instantiate a quantifier with a type that itself contains quantifiers. More recently, Dunfield et al. [2016] studied type inference for existential types in the context of GADT type inference, which, while still predicative, supported higher-rank types (i.e., quantifiers can occur anywhere inside of a type scheme). Eisenberg et al. [2021] propose a system which only permits types of the form `forall-exists`, but which permits projective elimination in the style of ML modules.

All of these papers are restricted to *predicative* quantification, where quantifiers can only be instantiated with monotypes. However, existential types in full System F are *impredicative*—that is, quantifiers can be instantiated with arbitrary types, specifically including types containing quantifiers.

Historically, inference for impredicative quantification has been neglected, due to results by Tiuryn et al. [1996] and Chrzaszcz [1998] which show that not only is full type inference for System F undecidable, but even that the subtyping relation induced by type instantiation is undecidable. However, in recent years, interest in impredicative inference has revived (for example, the work of Serrano et al. [2020]), with a focus on avoiding the undecidability barriers by doing *partial* type inference. That is, we no longer try to infer the type of any term in the language, but rather accept that the programmer will need to write annotations in cases where inference would be too difficult. Unfortunately, it is often difficult to give a specification for when the partial algorithm succeeds—for example, Eisenberg et al. [2021] observe that their algorithm lacks a declarative specification, and explain why existential types make this particularly difficult to do.

One especially well-behaved form of partial type inference is *local type inference*, introduced by Pierce et al. [2000]. In this approach, quantifiers in a polymorphic function are instantiated using only the type information available in the arguments at each call site. While this infers fewer quantifiers than global algorithms such as Damas-Milner can, it has the twin benefits that it is easy to give a mathematical specification to, and that failures of type inference are easily explained in terms of local features of the call site. In fact, many production programming languages such as C#

and Java use a form of local type inference, due to the ease of extending this style of inference to support object-oriented features.

In this paper, we extend the local type inference algorithm to a language with both universal and existential quantifiers, which can both be instantiated impredicatively. This combination of features broke a number of the invariants which traditional type inference algorithms depend on, and required us to invent new algorithms which combine both unification and (surprisingly) the dual—anti-unification.

In this paper, we extend the local type inference algorithm to a language with both universal and existential quantifiers, which can both be instantiated impredicatively. This extended system is referred to as F_{\exists}^{\pm} . The combination of features in F_{\exists}^{\pm} broke a number of the invariants which traditional type inference algorithms depend on and required us to invent new algorithms which combine both unification and (surprisingly) anti-unification.

To research the type inference in impredicative polymorphic systems in the presence of existential types, we make the following **contributions**.

Type Inference Algorithm (section 4). We give a local type inference algorithm which supports both first-class existential and universal quantifiers, both of which can be instantiated impredicatively. To evade the undecidability results surrounding type inference for System F, we introduce F_{\exists}^{\pm} —a variant of Call-By-Push-Value [Levy 2006], which lets us formulate a subtyping relation which is still decidable.

Declarative Specification (sections 3 and 5). We give a declarative type system (again, based on Call-By-Push-Value) to serve as a specification of our algorithm, and we prove our algorithm is sound and complete with respect to the declarative type system. The specification makes it easy to see that all type applications (for \forall -elimination) and all packs (for \exists -introduction) are inferred.

Employment of the Anti-unification (section 4.7). The introduction of impredicative existentials breaks some of the fundamental invariants of HM-style type inference. As a result, it needs to mix unification with *anti-unification*—a dual operation which has been used in the literature on term rewriting systems, but which has never appeared in the context of polymorphic type inference.

Extrapolation of the Inference Framework to Other Type Systems (section 6). We explore the design space further to show how the same type inference scheme can be applied to work with different type systems. In particular, (i) in combination with *bidirectional* typechecking—the approach used in the original local type inference paper of Pierce et al. [2000] that minimizes the number of needed annotations; (ii) in systems with different subtyping rules, such as the one used in Zhao et al. [2022]—a system that permits explicit type applications; (iii) in systems with *bounded* polymorphic quantifiers. That is, one can only instantiate a quantifier with a type that is a subtype of a given bound.

Ilya: Is the following needed since we have the contributions? I just wanted to introduce the appendix somewhere...

Structure of this Paper. In section 2, we discuss the language of the F_{\exists}^{\pm} system, providing examples, outlining limitations, and summarizing the main ideas behind the algorithm. Section 3 presents a formal definition of the declarative specification and its key properties. In section 4, the type inference algorithm is introduced, detailing a hierarchy of subroutines; upon perusing this section, one could *implement* the type inference algorithm. Section 5 overviews the algorithm's correctness proof, explaining its structure and how the declarative system's features make it decidable. Section 6 covers potential F_{\exists}^{\pm} extensions and future work directions. Finally, ?? contains

comprehensive formal definitions and proofs of lemmas and theorems used in the paper; in instances where the full definitions might not be of immediate interest, we direct the reader to the appendix.

2 OVERVIEW

In the presence of impredicative polymorphism, the standard System F has undecidable subtyping and type inference. It means that any inference algorithm is incomplete: it only infers the types for a *fragment* of System F.

To tackle this issue, we restrict System F by employing the ideas from the Call-By-Push-Value system [Levy 2006]. We introduce a new language, F_{\exists}^{\pm} , which stratifies the syntax of System F into positive and negative parts. It provides the structure that allows us to restrict the typing specification to make it decidable.

Every term and type of F_{\exists}^{\pm} has either *positive* or *negative* polarity. The type variables are annotated with their polarities (e.g. α^+ or β^-), and the types can change polarity via the shift operators \uparrow (positive to negative) and \downarrow (negative to positive).

Positive expressions correspond to *values* in the Call-By-Push-Value system, and negative expressions correspond to *computations*. The difference between the computations and values is intuitively described in the motto of Call-By-Push-Value: “a value is, a computation does”. In particular, an argument of a function is always a value (i.e., positive), and the resulting type, as well as the function itself, is a computation (i.e., negative). Analogously, the polymorphic \forall quantifies a computation over positive types.

2.1 Examples

Generally, any term and type of System F can be embedded into F_{\exists}^{\pm} in multiple ways. This is because at any given point, one can opt for either the positive or negative polarity of the subterm. We will delve into the consistent term-level and type-level embeddings (i.e., preserving typing) in more detail in section 3.4. Figure 1 showcases the polarized versions of the standard System F terms. Next, we will discuss several examples (using the terms from fig. 1) to illustrate what can and cannot be inferred in F_{\exists}^{\pm} .

$$\begin{aligned} \text{map} &: \downarrow(\forall \alpha^+. \forall \beta^+. \downarrow(\alpha^+ \rightarrow \uparrow \beta^+) \rightarrow [\alpha^+] \rightarrow \uparrow[\beta^+]) \\ \text{len} &: \downarrow(\forall \alpha^+. [\alpha^+] \rightarrow \uparrow \text{Int}) \\ \text{choose} &: \downarrow(\forall \alpha^+. \alpha^+ \rightarrow \alpha^+ \rightarrow \uparrow \alpha^+) \\ \text{id} &: \downarrow(\forall \alpha^+. \alpha^+ \rightarrow \uparrow \alpha^+) \\ \text{auto} &: \downarrow(\downarrow(\forall \alpha^+. \alpha^+ \rightarrow \uparrow \alpha^+) \rightarrow (\forall \alpha^+. \alpha^+ \rightarrow \uparrow \alpha^+)) \end{aligned}$$

Fig. 1. Polarization of standard System F terms

Application of a Polymorphic Function. A polymorphic function can be applied to a value of a concrete type. In this case, the quantified variables are instantiated to make the application type-check. For instance, the polymorphic `len` can be applied to a list of any positive type, and the resulting type will be inferred by the algorithm:

$$\cdot; \Gamma \vdash \text{let } x = \text{len}([1, 2, 3]); \text{return } x: \uparrow \text{Int}$$

Polymorphic Self-application. Unlike predicative systems, F_{\exists}^{\pm} permits self-application of polymorphic functions. For instance, `id` can be applied to itself resulting in a suspended polymorphic

Negative declarative types

$$N, M, K ::= \alpha^- \mid \uparrow P \mid P \rightarrow N \mid \forall \alpha^+. N$$

Positive declarative types

$$P, Q, R ::= \alpha^+ \mid \downarrow N \mid \exists \alpha^-. P$$
Fig. 2. Declarative Types of F_{\exists}^{\pm}

computation of type $\downarrow(\forall \alpha^+. \alpha^+ \rightarrow \uparrow \alpha^+)$, which, in turn, can be applied to itself:

$$;\Gamma \vdash \text{let } x = \text{id}(\text{id}); \text{let } y = x(x); \text{return } y: \uparrow \downarrow(\forall \alpha^+. \alpha^+ \rightarrow \uparrow \alpha^+)$$

Non-example: Polymorphic Arguments. The focus of our system is the *local* type inference [Pierce et al. 2000], which has certain limitations. In particular, the instantiation does not happen when the polymorphic term is in the argument position. For example, one could expect the following to be inferable: $;\Gamma \vdash \text{let } x = \text{map}(\text{id}, [2, 3, 9]); \text{return } x: \uparrow[\text{Int}]$. However, to infer this, id must be instantiated to $\text{Int} \rightarrow \uparrow \text{Int}$ which requires the information to be propagated from the neighboring branch of the syntax tree $([2, 3, 9])$, i.e., bypass the locality.

The annotated version of the same term infers the type successfully, since all the polymorphic variables of the arguments are instantiated:

$$;\Gamma \vdash \text{let } x = \text{map}((\text{id} : \downarrow(\text{Int} \rightarrow \uparrow \text{Int})), [2, 3, 9]); \text{return } x: \uparrow[\text{Int}]$$

Inferring the Common Supertype. The function `choose` takes two arguments of the same polymorphic type. To apply it to two values of different (concrete) types, the inference algorithm must find the minimal type they can be coerced to—the least common supertype.

Let us assume that in Γ , we have two identity-like functions $\text{id}_M : \downarrow(\downarrow M \rightarrow M)$ and $\text{id}_N : \downarrow(\downarrow N \rightarrow N)$. Then their least common supertype—the existential $\exists y^-. \downarrow(\downarrow y^- \rightarrow y^-)$ is inferable:

$$;\Gamma \vdash \text{let } x = \text{choose}(\text{id}_N, \text{id}_M); \text{return } x: \uparrow \exists y^-. \downarrow(\downarrow y^- \rightarrow y^-)$$

Inferring a general type. As mentioned before, the local inference does not instantiate the polymorphic types in the argument position. Because of that, the application $\text{let } x = \text{choose}(\text{id}, \text{auto}); \text{return } x$ cannot infer $\uparrow \downarrow(\downarrow(\forall \alpha^+. \alpha^+ \rightarrow \uparrow \alpha^+) \rightarrow (\forall \alpha^+. \alpha^+ \rightarrow \uparrow \alpha^+))$. However, the inference succeeds and infers a less informative type:

$$;\Gamma \vdash \text{let } x = \text{choose}(\text{id}, \text{auto}); \text{return } x: \uparrow \exists y^-. \downarrow y^-$$

2.2 The Language of Types

The types of F_{\exists}^{\pm} are given in fig. 2. They are stratified into two syntactic categories (polarities): positive and negative, similar to the values and computations in the Call-By-Push-Value system.

- ± α^- and α^+ denote negative and positive type variables, respectively. The variables are either introduced by the corresponding quantifiers (\forall and \exists , respectively) or declared in the type context;
- a function $P \rightarrow N$ takes a *positive* input P and returns a *negative* output N ; the function type itself is *negative*, so the ‘arrow’ operator is right-associative: $P \rightarrow Q \rightarrow N$ is equivalent to $P \rightarrow (Q \rightarrow N)$;
- similarly to functions, a polymorphic abstraction $\forall \alpha^+. N$ quantifies a negative type N over a list of positive type variables α^+ ;
- + dually, $\exists \alpha^-. P$ binds negative variables in a positive type P ;
- ± an upshift $\uparrow P$ and a downshift $\downarrow N$ change the polarity of a type from positive to negative and vice versa. At the level of terms, the constructors for shift types are `return v` and `{c}` (thunked computation), respectively.

Computation Terms

$$c, d ::= (c : N) \mid \lambda x : P. c \mid \Lambda \alpha^+. c \mid \text{return } v \mid \text{let } x = v; c \mid \text{let } x : P = c; c' \mid \\ \text{let } x : P = v(\vec{v}); c \mid \text{let } x = v(\vec{v}); c \mid \text{let}^\exists(\vec{\alpha}^+, x) = v; c$$

Value Terms

$$v, w ::= x \mid \{c\} \mid (v : P)$$

Fig. 3. Declarative Terms of F_{\exists}^{\pm}

Syntactic Conventions. We always consider terms and types up to alpha-equivalence, and capture-avoiding substitution. Besides, we assume the quantification is associative, for example, $\forall \alpha^+. \forall \beta^+. N$ and $\forall \alpha^+, \beta^+. N$ represent the same type.

Type Context and Type Well-formedness. In the construction of F_{\exists}^{\pm} , a type context (denoted as Θ) is represented as a set of positive and negative type variables and it is used to assert the well-formedness of types. The well-formedness (or well-scopeness) of a type is denoted as $\Theta \vdash P$ and $\Theta \vdash N$ and it asserts that all type variables are either bound by a quantifier (\forall and \exists) or declared in the context Θ . The well-formedness checking is an *algorithmic* procedure. We represent it as a system of inference rules, that correspond to a recursive algorithm taking the context and the type as input. For brevity, we omit these rules, as they are standard binder scoping rules.

2.3 The Language of Terms

In fig. 3, we define the language of terms of F_{\exists}^{\pm} . The language is a combination of System F and Call-By-Push-Value constructs: the terms are stratified into values and computations; there are return and thunk constructors, as well as several forms of let bindings to sequentialize the computation/application with or without type annotation;

- + x denotes a term variable. Following the Call-By-Push-Value stratification, we only have *positive* (value) term variables;
- + $\{c\}$ is a value corresponding to a thunked or suspended computation;
- + $(c : N)$ and $(v : P)$ allow one to annotate positive and negative terms;
- $\text{return } v$ is a pure computation, returning a value;
- $\lambda x : P. c$ and $\Lambda \alpha^+. c$ are standard lambda abstractions. Notice that we require the type annotation for the argument of λ ;
- $\text{let } x = v; c$ is a standard let, binding a value v to a variable x in a computation c ;
- computational let form $\text{let } x : P = c; c'$ operates similarly to a monadic bind. It binds the result of a computation $c : \uparrow P$ to a variable $x : P$, and continues with a computation c' ;
- applicative let forms $\text{let } x : P = v(\vec{v}); c$ and $\text{let } x = v(\vec{v}); c$ can be viewed as a special case of the computational let, when the first computation is a function application $v(\vec{v})$. They take a suspended computation v , apply it to a list of arguments, bind the (pure) result to a variable x , and continues with a computation c . If the resulting type of the application is unique, one can omit the type annotation, as in the second form: it will be inferred by the algorithm;
- $\text{let}^\exists(\vec{\alpha}^+, x) = v; c$ is the standard eliminator of an existential type (unpack): expecting v to be of an existential type, it binds the packed negative types to a list of variables $\vec{\alpha}^+$, binds the body of the existential to x , and continues with a computation c .

Missing constructors. Notice that the language does not have first-class applications: their role is played by the applicative let forms, binding the result of a *fully applied* function to a variable. Also notice that the language does not have a type application (i.e. the eliminator of \forall) and dually, it

does not have pack (i.e. the constructor of \exists). This is because the instantiation of polymorphic and existential types is inferred by the algorithm. In section 6, we discuss the way to modify the system to introduce *explicit* type applications.

2.4 The Key Ideas of the Algorithm

The inference algorithm for F_{\exists}^{\pm} is *local*: it has a limited scope of the inference information and does not propagate it between far-apart branches of the syntax tree. Nevertheless, many difficulties appear in the inference of the combination of polymorphic and existential types. Next, we discuss the most challenging of these difficulties and how they are addressed in the algorithm.

Subtyping. The inference algorithm's complexity mainly lies in the subtyping relation, which is, in particular, responsible for polymorphic instantiation. It's this subtyping relation that enables us to apply polymorphic terms in situations where their specific instantiations are anticipated.

To check whether one type is a subtype of another, the algorithm introduces *algorithmic* type variables $(\widehat{\alpha}^{\pm}, \widehat{\beta}^{\pm}, \dots)$ —the ‘placeholders’ representing the polymorphic variables whose instantiation is postponed. This way, the problem of $\forall \alpha^+. \uparrow \alpha^+ \leq \uparrow \text{Int}$ becomes the problem of finding $\widehat{\alpha}^+$ such that $\uparrow \widehat{\alpha}^+ \leq \uparrow \text{Int}$.

Equivalence. The system introduces another complexity: the equivalence induced by subtyping $(\Theta \vdash N \simeq^{\leq} M \stackrel{\text{def}}{\iff} \Theta \vdash N \leq M \text{ and } \Theta \vdash M \leq N)$ is not straightforward. This relation extends beyond just alpha-equivalent types. As the neighboring polymorphic quantifiers can be instantiated in arbitrary order, the equivalence relation must permit the *permutations* of the quantifiers. For example, we have $\cdot \vdash \forall \alpha^+. \forall \beta^+. N \simeq^{\leq} \forall \beta^+. \forall \alpha^+. N$. Similarly, the equivalence permits the *removal/addition* of the unused quantifiers: $\cdot \vdash \forall \alpha^+. \uparrow \text{Int} \simeq^{\leq} \uparrow \text{Int}$.

To handle the complexity of non-trivial equivalence, we employ a type normalization algorithm. This approach simplifies mutual subtyping to a more manageable alpha-equivalence by eliminating unused quantifiers and reordering the rest into a canonical permutation, effectively replacing subtyping-induced equivalence on terms with alpha-equivalence on their normal forms.

Least Upper Bound. Eventually, the subtyping algorithm reduces the initial subtyping problem to a set of constraints that need to be resolved. At this moment, another problem arises: as one variable can occur in multiple constraints, the resolution becomes non-trivial. In particular, to resolve the conjunction of $\widehat{\alpha}^+ \geq P$ and $\widehat{\alpha}^+ \geq Q$, we would need to find the least upper bound of P and Q .

In many type systems, the least upper bound boils down to the trivial syntactic equality. Because of that, the set of constraints is typically resolved using *unification*. However, this problem is trickier in the presence of existential types. For example, in F_{\exists}^{\pm} , the least upper bound of $\downarrow \uparrow \text{Int}$ and $\downarrow \uparrow \text{Bool}$ is $\exists \alpha^-. \downarrow \alpha^-$, because the quantified α^- can be instantiated to either $\uparrow \text{Int}$ or $\uparrow \text{Bool}$.

To find this least upper bound, we use *anti-unification*—the dual of unification. While unification finds *the most general instance* of two given patterns, the anti-unification finds *the most detailed* (under the restrictions of the system) *pattern* that matches two given types. In the example above, the anti-unifier of $\downarrow \uparrow \text{Int}$ and $\downarrow \uparrow \text{Bool}$ is $\downarrow \widehat{\alpha}^-$, and in F_{\exists}^{\pm} , this unifier is the most specific (in particular because $\downarrow \widehat{\alpha}^+$ is not considered as only *negative* variables are existentially quantified).

Impredicativity. Another difficulty that hinders the constraint resolution arises from the impredicativity of the system—the ability to quantify over the types that are themselves polymorphic. In the impredicative polymorphic system (even with only \forall -quantifiers), the naturally formulated subtyping is undecidable, which follows from the undecidability of second-order unification. All the more the undecidability is evident in the presence of existential types and subtyping constraints.

<div style="display: flex; justify-content: space-between; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px;">$\Theta \vdash N \leq M$</div> <div>Negative subtyping</div> </div> $\frac{}{\Theta \vdash \alpha^- \leq \alpha^-} \text{ (VAR}_{\leq}^{\text{--}})$ $\frac{\Theta \vdash P \simeq^{\leq} Q}{\Theta \vdash \uparrow P \leq \uparrow Q} \text{ (}\uparrow^{\leq}\text{)}$ $\frac{\Theta \vdash P \geq Q \quad \Theta \vdash N \leq M}{\Theta \vdash P \rightarrow N \leq Q \rightarrow M} \text{ (}\rightarrow^{\leq}\text{)}$ $\frac{\Theta, \vec{\beta}^+ \vdash \sigma : \vec{\alpha}^+ \quad \Theta, \vec{\beta}^+ \vdash [\sigma]N \leq M}{\Theta \vdash \forall \vec{\alpha}^+. N \leq \forall \vec{\beta}^+. M} \text{ (}\forall^{\leq}\text{)}$ <div style="display: flex; justify-content: space-between; align-items: center; padding-top: 10px;"> <div style="border: 1px solid black; padding: 2px 5px;">$\Theta \vdash N \simeq^{\leq} M$</div> <div>Negative equivalence</div> </div> $\frac{\Theta \vdash N \leq M \quad \Theta \vdash M \leq N}{\Theta \vdash N \simeq^{\leq} M} \text{ (}\simeq_{\leq}^{\text{--}}\text{)}$	<div style="display: flex; justify-content: space-between; align-items: center;"> <div style="border: 1px solid black; padding: 2px 5px;">$\Theta \vdash P \geq Q$</div> <div>Positive supertyping</div> </div> $\frac{}{\Theta \vdash \alpha^+ \geq \alpha^+} \text{ (VAR}_{\geq}^{\text{+}})$ $\frac{\Theta \vdash N \simeq^{\leq} M}{\Theta \vdash \downarrow N \geq \downarrow M} \text{ (}\downarrow^{\geq}\text{)}$ $\frac{\Theta, \vec{\beta}^- \vdash \sigma : \vec{\alpha}^- \quad \Theta, \vec{\beta}^- \vdash [\sigma]P \geq Q}{\Theta \vdash \exists \vec{\alpha}^-. P \geq \exists \vec{\beta}^-. Q} \text{ (}\exists^{\geq}\text{)}$ <div style="display: flex; justify-content: space-between; align-items: center; padding-top: 10px;"> <div style="border: 1px solid black; padding: 2px 5px;">$\Theta \vdash P \simeq^{\leq} Q$</div> <div>Positive equivalence</div> </div> $\frac{\Theta \vdash P \geq Q \quad \Theta \vdash Q \geq P}{\Theta \vdash P \simeq^{\leq} Q} \text{ (}\simeq_{\leq}^{\text{+}}\text{)}$
---	--

Fig. 4. Declarative Subtyping

For instance, in F_{\exists}^{\pm} the constraint $\widehat{\alpha}^+ : \geq \downarrow \uparrow \widehat{\alpha}^+$ would have a surprising impredicative solution $\widehat{\alpha}^+ = \exists \alpha^-. \downarrow \alpha^-$.

