

Local Type Inference for Polarised System F with Existentials

ANONYMOUS AUTHOR(S)

CHANGE!! This paper addresses the challenging problem of type inference for Impredicative System F with existential types, a critical aspect of many programming languages. While System F serves as the basis for type systems in numerous languages, existing type inference techniques for Impredicative System F are undecidable due to the presence of existential (\exists) and polymorphic (\forall) types. Consequently, current algorithms are often ad-hoc and sub-optimal. This paper presents novel contributions in the form of a local type inference algorithm for Impredicative System F with existential types. The algorithm introduces innovative techniques, such as a unique combination of unification and anti-unification, a full correctness proof, and the use of control structures inspired by Call-By-Push-Value. Additionally, the paper discusses a type inference framework that allows the algorithm to be applied to different type systems, offering insights into the under-researched area of impredicative existential 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), 32 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), the polymorphic type $\forall a. A$ has a big lambda $\lambda a. e$ 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 sophisticated enough to model things like parametricity and representation independence. However, System F by itself is unwieldy as a programming language. The fact that the universal type $\forall a. A$ 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 $\text{len} : \forall a. \text{List } a \rightarrow \text{int}$ 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 most famous of the algorithms for solving these constraints is the Damas-Hindley-Milner algorithm. The idea is that type instantiation induces a subtype ordering: the type $\forall a. A$ is a subtype of all of its instantiations. So we wish to be free to use the same function len at many different types

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, or 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>

such as $\text{List int} \rightarrow \text{int}$, $\text{List bool} \rightarrow \text{int}$, $\text{List (int} \times \text{bool)} \rightarrow \text{int}$, and so on. However, the subtype relation is nondeterministic: it tells us that whenever we see a polymorphic type $\forall a. A$, we know it is a subtype of *any* of its instantiations. To turn this into an algorithm, we have to actually make some choices, and DHM works by 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. 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, **leijen06** studied an extension of Damas-Milner inference in which type schemes contain alternating universal and existential quantifiers. Quantifiers still range over monotypes, and higher-rank polymorphism is not permitted. More recently, **dk19** 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). **existential-crisis** 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 (i.e., types without any occurrences of quantifiers). 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 do full type inference, 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 what the partial algorithm does – for example, **existential-crisis** 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 OO 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) anti-unification.

Contributions. Our contributions are as follows:

- 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.
- 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 work in a variant of call-by-push-value, which lets us formulate a subtyping relation which is still decidable.
- Our algorithm breaks some of the fundamental invariants of HM-style type inference. As a result, it needs to mix unification and anti-unification, uses the call-by-push-value structure to control function arities and how quantifiers can be instantiated.
- The original local type inference paper combined local type inference with bidirectional typechecking to minimize the number of needed annotations, but we show how existential types complicate the integration of bidirectionality with local type inference, and we explore the design space to show how the same scheme could be applied to work with different type systems.

Neel: We need to explain what local type inference is, and how it is and isn't related to bidirectional typechecking.

2 OVERVIEW

2.1 What types do we infer?

2.2 The Language of Types

The types of $F^{\pm}\exists$ are given in fig. 1. They are stratified into two syntactic categories (polarities): positive and negative, similarly to the Call-By-Push-Value system [Levy 2006]. The negative types represent computations, and the positive types represent values:

- α^- is a negative type variable, which can be taken from a context or introduced by \exists .
- a function $P \rightarrow N$ takes a value as input and returns a computation;
- a polymorphic abstraction $\forall \vec{\alpha}^+. N$ quantifies a computation over a list of positive type variables $\vec{\alpha}^+$. The polarities are chosen to follow the definition of functions.
- a shift $\uparrow P$ allows a value to be used as a computation, which at the term level corresponds to a pure computation **return** v .
- + α^+ is a positive type variable, taken from a context or introduced by \forall .
- + $\exists \vec{\alpha}^+. P$, symmetrically to \forall , binds negative variables in a positive type P .
- + a shift $\downarrow N$, symmetrically to the up-shift, thunk a computation, which at the term level corresponds to $\{c\}$.

Fig. 1. Declarative Types of $F^{\pm}\exists$

Definitional Equalities. For simplicity, we assume that alpha-equivalent terms are equal. This way, we assume that substitutions do not capture bound variables. Besides, we equate $\forall \vec{\alpha}^+. \forall \vec{\beta}^+. N$

with $\forall \alpha^+, \beta^+. N$, as well as $\exists \alpha^-. \exists \beta^-. P$ with $\exists \alpha^-, \beta^-. P$, and lift these equations transitively and congruently to the whole system.