To maintain the decidability in the presence of impredicativity, we carefully restrict the system using the polarity stratification. One of the main constraints that we put on subtyping is the invariance of the shift operators: the subtyping $\uparrow P \leq \uparrow Q$ is only allowed if $Q \geq P$ and $P \geq Q$. These constraints do not trivialize the system: we still have non-trivial upper bounds, and thus, need to use anti-unification to find them. However, it prevents such examples as $\widehat{\alpha}^+ : \geq \downarrow \uparrow \widehat{\alpha}^+$, and as we prove, makes the constraint resolution decidable.

3 DECLARATIVE SYSTEM

In this section, we present the declarative system of F_{\exists}^{\pm} , which serves as a specification of the type inference algorithm. The declarative system is represented as a set of inference rules and consists of two main subsystems: subtyping and type inference. First, we present the declarative subtyping rules specifying when one type is a subtype of another; next, we discuss the equivalence relation induced by mutual subtyping; finally, we present the type inference rules, that refer to the subtyping and equivalence relation. We conclude this section by discussing the relation between the proposed type system and the standard System F.

3.1 Subtyping

The inference rules representing declarative subtyping are shown in fig. 4. Let us discuss them in more detail.

Quantifiers. Symmetric rules (\forall^{\leq}) and (\exists^{\geq}) specify subtyping between top-level quantified types. Usually, the polymorphic subtyping is represented by two rules introducing quantifiers to the left and to the right-hand side of subtyping. For conciseness of representation, we compose these rules into one. First, our rule extends context Θ with the quantified variables from the right-hand side ($\vec{\beta}^+$ or $\vec{\beta}^-$), as these variables must remain abstract. Second, it verifies that the left-hand side

quantifiers ($\vec{\alpha}^+$ or $\vec{\alpha}^-$) can be instantiated to continue subtyping recursively. This is the point of non-determinism as the instantiation can be chosen arbitrarily

The instantiation of quantifiers is modeled by substitution σ . The notation $\Theta_2 \vdash \sigma : \Theta_1$ specifies its domain and range. For instance, $\Theta, \vec{\beta}^+ \vdash \sigma : \vec{\alpha}^+$ means that σ maps the variables from $\vec{\alpha}^+$ to (positive) types well-formed in $\Theta, \vec{\beta}^+$. This way, application $[\sigma]N$ instantiates (replaces) every α_i^- in N with $\sigma(\alpha_i^-)$.

Invariant Shifts. As mentioned above, an important restriction that we put on the subtyping system is that subtyping of shifted types requires their equivalence, as shown in (\Downarrow^\geq) and (\Uparrow^\leq) . The reason for this is that if both of these rules were relaxed to the covariant form, the subtyping relation would become equivalent to the standard subtyping of System F, which is undecidable [Tiuryn et al. 1996]. However, after certain changes (\Uparrow^\leq) can be relaxed to the covariant form, thereby increasing the expressiveness of the system. These changes are discussed in section 6.2.

Functions. Standardly, subtyping of function types is covariant in the return type and contravariant in the argument type.

Variables. Subtyping of variables is defined reflexively, which is enough to ensure the reflexivity of subtyping in general. The algorithm—specifically the least upper bound procedure—will use the fact that the subtypes of a variable coincide with its supertypes (property 2), which however does not hold for arbitrary types.

3.1.1 Properties of Declarative Subtyping. A property that is important for the subtyping algorithm, in particular for the type *upgrade* procedure (section 4.6), is the preservation of free variables by subtyping. Informally, it says that the free variables of a positive type cannot disappear in its subtypes, and the free variables of a negative type cannot disappear in its supertypes.

Property 1 (Subtyping preserves free variables). *Let us assume that all the mentioned types are well-formed in Θ . Then $\Theta \vdash N_1 \leq N_2$ implies $\text{fv}(N_1) \subseteq \text{fv}(N_2)$, and $\Theta \vdash P_1 \geq P_2$ implies $\text{fv}(P_1) \subseteq \text{fv}(P_2)$.*

Property 2 (Variable subtyping is trivial). *A subtype or a supertype of a variable is equivalent to the variable itself:*

$$\begin{array}{llll} - & \Theta \vdash \alpha^- \leq N & \iff & \Theta \vdash N \leq \alpha^- & \iff & N = \forall \vec{\beta}^+. \alpha^- \\ + & \Theta \vdash \alpha^+ \geq P & \iff & \Theta \vdash P \geq \alpha^+ & \iff & P = \exists \vec{\beta}^-. \alpha^+ \end{array}$$

Another property that we extensively use is that subtyping is reflexive and transitive, and agrees with substitution.

Property 3 (Subtyping forms a preorder). *For a fixed context Θ , the negative subtyping relation $\Theta \vdash N_1 \leq N_2$ and the positive subtyping relation $\Theta \vdash P_1 \geq P_2$ are reflexive and transitive on types well-formed in Θ .*

Property 4 (Subtyping agrees with substitution). *Suppose that σ is a substitution such that $\Theta_2 \vdash \sigma : \Theta_1$. Then*

$$\begin{array}{l} - \quad \Theta_1 \vdash N \leq M \text{ implies } \Theta_2 \vdash [\sigma]N \leq [\sigma]M, \text{ and} \\ + \quad \Theta_1 \vdash P \geq Q \text{ implies } \Theta_2 \vdash [\sigma]P \geq [\sigma]Q. \end{array}$$

Moreover, any two *positive* types have a least upper bound, which makes positive subtyping a semilattice. The positive least upper bound can be found algorithmically, which we will discuss in the next section.

Property 5 (Positive least upper bound exists). Suppose that P_1 and P_2 are positive types well-formed in Θ . Then there exists a least common supertype—a type P such that

- $\Theta \vdash P \geq P_1$ and $\Theta \vdash P \geq P_2$, and
- for any Q such that $\Theta \vdash Q \geq P_1$ and $\Theta \vdash Q \geq P_2$, $\Theta \vdash Q \geq P$.

Negative greatest lower bound might not exist. The symmetric construction—the greatest lower bound of two negative types—does not always exist, as the following counterexample shows. Consider the following types:

- N and Q are arbitrary closed types,
- P, P_1 , and P_2 are non-equivalent closed types such that $\cdot \vdash P_1 \geq P$ and $\cdot \vdash P_2 \geq P$, and none of the types is equivalent to Q .

What is the greatest common subtype of $Q \rightarrow \downarrow\uparrow Q \rightarrow \downarrow\uparrow Q \rightarrow N$ and $P \rightarrow \downarrow\uparrow P_1 \rightarrow \downarrow\uparrow P_2 \rightarrow N$? The type $\forall\alpha^+, \beta^+, \gamma^+. \alpha^+ \rightarrow \downarrow\uparrow\beta^+ \rightarrow \downarrow\uparrow\gamma^+ \rightarrow N$ is a common subtype, however, it is not the greatest one, as it is too general.

One can find two greater candidates: $M_1 = \forall\alpha^+, \beta^+. \alpha^+ \rightarrow \downarrow\uparrow\alpha^+ \rightarrow \downarrow\uparrow\beta^+ \rightarrow N$ and $M_2 = \forall\alpha^+, \beta^+. \beta^+ \rightarrow \downarrow\uparrow\alpha^+ \rightarrow \downarrow\uparrow\beta^+ \rightarrow N$. Instantiating α^+ and β^+ with Q ensures that both of these types are subtypes of $Q \rightarrow \downarrow\uparrow Q \rightarrow \downarrow\uparrow Q \rightarrow N$; instantiating α^+ with P_1 and β^+ with P_2 demonstrates the subtyping with $P \rightarrow \downarrow\uparrow P_1 \rightarrow \downarrow\uparrow P_2 \rightarrow N$, as P is a subtype of both P_1 and P_2 .

By analyzing the inference rules, we can prove that both M_1 and M_2 are *maximal* common subtypes, i.e., there is no common subtype that is greater than them. However, M_1 and M_2 are not equivalent, which means that none of them is the greatest.

3.2 Equivalence and Normalization

The subtyping relation forms a preorder on types, and thus, it induces an equivalence relation also known as bicoercibility [Tiuryn 1995]. The declarative specification of subtyping must be defined up to this equivalence. Moreover, the algorithms we use must withstand changes in input types within the equivalence class. To deal with non-trivial equivalence, we use normalization—a function that uniformly selects a representative of the equivalence class.

Using normalization gives us two benefits: (i) we do not need to modify significantly standard operations such as unification to withstand non-trivial equivalence, and (ii) if the subtyping (and thus, the equivalence) changes, we only need to modify the normalization function, while the rest of the algorithm remains the same.

In our system, equivalence is richer than equality. Specifically, while staying within one equivalence class, one can change the type:

- introduce and remove redundant quantifiers. For example, $\forall\alpha^+, \beta^+. \uparrow\alpha^+$ is equivalent but not equal to $\forall\alpha^+. \uparrow\alpha^+$;
- reorder adjacent quantifiers. For example, $\forall\alpha^+, \beta^+. \alpha^+ \rightarrow \beta^+ \rightarrow \gamma^-$ is equivalent but not equal to $\forall\alpha^+, \beta^+. \beta^+ \rightarrow \alpha^+ \rightarrow \gamma^-$;
- make the transformations (i) and (ii) at any position in the type.

It turns out that the transformations (i-iii) are complete, in the sense that they generate the whole equivalence class. This way, to normalize the type, one must

- remove the redundant quantifiers,
- reorder the quantifiers to the canonical order,
- do the procedures (i) and (ii) recursively on the subterms.

The normalization algorithm is shown in fig. 5. The steps (i-ii) are implemented by the *ordering* function ‘ord vars in N ’ and ‘ord vars in P ’. For a set of variables $vars$, and a type, it returns a list of

variables from $vars$ that occur in the type in the *order of their first occurrence*. For brevity, we omit the formal definition of ordering referring the reader to the appendix (??).

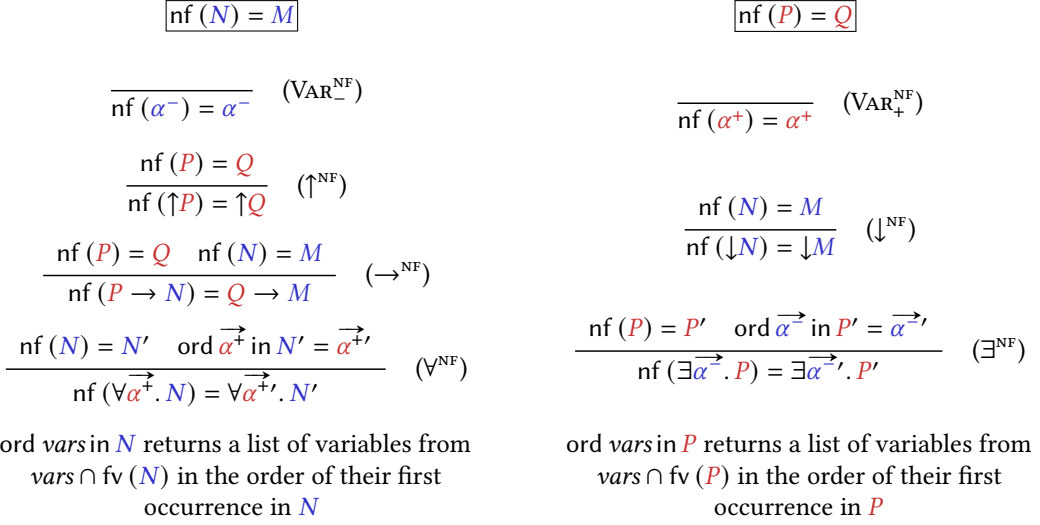


Fig. 5. Type Normalization Procedure

For the normalization procedure, we prove soundness and completeness w.r.t. the equivalence relation.

Property 6 (Correctness of normalization). *Assuming all types are well-formed in Θ ,*

- $\Theta \vdash N \approx^{\leq} M$ is equivalent to $\text{nf}(N) = \text{nf}(M)$, and
- + $\Theta \vdash P \approx^{\leq} Q$ is equivalent to $\text{nf}(P) = \text{nf}(Q)$.

3.3 Typing

The declarative specification of the type inference is given in fig. 6. The positive typing judgment $\Theta; \Gamma \vdash v : P$ is read as “under the type context Θ and variable context Γ , the term v is allowed to infer type P ”, where Γ —the variable context—is defined standardly as a set of pairs of the form $x : P$. The negative typing judgment is read similarly.

The *Application typing* judgment infers the type of the application of a function to a list of arguments. It has form $\Theta; \Gamma \vdash N \bullet \vec{v} \Rightarrow M$, which reads “under the type context Θ and variable context Γ , the application of a function of type N to the list of arguments \vec{v} is allowed to infer type M ”.

Let us discuss the rules of the declarative system in more detail.

Variables. Rule $(\text{VAR}^{\text{INF}})$ allows us to infer the type of a variable from the context. In literature, one can find another version of this rule, that enables inferring a type *equivalent* to the type from the context. In our case, the inference of equivalent types is admissible in general case by (\approx_+^{INF}) .

Annotations. The annotation rules $(\text{ANN}_-^{\text{INF}})$ and $(\text{ANN}_+^{\text{INF}})$ use subtyping. The annotation is only valid if the inferred type is a subtype of the annotation type.

Abstractions. The typing of lambda abstraction is standard. Rule (λ^{INF}) first checks that the given type annotating the argument is well-formed, and then infers the type of the body in the extended

$\boxed{\Theta; \Gamma \vdash c: N}$	Negative typing
$\frac{\Theta \vdash P \quad \Theta; \Gamma, x: P \vdash c: N}{\Theta; \Gamma \vdash \lambda x: P. c: P \rightarrow N} \quad (\lambda^{\text{INF}})$	$\frac{\Theta; \Gamma \vdash v: \downarrow M \quad \Theta; \Gamma \vdash M \bullet \vec{v} \Rightarrow \uparrow Q \text{ principal} \quad \Theta; \Gamma, x: Q \vdash c: N}{\Theta; \Gamma \vdash \text{let } x = v(\vec{v}); c: N} \quad (\text{LET}_{@}^{\text{INF}})$
$\frac{\Theta, \alpha^+; \Gamma \vdash c: N}{\Theta; \Gamma \vdash \Lambda \alpha^+. c: \forall \alpha^+. N} \quad (\Lambda^{\text{INF}})$	
$\frac{\Theta; \Gamma \vdash v: P}{\Theta; \Gamma \vdash \text{return } v: \uparrow P} \quad (\text{RET}^{\text{INF}})$	$\frac{\Theta \vdash P \quad \Theta; \Gamma \vdash v: \downarrow M \quad \Theta; \Gamma \vdash M \bullet \vec{v} \Rightarrow M' \quad \Theta \vdash M' \leq \uparrow P \quad \Theta; \Gamma, x: P \vdash c: N}{\Theta; \Gamma \vdash \text{let } x: P = v(\vec{v}); c: N} \quad (\text{LET}_{:@}^{\text{INF}})$
$\frac{\Theta; \Gamma \vdash v: P \quad \Theta; \Gamma, x: P \vdash c: N}{\Theta; \Gamma \vdash \text{let } x = v; c: N} \quad (\text{LET}^{\text{INF}})$	
$\frac{\Theta \vdash P \quad \Theta; \Gamma \vdash c: M \quad \Theta \vdash M \leq \uparrow P \quad \Theta; \Gamma, x: P \vdash c': N}{\Theta; \Gamma \vdash \text{let } x: P = c; c': N} \quad (\text{LET}_c^{\text{INF}})$	$\frac{\Theta; \Gamma \vdash v: \exists \vec{\alpha}^-. P \quad \text{nf}(\exists \vec{\alpha}^-. P) = \exists \vec{\alpha}^-. P \quad \Theta, \vec{\alpha}^-; \Gamma, x: P \vdash c: N \quad \Theta \vdash N}{\Theta; \Gamma \vdash \text{let}^{\exists}(\vec{\alpha}^-, x) = v; c: N} \quad (\text{LET}_{\exists}^{\text{INF}})$
$\frac{\Theta; \Gamma \vdash c: N \quad \Theta \vdash N \simeq^{\leq} N'}{\Theta; \Gamma \vdash c: N'} \quad (\simeq_-^{\text{INF}})$	$\frac{\Theta \vdash M \quad \Theta; \Gamma \vdash c: N \quad \Theta \vdash N \leq M}{\Theta; \Gamma \vdash (c: M): M} \quad (\text{ANN}_{-}^{\text{INF}})$
$\boxed{\Theta; \Gamma \vdash v: P}$	Positive typing
$\frac{x: P \in \Gamma}{\Theta; \Gamma \vdash x: P} \quad (\text{VAR}^{\text{INF}})$	$\frac{\Theta \vdash Q \quad \Theta; \Gamma \vdash v: P \quad \Theta \vdash Q \geq P}{\Theta; \Gamma \vdash (v: Q): Q} \quad (\text{ANN}_{+}^{\text{INF}})$
$\frac{\Theta; \Gamma \vdash c: N}{\Theta; \Gamma \vdash \{c\}: \downarrow N} \quad (\{\}^{\text{INF}})$	$\frac{\Theta; \Gamma \vdash v: P \quad \Theta \vdash P \simeq^{\leq} P'}{\Theta; \Gamma \vdash v: P'} \quad (\simeq_{+}^{\text{INF}})$
$\boxed{\Theta; \Gamma \vdash N \bullet \vec{v} \Rightarrow M}$	Application typing
$\frac{\Theta \vdash N \simeq^{\leq} N'}{\Theta; \Gamma \vdash N \bullet \cdot \Rightarrow N'} \quad (\emptyset_{\bullet \Rightarrow}^{\text{INF}})$	$\frac{\vec{v} \neq \vec{\alpha}^+ \neq \cdot \quad \Theta \vdash \sigma: \vec{\alpha}^+ \quad \Theta; \Gamma \vdash [\sigma] N \bullet \vec{v} \Rightarrow M}{\Theta; \Gamma \vdash \forall \vec{\alpha}^+. N \bullet \vec{v} \Rightarrow M} \quad (\forall_{\bullet \Rightarrow}^{\text{INF}})$
$\frac{\Theta; \Gamma \vdash v: P \quad \Theta \vdash Q \geq P \quad \Theta; \Gamma \vdash N \bullet \vec{v} \Rightarrow M}{\Theta; \Gamma \vdash Q \rightarrow N \bullet v, \vec{v} \Rightarrow M} \quad (\rightarrow_{\bullet \Rightarrow}^{\text{INF}})$	

Fig. 6. Declarative Inference

context. As a result, it returns an arrow type of function from the annotated type of the argument to the type of the body. Rule (Λ^{INF}) infers polymorphic \forall -type. It extends the type context with the quantifying variable α^+ and infers the type of the body. As a result, it returns a polymorphic type quantifying the abstracted variable α^+ over the type of the body.

Return and thunk. Rules $(\text{RET}^{\text{INF}})$ and $(\{\}^{\text{INF}})$ simply add the corresponding shift constructors to the type of the body.

Unpack. Rule $(\text{LET}_{\exists}^{\text{INF}})$ types elimination of \exists . First, it infers the normalized type of the existential package. The normalization is required to fix the order of the quantifying variables to bind them. Alternative approaches that do not require normalization will be discussed in section 6.1. After the

bind, the rule infers the type of the body and checks that it does not use the bound variables so that they do not escape the scope.

Applicative let binders. Rules $(\text{LET}_{@}^{\text{INF}})$ and $(\text{LET}_{@}^{\text{INF}})$ infer the type of the applicative let binders. Both of them infer the type of the head v and invoke the application typing to infer the type of the application before recursing on the body of the let binder.

The former rule infers the type of an *unannotated* let binder, and thus it requires the resulting type of application to be *the principal type*, so that the type we assign to the bound variable x is determined. In this context, principality means minimality. In other words, $\Theta; \Gamma \vdash M \bullet \vec{v} \Rightarrow \uparrow Q$ principal means that any other type Q' inferable for the application (i.e., $\Theta; \Gamma \vdash M \bullet \vec{v} \Rightarrow \uparrow Q'$) is greater than the principal type Q , i.e., $\Theta \vdash Q' \geq Q$.

The latter rule is for the *annotated* binder, and thus, the type of the bound x is given, however, the rule must check that this type is a supertype of the inferred type of the application. This check is done by invoking the subtyping judgment $\Theta \vdash \uparrow Q \leq \uparrow P$. This judgment is more restrictive than checking bare $\Theta \vdash P \geq Q$, however, it is necessary to make the algorithm complete as it allows us to preserve certain invariants throughout the algorithm to ensure the decidability of the produced unification task. In section 5.4, we discuss these invariants in more detail, and in section 6.2 we suggest a way to relax them.

Typing up to equivalence. As discussed in section 3.2, the subtyping, as a preorder, induces a non-trivial equivalence relation on types. The system must not distinguish between equivalent types, and thus, type inference must be defined up to equivalence. For this purpose, we use rules $(\simeq_{+}^{\text{INF}})$ and $(\simeq_{-}^{\text{INF}})$. They allow one to replace the inferred type with an equivalent one.

Application to an empty list of arguments. The base case of the application type inference is represented by rule $(\emptyset_{\bullet \Rightarrow}^{\text{INF}})$. If the head of the type N is applied to no arguments, the type of the result is allowed to be N or any equivalent type. We need to relax this rule up to equivalence to ensure the corresponding property globally: the inferred application type can be replaced with an equivalent one. Alternatively, we could have added a separate rule similar to $(\simeq_{+}^{\text{INF}})$, however, the local relaxation is sufficient to prove the global property.

Application of a polymorphic type \forall . The complexity of the system lives in the rules whose output type is not immediately defined by their input and the output of their premises (also known as not mode-correct [Dunfield et al. 2020]). In our typing system, $(\forall_{\bullet \Rightarrow}^{\text{INF}})$ is such a rule: the instantiation of the quantifying variables is not known a priori. The algorithm we present in section 4 delays this instantiation until more information about it (in particular, typing constraints) is collected.

To ensure the priority of application between this rule and $(\emptyset_{\bullet \Rightarrow}^{\text{INF}})$, we also check that the list of arguments is not empty.

Application of an arrow type. Another application rule is $(\rightarrow_{\bullet \Rightarrow}^{\text{INF}})$. This is where the subtyping is used to check that the type of the argument is convertible to (a subtype of) the type of the function parameter. In the algorithm (section 4), this subtyping check will provide the constraints we need to resolve the delayed instantiations of the quantifying variables.

An important property that the declarative system has is that the declarative specification is correctly defined for equivalence classes.

Property 7 (Declarative typing is defined up to equivalence). *Let us assume that $\Theta \vdash \Gamma_1 \simeq^{\leq} \Gamma_2$, i.e., the corresponding types assigned by Γ_1 and Γ_2 are equivalent in Θ . Also, let us assume that $\Theta \vdash N_1 \simeq^{\leq} N_2$, $\Theta \vdash P_1 \simeq^{\leq} P_2$, and $\Theta \vdash M_1 \simeq^{\leq} M_2$. Then*

- $\Theta; \Gamma_1 \vdash c: N_1$ holds if and only if $\Theta; \Gamma_2 \vdash c: N_2$,

- + $\Theta; \Gamma_1 \vdash v: P_1$ holds if and only if $\Theta; \Gamma_2 \vdash v: P_2$, and
- $\Theta; \Gamma_1 \vdash N_1 \bullet \vec{v} \Rightarrow M_1$ holds if and only if $\Theta; \Gamma_2 \vdash N_2 \bullet \vec{v} \Rightarrow M_2$.

3.4 Relation to System F

To establish a correspondence between F_{\exists}^{\pm} and standard unpolarized System F, we present a translation in both ways: the polarization from System F to F_{\exists}^{\pm} and the depolarization from F_{\exists}^{\pm} to System F.

Type-level Translation. The type depolarization (fig. 7) simply forgets the polarization structure of types: it removes the shift operators and the polarities of the free type variables.

The type polarization (fig. 8) is more involved, as there are multiple ways to polarize a type: for instance, a variable α can be polarized to α^+ or α^- . The choice of polarization affects the execution strategy of the program. **Ilya: How?** We chose the positive polarization: every System F type is translated to a *positive* type of F_{\exists}^{\pm} . This choice makes it smoother to lift the translation to the term level, which will be disused next.

The type-level translation is naturally lifted to the contexts: $|\Theta|$ is the context Θ with every type depolarized, and symmetrically, $|\Theta|$ is the context Θ with every type polarized. The context translation is used to formulate the soundness of the term-level translation, which we discuss next.

$$\begin{array}{c}
 \boxed{|P|} \qquad \qquad \qquad \boxed{|N|} \\
 \\
 |\alpha^+| \equiv \alpha \qquad \qquad |\alpha^-| \equiv \alpha \\
 |\downarrow N| \equiv |N| \qquad \qquad |\uparrow P| \equiv |P| \\
 |\exists \alpha^{\rightarrow}. P| \equiv \exists \vec{\alpha}. |P| \qquad |\forall \alpha^+. N| \equiv \forall \alpha. |N| \\
 |P \rightarrow N| \equiv |P| \rightarrow |N|
 \end{array}$$

Fig. 7. Type Depolarization ¹

$$\begin{array}{c}
 \boxed{|T|} \\
 \\
 |\alpha| \equiv \alpha^+ \\
 |\forall \alpha. T| \equiv \downarrow (\forall \alpha^+. \uparrow |T|) \\
 |A \rightarrow B| \equiv \downarrow (|A| \rightarrow \uparrow |B|)
 \end{array}$$

Fig. 8. Type Polarization

Term-level Translation. The term-level translations between the systems are more complex. They are defined as *elaborations* of the typing derivations. In other words, the input of the translation is not a term, but a whole typing derivation in the corresponding system. For example, the subtyping elaboration $\Theta \vdash N \leq M \rightsquigarrow t$ constructs a System F function t of type $|N| \rightarrow |M|$ that witnesses the given F_{\exists}^{\pm} -subtyping $\Theta \vdash N \leq M$; a typing elaboration $\Theta; \Gamma \vdash v: P \rightsquigarrow t$ builds a term t —a System F counterpart of v —of type $|P|$.

For brevity, we omit the formal definition of the term-level elaboration, referring the reader to the appendix (??). We define the elaboration in such a way that it preserves the soundness invariants given in ??.

The least obvious aspect of the elaboration soundness is the application typing. The elaboration judgment $\Theta; \Gamma \vdash N \bullet \vec{v} \Rightarrow M \rightsquigarrow e; \vec{t}$ outputs two things: (i) \vec{t} —the System F counterpart of the argument list \vec{v} , and (ii) e —a System F term that is used for the final cast of the application $((x t_1) \cdots t_n)$ to the resulting type $|M|$. This final cast is needed to reflect the F_{\exists}^{\pm} feature of $(\Theta_{\bullet}^{\text{INF}})$ that permits the equivalent conversion of the resulting type $\Theta \vdash N \simeq^{\epsilon} N'$, as F_{\exists}^{\pm} -equivalent types might have different System F representations (i.e., $|N| \neq |N'|$).

4 THE ALGORITHM

In this section, we present the algorithmization of the declarative system described above. The algorithmic system follows the structure of the declarative specification closely. First, it is also given by a set of inference rules, which, however, are mode-correct ([Dunfield et al. 2020]), i.e., the

¹Here the System F existential $\exists \vec{\alpha}. T$ is a syntax sugar for its standard encoding $\forall \beta. (\forall \vec{\alpha}. (T \rightarrow \beta)) \rightarrow \beta$

	Elaboration Judgment	The Soundness Property
Negative Subtyping	$\Theta \vdash N \leq M \rightsquigarrow t$	$ \Theta ; \cdot \vdash t: N \rightarrow M $
Positive Subtyping	$\Theta \vdash P \geq Q \rightsquigarrow t$	$ \Theta ; \cdot \vdash t: Q \rightarrow P $
Negative Typing	$\Theta; \Gamma \vdash c: N \rightsquigarrow t$	$ \Theta ; \Gamma \vdash t: N $
Positive Typing	$\Theta; \Gamma \vdash v: P \rightsquigarrow t$	$ \Theta ; \Gamma \vdash t: P $
Application Typing	$\Theta; \Gamma \vdash N \bullet \vec{v} \Rightarrow M \rightsquigarrow e; \vec{t}$	$ \Theta ; \Gamma , x: N \vdash e((x t_1) \cdots t_n): M $
System F Typing	$\Theta; \Gamma \vdash t: T \rightsquigarrow^\pm c$	$ \Theta ; \Gamma \vdash c: \uparrow \downarrow T$

Table 1. Soundness of the Elaboration: the elaboration judgment on the left implies the typing judgment on the right

Negative Algorithmic Variables

$\widehat{\alpha}^-, \widehat{\beta}^-, \widehat{\gamma}^-, \dots$

Positive Algorithmic Variables

$\widehat{\alpha}^+, \widehat{\beta}^+, \widehat{\gamma}^+, \dots$

Negative Algorithmic Types

$N, M ::= \dots \mid \widehat{\alpha}^-$

Positive Algorithmic Types

$P, Q ::= \dots \mid \widehat{\alpha}^+$

Algorithmic Type Context

$\widehat{\Theta} ::= \{\widehat{\alpha}_1^\pm, \dots, \widehat{\alpha}_n^\pm\}$ where $\widehat{\alpha}_1^\pm, \dots, \widehat{\alpha}_n^\pm$ are pairwise distinct

Instantiation Context

$\Xi ::= \{\widehat{\alpha}_1^\pm \{\Theta_1\}, \dots, \widehat{\alpha}_n^\pm \{\Theta_n\}\}$ where $\widehat{\alpha}_1^\pm, \dots, \widehat{\alpha}_n^\pm$ are pairwise distinct

Fig. 9. Algorithmic Syntax

output of each rule is always uniquely defined by its input. And second, most of the declarative rules (except for the rules (\simeq_+^{INF}) and (\simeq_-^{INF})) have a unique algorithmic counterpart, which simplifies reasoning about the algorithm and its correctness proofs.

4.1 Algorithmic Syntax

First, let us discuss the syntax of the algorithmic system.

Algorithmic Variables. To design a mode-correct inference system, we slightly modify the language we operate on. The entities (terms, types, contexts) that the algorithm manipulates we call *algorithmic*. They extend the previously defined declarative terms and types by adding *algorithmic type variables* (also known as unification variables). The algorithmic variables represent unknown types, which cannot be inferred immediately but are promised to be instantiated as the algorithm proceeds.

We denote algorithmic variables as $\widehat{\alpha}^+, \widehat{\beta}^-, \dots$ to distinguish them from normal variables α^+, β^- . In a few places, we replace the quantified variables $\vec{\alpha}^+$ with their algorithmic counterpart $\vec{\widehat{\alpha}}^+$. The procedure of replacing declarative variables with algorithmic ones we call *algorithmization* and denote as $\vec{\alpha}^+ / \vec{\widehat{\alpha}}^+$ and $\vec{\alpha}^- / \vec{\widehat{\alpha}}^-$.

Algorithmic Types. The syntax of algorithmic types extends the declarative syntax by adding algorithmic variables as new terminals. We add positive algorithmic variables $\widehat{\alpha}^+$ to the syntax of positive types, and negative algorithmic variables $\widehat{\alpha}^-$ to the syntax of negative types. All the

$$\begin{array}{ccc}
\vdots & & \vdots \\
\frac{\widehat{\alpha}^+ \in \widehat{\Theta}}{\Theta; \widehat{\Theta} \vdash \widehat{\alpha}^+} & (\text{UVar}_+^{\text{WF}}) & \frac{\widehat{\alpha}^- \in \widehat{\Theta}}{\Theta; \widehat{\Theta} \vdash \widehat{\alpha}^-} & (\text{UVar}_-^{\text{WF}})
\end{array}$$

Fig. 10. Well-formedness of Algorithmic Types

$$\begin{array}{ccc}
\vdots & & \vdots \\
\overline{\text{nf}(\widehat{\alpha}^+) = \widehat{\alpha}^+} & (\text{UVar}_+^{\text{NF}}) & \overline{\text{nf}(\widehat{\alpha}^-) = \widehat{\alpha}^-} & (\text{UVar}_-^{\text{NF}})
\end{array}$$

Fig. 11. Normalization of Algorithmic Types

constructors of the system can be applied to *algorithmic* types, however, algorithmic variables cannot be abstracted by the quantifiers \forall and \exists .

Algorithmic Contexts $\widehat{\Theta}$ and Well-formedness. To specify when algorithmic types are well-formed, we define algorithmic contexts $\widehat{\Theta}$ as sets of algorithmic variables. Then $\Theta; \widehat{\Theta} \vdash P$ and $\Theta; \widehat{\Theta} \vdash N$ represent the well-formedness judgment of algorithmic terms defined as expected. Informally, they check that all free declarative variables are in Θ , and all free algorithmic variables are in $\widehat{\Theta}$. Most of the rules are inherited from the well-formedness of *declarative* types: the declarative variables are checked to belong to the context Θ , the quantifiers extend the context, type constructors are well-formed congruently. As we extend the syntax with algorithmic variables, we also add two base-case rules for them: $(\text{UVar}_+^{\text{WF}})$ and $(\text{UVar}_-^{\text{WF}})$ (see fig. 10).

Instantiation Context Ξ . When one instantiates an algorithmic variable, one may only use type variables *available in its scope*. As such, each algorithmic variable must remember the context at the moment when it was introduced. In our algorithm, this information is represented by an *instantiation context* Ξ —a set of pairs associating algorithmic variables and declarative contexts.

Algorithmic Normalization. Similarly to well-formedness, the normalization of algorithmic types is defined by extending the declarative definition (fig. 5) with the algorithmic variables. To the rules repeating the declarative normalization, we add rules saying that normalization is trivial on algorithmic variables.

4.2 Type Constraints

Throughout the algorithm's operation, it gathers information about the algorithmic type variables, represented as *constraints*. These constraints in our system can be either *subtyping constraints* or *unification constraints*. As preserved by the algorithm, each subtyping constraint has a positive shape $\widehat{\alpha}^+ : \geq P$, meaning it confines a positive algorithmic variable to be the supertype of a positive declarative type. Unification constraints can take a positive $\widehat{\alpha}^+ : \simeq P$ or negative $\widehat{\alpha}^- : \simeq N$ form, but algorithmic type variables cannot occur on their right-hand side. The constraints set is denoted as C , and we presume that each algorithmic variable can be restricted by at most one constraint.

We define UC separately as a set solely containing *unification* constraints for simpler algorithm representation. The unification algorithm, which we use as a subroutine of the subtyping algorithm, can only produce unification constraints. A resolution of unification constraints is simpler than that

of a general constraint set. This way, this separation allows us to better decompose the algorithm's structure, thus simplifying the inductive proofs.

Constraint Entry		Unification Constraint Entry	
e	$::=$	ue	$::=$
	$\mid \hat{a}^+ \approx P$		$\mid \hat{a}^+ \approx P$
	$\mid \hat{a}^- \approx N$		$\mid \hat{a}^- \approx N$
	$\mid \hat{a}^+ \geq P$		
Constraint Set		Unification Constraint Set	
C	$::= \{e_1, \dots, e_n\}$	UC	$::= \{ue_1, \dots, ue_n\}$

Fig. 12. Constraint Entries and Sets

Auxiliary Functions. We define $\text{dom}(C)$ —the domain of constraint set C —as a set of algorithmic variables that it restricts. Similarly, we define $\text{dom}(\Xi)$ —the domain of instantiation context—as a set of algorithmic variables that Ξ associates with their contexts. We write $\Xi(\hat{a}^\pm)$ to denote the declarative context associated with \hat{a}^\pm in Ξ .

4.3 Subtyping Algorithm

For convenience and scalability, we decompose the subtyping algorithm into several procedures. Figure 13 shows these procedures and the dependencies between them: arrows denote the invocation of one procedure from another.

Some of the procedures (in particular, the unification and the anti-unification) assume that the input types are normalized. Therefore, we call the normalization procedure before invoking them, indicating this by the ‘nf’ annotation on the arrows in fig. 13. Alternatively, one could normalize the input types in the very beginning, as the normalization is preserved throughout the algorithm. However, we delay the normalization to the places where it is required to show that normalization is needed only at these stages to maintain consistent invariants.

In the remainder of this section, we will delve into each of these procedures in detail, following the top-down order of the dependency graph. First, we present the subtyping algorithm itself.

As an input, the subtyping algorithm takes a type context Θ , an instantiation context Ξ , and two types of the corresponding polarity: N and M for the negative subtyping, and P and Q for the positive subtyping. We assume the second type (M and Q) to be declarative (with no algorithmic variables) and well-formed in Θ , but the first type (N and P) may contain algorithmic variables, whose instantiation contexts are specified by Ξ .

Notice that the shape of the input types uniquely determines the applied subtyping rule. If the subtyping is successful, it returns a set of constraints C restricting the algorithmic variables of the first type. If the subtyping does not hold, there will be no inference tree with such inputs.

The rules of the subtyping algorithm bijectively correspond to the rules of the declarative system. Let us discuss them in detail.

Variables. Rules (VAR^\leq) and (VAR^\geq) say that if both of the input types are equal declarative variables, they are subtypes of each other, with no constraints (as there are no algorithmic variables).

Shifts. Rules (\downarrow^\geq) and (\uparrow^\leq) cover the downshift and the upshift cases, respectively. If the input types are constructed by shifts, then the subtyping can only hold if they are equivalent. This way, the algorithm must find the instantiations of the algorithmic variables on the left-hand side such that these instantiations make the left-hand side and the right-hand side equivalent. For this purpose,

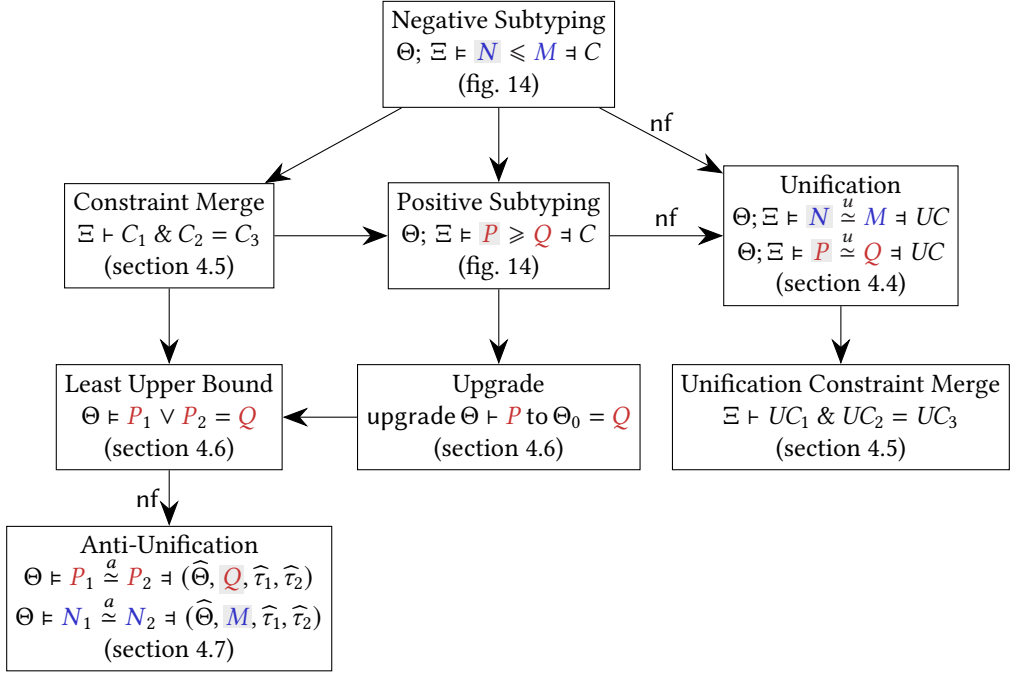


Fig. 13. Dependency graph of the subtyping algorithm

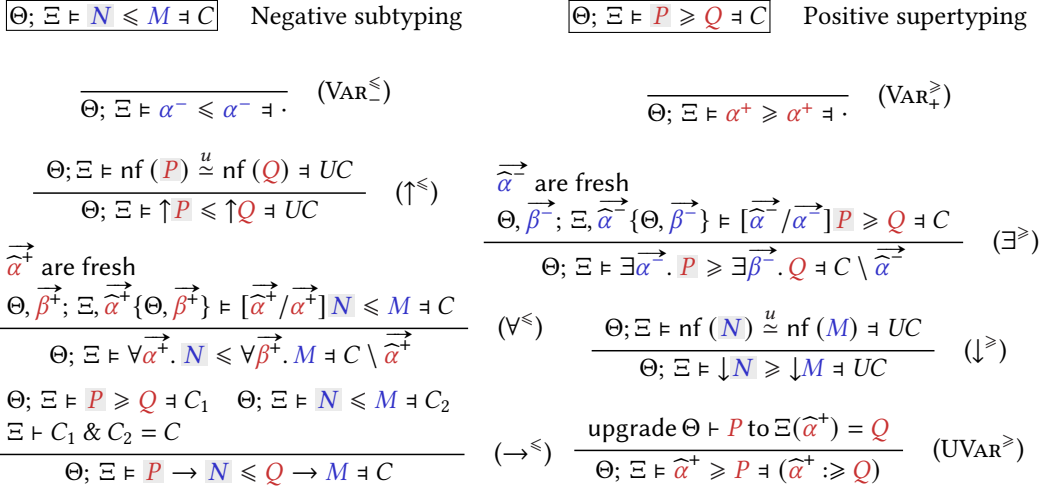


Fig. 14. Subtyping Algorithm

the algorithm invokes the unification procedure (section 4.4) preceded by the normalization of the input types. It returns the resulting constraints given by the unification algorithm.

Quantifiers. Rules (\forall^{\leq}) and (\exists^{\geq}) are symmetric. Declaratively, the quantified variables on the left-hand side must be instantiated with types, which are not known beforehand. We address this

problem by algorithmization of the quantified variables (see section 4.1). The rule introduces fresh algorithmic variables $\vec{\alpha}^+$ or $\vec{\alpha}^-$, puts them into the instantiation context Ξ (specifying that they must be instantiated in the extended context $\Theta, \vec{\beta}^+$ or $\Theta, \vec{\beta}^-$) and substitute the quantified variables for them in the input type.

After algorithmization of the quantified variables, the algorithm proceeds with the recursive call, returning constraints C . As the output, the algorithm removes the freshly introduced algorithmic variables from the instantiation context. This operation is sound: it is guaranteed that C always has a solution, but the specific instantiation of the freshly introduced algorithmic variables is not important, as they do not occur in the input types.

Functions. To infer the subtyping of the function types, the algorithm makes two calls: (i) a recursive call ensuring the subtyping of the result types, and (ii) a call to positive subtyping (or rather super-typing) on the argument types. The resulting constraints are merged using a special procedure defined in section 4.5 and returned as the output.

Algorithmic variables. If one of the sides of the subtyping is a unification variable, the algorithm creates a new constraint. Because the right-hand side of the subtyping is always declarative, it is only the left-hand side that can be a unification variable. Moreover, another invariant we preserve prevents the negative algorithmic variables from occurring in types during the negative subtyping algorithm. It means that the only possible form of the subtyping here is $\vec{\alpha}^+ \geq P$, which is covered by (UVar[≥]).

The potential problem here is that the type P might be not well-formed in the instantiation context required for $\vec{\alpha}^+$ by Ξ because this context might be smaller than the current context Θ . As we wish the resulting constraint set to be sound w.r.t. Ξ , we cannot simply put $\vec{\alpha}^+ : \geq P$ into the output. Prior to that, we update the type P to its lowest supertype Q well-formed in $\Xi(\vec{\alpha}^+)$. It is done by the *upgrade* procedure, which we discuss in detail in section 4.6.

To summarize, the subtyping algorithm uses the following additional subroutines: (i) rules (\downarrow^{\geq}) and (\uparrow^{\leq}) invoke the *unification* algorithm to equate the input types; (ii) rule (\rightarrow^{\leq}) merges the constraints produced by the recursive calls on the result and the argument types; and (iii) rule (UVar[≥]) upgrades the input type to its least supertype well-formed in the context required by the algorithmic variable. The following sections discuss these additional procedures in detail.

4.4 Unification

As an input, the unification context takes a type context Θ , an instantiation context Ξ , and two types of the required polarity: N and M for the negative unification, and P and Q for the positive unification. It is assumed that only the left-hand side type may contain algorithmic variables. This way, the left-hand side is well-formed as an algorithmic type in Θ and Ξ , whereas the right-hand side is well-formed declaratively in Θ .

Since only the left-hand side may contain algorithmic variables that the unification instantiates, we could have called this procedure *matching*. However, in section 6.2, we will discuss several modifications of the type system, where this invariant is not preserved, and therefore, this procedure requires general first-order pattern unification [Miller 1991].

As the output, the unification algorithm returns the *weakest* set of unification constraints UC such that any instantiation satisfying these constraints unifies the input types.

The algorithm works as one might expect: if both sides are formed by constructors, it is required that the constructors are the same, and the types unify recursively. If one of the sides is a unification

<div style="border: 1px solid black; padding: 2px; display: inline-block; margin-bottom: 5px;"> $\Theta; \Xi \models \mathbf{N} \stackrel{u}{\simeq} \mathbf{M} \models UC$ </div> <p style="margin: 0;">Negative unification</p> $ \begin{array}{c} \frac{}{\Theta; \Xi \models \alpha^- \stackrel{u}{\simeq} \alpha^- \models \cdot} \quad (\text{VAR}_{-}^u) \\ \\ \frac{\Theta; \Xi \models \mathbf{P} \stackrel{u}{\simeq} \mathbf{Q} \models UC}{\Theta; \Xi \models \uparrow \mathbf{P} \stackrel{u}{\simeq} \uparrow \mathbf{Q} \models UC} \quad (\uparrow^u) \\ \\ \frac{\Theta; \Xi \models \mathbf{P} \stackrel{u}{\simeq} \mathbf{Q} \models UC_1 \quad \Theta; \Xi \models \mathbf{N} \stackrel{u}{\simeq} \mathbf{M} \models UC_2}{\Theta; \Xi \models \mathbf{P} \rightarrow \mathbf{N} \stackrel{u}{\simeq} \mathbf{Q} \rightarrow \mathbf{M} \models UC_1 \ \& \ UC_2} \quad (\rightarrow^u) \\ \\ \frac{\Theta, \vec{\alpha}^+; \Xi \models \mathbf{N} \stackrel{u}{\simeq} \mathbf{M} \models UC}{\Theta; \Xi \models \forall \vec{\alpha}^+. \mathbf{N} \stackrel{u}{\simeq} \forall \vec{\alpha}^+. \mathbf{M} \models UC} \quad (\forall^u) \\ \\ \frac{\Xi(\widehat{\alpha}^-) \vdash \mathbf{N}}{\Theta; \Xi \models \widehat{\alpha}^- \stackrel{u}{\simeq} \mathbf{N} \models (\widehat{\alpha}^- \simeq \mathbf{N})} \quad (\text{UVar}_{-}^u) \end{array} $	<div style="border: 1px solid black; padding: 2px; display: inline-block; margin-bottom: 5px;"> $\Theta; \Xi \models \mathbf{P} \stackrel{u}{\simeq} \mathbf{Q} \models UC$ </div> <p style="margin: 0;">Positive unification</p> $ \begin{array}{c} \frac{}{\Theta; \Xi \models \alpha^+ \stackrel{u}{\simeq} \alpha^+ \models \cdot} \quad (\text{VAR}_{+}^u) \\ \\ \frac{\Theta; \Xi \models \mathbf{N} \stackrel{u}{\simeq} \mathbf{M} \models UC}{\Theta; \Xi \models \downarrow \mathbf{N} \stackrel{u}{\simeq} \downarrow \mathbf{M} \models UC} \quad (\downarrow^u) \\ \\ \frac{\Theta, \vec{\alpha}^+; \Xi \models \mathbf{P} \stackrel{u}{\simeq} \mathbf{Q} \models UC}{\Theta; \Xi \models \exists \vec{\alpha}^+. \mathbf{P} \stackrel{u}{\simeq} \exists \vec{\alpha}^+. \mathbf{Q} \models UC} \quad (\exists^u) \\ \\ \frac{\Xi(\widehat{\alpha}^+) \vdash \mathbf{P}}{\Theta; \Xi \models \widehat{\alpha}^+ \stackrel{u}{\simeq} \mathbf{P} \models (\widehat{\alpha}^+ \simeq \mathbf{P})} \quad (\text{UVar}_{+}^u) \end{array} $
---	---

Fig. 15. Unification Algorithm

variable (in our case it can only be the left-hand side), we create a new unification constraint restricting it to be equal to the other side. Let us discuss the rules that implement this strategy.

Variables. The variable rules (VAR_{-}^u) and (VAR_{+}^u) are trivial: as the input types do not have algorithmic variables, and are already equal, the unification returns no constraints.

Shifts. The shift rules (\downarrow^u) and (\uparrow^u) require the two input types to be formed by the same shift constructor. They remove this constructor, unify the types recursively, and return the resulting set of constraints.

Quantifiers. Similarly, the quantifier rules (\forall^u) and (\exists^u) require the quantifier variables on the left-hand side and the right-hand side to be the same. This requirement is complete because we assume the input types of the unification to be normalized, and thus, the equivalence implies alpha-equivalence. In the implementation of this rule, an alpha-renaming might be needed to ensure that the quantified variables are the same, however, we omit it for brevity.

Functions. Rule (\rightarrow^u) unifies two functional types. First, it unifies their argument types and their result types recursively. Then it merges the resulting constraints using the procedure described in section 4.5.

Notice that the resulting constraints can only have *unification* entries. It means that they can be merged in a simpler way than general constraints. In particular, the merging procedure does not call any of the subtyping subroutines but rather simply checks the matching constraint entries for equality.

Algorithmic variable. Finally, if the left-hand side of the unification is an algorithmic variable, (VAR_{-}^u) or (VAR_{+}^u) is applied. It simply checks that the right-hand side type is well-formed in the required instantiation context, and returns a newly created constraint restricting the variable to be equal to the right-hand side type.

As one can see, the unification procedure is standard; the only peculiarity (although it is common for type inference) is that it makes sure that the resulting instantiations agree with the input instantiation context Ξ . As a subroutine, the unification algorithm only uses the (unification) constraint merge procedure and the well-formedness checking.

4.5 Constraint Merge

In this section, we discuss the constraint merging procedure. It allows one to combine two constraint sets into one. A simple union of two constraint sets is not sufficient, since the resulting set must not contain two entries restricting the same algorithmic variable—we call such entries *matching*.

The matching entries must be combined into *one* constraint entry, that would represent their conjunction. This way, to merge two constraint sets, we unite the entries of two sets and then merge the matching pairs.

Merging matching constraint entries. Two *matching* entries formed in the same context Θ can be merged as shown in fig. 16. Suppose that e_1 and e_2 are input entries. The result of the merge $e_1 \& e_2$ is *the weakest entry which implies both e_1 and e_2* .

$\boxed{\Theta \vdash e_1 \& e_2 = e_3}$ Subtyping Constraint Entry Merge

$$\begin{array}{c}
 \frac{\Theta \vdash P_1 \vee P_2 = Q}{\Theta \vdash (\widehat{\alpha}^+ : \geq P_1) \& (\widehat{\alpha}^+ : \geq P_2) = (\widehat{\alpha}^+ : \geq Q)} \quad (\geq \&^+ \geq) \\
 \\
 \frac{\Theta; \cdot \vdash P \geq Q \ni \cdot}{\Theta \vdash (\widehat{\alpha}^+ : \simeq P) \& (\widehat{\alpha}^+ : \geq Q) = (\widehat{\alpha}^+ : \simeq P)} \quad (\simeq \&^+ \geq) \\
 \\
 \frac{\Theta; \cdot \vdash Q \geq P \ni \cdot}{\Theta \vdash (\widehat{\alpha}^+ : \geq P) \& (\widehat{\alpha}^+ : \simeq Q) = (\widehat{\alpha}^+ : \simeq Q)} \quad (\geq \&^+ \simeq) \\
 \\
 \frac{\text{nf}(P) = \text{nf}(P')}{\Theta \vdash (\widehat{\alpha}^+ : \simeq P) \& (\widehat{\alpha}^+ : \simeq P') = (\widehat{\alpha}^+ : \simeq P)} \quad (\simeq \&^+ \simeq) \\
 \\
 \frac{\text{nf}(N) = \text{nf}(N')}{\Theta \vdash (\widehat{\alpha}^- : \simeq N) \& (\widehat{\alpha}^- : \simeq N') = (\widehat{\alpha}^- : \simeq N)} \quad (\simeq \&^- \simeq)
 \end{array}$$

Fig. 16. Merge of Matching Constraint Entries

Suppose that one of the input entries, say e_1 , is a *unification* constraint entry. Then the resulting entry e_1 must coincide with it (up-to-equivalence), and thus, it is only required to check that e_2 is implied by e_1 . We consider two options:

- (i) if e_2 is also a *unification* entry, then the types on the right-hand side of e_1 and e_2 must be equivalent, as given by rules $(\simeq \&^+ \simeq)$ and $(\simeq \&^- \simeq)$;
- (ii) if e_2 is a *supertype* constraint entry $\widehat{\alpha}^+ : \geq P$, the algorithm must check that the type assigned by e_1 is a supertype of P . The corresponding symmetric rules are $(\geq \&^+ \simeq)$ and $(\simeq \&^+ \geq)$. In the premises of these rules, P is the same as P below, and Q is the same as Q below but relaxed to an algorithmic type.

If both input entries are supertype constraints: $\widehat{\alpha}^+ : \geq P$ and $\widehat{\alpha}^+ : \geq Q$, then their conjunction is $\widehat{\alpha}^+ : \geq P \vee Q$, as given by $(\geq \&^+ \geq)$. The least upper bound— $P \vee Q$ —is the least supertype of both P

and Q , and this way, $\widehat{\alpha}^+ : \geq P \vee Q$ is the weakest constraint entry that implies $\widehat{\alpha}^+ : \geq P$ and $\widehat{\alpha}^+ : \geq Q$. The algorithm finding the least upper bound is discussed in section 4.6.

Merging constraint sets. The algorithm for merging constraint sets is shown in fig. 17. As discussed, the result of merge C_1 and C_2 consists of three parts: (i) the entries of C_1 that do not match any entry of C_2 ; (ii) the entries of C_2 that do not match any entry of C_1 ; and (iii) the merge (fig. 16) of matching entries.

Suppose that $\Xi \vdash C_1$ and $\Xi \vdash C_2$.

Then $\Xi \vdash C_1 \& C_2 = C$ defines a set of constraints C such that $e \in C$ iff either:

- $e \in C_1$ and there is no matching $e' \in C_2$; or
- $e \in C_2$ and there is no matching $e' \in C_1$; or
- $\Xi(\widehat{\alpha}^+) \vdash e_1 \& e_2 = e$ for some $e_1 \in C_1$ and $e_2 \in C_2$ such that e_1 and e_2 both restrict variable $\widehat{\alpha}^+$.

Fig. 17. Constraint Merge

As shown in fig. 16, the merging procedure relies substantially on the least upper bound algorithm. In the next section, we discuss this algorithm in detail, together with the upgrade procedure, selecting the least supertype well-formed in a given context.

4.6 Type Upgrade and the Least Upper Bounds

Both type upgrade and the least upper bound algorithms are used to find a minimal supertype under certain conditions. For a given type P well-formed in Θ , the *upgrade* operation finds the least among those supertypes of P that are well-formed in a smaller context $\Theta_0 \subseteq \Theta$. For given two types P_1 and P_2 well-formed in Θ , the *least upper bound* operation finds the least among common supertypes of P_1 and P_2 well-formed in Θ . These algorithms are shown in fig. 18.

<div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;"> $\text{upgrade } \Theta \vdash P \text{ to } \Theta_0 = Q$ </div> <div style="display: flex; justify-content: space-between;"> <div style="width: 45%;"> <p style="text-align: center;">Type Upgrade</p> $\begin{array}{c} \Theta = \Theta_0, \vec{\alpha}^\pm \\ \vec{\beta}^\pm \text{ are fresh } \quad \vec{\gamma}^\pm \text{ are fresh} \\ \Theta_0, \vec{\beta}^\pm, \vec{\gamma}^\pm \vdash [\vec{\beta}^\pm / \vec{\alpha}^\pm] P \vee [\vec{\gamma}^\pm / \vec{\alpha}^\pm] P = Q \\ \hline \text{upgrade } \Theta \vdash P \text{ to } \Theta_0 = Q \end{array}$ <p style="text-align: right;">(UPG)</p> </div> <div style="width: 45%;"> <p style="text-align: center;">Least Upper Bound</p> $\begin{array}{c} \frac{\Theta, \vec{\alpha}^\pm, \vec{\beta}^\pm \vdash P_1 \vee P_2 = Q}{\Theta \vdash \exists \vec{\alpha}^\pm. P_1 \vee \exists \vec{\beta}^\pm. P_2 = Q} \quad (\exists^\vee) \\ \frac{}{\Theta \vdash \alpha^+ \vee \alpha^+ = \alpha^+} \quad (\text{VAR}^\vee) \\ \frac{\Theta \vdash \text{nf}(\downarrow N) \stackrel{a}{\approx} \text{nf}(\downarrow M) \ni (\widehat{\Theta}, \vec{P}, \widehat{\tau}_1, \widehat{\tau}_2)}{\Theta \vdash \downarrow N \vee \downarrow M = \exists \vec{\alpha}^\pm. [\vec{\alpha}^\pm / \widehat{\Theta}] \vec{P}} \quad (\downarrow^\vee) \end{array}$ </div> </div>	<div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;"> $\Theta \vdash P_1 \vee P_2 = Q$ </div> <p style="text-align: center;">Least Upper Bound</p>
--	---

Fig. 18. Type Upgrade and Least Upper Bound Algorithms

The Type Upgrade. The type upgrade algorithm uses the least upper bound algorithm as a subroutine. It exploits the idea that the free variables of a positive type Q cannot disappear in its subtypes (see property 1). It means that if a type P has free variables not occurring in P' , then any common supertype of P and P' must not contain these variables either. This way, any supertype of P not containing certain variables $\vec{\alpha}^\pm$ must also be a supertype of $P' = [\vec{\beta}^\pm / \vec{\alpha}^\pm] P$, where $\vec{\beta}^\pm$ are fresh; and vice versa: any common supertype of P and P' does not contain $\vec{\alpha}^\pm$ nor $\vec{\beta}^\pm$.

This way, to find the least supertype of P well-formed in $\Theta_0 = \Theta \setminus \vec{\alpha}^\pm$ (i.e., not containing $\vec{\alpha}^\pm$), we can do the following. First, construct a new type P' by renaming $\vec{\alpha}^\pm$ in P to fresh $\vec{\beta}^\pm$, and second, find the *least upper bound* of P and P' in the appropriate context. However, for reasons of symmetry, in rule (UPG) we employ a different but equivalent approach: we create *two* types P_1 and P_2 constructed by renaming $\vec{\alpha}^\pm$ in P to fresh disjoint variables $\vec{\beta}^\pm$ and $\vec{\gamma}^\pm$ respectively, and then find the least upper bound of P_1 and P_2 .

The Least Upper Bound. The Least Upper Bound algorithm we use operates on *positive* types. This way, the inference rules of the algorithm analyze the three possible shapes of the input types: a variable type, an existential type, and a shifted computation.

Rule (\exists^\vee) covers the case when at least one of the input types is an existential type. In this case, we can simply move the existential quantifiers from both sides to the context, and make a tail-recursive call. However, it is important to make sure that the quantified variables $\vec{\alpha}^\pm$ and $\vec{\beta}^\pm$ are disjoint (i.e., alpha-renaming might be required in the implementation).

Rule (VAR^\vee) applies when both sides are variables. In this case, the common supertype only exists if these variables are the same. And if they are, the common supertypes must be equivalent to this variable.

Rule (\downarrow^\vee) is the most interesting. If both sides are not quantified, and one of the sides is a shift, so must be the other side. However, the set of common upper bounds is not trivial in this case. For example, $\downarrow(\beta^+ \rightarrow \gamma_1^-)$ and $\downarrow(\beta^+ \rightarrow \gamma_2^-)$ have two non-equivalent common supertypes: $\exists \alpha^-. \downarrow \alpha^-$ (by instantiating α^- with $\beta^+ \rightarrow \gamma_1^-$ and $\beta^+ \rightarrow \gamma_2^-$ respectively) and $\exists \alpha^-. \downarrow(\beta^+ \rightarrow \alpha^-)$ (by instantiating α^- with γ_1^- and γ_2^- respectively). As one can see, the second supertype $\exists \alpha^-. \downarrow(\beta^+ \rightarrow \alpha^-)$ is the least among them because it abstracts over a ‘deeper’ negative subexpression.

In general, we must (i) find the most detailed pattern (a type with ‘holes’ at negative positions) that matches both sides, and (ii) abstract over the ‘holes’ by existential quantifiers. The algorithm that finds the most detailed common pattern is called *anti-unification*. As output, it returns $(\hat{\Theta}, \hat{P}, \hat{\tau}_1, \hat{\tau}_2)$, where important for us is \hat{P} —the pattern and $\hat{\Theta}$ —the set of ‘holes’ represented by negative algorithmic variables. We discuss the anti-unification algorithm in detail in the following section.

4.7 Anti-Unification

The anti-unification algorithm [Plotkin 1970; Reynolds 1970], is a procedure dual to unification. For two given (potentially different) expressions, it finds the most specific generalizer—the most detailed pattern that matches both of the input expressions. As evidence, it can also return two substitutions that instantiate the ‘holes’ of the pattern to the input expressions.

In our case, we have to be more demanding on the anti-unification algorithm. Since we use it to construct an existential type, whose (negative) quantified variables can only be instantiated with negative types, we must make sure that the pattern has ‘holes’ only at negative positions. Moreover, we must make sure that the resulting substitutions for the ‘holes’ are well-formed in the initial context, and do not contain variables bound later. For example, the anti-unification of $N_1 = \forall \beta^+. \alpha_1^+ \rightarrow \uparrow \beta^+$ and $N_2 = \forall \beta^+. \alpha_2^+ \rightarrow \uparrow \beta^+$ must result in a ‘hole’, which we model as an algorithmic type variable $\hat{\gamma}^-$, with a pair of substitutions $\hat{\gamma}^- \mapsto N_1$ and $\hat{\gamma}^- \mapsto N_2$. But it cannot be more specific such as $\forall \beta^+. \hat{\gamma}^+ \rightarrow \uparrow \beta^+$ (since the hole cannot be positive) or $\forall \beta^+. \hat{\gamma}^-$ (since the instantiation cannot capture the bound variable β^+).

The algorithm that finds the most specific generalizer of two types under required conditions is given in fig. 19. It consists of two mutually recursive procedures: the positive and the negative anti-unification. As the positive and the negative anti-unification procedures are symmetric in their interface, let us discuss how to read the positive judgment.

The positive anti-unification judgment has form $\Theta \models P_1 \stackrel{a}{\simeq} P_2 \models (\widehat{\Theta}, \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2)$. As an input, it takes a context Θ , in which the ‘holes’ instantiations must be well-formed, and two positive types: P_1 and P_2 ; it returns a tuple of four components: $\widehat{\Theta}$ —a set of ‘holes’ represented by negative algorithmic variables, \underline{Q} —a pattern represented as a positive algorithmic type, whose algorithmic variables are in $\widehat{\Theta}$, and two substitutions $\widehat{\tau}_1$ and $\widehat{\tau}_2$ instantiating the variables from $\widehat{\Theta}$ such that $[\widehat{\tau}_1] \underline{Q} = P_1$ and $[\widehat{\tau}_2] \underline{Q} = P_2$.

$$\begin{array}{c}
 \boxed{\Theta \models P_1 \stackrel{a}{\simeq} P_2 \models (\widehat{\Theta}, \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2)} \\
 \\
 \frac{}{\Theta \models \alpha^+ \stackrel{a}{\simeq} \alpha^+ \models (\cdot, \alpha^+, \cdot, \cdot)} \quad (\text{VAR}_+^a) \\
 \\
 \frac{\Theta \models N_1 \stackrel{a}{\simeq} N_2 \models (\widehat{\Theta}, \underline{M}, \widehat{\tau}_1, \widehat{\tau}_2)}{\Theta \models \downarrow N_1 \stackrel{a}{\simeq} \downarrow N_2 \models (\widehat{\Theta}, \downarrow \underline{M}, \widehat{\tau}_1, \widehat{\tau}_2)} \quad (\downarrow^a) \\
 \\
 \frac{\begin{array}{c} \overrightarrow{\alpha} \cap \Theta = \emptyset \quad \Theta \models P_1 \stackrel{a}{\simeq} P_2 \models (\widehat{\Theta}, \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2) \\ \Theta \models \exists \overrightarrow{\alpha}. P_1 \stackrel{a}{\simeq} \exists \overrightarrow{\alpha}. P_2 \models (\widehat{\Theta}, \exists \overrightarrow{\alpha}. \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2) \end{array}}{\Theta \models \exists \overrightarrow{\alpha}. P_1 \stackrel{a}{\simeq} \exists \overrightarrow{\alpha}. P_2 \models (\widehat{\Theta}, \exists \overrightarrow{\alpha}. \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2)} \quad (\exists^a) \\
 \\
 \boxed{\Theta \models N_1 \stackrel{a}{\simeq} N_2 \models (\widehat{\Theta}, \underline{M}, \widehat{\tau}_1, \widehat{\tau}_2)} \\
 \\
 \frac{}{\Theta \models \alpha^- \stackrel{a}{\simeq} \alpha^- \models (\cdot, \alpha^-, \cdot, \cdot)} \quad (\text{VAR}_-^a) \\
 \\
 \frac{\Theta \models P_1 \stackrel{a}{\simeq} P_2 \models (\widehat{\Theta}, \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2)}{\Theta \models \uparrow P_1 \stackrel{a}{\simeq} \uparrow P_2 \models (\widehat{\Theta}, \uparrow \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2)} \quad (\uparrow^a) \\
 \\
 \frac{\begin{array}{c} \overrightarrow{\alpha^+} \cap \Theta = \emptyset \quad \Theta \models N_1 \stackrel{a}{\simeq} N_2 \models (\widehat{\Theta}, \underline{M}, \widehat{\tau}_1, \widehat{\tau}_2) \\ \Theta \models \forall \overrightarrow{\alpha^+}. N_1 \stackrel{a}{\simeq} \forall \overrightarrow{\alpha^+}. N_2 \models (\widehat{\Theta}, \forall \overrightarrow{\alpha^+}. \underline{M}, \widehat{\tau}_1, \widehat{\tau}_2) \end{array}}{\Theta \models \forall \overrightarrow{\alpha^+}. N_1 \stackrel{a}{\simeq} \forall \overrightarrow{\alpha^+}. N_2 \models (\widehat{\Theta}, \forall \overrightarrow{\alpha^+}. \underline{M}, \widehat{\tau}_1, \widehat{\tau}_2)} \quad (\forall^a) \\
 \\
 \frac{\Theta \models P_1 \stackrel{a}{\simeq} P_2 \models (\widehat{\Theta}_1, \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2) \quad \Theta \models N_1 \stackrel{a}{\simeq} N_2 \models (\widehat{\Theta}_2, \underline{M}, \widehat{\tau}_1', \widehat{\tau}_2')}{\Theta \models P_1 \rightarrow N_1 \stackrel{a}{\simeq} P_2 \rightarrow N_2 \models (\widehat{\Theta}_1 \cup \widehat{\Theta}_2, \underline{Q} \rightarrow \underline{M}, \widehat{\tau}_1 \cup \widehat{\tau}_1', \widehat{\tau}_2 \cup \widehat{\tau}_2')} \quad (\rightarrow^a) \\
 \\
 \frac{\text{if other rules are not applicable} \quad \Theta \vdash N \quad \Theta \vdash M}{\Theta \models N \stackrel{a}{\simeq} M \models (\widehat{\alpha}_{\{N,M\}}^-, \widehat{\alpha}_{\{N,M\}}^-, (\widehat{\alpha}_{\{N,M\}}^- \mapsto N), (\widehat{\alpha}_{\{N,M\}}^- \mapsto M))} \quad (\text{AU})
 \end{array}$$

Fig. 19. Anti-Unification Algorithm

At the high level, the algorithm scheme follows the standard approach [Plotkin 1970] represented as a recursive procedure. Specifically, it follows two principles:

- (i) if the input terms start with the same constructor, we anti-unify the corresponding parts recursively and unite the results. This principle is followed by all the rules except (AU), which works as follows:
- (ii) if the first principle does not apply to the input terms N and M (for instance, if they have different outer constructors), the anti-unification algorithm returns a ‘hole’ such that one

substitution maps it to N and the other maps it to M . This ‘hole’ should have a name uniquely defined by the pair (N, M) , so that it automatically merges with other ‘holes’ mapped to the same pair of types, and thus, the initiality of the generalizer is ensured.

Let us discuss the specific rules of the algorithm in detail.

Variables. Rules (VAR_+^a) and (VAR_-^a) generalize two equal variables. In this case, the resulting pattern is the variable itself, and no ‘holes’ are needed.

Shifts. Rules (\downarrow^a) and (\uparrow^a) operate by congruence: they anti-unify the bodies of the shifts recursively and add the shift constructor back to the resulting pattern.

Quantifiers. Rules (\forall^a) and (\exists^a) are symmetric. They generalize two quantified types congruently, similarly to the shift rules. However, we also require that the quantified variables are fresh, and that the left-hand side variables are equal to the corresponding variables on the right-hand side. To ensure it, alpha-renaming might be required in the implementation.

Notice that the context Θ is *not* extended with the quantified variables. In this algorithm, Θ does not play the role of a current typing context, but rather a snapshot of a context at the moment of calling the anti-unification, i.e., the context in which the instantiations of the ‘holes’ must be well-formed.

Functions. Rule (\rightarrow^a) congruently generalizes two function types. An arrow type is the only binary constructor, and thus, it is the only rule where the union of the anti-unification results is substantial. The interesting is the case when the resulting generalization of the input types and the resulting generalization of the output types have ‘holes’ mapped to the same pair of types. In this case, the algorithm must merge the ‘holes’ into one. For example, the anti-unification of $\downarrow\alpha^- \rightarrow \alpha^-$ and $\downarrow\beta^- \rightarrow \beta^-$ must result in $\downarrow\gamma^- \rightarrow \gamma^-$, rather than $\downarrow\gamma_1^- \rightarrow \gamma_2^-$.

In our representation of the anti-unification algorithm, this ‘merge’ happens automatically: following the rule (AU), the name of the ‘hole’ is uniquely defined by the pair of types it is mapped to. Specifically, when anti-unifying $\downarrow\alpha^- \rightarrow \alpha^-$ and $\downarrow\beta^- \rightarrow \beta^-$ our algorithm returns $\downarrow\hat{\alpha}_{\{\alpha^-, \beta^-\}}^- \rightarrow \hat{\alpha}_{\{\alpha^-, \beta^-\}}^-$, that is a renaming of $\downarrow\gamma^- \rightarrow \gamma^-$.

This way, as the output the rule returns the following tuple:

- $\hat{\Theta}_1 \cup \hat{\Theta}_2$ —a simple union of the sets of ‘holes’ returned from by the recursive calls,
- $\hat{Q} \rightarrow \hat{M}$ —the resulting pattern constructed from the patterns returned recursively.
- $\hat{\tau}_1 \cup \hat{\tau}_1'$ and $\hat{\tau}_2 \cup \hat{\tau}_2'$ — a union (in a relational sense) of the substitutions returned by the recursive calls. It is worth noting that the union is well-defined because the result of the substitution on a ‘hole’ is determined by the name of the ‘hole’.

The Anti-Unification Rule. Rule (AU) is the base case of the anti-unification algorithm. If the congruent rules are not applicable, it means that the input types have a substantially different structure, and thus, the only option is to create a ‘hole’. There are three important aspects of this rule that we would like to discuss.

First, as mentioned earlier, the freshly created ‘hole’ has a name that is uniquely defined by the pair of input types. It is ensured by the following invariant: all the ‘holes’ in the algorithm have name $\hat{\alpha}^-$ indexed by the pair of negative types it is mapped to. This way, the returning set of ‘holes’ is a singleton set $\{\hat{\alpha}_{\{N, M\}}^-\}$; the resulting pattern is the ‘hole’ $\hat{\alpha}_{\{N, M\}}^-$, and the mappings simply send it to the corresponding types: $\hat{\alpha}_{\{N, M\}}^- \mapsto N$ and $\hat{\alpha}_{\{N, M\}}^- \mapsto M$.

Second, this rule is only applicable to negative types, moreover, the input types are checked to be well-formed in the outer context Θ . This is required by the usage of anti-unification: we call it to build an existential type that would be an upper bound of two input types via abstracting some of

their subexpressions under existential quantifiers. The existentials quantify over *negative* variables, and they must be instantiated in the context available at that moment.

Third, the rule is only applicable if all other rules fail. Notice that it could happen even when the input types have matching constructors. For example, the generalizer of $\uparrow\alpha^+$ and $\uparrow\beta^+$ is $\widehat{\gamma}^-$ (with mappings $\widehat{\gamma}^- \mapsto \uparrow\alpha^+$ and $\widehat{\gamma}^- \mapsto \uparrow\beta^+$), rather than $\uparrow\gamma^+$. This way, the algorithm must try to apply the congruent rules first, and only if they fail, apply (AU). This principle makes the inference system not syntax-directed: it is not known a priori (before the recursive call) whether the corresponding congruent rule or rule (AU) will be applied.

4.8 Type Inference

Finally, we present the type inference algorithm. Similarly to the subtyping algorithm, it structurally corresponds to the declarative inference specification, meaning that most of the algorithmic rules have declarative counterparts, with respect to which they are sound and complete.

This way, the inference algorithm also consists of three mutually recursive procedures: the positive type inference, the negative type inference, and the application type inference. As subroutines, the inference algorithm uses subtyping, constraint merge, and normalization. The corresponding graph is shown in fig. 20.

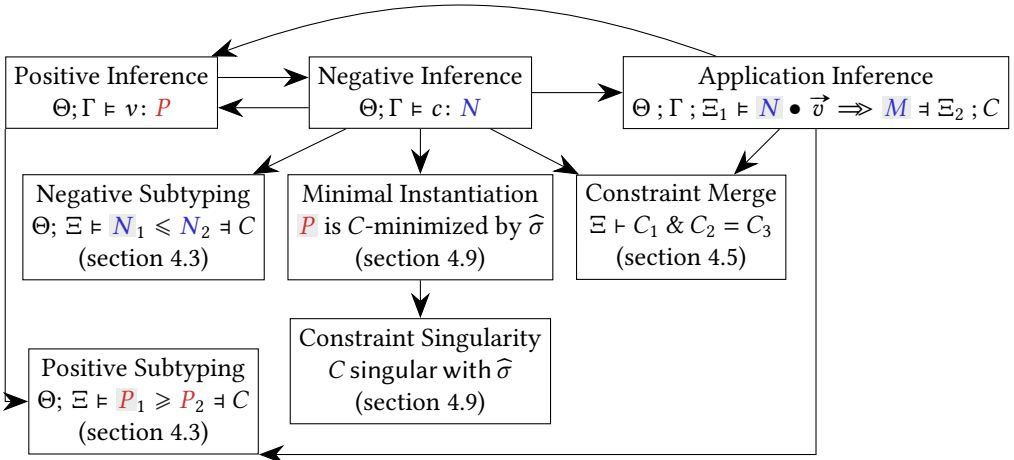


Fig. 20. Dependency graph of the typing algorithm

The positive and the negative type inference judgments have symmetric forms: $\Theta; \Gamma \vdash v : P$ and $\Theta; \Gamma \vdash c : N$. Both of these algorithms take as an input typing context Θ , a variable context Γ , and a term (a value or a computation) taking its type variables from Θ , and term variables from Γ . As an output, they return a type of the given term, which we guarantee to be normalized.

The application type inference judgment has form $\Theta; \Gamma; \Xi_1 \vdash N \bullet \vec{v} \Rightarrow M \vdash \Xi_2; C$. As an input, it takes three contexts: typing context Θ , a variable context Γ , and an instantiation context Ξ_1 . It also takes a head type N and a list of arguments (terms) \vec{v} the head is applied to. The head may contain algorithmic variables specified by Ξ_1 , in other words, $\Theta; \text{dom}(\Xi_1) \vdash N$. As a result, the application inference judgment returns M —a normalized type of the result of the application. Type M may contain new algorithmic variables, and thus, the judgment also returns Ξ_2 —an updated instantiation context and C —a set of subtyping constraints. Together Ξ_2 and C specify how the algorithmic variables must be instantiated.

$\boxed{\Theta; \Gamma \vDash v : P}$ Positive typing

$$\frac{x : P \in \Gamma}{\Theta; \Gamma \vDash x : \text{nf}(P)} \quad (\text{VAR}^{\text{INF}}) \quad \frac{\Theta \vdash Q \quad \Theta; \Gamma \vDash v : P}{\Theta; \Gamma \vDash (v : Q) : \text{nf}(Q)} \quad (\text{ANN}_+^{\text{INF}}) \quad \frac{\Theta; \Gamma \vDash c : N}{\Theta; \Gamma \vDash \{c\} : \downarrow N} \quad (\{\}^{\text{INF}})$$

$\boxed{\Theta; \Gamma \vDash c : N}$ Negative typing

$$\frac{\Theta \vdash M \quad \Theta; \Gamma \vDash c : N}{\Theta; \Gamma \vDash (c : M) : \text{nf}(M)} \quad (\text{ANN}_-^{\text{INF}}) \quad \frac{\Theta; \Gamma \vDash v : P \quad \Theta; \Gamma, x : P \vDash c : N}{\Theta; \Gamma \vDash \text{let } x = v; c : N} \quad (\text{LET}^{\text{INF}})$$

$$\frac{\Theta \vdash P \quad \Theta; \Gamma, x : P \vDash c : N}{\Theta; \Gamma \vDash \lambda x : P. c : \text{nf}(P \rightarrow N)} \quad (\lambda^{\text{INF}}) \quad \frac{\Theta \vdash P \quad \Theta; \Gamma \vDash c : M}{\Theta; \Gamma \vDash M \leq \uparrow P \dashv \cdot} \quad \frac{\Theta; \Gamma, x : P \vDash c' : N}{\Theta; \Gamma \vDash \text{let } x : P = c; c' : N} \quad (\text{LET}_c^{\text{INF}})$$

$$\frac{\Theta, \alpha^+; \Gamma \vDash c : N}{\Theta; \Gamma \vDash \Lambda \alpha^+. c : \text{nf}(\forall \alpha^+. N)} \quad (\Lambda^{\text{INF}}) \quad \frac{\Theta; \Gamma \vDash v : \exists \alpha^-. P}{\Theta, \alpha^-; \Gamma, x : P \vDash c : N} \quad \frac{\Theta \vdash N}{\Theta; \Gamma \vDash \text{let}^\exists(\alpha^-, x) = v; c : N} \quad (\text{LET}_\exists^{\text{INF}})$$

$$\frac{\Theta; \Gamma \vDash v : P}{\Theta; \Gamma \vDash \text{return } v : \uparrow P} \quad (\text{RET}^{\text{INF}})$$

$$\frac{\Theta \vdash P \quad \Theta; \Gamma \vDash v : \downarrow M \quad \Theta; \Gamma; \cdot \vDash M \bullet \vec{v} \Rightarrow M' \dashv \Xi; C_1}{\Theta; \Xi \vDash M' \leq \uparrow P \dashv C_2 \quad \Xi \vdash C_1 \& C_2 = C \quad \Theta; \Gamma, x : P \vDash c : N} \quad (\text{LET}_{\text{@@}}^{\text{INF}})$$

$$\frac{\Theta; \Gamma \vDash v : \downarrow M \quad \Theta; \Gamma; \cdot \vDash M \bullet \vec{v} \Rightarrow \uparrow Q \dashv \Xi; C}{Q \text{ is } C\text{-minimized by } \widehat{\sigma} \quad \Theta; \Gamma, x : [\widehat{\sigma}] Q \vDash c : N} \quad (\text{LET}_{\text{@@}}^{\text{INF}})$$

$\boxed{\Theta; \Gamma; \Xi_1 \vDash N \bullet \vec{v} \Rightarrow M \dashv \Xi_2; C}$ Application typing

$$\frac{}{\Theta; \Gamma; \Xi \vDash N \bullet \Rightarrow \text{nf}(N) \dashv \Xi; \cdot} \quad (\emptyset^{\text{INF}}) \quad \frac{\Theta; \Gamma \vDash v : P \quad \Theta; \Xi \vDash Q \geq P \dashv C_1}{\Theta; \Gamma; \Xi \vDash N \bullet \vec{v} \Rightarrow M \dashv \Xi'; C_2} \quad \frac{\Xi \vdash C_1 \& C_2 = C}{\Theta; \Gamma; \Xi \vDash Q \rightarrow N \bullet v, \vec{v} \Rightarrow M \dashv \Xi'; C} \quad (\rightarrow^{\text{INF}})$$

$$\frac{\Theta; \Gamma; \Xi, \vec{\alpha}^+ \{ \Theta \} \vDash [\vec{\alpha}^+ / \alpha^+] N \bullet \vec{v} \Rightarrow M \dashv \Xi'; C}{\vec{\alpha}^+ \text{ are fresh} \quad \vec{v} \neq \vec{\alpha}^+ \neq \cdot} \quad (\forall^{\text{INF}})$$

$$\Theta; \Gamma; \Xi \vDash \forall \alpha^+. N \bullet \vec{v} \Rightarrow M \dashv \Xi'; C|_{\text{fav}(N) \cup \text{fav}(M)}$$

Fig. 21. Algorithmic Type Inference

The inference rules are shown in fig. 21. Next, we discuss them in detail.

Variables. Rule (VAR^{INF}) infers the type of a positive variable by looking it up in the term variable context and normalizing the result.

Annotations. Rules ($\text{ANN}_+^{\text{INF}}$) and ($\text{ANN}_-^{\text{INF}}$) are symmetric. First, they check that the annotated type is well-formed in the given context Θ . Then they make a recursive call to infer the type of annotated expression, check that the inferred type is a subtype of the annotation, and return the normalized annotation.

Abstractions. Rule (λ^{INF}) infers the type of a lambda abstraction. It checks the well-formedness of the annotation P , makes a recursive call to infer the type of the body in the extended context, and returns the corresponding arrow type. Since the annotation P is allowed to be non-normalized, the rule also normalizes the resulting type.

Rule (Λ^{INF}) infers the type of a big lambda. Similarly to the previous case, it makes a recursive call to infer the type of the body in the extended *type* context. After that, it returns the corresponding universal type. It is also required to normalize the result. For instance, if α^+ does not occur in the body of the lambda, the corresponding \forall will be removed.

Return and Thunk. Rules ($\{\}^{\text{INF}}$) and (RET^{INF}) are similar to the declarative rules: they make a recursive call to type the body of the thunk or the return expression and put the shift on top of the result.

Unpack. Rule ($\text{LET}_{\exists}^{\text{INF}}$) allows one to unpack an existential type. First, it infers the existential type $\exists \vec{\alpha}. P$ of the value being unpacked, and since the type is guaranteed to be normalized, binds the quantified variables with $\vec{\alpha}$. Then it infers the type of the body in the appropriately extended context and checks that the inferred type does not depend on $\vec{\alpha}$ by checking well-formedness $\Theta \vdash N$.

Let Binders. Rule (LET^{INF}) represents the type inference of a standard let binder. It infers the type of the bound value v , and makes a recursive call to infer the type of the body in the extended context.

Rule ($\text{LET}_{@}^{\text{INF}}$) infers a type of *annotated* applicative let binder. First, it infers the type of the head of the application, ensuring that it is a *thunked computation* $\downarrow M$. After that, it makes a recursive call to the application inference procedure, returning an algorithmic type M' , that must be a subtype of the annotation $\uparrow P$.

Then premise $\Theta; \Xi \models M' \leq \uparrow P \models C_2$ together with $\Xi \vdash C_1 \ \& \ C_2 = C$ check that M' can be instantiated to the annotated type $\uparrow P$, and if it is, the algorithm infers the type of the body in the extended context, and returns it as the result.

Rule ($\text{LET}_{@}^{\text{INF}}$) works similarly to ($\text{LET}_{@}^{\text{INF}}$). However, since no annotation is given, the algorithm must ensure that the inferred Q has the ‘canonical’ minimal instantiation. To find it, it makes a call to the minimal instantiation algorithm (section 4.9) that finds the substitution that satisfies the inferred constraints C and instantiates Q to the minimal (among other such instantiations) type $[\sigma] Q$.

Application to an Empty List of Arguments. Rule ($\emptyset_{\bullet \Rightarrow}^{\text{INF}}$) is the base case of application inference. If the list of applied arguments is empty, the inferred type is the type of the head, and the algorithm returns it after normalizing.

Application of a Polymorphic Type \forall . Rule ($\forall_{\bullet \Rightarrow}^{\text{INF}}$), analogously to the declarative case, is the rule ensuring the implicit elimination of the universal quantifiers. This is the place where the

algorithmic variables are introduced. The algorithm simply replaces the quantified variables $\vec{\alpha}^+$ with fresh algorithmic variables $\vec{\alpha}^+$, and makes a recursive call in the extended context.

To ensure that this step does not cause infinite recursion, we also check that the head type has at least one \forall -quantifier. Also, to force the algorithm to apply rule $(\emptyset_{\bullet \Rightarrow}^{\text{INF}})$ when there are no arguments, we require $\vec{v} \neq \emptyset$.

Application of an Arrow Type. Rule $(\rightarrow_{\bullet \Rightarrow}^{\text{INF}})$ is the main rule of algorithmic application inference. It is applied when the head has an arrow type $Q \rightarrow N$. First, it infers the type of the first argument v , and then, calling the algorithmic subtyping, finds C_1 —the minimal constraint ensuring that Q is a supertype of the type of v . Then it makes a recursive call applying N to the rest of the arguments and merges the resulting constraint with C_1 .

4.9 Minimal Instantiation and Constraint Singularity

Multiple types M can be inferred for a type application: $\Theta; \Gamma \vdash N \bullet \vec{v} \Rightarrow M$ but only one *principal* type should be chosen for a variable in an unannotated let binder. Declaratively, we require the principal type P to be minimal among other all types P' that can be inferred for the application $\Theta; \Gamma \vdash N \bullet \vec{v} \Rightarrow \uparrow P'$. Algorithmically, the inference returns an *algorithmic* type P together with a set of necessary and sufficient constraints C that any instantiation of P must satisfy. This way, to guarantee principality, we must minimize the instantiation of P , while ensuring that it satisfies the given constraints C .

The minimal instantiation judgment ' P is C -minimized by $\hat{\sigma}$ ' means that $\hat{\sigma}$ is the minimum among all the substitutions satisfying C , that is $\hat{\sigma}$ instantiates P to a subtype of any other instantiation of P . The algorithm is given in fig. 22. First, it removes the existential quantifiers by (\exists^{MIN}) . Then it considers two cases. If the type is an algorithmic variable $\hat{\alpha}^+$ restricted by a constraint $(\hat{\alpha}^+ \geq P) \in C$, its minimal instantiation is the type $\text{nf}(P)$, and $(\text{UVar}^{\text{MIN}})$ is applied. Otherwise, $(\text{Sing}^{\text{MIN}})$ is applied, meaning that all the instantiations of the type must be equivalent. The algorithm ensures that by two premises: (i) checking that all the algorithmic variables of the type are restricted by the constraint set; and (ii) building the substitution that satisfies the constraint set on these variables (and simultaneously checking that such substitution is unique up-to-equivalence) using the constraint singularity algorithm.

$$\begin{array}{c}
 \frac{(\hat{\alpha}^+ \geq P) \in C}{\hat{\alpha}^+ \text{ is } C\text{-minimized by } (\text{nf}(P)/\hat{\alpha}^+)} \quad (\text{UVar}^{\text{MIN}}) \qquad \frac{\text{fav}(P) \subseteq \text{dom}(C) \quad C|_{\text{fav}(P)} \text{ singular with } \hat{\sigma}}{P \text{ is } C\text{-minimized by } \hat{\sigma}} \quad (\text{Sing}^{\text{MIN}}) \\
 \frac{P \text{ is } C\text{-minimized by } \hat{\sigma}}{\exists \vec{\alpha}^+. P \text{ is } C\text{-minimized by } \hat{\sigma}} \quad (\exists^{\text{MIN}})
 \end{array}$$

Fig. 22. Minimal Instantiation

The singularity algorithm performs two tasks: it checks that a *constraint set* has a single substitution satisfying it, and if it does, it builds this substitution.

To implement the singularity algorithm, we define a partial function C singular with $\hat{\sigma}$, taking a subtyping constraint C as an argument and returning a substitution $\hat{\sigma}$ —the only solution of C . The constraint C consists of constraint entries. Thus, we define the singularity using the singularity of constraint entries: each entry of C must have a unique instantiation and the resulting substitution $\hat{\sigma}$ must be the combination of these instantiations (see fig. 23).

For a constraint set C , we write C singular with $\hat{\sigma}$ if and only if:

- for any *positive* constraint entry $e \in C$, there exists P such that e singular with P , and for the variable $\hat{\beta}^+$ restricted by e , $[\hat{\sigma}]\hat{\beta}^+ = P$;
- the symmetric property holds for all *negative* $e \in C$; and
- for any $\hat{\alpha}^\pm \notin \text{dom}(C)$, $[\hat{\sigma}]\hat{\alpha}^\pm = \hat{\alpha}^\pm$.

Fig. 23. Constraint Singularity

The singularity of constraint *entries* is defined in fig. 24. The *equivalence* entries are always singular, as the only possible type satisfying them is the one given in the constraint itself, which is reflected in (\simeq_-^{SING}) and (\simeq_+^{SING}) . The *subtyping* constraints are trickier. As will be discussed in section 5.3, variables (and equivalent types such as $\exists \vec{\alpha}^-. \alpha^+$) do not have *proper* supertypes, and thus, the constraints of a form $\hat{\alpha}^+ :\geq \exists \vec{\alpha}^-. \alpha^+$ are singular with the only possible normalized solution α^+ ($:\geq \alpha^{\text{SING}}$). However, if the body of the existential type is guarded by a shift $\exists \vec{\alpha}^-. \downarrow N$, it is singular if and only if N is equivalent to some $\alpha_i^- \in \vec{\alpha}^-$ bound by the quantifier, see the rule $(:\geq \downarrow^{\text{SING}})$.

$$\begin{array}{c}
 \frac{}{\hat{\alpha}^+ :\simeq P \text{ singular with nf } (P)} \quad (\simeq_+^{\text{SING}}) \qquad \frac{}{\hat{\alpha}^+ :\geq \exists \vec{\alpha}^-. \alpha^+ \text{ singular with } \alpha^+} \quad (:\geq \alpha^{\text{SING}}) \\
 \frac{}{\hat{\alpha}^- :\simeq N \text{ singular with nf } (N)} \quad (\simeq_-^{\text{SING}}) \qquad \frac{\text{nf } (N) = \alpha_i^-}{\hat{\alpha}^+ :\geq \exists \vec{\alpha}^-. \downarrow N \text{ singular with } \exists \alpha^-. \downarrow \alpha^-} \quad (:\geq \downarrow^{\text{SING}})
 \end{array}$$

Fig. 24. Singular Constraint Entry

5 ALGORITHM CORRECTNESS

The central results ensuring the correctness of the inference algorithm are its soundness and completeness with respect to the declarative specification. The soundness means the algorithm will always produce a typing *allowed* by the declarative system; Dually, the completeness says that once a term has some type declaratively, the inference algorithm succeeds.

The formulation of soundness and completeness of *application inference* is more complex. Both of them assume that the input head type N is free from *negative* algorithmic variables (it is achieved by polarization invariants preserved by the inference rules). The soundness states that the output of the algorithm— M and C is viable. Specifically, that the constraint set C provides a sufficient set of restrictions that a substitution $\hat{\sigma}$ must satisfy to ensure the *declarative* inference of the output type M , that is $\Theta; \Gamma \vdash [\hat{\sigma}]N \bullet \vec{v} \Rightarrow [\hat{\sigma}]M$.

The application inference completeness means that if there exists a substitution $\hat{\sigma}$ and the resulting type M ensuring the declarative inference $\Theta; \Gamma \vdash [\hat{\sigma}]N \bullet \vec{v} \Rightarrow M$ then the algorithm succeeds and gives the most general result M_0 and C_0 . The property of ‘being the most general’ is specified in pt. (2). Intuitively, it means that for any other solution—substitution $\hat{\sigma}$ and the resulting type M , if it ensures the declarative inference, then $\hat{\sigma}$ can be extended in a C_0 -complying way to equate M_0 with M .

Theorem (Soundness of Typing). *Suppose that $\Theta \vdash \Gamma$. Then¹*

$$+ \Theta; \Gamma \models v : P \text{ implies } \Theta; \Gamma \vdash v : P,$$

¹The presented properties hold, but the actual inductive proof requires strengthening of the statement and the corresponding theorem is more involved. See the appendix for details.

- $\Theta; \Gamma \vdash c: N$ implies $\Theta; \Gamma \vdash c: N$,
- $\Theta; \Gamma; \Xi \vdash N \bullet \vec{v} \Rightarrow M \triangleq \Xi'; C$ implies $\Theta; \Gamma \vdash [\hat{\sigma}]N \bullet \vec{v} \Rightarrow [\hat{\sigma}]M$, for any instantiation of $\hat{\sigma}$ satisfying constraints C , assuming $\Theta \vdash^{\supset} \Xi$ and $\Theta; \text{dom}(\Xi) \vdash N$ and that N is free from negative algorithmic variables.

Theorem (Completeness of Typing). *Suppose that $\Theta \vdash \Gamma$. Then¹*

- + $\Theta; \Gamma \vdash v: P$ implies $\Theta; \Gamma \vdash v: \text{nf}(P)$,
- $\Theta; \Gamma \vdash c: N$ implies $\Theta; \Gamma \vdash c: \text{nf}(N)$,
- If $\Theta; \Gamma \vdash [\hat{\sigma}]N \bullet \vec{v} \Rightarrow M$ where (1) $\Theta \vdash^{\supset} \Xi$, (2) $\Theta \vdash M$, (3) $\Theta; \text{dom}(\Xi) \vdash N$ (free from negative algorithmic variables), and (4) $\Xi \vdash \hat{\sigma} : \text{fav}(N)$, then there exist M_0, Ξ_0 , and C_0 such that
 - (1) $\Theta; \Gamma; \Xi \vdash N \bullet \vec{v} \Rightarrow M_0 \triangleq \Xi_0; C_0$ and
 - (2) for any other $\hat{\sigma}$ and M (where $\Xi \vdash \hat{\sigma} : \text{fav}(N)$ and $\Theta \vdash M$) such that $\Theta; \Gamma \vdash [\hat{\sigma}]N \bullet \vec{v} \Rightarrow M$, there exists $\hat{\sigma}'$ such that (a) $\Xi_0 \vdash \hat{\sigma}' : \text{fav}N \cup \text{fav}M_0$ and $\Xi_0 \vdash \hat{\sigma}' : C_0$, (b) $\Xi \vdash \hat{\sigma}' \simeq^{\leq} \hat{\sigma} : \text{fav}N$, and (c) $\Theta \vdash [\hat{\sigma}']M_0 \simeq^{\leq} M$.

The proof of soundness and completeness result is done gradually for all the subroutines, following the structure of the algorithm (figs. 13 and 20) bottom-up. Next, we discuss the main of these results.

5.1 Normalization

The point of type normalization is factoring out non-trivial equivalence by selecting a representative from each equivalence class. This way, the correctness of normalization means that checking for equivalence of two types is the same as checking for equality of their normal forms.

Lemma (Normalization Correctness). *Assuming all types are well-formed in Θ , we have $\Theta \vdash N \simeq^{\leq} M \iff \text{nf}(N) = \text{nf}(M)$ and $\Theta \vdash P \simeq^{\leq} Q \iff \text{nf}(P) = \text{nf}(Q)$.*

To prove the correctness of normalization, we introduce an *intermediate* relation on types—*declarative equivalence* (the notation is $N \simeq^D M$ and $P \simeq^D Q$). In contrast to $\Theta \vdash N \simeq^{\leq} M$ (which means mutual subtyping), $N \simeq^D M$ does not depend on subtyping judgments, but explicitly allows quantifier reordering and redundant quantifier removal. Then the statement $\Theta \vdash N \simeq^{\leq} M \iff \text{nf}(N) = \text{nf}(M)$ (as well as its positive counterpart) is split into two steps: $\Theta \vdash N \simeq^{\leq} M \iff N \simeq^D M \iff \text{nf}(N) = \text{nf}(M)$.

5.2 Anti-Unification

The anti-unifier of the two types is the most specific pattern that matches against both of them. In our setting, the anti-unifiers are restricted further: first, the pattern might only contain placeholders at *negative* positions (because eventually, the placeholders become variables abstracted by the existential quantifier); second, the anti-unification is parametrized with a context Θ , in which the pattern instantiations must be well-formed.

This way, the correctness properties of the anti-unification algorithm are refined accordingly. The soundness of anti-unification not only ensures that the resulting pattern matches with the input types, but also that the pattern instantiations are well-formed in the corresponding context, and that all the ‘placeholder’ variables are negative. In turn, completeness states that if there exists a solution satisfying the soundness criteria, then the algorithm succeeds.

The correctness properties are formulated by the following lemmas. For brevity, we only provide the statements for the positive case, since the negative case is symmetric.

Lemma (Soundness of (Positive) Anti-Unification). *Assuming P_1 and P_2 are normalized, if $\Theta \models P_1 \stackrel{a}{\approx} P_2 \models (\hat{\Theta}, \underline{Q}, \hat{\tau}_1, \hat{\tau}_2)$ then (1) $\Theta; \hat{\Theta} \vdash \underline{Q}$, (2) $\Theta; \cdot \vdash \hat{\tau}_i : \hat{\Theta}$ for $i \in \{1, 2\}$ are anti-unification substitutions (in particular, $\hat{\Theta}$ contains only negative algorithmic variables), and (3) $[\hat{\tau}_i] \underline{Q} = P_i$ for $i \in \{1, 2\}$.*

Lemma (Completeness of (Positive) Anti-Unification). *Assuming that P_1 and P_2 are normalized terms well-formed in Θ and there exist $(\hat{\Theta}', \underline{Q}', \hat{\tau}_1', \hat{\tau}_2')$ such that (1) $\Theta; \hat{\Theta}' \vdash \underline{Q}'$, (2) $\Theta; \cdot \vdash \hat{\tau}_i' : \hat{\Theta}'$ for $i \in \{1, 2\}$ are anti-unification substitutions, and (3) $[\hat{\tau}_i'] \underline{Q}' = P_i$ for $i \in \{1, 2\}$.*

Then the anti-unification algorithm terminates, that is there exists $(\hat{\Theta}, \underline{Q}, \hat{\tau}_1, \hat{\tau}_2)$ such that $\Theta \models P_1 \stackrel{a}{\approx} P_2 \models (\hat{\Theta}, \underline{Q}, \hat{\tau}_1, \hat{\tau}_2)$.

The anti-unification is used to find the least upper bound (LUB). To guarantee that the result of the LUB algorithm is indeed the least, we also need to show the corresponding property of anti-unification. We prove that the anti-unifier that the algorithm provides is the most specific (or the most ‘detailed’). Specifically, it means that any other sound anti-unification solution can be ‘refined’ to the result of the algorithm. The ‘refinement’ is represented as an instantiation of the anti-unifier—a substitution $\Theta; \hat{\Theta}_2 \vdash \hat{\rho} : \hat{\Theta}_1$ replacing the placeholders $\hat{\Theta}_1$ with types that themselves can contain placeholders from $\hat{\Theta}_2$.

Lemma (Initiality of Anti-Unification). *Assume that P_1 and P_2 are normalized types well-formed in Θ , and the anti-unification algorithm succeeds: $\Theta \models P_1 \stackrel{a}{\approx} P_2 \models (\hat{\Theta}, \underline{Q}, \hat{\tau}_1, \hat{\tau}_2)$. Then $(\hat{\Theta}, \underline{Q}, \hat{\tau}_1, \hat{\tau}_2)$ is more specific than any other sound anti-unifier $(\hat{\Theta}', \underline{Q}', \hat{\tau}_1', \hat{\tau}_2')$, i.e., if (1) $\Theta; \hat{\Theta}' \vdash \underline{Q}'$, (2) $\Theta; \cdot \vdash \hat{\tau}_i' : \hat{\Theta}'$ for $i \in \{1, 2\}$, and (3) $[\hat{\tau}_i'] \underline{Q}' = P_i$ for $i \in \{1, 2\}$ then there exists a ‘refining’ substitution $\hat{\rho}$ such that $\Theta; \hat{\Theta} \vdash \hat{\rho} : (\hat{\Theta}'|_{\text{fav } \underline{Q}'})$ and $[\hat{\rho}] \underline{Q}' = \underline{Q}$.*

To prove the correctness properties of the anti-unification algorithm, one extra observation is essential. The algorithm relies on the fact that in the resulting tuple $(\hat{\Theta}, \underline{Q}, \hat{\tau}_1, \hat{\tau}_2)$, there are no two different placeholders $\hat{\beta}^-$ mapped to the same pair of types by $\hat{\tau}_1$ and $\hat{\tau}_2$. This is important to guarantee that, for example, the anti-unifier of $\downarrow \uparrow \text{Int} \rightarrow \uparrow \text{Int}$ and $\downarrow \uparrow \text{Bool} \rightarrow \uparrow \text{Bool}$ is $\downarrow \hat{\alpha}^- \rightarrow \hat{\alpha}^-$, but not (less specific) $\downarrow \hat{\alpha}^- \rightarrow \hat{\beta}^-$. To preserve this property, we arrange the algorithm in such a way that the name of the placeholder is determined by the types it is mapped to. The following lemma specifies this observation.

Lemma (Uniqueness of Anti-unification Variable Names). *Names of the anti-unification variables are uniquely defined by the types they are mapped to by the resulting substitutions. Assuming P_1 and P_2 are normalized, if $\Theta \models P_1 \stackrel{a}{\approx} P_2 \models (\hat{\Theta}, \underline{Q}, \hat{\tau}_1, \hat{\tau}_2)$ then for any $\hat{\beta}^- \in \hat{\Theta}$, $\hat{\beta}^- = \hat{\alpha}^-_{\{[\hat{\tau}_1] \hat{\beta}^-, [\hat{\tau}_2] \hat{\beta}^-\}}$.*

5.3 Least Upper Bound and Upgrade

The Least Upper Bound algorithm finds the least type that is a supertype of two given types. The *soundness* means that the returned type is indeed a supertype of the given ones; the *completeness* means that the algorithm succeeds if the least upper bound exists; and the *initiality* means that the returned type is the least among common supertypes.

Lemma (Least Upper Bound Soundness). *For types $\Theta \vdash P_1$, and $\Theta \vdash P_2$, if $\Theta \models P_1 \vee P_2 = \underline{Q}$ then*

- (i) $\Theta \vdash \underline{Q}$
- (ii) $\Theta \vdash \underline{Q} \geq P_1$ and $\Theta \vdash \underline{Q} \geq P_2$

Lemma (Least Upper Bound Completeness and Initiality). *For types $\Theta \vdash P_1$, $\Theta \vdash P_2$, and $\Theta \vdash \underline{Q}$ such that $\Theta \vdash \underline{Q} \geq P_1$ and $\Theta \vdash \underline{Q} \geq P_2$, there exists \underline{Q}' s.t. $\Theta \models P_1 \vee P_2 = \underline{Q}'$ and $\Theta \vdash \underline{Q} \geq \underline{Q}'$.*

The key observation that allows us to prove these properties is the characterization of positive supertypes. The following lemma justifies the cases of the Least Upper Bound algorithm (section 4.6). In particular, it establishes the correspondence between the upper bounds of shifted types $\downarrow M$ and patterns fitting M (represented by existential types), which explains the usage of anti-unification as a way to find a common pattern.

Lemma (Characterization of Normalized Supertypes). *For a normalized positive type P well-formed in Θ , the set of normalized Θ -formed supertypes of P is the following:*

- if P is a variable β^+ , its only normalized supertype is β^+ itself;
- if P is an existential type $\exists \beta^-. P'$ then its Θ -formed supertypes are the $(\Theta, \vec{\beta}^-)$ -formed supertypes of P' not using β^- ;
- if P is a downshift type $\downarrow M$, its supertypes have form $\exists \vec{\alpha}^-. \downarrow M'$ such that there exists a Θ -formed instantiation of $\vec{\alpha}^-$ in $\downarrow M'$ that makes $\downarrow M'$ equal to $\downarrow M$, i.e., $[\vec{N}/\vec{\alpha}^-] \downarrow M' = \downarrow M$.

Similarly to the Least Upper Bound algorithm, the Upgrade finds the least type among upper bounds (this time the set of considered upper bounds consists of supertypes well-formed in a *smaller* context). This way, we also use the supertype characterization to prove the following properties of the Upgrade algorithm.

Lemma (Upgrade Soundness). *Assuming P is well-formed in $\Theta = \Theta_0, \vec{\alpha}^\pm$, if upgrade $\Theta \vdash P$ to $\Theta_0 = Q$ then (1) $\Theta_0 \vdash Q$ (2) $\Theta \vdash Q \geq P$*

Lemma (Upgrade Completeness). *Assuming P is well-formed in $\Theta = \Theta_0, \vec{\alpha}^\pm$, for any Q' such that Q' is a Θ_0 -formed upper bound of P , i.e., (1) $\Theta_0 \vdash Q'$ and (2) $\Theta \vdash Q' \geq P$, there exists Q such that (upgrade $\Theta \vdash P$ to $\Theta_0 = Q$) and $\Theta_0 \vdash Q' \geq Q$.*

5.4 Subtyping

As for other properties, the correctness of subtyping means that the algorithm produces a valid result (soundness) whenever it exists (completeness). In fig. 14, one can see that the positive and negative subtyping relations are not *mutually* recursive: negative subtyping algorithm uses the positive subtyping, but not vice versa. Because of that, the inductive proofs of the *positive* subtyping correctness are done independently.

The soundness of positive subtyping states that the output constraint C provides a sufficient set of restrictions. In other words, any substitution satisfying C , ensures the desired declarative subtyping: $\Theta \vdash [\widehat{\sigma}] P \geq Q$. Notice that the soundness requires that only the left-hand side input type (P) can have algorithmic variables. This is one of the invariants of the algorithm that significantly simplifies the unification and constraint resolution.

Lemma (Soundness of the Positive Subtyping). *If $\Theta \vdash^\geq \Xi$ and $\Theta \vdash Q$ and $\Theta; \text{dom}(\Xi) \vdash P$ and $\Theta; \Xi \models P \geq Q \models C$, then $\Xi \vdash C : \text{fav } P$ and for any $\widehat{\sigma}$ such that $\Xi \vdash \widehat{\sigma} : C$, we have $\Theta \vdash [\widehat{\sigma}] P \geq Q$.*

The completeness of subtyping says that if the substitution ensuring the declarative subtyping exists, then the algorithm terminates.

Lemma (Completeness of the Positive Subtyping). *Suppose that $\Theta \vdash^\geq \Xi$, $\Theta \vdash Q$ and $\Theta; \text{dom}(\Xi) \vdash P$. Then if there exists $\widehat{\sigma}$ such that $\Xi \vdash \widehat{\sigma} : \text{fav}(P)$ and $\Theta \vdash [\widehat{\sigma}] P \geq Q$, then the subtyping algorithm succeeds: $\Theta; \Xi \models P \geq Q \models C$.*

After the correctness properties of the positive subtyping are established, they are used to prove the correctness of the negative subtyping. The soundness is formulated symmetrical to the positive case, however, the completeness requires an additional invariant to be preserved. The algorithmic

input type N must be free from *negative* algorithmic variables. In particular, it ensures that the constraint restricting a *negative* algorithmic variable will never be generated, and thus, we do not need the Greatest Common Subtype procedure to resolve the constraints.

Lemma (Completeness of the Negative Subtyping). *Suppose that $\Theta \vdash^{\supset} \Xi$ and $\Theta \vdash M$ and $\Theta ; \text{dom}(\Xi) \vdash N$ and N does not contain negative unification variables ($\widehat{\alpha}^- \notin \text{fav} N$). Then for any $\Xi \vdash \widehat{\sigma} : \text{fav}(N)$ such that $\Theta \vdash [\widehat{\sigma}]N \leq M$, the subtyping algorithm succeeds: $\Theta ; \Xi \vdash N \leq M \ni C$.*

5.5 Minimal Instantiation and Singularity

Minimal Instantiation. Algorithmic typing relies on the *minimal instantiation* procedure. For a given positive type P and a set of constraints C it returns a substitution $\Xi \vdash \widehat{\sigma} : C$ such that $[\widehat{\sigma}]P$ is a subtype of all the other instantiations of P satisfying C . This way, the correctness of minimal instantiation is formulated as the following two lemmas.

Lemma (Soundness of Minimal Instantiation). *Suppose that $\Theta \vdash^{\supset} \Xi$, $\Xi \vdash C$, and $\Theta ; \text{dom}(\Xi) \vdash P$. Then P is C -minimized by $\widehat{\sigma}$ implies that $\Xi \vdash \widehat{\sigma} : \text{fav} P$ is a normalized substitution satisfying C (i.e., $\Xi \vdash \widehat{\sigma} : C$) and for any other substitution $\Xi \vdash \widehat{\sigma}' : \text{fav} P$ satisfying C , $\Theta \vdash [\widehat{\sigma}']P \geq [\widehat{\sigma}]P$ holds*

Lemma (Completeness of Minimal Instantiation). *Suppose that $\Theta \vdash^{\supset} \Xi$, $\Xi \vdash C$, $\Theta ; \text{dom}(\Xi) \vdash P$, and there exists $\Xi \vdash \widehat{\sigma} : \text{fav} P$ satisfying C ($\Xi \vdash \widehat{\sigma} : C$) such that for any other $\Xi \vdash \widehat{\sigma}' : \text{fav} P$ satisfying C , we have $\Theta \vdash [\widehat{\sigma}']P \geq [\widehat{\sigma}]P$. Then the minimal instantiation procedure succeeds, specifically, P is C -minimized by $\text{nf}(\widehat{\sigma})$ holds.*

The soundness is rather straightforwardly proved by induction on the judgment inference tree, relying on the soundness of the singularity procedure. The proof of completeness is done by induction on the structure of the type P . It uses the fact that $\Theta \vdash \exists \vec{\alpha}^{\supset} . [\widehat{\sigma}']P \geq \exists \vec{\alpha}^{\supset} . [\widehat{\sigma}]P$ implies $\Theta, \vec{\alpha}^{\supset} \vdash [\widehat{\sigma}']P \geq [\widehat{\sigma}]P$.

Singularity. The soundness of singularity states that C singular with $\widehat{\sigma}$ implies that any substitution satisfying C is equivalent to $\widehat{\sigma}$ on the domain. The completeness states that if all the C -compliant substitutions are equivalent, then the singularity procedure succeeds.

Lemma (Soundness of Singularity). *Suppose $\Xi \vdash C : \widehat{\Theta}$, and C singular with $\widehat{\sigma}$. Then $\Xi \vdash \widehat{\sigma} : \widehat{\Theta}$, $\Xi \vdash \widehat{\sigma} : C$ and for any $\widehat{\sigma}'$ such that $\Xi \vdash \widehat{\sigma} : C$, $\Xi \vdash \widehat{\sigma}' \simeq^{\leq} \widehat{\sigma} : \widehat{\Theta}$.*

Lemma (Completeness of Singularity). *For a given $\Xi \vdash C$, suppose that all the substitutions satisfying C are equivalent on $\text{dom}(C)$. In other words, suppose that there exists $\Xi \vdash \widehat{\sigma}_1 : \text{dom}(C)$ such that for any $\Xi \vdash \widehat{\sigma} : \text{dom}(C)$, $\Xi \vdash \widehat{\sigma} : C$ implies $\Xi \vdash \widehat{\sigma} \simeq^{\leq} \widehat{\sigma}_1 : \text{dom}(C)$. Then¹ C singular with $\widehat{\sigma}_0$ for some $\widehat{\sigma}_0$.*

5.6 Typing

Finally, we discuss the proofs of the soundness and completeness of type inference algorithm that we stated at the beginning of this section. There are three subtleties that we will cover that are important for the proof to go through: the determinacy of the algorithm, the mutual dependence of the soundness and completeness proofs, and the non-trivial inductive metric that we use.

Determinacy. One of the properties that our proof relies on is the determinacy of the typing algorithm: the output (the inferred type) is uniquely determined by the input (the term and the contexts). Determinacy is not hard to demonstrate by structural induction: in every algorithmic inference system, only one inference rule can be applied for a given input. However, we need to prove it for *every* subroutine of the algorithm. Ultimately, it requires the determinacy of such procedures as *generation of fresh variables*, which is easy to ensure, but must be taken into account in the implementation.

Lemma (Determinacy of the Typing Algorithm). *Suppose that $\Theta \vdash \Gamma$ and $\Theta \vdash^2 \Xi$. Then*

- + *If $\Theta; \Gamma \vdash v: P$ and $\Theta; \Gamma \vdash v: P'$ then $P = P'$.*
- *If $\Theta; \Gamma \vdash c: N$ and $\Theta; \Gamma \vdash c: N'$ then $N = N'$.*
- *If $\Theta; \Gamma; \Xi \vdash N \bullet \vec{v} \Rightarrow M \dashv \Xi'; C$ and $\Theta; \Gamma; \Xi \vdash N \bullet \vec{v} \Rightarrow M' \dashv \Xi'; C'$ then $M = M'$, $\Xi = \Xi'$, and $C = C'$.*

Mutuality of the Soundness and Completeness Proofs. Typically in our inductive proofs, the soundness is proven before completeness, as the completeness requires certain properties of the premise subtrees. However, in the case of the typing algorithm, the soundness and completeness proofs cannot be separated: the inductive proof of one requires another, and vice versa.

The soundness proof can be viewed as a mapping from an algorithmic tree to a declarative one. We show that each algorithmic inference rule can be transformed into the corresponding declarative rule, as long as the premises are transformed accordingly, and apply the induction principle.

The soundness requires completeness in the case of $(\text{LET}_{\Theta}^{\text{INF}})$. To prove the soundness, in other words, to transform this rule into its declarative counterpart $(\text{LET}_{\Theta}^{\text{INF}})$, one needs to prove *uniqueness* of the inferred application type $\uparrow Q$. To show the uniqueness, we need to conclude that any other *declarative* tree (which is given by the induction hypothesis) infers an equivalent type. The only way to do so is by applying the soundness of singularity (section 5.5), however, first, the declarative tree must be converted to the algorithmic one (by completeness!). This way, the soundness and completeness proofs are mutually dependent.

Inductive Metric. The soundness and completeness lemmas are proven by *mutual* induction. Since the declarative and the algorithmic systems do not depend on each other, we must introduce a uniform *metric* on which the induction is conducted. We define the metric gradually, starting from the auxiliary function—the size of a judgment $\text{size}(J)$.

The size of *declarative* and *algorithmic* judgments is defined as a pair: the first component is the syntactic size of the terms used in the judgment. The second component depends on the kind of judgment. For regular type inference judgments (such as $\Theta; \Gamma \vdash v: P$ or $\Theta; \Gamma \vdash c: N$), it is always zero. For application inference judgments ($\Theta; \Gamma \vdash N \bullet \vec{v} \Rightarrow M$ or $\Theta; \Gamma; \Xi \vdash N \bullet \vec{v} \Rightarrow M \dashv \Xi'; C$), it is equal to the number of prenex quantifiers of the head type N . We need this adjustment to ensure the monotonicity of the metric, since the rules $(\forall_{\bullet}^{\text{INF}})$ and $(\forall_{\Rightarrow}^{\text{INF}})$ only reduce the quantifiers in the head type but not change the list of arguments.

Definition (Judgement Size). *For a declarative or an algorithmic typing judgment J , we define a metric $\text{size}(J)$ as a pair of integers in the following way:*

- + $\text{size}(\Theta; \Gamma \vdash v: P) = (\text{size}(v), 0);$ + $\text{size}(\Theta; \Gamma \vdash v: P) = (\text{size}(v), 0);$
- $\text{size}(\Theta; \Gamma \vdash c: N) = (\text{size}(c), 0);$ – $\text{size}(\Theta; \Gamma \vdash c: N) = (\text{size}(c), 0);$
- $\text{size}(\Theta; \Gamma \vdash N \bullet \vec{v} \Rightarrow M) = (\text{size}(\vec{v}), \text{npq}(N));$ • $\text{size}(\Theta; \Gamma; \Xi \vdash N \bullet \vec{v} \Rightarrow M \dashv \Xi'; C) = (\text{size}(\vec{v}), \text{npq}(N)).$

Here $\text{size}(v)$ and $\text{size}(c)$ is the size of the syntax tree of the term, and $\text{size}(\vec{v})$ is the sum of sizes of the terms in \vec{v} ; and $\text{npq}(N)$ and $\text{npq}(P)$ represent the number of prenex quantifiers, i.e.,

$$\begin{aligned} &+ \text{npq}(\exists \vec{\alpha}^{\rightarrow}. P) = |\vec{\alpha}^{\rightarrow}|, \text{ if } P \neq \exists \vec{\beta}^{\rightarrow}. P', \\ &- \text{npq}(\forall \vec{\alpha}^{\rightarrow}. N) = |\vec{\alpha}^{\rightarrow}|, \text{ if } N \neq \forall \vec{\beta}^{\rightarrow}. N'. \end{aligned}$$

Notice that for *algorithmic* inference system, $\text{size}(J)$ decreases in all the inductive steps, i.e., for each inference rule, the size of the premise judgments is strictly less than the size of the conclusion. However, the *declarative* inference system has rules $(\simeq_{-}^{\text{INF}})$ and $(\simeq_{+}^{\text{INF}})$, that ‘step to’ an equivalent type, and thus, technically, might keep the judgment unchanged altogether.

To deal with this issue, we introduce the metric on the entire *inference trees* rather than on judgments, and plug into this metric the parameter that certainly decreases in rules (\simeq_-^{INF}) and (\simeq_+^{INF}) —the number of these rules in the inference tree. We denote this number as $\text{eq_nodes}(T)$. Then the final metric is defined as a pair in the following way.

Definition (Inference Tree Metric). *For a tree T , inferring a declarative or an algorithmic judgement J , we define $\text{metric}(T)$ as follows:*

$$\text{metric}(T) = \begin{cases} (\text{size}(J), \text{eq_nodes}(T)) & \text{if } J \text{ represents a declarative judgement,} \\ (\text{size}(J), 0) & \text{if } J \text{ represents an algorithmic judgement.} \end{cases}$$

Here $\text{eq_nodes}(T)$ is the number of nodes in T labeled with (\simeq_+^{INF}) or (\simeq_-^{INF}) .

This metric is suitable for mutual induction on the soundness and completeness of the typing algorithm. First, it is monotonous w.r.t. the inference rules, and this way, we can always apply the induction hypothesis to premises of each rule. Second, the induction hypothesis is powerful enough, so we can use the completeness of the algorithm in the soundness proof, where required. For instance, to prove the soundness of typing in case of $\Theta; \Gamma \vdash \text{let } x = v(\vec{v}); c' : N$, we can assume, that *completeness* holds for a term of shape $\Theta; \Gamma \vdash M \bullet \vec{v} \Rightarrow K$, since $\text{size}(\text{args}) < \text{size}(\text{let } x = v(\text{args}); c')$. This is exactly what allows us to deal with the case of $(\text{LET}_{\oplus}^{\text{INF}})$, because then we can conclude that the inferred type (of a declarative judgment $\Theta; \Gamma \vdash M \bullet \vec{v} \Rightarrow \uparrow[\hat{\sigma}] Q$ constructed by the induction hypothesis) is unique.

6 EXTENSIONS AND FUTURE WORK

In this section, we discuss several extensions to the system and the algorithm. Some of them can be incorporated into the system immediately, while others are beyond the scope of this work, and thus are left for future research.

6.1 Explicit Type Application

In our system, all type applications are inferred implicitly: the algorithm automatically instantiates the variables abstracted by \forall . The implicit type application can be added to the declarative system by the following rule:

$$\frac{\Theta; \Gamma \vdash c : \forall \alpha^+. N}{\Theta; \Gamma \vdash c\{\vec{P}\} : [\vec{P}/\alpha^+]N} \quad (TApp^{\text{INF}})$$

However, this rule alone would cause ambiguity. The declarative system does not fix the order of the quantifiers, which means that $\forall \alpha^+. \forall \beta^+. N$ and $\forall \beta^+. \forall \alpha^+. N$ can be inferred as a type of c interchangeably. But then explicit instantiation ($c\{P\}$) would be ambiguous, as it is unclear whether α^+ or β^+ should be instantiated with P .

Solution 1: Declarative Normalization. One way to resolve this ambiguity is to fix the order of quantifiers. The algorithm already performs the ordering of the quantifiers in the normalization procedure. This way, we could require the inferred type to be normalized to specify the order of the quantifiers:

$$\frac{\Theta; \Gamma \vdash c : \forall \alpha^+. N \quad \text{nf}(\forall \alpha^+. N) = \forall \alpha^+. N}{\Theta; \Gamma \vdash c\{\vec{P}\} : [\vec{P}/\alpha^+]N} \quad (TApp^{\text{INF}})$$

The drawback of this approach is that it would cause the ‘leakage’ of the internal algorithmic concept of type normalization into the ‘surface’ declarative system.

Solution 2: Elementary Type Inference. An alternative approach to provide explicit type application was proposed by [Zhao et al. 2022]. In this work, the subtyping relation is restricted in such a way that $\forall \alpha^+. \forall \beta^+. N$ and $\forall \beta^+. \forall \alpha^+. N$ are *not* mutual subtypes (as long as $\alpha^+ \in \text{fv}(N)$ and $\beta^+ \in \text{fv}(N)$). It implies that the order of the quantifiers of the inferred type is unique, and thus, the explicit type application is unambiguous.

These restrictions can be incorporated into our system by replacing the polymorphic subtyping rules (\forall^\leq) and (\exists^\geq) with the following stronger versions:

$$\frac{\Theta, \vec{\beta}^+ \vdash \sigma : \vec{\alpha}^+ \quad \Theta, \vec{\beta}^+ \vdash [\sigma]N \leq M}{\Theta \vdash \forall \vec{\beta}^+, \vec{\alpha}^+. N \leq \forall \vec{\beta}^+. M} \quad (\text{EV}^\leq) \quad \frac{\Theta, \vec{\beta}^+ \vdash \sigma : \vec{\alpha}^+ \quad \Theta, \vec{\beta}^+ \vdash [\sigma]P \geq Q}{\Theta \vdash \exists \vec{\beta}^+, \vec{\alpha}^+. P \geq \exists \vec{\beta}^+. Q} \quad (\text{E}\exists^\geq)$$

According to these rules, if two polymorphic (\forall or \exists) types are subtypes of each other, they must have the same top-level quantifiers. Moreover, the equivalence on types (mutual subtyping) degenerates to *equality* up to alpha-conversion.

To accommodate these changes in the *algorithm*, it suffices to (i) replace the normalization procedure with identity: $\text{nf}(N) \stackrel{\text{def}}{=} N$, $\text{nf}(P) \stackrel{\text{def}}{=} P$, since the new equivalence classes consist of singletons; (ii) modify the least upper bound polymorphic rule (\exists^\vee) so that it requires the quantifiers to be equal (and performs alpha-conversion if necessary):

$$\frac{\Theta, \vec{\alpha}^+ \models P_1 \vee P_2 = Q \quad \vec{\alpha}^+ \subseteq \text{fv } Q}{\Theta \models \exists \vec{\alpha}^+. P_1 \vee \exists \vec{\alpha}^+. P_2 = Q} \quad (\text{E}\exists^\vee)$$

(iii) replace the subtyping polymorphic rule (\forall^\leq) by the following rule:

$$\frac{\vec{\alpha}^+ \text{ are fresh} \quad \Theta, \vec{\beta}^+, \Xi, \vec{\alpha}^+ \{ \Theta, \vec{\beta}^+ \} \models [\vec{\alpha}^+ / \alpha^+] N \leq M \vdash C}{\Theta; \Xi \models \forall \vec{\beta}^+, \vec{\alpha}^+. N \leq \forall \vec{\beta}^+. M \vdash C} \quad (\text{EV}^\leq)$$

and (iv) update the existential rule (\exists^\geq) symmetrically.

After these changes, the rule ($\text{TA}pp^{\text{INF}}$) and its algorithmic counterpart can be used to infer the type of $c\{P\}$.

The elementary type inference restricts the expressiveness of the subtyping perhaps too much. As mentioned, it forbids the quantifier reordering; besides, as soon as the right-hand side of the subtyping relation is polymorphic, it restricts the instantiation of the left-hand side quantifiers, which *disallows*, for example, $\cdot \vdash \forall \alpha^+. \uparrow \alpha^+ \leq \forall \alpha^+. \uparrow \downarrow \alpha^+$. On the other hand, the elementary subtyping system can be smoothly extended with bounded quantification, which we discuss later.

Solution 3: Labeled Quantifiers. A compromise solution that resolves the ambiguity of explicit the instantiation without fixing the order of quantifiers is using *labeled* quantifiers. Let us discuss how to introduce them to the *negative* part of the system, since the positive subsystem is symmetric.

Each type abstraction $\Lambda l \triangleright \alpha^+. c$ must be annotated with a label l . This label propagates to the inferred polymorphic type. This way, each polymorphic quantifier $\forall l_i \triangleright \alpha_i^{+I}. N$ is annotated with a label l_i whose index is taken from the *list* of labels I associated with the group of quantifiers. To instantiate a quantifying variable, one refers to it by its label: $c\{l_j \triangleright P\}$, which is expressed by the following typing rules:

$$\frac{\Theta, \alpha^+; \Gamma \vdash c : N}{\Theta; \Gamma \vdash \Lambda l \triangleright \alpha^+. c : \forall l \triangleright \alpha^+. N} \quad (\Lambda^{\text{INF}}) \quad \frac{\Theta; \Gamma \vdash c : \forall l_i \triangleright \alpha_i^{+I}. N}{\Theta; \Gamma \vdash c\{l_j \triangleright P\} : \forall l_i \triangleright \alpha_i^{+I \setminus \{j\}}. [P / \alpha_j^+] N} \quad (\text{TA}pp^{\text{INF}})$$

The subtyping rule allows the quantifiers to have an arbitrary order. However, similarly to the elementary subtyping, each of the right-hand side quantifiers must have a matching left-hand side quantifier with the same label. These matching quantifiers synchronously removed (in other words,

each variable is abstractly instantiated to itself), and the remaining quantifiers are instantiated as usual by a substitution σ .

$$\frac{\{I\} \subseteq \{J\} \quad \Theta, \vec{\alpha}_i^{+I} \vdash \sigma : \vec{\alpha}_j^{+J} \setminus \{I\} \quad \Theta, \vec{\alpha}_i^{+I} \vdash [\sigma]N \leq M}{\Theta \vdash \forall l_j \triangleright \vec{\alpha}_j^{+J}. N \leq \forall l_i \triangleright \vec{\alpha}_i^{+I}. M} \quad (\forall^{\leq})$$

The equivalence (mutual subtyping) in this system allows the quantifiers within the same group to be reordered together with their labels. For instance, $\forall l \triangleright \alpha^+. \forall m \triangleright \beta^+. N$ and $\forall m \triangleright \beta^+. \forall l \triangleright \alpha^+. N$ are subtypes of each other. However, the right-hand side quantifiers are removed synchronously with the matching left-hand side quantifiers, which means that $\cdot \vdash \forall l \triangleright \alpha^+. \uparrow \alpha^+ \leq \forall l \triangleright \alpha^+. \uparrow \downarrow \uparrow \alpha^+$ is still not allowed.

6.2 Weakening of the Subtyping Invariant

In order to make the subtyping relation decidable, we made the subtyping of shifts *invariant*, which manifests in the following rules:

$$\frac{\Theta \vdash P \simeq^{\leq} Q}{\Theta \vdash \uparrow P \leq \uparrow Q} \quad (\uparrow^{\leq}) \qquad \frac{\Theta \vdash N \simeq^{\leq} M}{\Theta \vdash \downarrow N \geq \downarrow M} \quad (\downarrow^{\geq})$$

To make the system more expressive, we can relax these rules. Although making both rules covariant would make the system undecidable, it is possible to make the up-shift subtyping covariant while keeping the down-shift invariant:

$$\frac{\Theta \vdash Q \geq P}{\Theta \vdash \uparrow P \leq \uparrow Q} \quad (\uparrow_{\text{RLX}}^{\leq})$$

However, this change requires an update to the algorithm. As one can expect, the corresponding algorithmic rule (\uparrow^{\leq}) must be changed accordingly:

$$\frac{\Theta; \Xi \models Q \leq P \models C}{\Theta; \Xi \models \uparrow P \leq \uparrow Q \models C} \quad (\uparrow_{\text{RLX}}^{\leq})$$

Also, notice that now the algorithmic variables can occur on *both* sides of the positive subtyping relation, and thus, a more sophisticated constraint solver is needed. To see what kind of constraints will be generated, let us keep the other algorithmic rules unchanged for a moment and make several observations.

First, notice that the negative quantification variables still only occur at *invariant* positions: the negative variables are generated from the existential quantifiers $\exists \alpha^-. P$, and the switch from a positive to a negative type in the syntax path to α^- in P is ‘guarded’ by *invariant* down-shift. It means that in the algorithm, the negative constraint entries can only be equivalence entries ($\widehat{\alpha}^- := N$) but not subtyping entries.

Second, notice that the positive subtyping is still not dependent on the negative subtyping. It means, that any leaf-to-root path in the inference tree is monotonous: in the first phase, the negative subtyping rules are applied, and in the second phase, only the positive ones.

The invariant that only the left-hand side of the judgment is algorithmic is not preserved. However, the only rule violating this invariant is ($\uparrow_{\text{RLX}}^{\leq}$), and it is placed at the interface between negative subtyping and positive subtyping. It means that the *negative* algorithmic variables (produced by \exists in the positive phase of the inference path) still only occur on the left-hand side of the judgment, and the positive algorithmic variables (produced by \forall in the first phase) can occur either on the left-hand side or on the right-hand side (if the switch from the negative to the positive phase is made by ($\uparrow_{\text{RLX}}^{\leq}$)) but not both. This way, the produced entries will be of the following form: (i) $\widehat{\alpha}^+ := P$, (ii) $\widehat{\alpha}^- := N$, (iii) $\widehat{\alpha}^+ := P$, and (iv) $\widehat{\alpha}^+ := P$ where P does not contain positive unification variables.

We do not present the details of the constraint resolution procedure or how it is integrated into the subtyping algorithm, leaving it for future work. However, we believe that the system consisting of the mentioned kinds of constraints is solvable by the following reasons. (i) The equivalence entries $\widehat{\alpha}^+ \approx P$ and $\widehat{\alpha}^- \approx N$ are resolved by standard first-order pattern unification techniques. (ii) The resolution of supertyping entries ($\widehat{\alpha}^+ \geq P$) is discussed in the original algorithm (section 4.5). (iii) The new kind of entries $\widehat{\alpha}^+ \leq P$ can be reduced by case analysis. If P is a variable β^+ or a shifted computation $\downarrow N$ then $\widehat{\alpha}^+ \leq P$ is equivalent to $\widehat{\alpha}^+ \approx P$; otherwise, $\widehat{\alpha}^+ \leq \exists \vec{\alpha}^-. \downarrow N$ is equivalently replaced by $\widehat{\alpha}^+ \approx [\vec{\alpha}^- / \alpha^-] \downarrow N$ for freshly created $\vec{\alpha}^-$.

This way, the algorithm can be extended to support the *covariant* up-shift subtyping. This, for example, enriches the relation with such subtypes as $\cdot \vdash \forall \alpha^+. \alpha^+ \rightarrow \uparrow \alpha^+ \leq \downarrow \text{Int} \rightarrow \uparrow \exists \beta^-. \downarrow \beta^-$, which does not hold in the original system. Moreover, this extension also increases the expressiveness of the *type inference*. Namely, when inferring a type of an annotated let binding $\text{let } x : P = v(\vec{v}); c$, we require $\uparrow P$ to be a supertype of $\uparrow Q$ —the resulting type of the application $v(\vec{v})$. Now, it is equivalent to the requirement that P is a supertype of Q (without the up-shift), which is more permissive.

In addition, the enhancement of the unification algorithm described above makes possible another improvement to the system—bounded quantification.

6.3 Bounded Quantification

After the weakening of subtyping invariants, it is possible to smoothly extend the *elementary* version of the system with bounded quantification. In particular, we can add upper bounds to polymorphic \forall -quantifiers: $\forall \alpha^+ \leq P. N$.

The declarative subtyping rules for bounded quantification are as expected: the instantiation of quantified variables must satisfy the corresponding bounds ($\Theta, \alpha^+ \vdash Q_i \geq [\sigma] \beta_i^+$ holds for each i); and the right-hand side quantifiers must be more restrictive than the matching left-hand side quantifiers ($\Theta \vdash P_i \geq P'_i$ for each i):

$$\frac{\overrightarrow{\Theta \vdash P_i \geq P'_i} \quad \overrightarrow{\Theta, \alpha^+ \vdash \sigma : \beta^+} \quad \overrightarrow{\Theta, \alpha^+ \vdash Q_i \geq [\sigma] \beta_i^+} \quad \overrightarrow{\Theta, \alpha^+ \vdash [\sigma] N \leq M}}{\Theta \vdash \forall \alpha^+ \leq P. \forall \beta^+ \leq Q. N \leq \forall \alpha^+ \leq P'. M} \quad (\text{EV}^{\leq})$$

The corresponding algorithmic rule, as before, replaces the polymorphic variables $\beta_1 \dots$ with fresh algorithmic variables $\widehat{\beta}_1^+ \dots$, and merges the generated constraints with the constraints given in the bounds $\widehat{\beta}_i^+ \leq Q_i$:

$$\frac{\overrightarrow{\Theta; \cdot \vdash P_i \geq P'_i} \quad \overrightarrow{\widehat{\beta}_i^+ \text{ are fresh}} \quad \overrightarrow{\Theta, \alpha^+; \Xi, \widehat{\beta}^+ \{ \Theta, \alpha^+ \} \vdash N \leq M \Rightarrow C_0} \quad \overrightarrow{\Xi \vdash C_0 \& (\widehat{\beta}_i^+ \leq Q_i) = C}}{\Theta; \Xi \vdash \forall \alpha^+ \leq P. \forall \beta^+ \leq Q. N \leq \forall \alpha^+ \leq P'. M \Rightarrow C} \quad (\text{EV}^{\leq})$$

Notice that the algorithmic rule (EV^{\leq}) requires the bounding types \vec{P} and \vec{Q} to be *declarative*. Because of that, the generated constraints have shape $\widehat{\beta}_i^+ \leq Q_i$, i.e., the bounding types do not contain algorithmic variables. This kind of constraints is covered by the resolution procedure described in section 6.2, and thus, (EV^{\leq}) fits into the algorithmic system without changes to the constraint solver.

However, in order for this rule to be complete, the declarative system must be restricted as well. To guarantee that no algorithmic variable is introduced in the bounding types, we need to forbid the quantifying variables to occur in the bounding types in the declarative system. For that purpose, we must include this requirement in the type well-formedness (here $\text{fbv}N$ denotes the set of variables

freely occurring in the quantifiers of N at any depth):

$$\frac{\overrightarrow{\Theta \vdash P_i} \quad \Theta, \overrightarrow{\alpha^+} \vdash N \quad \overrightarrow{\alpha^+} \cap \text{fbv} N = \emptyset}{\Theta \vdash \forall \alpha^+ \leq P. N} \quad (\forall^{\text{WF}})$$

For brevity, we do not discuss in detail the combination of bounded and unbounded quantification in one system. However, this update is straightforward: the unbounded quantifiers are treated as the weakest bound, which is trivially satisfied. Analogously, it is possible to extend the system with lower bound \forall -quantifiers: $\overrightarrow{\forall \alpha^+ \geq P. N}$. Then the generated constraints will change to $\widehat{\beta}_i^+ : \geq Q_i$, which is also covered by the resolution procedure. However, we do not consider bounded *existential* quantification, because in this case, the constraint resolution would require us to find the *greatest lower bound* of two negative types, which is not well-defined (see section 4.5).

6.4 Bidirectionalization

The algorithm we provide requires that all lambda functions are annotated. This restriction significantly simplifies the type inference by making the terms uniquely define their types. However, it leads to certain redundancies in typing, in particular, the type of a lambda expression cannot be inferred from the context it is used in: $((\lambda x. \text{return } x) : (\text{Int} \rightarrow \text{Int}))$ is not allowed.

The well-known way to incorporate this expressiveness into the system is to make the typing bidirectional [Dunfield et al. 2020]. The idea is to split the typing judgment into two kinds:

- (i) *Synthesis* judgments are used when the type of a term can be *inferred* based exclusively on the term itself. Syntactically, we denote synthesizing judgments as $\Theta ; \Gamma \vdash c \Rightarrow N$.
- (ii) *Checking* judgments assume that the type of a term is *given* from the context, and it is required to *check* that the given type can be assigned to the term. In checking judgments $\Theta ; \Gamma \vdash c \Leftarrow N$, the type N is considered as an *input*.

To bidirectionalize the system, each typing rule must be oriented either to synthesis or to checking (or both). In particular, the *annotated* lambda-abstractions are synthesizing, and the *unannotated* ones are checking.

$$\frac{\Theta \vdash P \quad \Theta ; \Gamma, x : P \vdash c \Rightarrow N}{\Theta ; \Gamma \vdash \lambda x : P. c \Rightarrow P \rightarrow N} \quad (\lambda^{\Leftarrow}) \qquad \frac{\Theta \vdash P \quad \Theta ; \Gamma, x : P \vdash c \Leftarrow N}{\Theta ; \Gamma \vdash \lambda x. c \Leftarrow P \rightarrow N} \quad (\lambda^{\Leftarrow})$$

As common for bidirectional systems with subtyping, we would also need to introduce the *subsumption* rules. They allow us to ‘forget’ the information about the type by switching to the synthesis mode, as long as the synthesized type is more polymorphic than the checked one:

$$\frac{\Theta ; \Gamma \vdash c \Rightarrow N \quad \Theta \vdash N \leq M}{\Theta ; \Gamma \vdash c \Leftarrow M} \quad (\text{SUB}_-) \qquad \frac{\Theta ; \Gamma \vdash v \Rightarrow P \quad \Theta \vdash Q \geq P}{\Theta ; \Gamma \vdash v \Leftarrow Q} \quad (\text{SUB}_+)$$

Most of the original rules will have two bidirectional counterparts, one for each mode. The premises of the rules are oriented with respect to the conclusion. For instance, we need both checking and inferring versions for return v :

$$\frac{\Theta ; \Gamma \vdash v \Leftarrow P}{\Theta ; \Gamma \vdash \text{return } v \Leftarrow \uparrow P} \quad (\text{RET}^{\Leftarrow})$$

For some rules, however, it only makes sense to have one mode. As mentioned above, the unannotated lambda-abstraction is always checking. On the other hand, an *annotated* lambda, type-level lambda abstraction, or a variable inference is always synthesizing, since the required type information is already given in the term:

$$\frac{\Theta, \alpha^+ ; \Gamma \vdash c \Rightarrow N}{\Theta ; \Gamma \vdash \Lambda \alpha^+. c \Rightarrow \forall \alpha^+. N} \quad (\Lambda^{\Leftarrow}) \qquad \frac{x : P \in \Gamma}{\Theta ; \Gamma \vdash x \Rightarrow P} \quad (\text{VAR}^{\Rightarrow})$$

Finally, let us discuss the application inference part of the declarative system. Rules $(\emptyset_{\bullet \Rightarrow}^{\text{INF}})$ and $(\forall_{\bullet \Rightarrow}^{\text{INF}})$ are unchanged: they simply do not have typing premises to be oriented. However, the arrow application rule $(\rightarrow_{\bullet \Rightarrow}^{\text{INF}})$ infers the result of the application of an arrow $Q \rightarrow N$ to a list of arguments v, \vec{v} , and it must make sure that the first argument v is typeable with the expected Q . As we will discuss further, simply checking $\Theta; \Gamma \vdash v \Leftarrow Q$ turns out to be too powerful and potentially leads to undecidability. Instead, this check is approximated by the combination of $\Theta; \Gamma \vdash v \Rightarrow P$ and $\Theta \vdash Q \geq P$.

The Algorithm. The algorithmic system is bidirectionalized in a similar way. Each typing premise and the conclusion of the rule is oriented to either ‘checking’ (\Leftarrow) or ‘synthesizing’ (\Rightarrow) mode, with respect to the declarative system.

For brevity, we omit the details of the bidirectional algorithm. However, let us discuss one especially important rule—arrow application inference. An intuitive way to bidirectionalize this rule is the following:

$$\frac{\begin{array}{l} \Theta; \Gamma \vdash v \Leftarrow Q \Leftarrow C_1 \\ \Theta; \Gamma; \Xi \vdash N \bullet \vec{v} \Rightarrow M \Leftarrow \Xi'; C_2 \quad \Xi \vdash C_1 \& C_2 = C \end{array}}{\Theta; \Gamma; \Xi \vdash Q \rightarrow N \bullet v, \vec{v} \Rightarrow M \Leftarrow \Xi'; C} \quad (\rightarrow^{\bullet \Rightarrow})$$

However, this rule is too allowing. The judgement $\Theta; \Gamma \vdash v \Leftarrow Q \Leftarrow C_1$ requires us to check a term against an *algorithmic* type Q . Further, through the lambda function checking (λ^{\Leftarrow}) the algorithmic types infiltrate the type context and then through variable inference and subsumption, *both* sides of subtyping. This way, $(\rightarrow^{\bullet \Rightarrow})$ compromises the important invariant of the subtyping algorithm: now the same algorithmic variable can occur on *both* sides of the subtyping relation.

We believe that the relaxation of this invariant brings the constraint resolution too close to the second-order pattern unification, which is undecidable [Goldfarb 1981]. In particular, the constraint $(\hat{\alpha}^+ : \geq \downarrow \uparrow \hat{\alpha}^+)$ is solvable with $\hat{\alpha}^+ = \exists \beta^-. \downarrow \beta^-$, since by (\exists^{\geq}) , the existential β^- can be impredicatively instantiated to $\uparrow \exists \beta^-. \downarrow \beta^-$.

The constraints such as $(\hat{\alpha}^+ : \geq \downarrow \uparrow \hat{\alpha}^+)$ are not merely hypothetical. It occurs, for example, when we infer the type of the following application:

$$\cdot; \cdot \vdash \forall \alpha^+. \downarrow (\alpha^+ \rightarrow \uparrow \text{Int}) \rightarrow \downarrow \uparrow \alpha^+ \rightarrow \uparrow \text{Int} \bullet \{\lambda x. \lambda y. \text{let } y' = x(y); \text{return } y'\} \Rightarrow \uparrow \text{Int}$$

For x to be applicable to y , the type of $y - \downarrow \uparrow \alpha^+$ must be a subtype of α^+ —the type expected by x . This way, this inference is possible if and only if $(\hat{\alpha}^+ : \geq \downarrow \uparrow \hat{\alpha}^+)$ is solvable. Using a similar scheme, we can construct different examples whose resolution significantly relies on second-order pattern unification. Thus, we leave the resolution of this type of constraints beyond the scope of this work.

Instead, we strengthen $(\rightarrow^{\bullet \Rightarrow})$ in such a way that it never checks a term against an algorithmic type. Instead, type checking v against Q is replaced by a more restrictive premise—checking that v *synthesizes a subtype* of Q :

$$\frac{\begin{array}{l} \Theta; \Gamma \vdash v \Rightarrow P \quad \Theta; \Xi \vdash Q \geq P \Leftarrow C_1 \\ \Theta; \Gamma; \Xi \vdash N \bullet \vec{v} \Rightarrow M \Leftarrow \Xi'; C_2 \quad \Xi \vdash C_1 \& C_2 = C \end{array}}{\Theta; \Gamma; \Xi \vdash Q \rightarrow N \bullet v, \vec{v} \Rightarrow M \Leftarrow \Xi'; C} \quad (\rightarrow^{\bullet \Rightarrow})$$

As one can expect, the declarative rule is also changed accordingly. In terms of practicality, this change disallows implicit checking against a *polymorphic* type: $\cdot; \cdot \vdash \lambda x. \text{return } x \Leftarrow \forall \alpha^+. \alpha^+ \rightarrow \uparrow \alpha^+$ is not allowed, because it would require adding algorithmic $x : \hat{\alpha}^+$ to the context. However, once the instantiation is made explicit or the lambda is explicitly polymorphic, the checking is allowed: $\cdot; \cdot \vdash \lambda x. \text{return } x \Leftarrow \text{Int} \rightarrow \uparrow \text{Int}$ and $\cdot; \cdot \vdash \Lambda \alpha^+. \lambda x : \alpha^+. \text{return } x \Leftarrow \forall \alpha^+. \alpha^+ \rightarrow \uparrow \alpha^+$ are both valid.

7 CONCLUSION

We have presented a type inference algorithm for an impredicative polymorphic lambda calculus with existential types and subtyping. The system is designed in call-by-value paradigm, that allowed us to restrict it to a decidable fragment of the language, that we describe declaratively. We presented the inference algorithm and proved its soundness and completeness with respect to the specification. The algorithm combines unification—a standard type inference technique—with anti-unification—which has not been used in type inference before.

REFERENCES

- Jacek Chrzaszcz (1998). “Polymorphic Subtyping without Distributivity.” In: *Proceedings of the 23rd International Symposium on Mathematical Foundations of Computer Science*. MFCS ’98. Berlin, Heidelberg: Springer-Verlag, pp. 346–355.
- Luis Damas and Robin Milner (1982). “Principal type-schemes for functional programs.” In: *ACM-SIGACT Symposium on Principles of Programming Languages*. URL: <https://api.semanticscholar.org/CorpusID:11319320>.
- Jana Dunfield and Neel Krishnaswami (Nov. 2020). “Bidirectional Typing.” In: arXiv: 1908.05839.
- Jana Dunfield and Neelakantan R. Krishnaswami (2016). “Sound and complete bidirectional typechecking for higher-rank polymorphism with existentials and indexed types.” In: *Proceedings of the ACM on Programming Languages* 3, pp. 1–28. URL: <https://api.semanticscholar.org/CorpusID:19625612>.
- Richard A. Eisenberg, Guillaume Duboc, Stephanie Weirich, and Daniel K. Lee (2021). “An existential crisis resolved: type inference for first-class existential types.” In: *Proceedings of the ACM on Programming Languages* 5, pp. 1–29. URL: <https://api.semanticscholar.org/CorpusID:237205112>.
- Jean-Yves Girard (1971). “Une Extension de L’interprétation de Gödel à L’Analyse, et Son Application à L’Élimination Des Coupures Dans L’Analyse et La Théorie Des Types.” In: *Proceedings of the Second Scandinavian Logic Symposium*. Ed. by J.E. Fenstad. Vol. 63. Studies in Logic and the Foundations of Mathematics. North Holland: Elsevier, pp. 63–92. DOI: 10.1016/S0049-237X(08)70843-7.
- Warren D. Goldfarb (1981). “The Undecidability of the Second-Order Unification Problem.” In: *Theor. Comput. Sci.* 13, pp. 225–230. DOI: 10.1016/0304-3975(81)90040-2.
- Roger Hindley (1969). “The Principal Type-Scheme of an Object in Combinatory Logic.” In: *Transactions of the American Mathematical Society* 146, pp. 29–60. URL: <https://api.semanticscholar.org/CorpusID:7223700>.
- Konstantin Läuffer and Martin Odersky (1994). “Polymorphic type inference and abstract data types.” In: *ACM Trans. Program. Lang. Syst.* 16, pp. 1411–1430. URL: <https://api.semanticscholar.org/CorpusID:13189837>.
- Daan Leijen (2006). “First-class polymorphism with existential types.” In: URL: <https://api.semanticscholar.org/CorpusID:16392798>.
- Paul Blain Levy (Dec. 2006). “Call-by-Push-Value: Decomposing Call-by-Value and Call-by-Name.” In: *Higher-Order and Symbolic Computation* 19.4, pp. 377–414. DOI: 10.1007/s10990-006-0480-6.
- Dale Miller (1991). “A Logic Programming Language with Lambda-Abstraction, Function Variables, and Simple Unification.” In: *J. Log. Comput.* 1.4, pp. 497–536. DOI: 10.1093/logcom/1.4.497.
- Robin Milner (1978). “A Theory of Type Polymorphism in Programming.” In: *J. Comput. Syst. Sci.* 17, pp. 348–375. URL: <https://api.semanticscholar.org/CorpusID:388583>.
- Benjamin Pierce and David Turner (Jan. 2000). “Local Type Inference.” In: *ACM Transactions on Programming Languages and Systems* 22.1, pp. 1–44. DOI: 10.1145/345099.345100.
- Gordon D. Plotkin (1970). “A Note on Inductive Generalization.” In: *Machine Intelligence* 5.1, pp. 153–163.
- John C Reynolds (1970). “Transformational systems and the algebraic structure of atomic formulas.” In: vol. 5. 1, pp. 135–151.
- John C Reynolds (1974). “Towards a Theory of Type Structure.” In: *Programming Symposium*. Springer. Paris, France: Springer, pp. 408–425.
- Alejandro Serrano, Jurriaan Hage, Simon Peyton Jones, and Dimitrios Vytiniotis (Aug. 2020). “A Quick Look at Impredicativity.” In: *Proceedings of the ACM on Programming Languages* 4.ICFP, pp. 1–29. DOI: 10.1145/3408971.
- Jerzy Tiuryn (1995). “Equational Axiomatization of Bicoercibility for Polymorphic Types.” In: *Foundations of Software Technology and Theoretical Computer Science, 15th Conference, Bangalore, India, December 18–20, 1995, Proceedings*. Ed. by P. S. Thiagarajan. Vol. 1026. Lecture Notes in Computer Science. Springer, pp. 166–179. DOI: 10.1007/3-540-60692-0_47.
- Jerzy Tiuryn and Paweł Urzyczyn (1996). “The Subtyping Problem for Second-Order Types Is Undecidable.” In: *Proceedings of the 11th Annual IEEE Symposium on Logic in Computer Science*. LICS ’96. USA: IEEE Computer Society, p. 74.
- Jinxu Zhao and Bruno C. d. S. Oliveira (2022). “Elementary Type Inference.” In: *36th European Conference on Object-Oriented Programming, ECOOP 2022, June 6–10, 2022, Berlin, Germany*. Ed. by Karim Ali and Jan Vitek. Vol. 222. LIPIcs. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2:1–2:28. DOI: 10.4230/LIPICS.ECOOP.2022.2.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107