2.3 The Language of Terms

In fig. 2, we define the language of terms of $F^\pm \exists$. The language combines System F with the Call-By-Push-Value approach.

- + x denotes a (positive) term variable; **Ilya: why no negatives? Following CBPV**
- + $\{c\}$ is a value corresponding to a thunked or suspended computation;
- $\pm (c : N)$ and $(v : P)$ allow one to annotate positive and negative terms;
- **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 ;
- **let** $x = v ; c$ is a standard let, binding a value v to a variable x in a computation c ;
- Applicative let forms **let** $x : P = v(\vec{v}) ; c$ and **let** $x = v(\vec{v}) ; c$ operate similarly to the bind of a monad: they take a suspended computation v , apply it to a list of arguments, bind the result (which is expected to be pure) to a variable x , and continue 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;
- **let $^\exists$** ($\vec{\alpha^-}, x$) = $v ; c$ is the standard unpack of an existential type: expecting v to be 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.

Computation Terms

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

Value Terms

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

Fig. 2. Declarative Terms of $F^\pm \exists$

2.4 The key ideas of the algorithm

3 DECLARATIVE SYSTEM

The declarative system serves as a specification of the type inference algorithm. It consists of two main parts: the subtyping and the type inference.

3.1 Subtyping

It is represented by a set of inference rules shown in fig. 3.

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

Fig. 3. Declarative Subtyping

Quantifiers. Symmetric rules (\forall^{\leqslant}) and (\exists^{\geqslant}) specify the 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 the subtyping. For conciseness of representation, we compose these rules into one. First, our rule extends context T 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.

The instantiation of quantifiers is modeled by substitution σ . The notation $T_2 \vdash \sigma : T_1$ specifies its domain and range. For instance, $T, \vec{\beta}^+ \vdash \sigma : \vec{\alpha}^+$ means that σ maps the variables from $\vec{\alpha}^+$ to (positive) types well-formed in $T, \vec{\beta}^+$. This way, application $[\sigma]N$ instantiates (replaces) every α_i^- in N with $\sigma(\alpha_i^-)$.

Invariant Shifts. An important restriction that we put on the subtyping system is that the subtyping on shifted types requires their equivalence, as shown in (\downarrow^{\geqslant}) and (\uparrow^{\leqslant}). Relaxing both of these invariants make the system equivalent to System F, and thus, undecidable. However, after certain changes (\uparrow^{\leqslant}) can be relaxed to the covariant form, as we will discuss in ??.

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

Variables. The subtyping of variables is defined reflexively, which is enough to ensure the reflexivity of subtyping in general. The algorithm will use the fact that the subtypes of a variable coincide with its supertypes, which however is not true for an arbitrary type.

3.1.1 Properties of the Declarative Subtyping. A property that is important for the subtyping algorithm, in particular for the type *upgrade* procedure ($??$), 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 T . Then $T \vdash N_1 \leq N_2$ implies $\text{fv}(N_1) \subseteq \text{fv}(N_2)$, and $T \vdash P_1 \geq P_2$ implies $\text{fv}(P_1) \subseteq \text{fv}(P_2)$.*

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

Property 2 (Subtyping forms a preorder). *Let us say that two types N_1 and N_2 are in the subtyping relation if there exists a context T such that $T \vdash N_1 \leq N_2$; symmetrically, two types P_1 and P_2 are in the subtyping relation if there exists T such that $T \vdash P_2 \geq P_1$. Then the subtyping relation defined this way is reflexive and transitive.*

Property 3 (Subtyping agrees with substitution). *Suppose that σ is a substitution such that $T_2 \vdash \sigma : T_1$. Then*

- $T_1 \vdash N \leq M$ implies $T_2 \vdash [\sigma]N \leq [\sigma]M$, and
- + $T_1 \vdash P \geq Q$ implies $T_2 \vdash [\sigma]P \geq [\sigma]Q$.

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

Property 4 (Positive Least Upper Bound exists). *Suppose that P_1 and P_2 are positive types well-formed in T . Then there exists the least common supertype—a type P such that*

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

Negative Greatest Lower Bound does not exist. However, the symmetric construction—the greatest lower bound of two negative types does not always exist. Let us consider the following counterexample. Let us 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 $P \rightarrow \downarrow\uparrow P_1 \rightarrow \downarrow\uparrow P_2 \rightarrow N$? One of the common subtypes is $\forall \alpha^+, \beta^+, \gamma^+. \alpha^+ \rightarrow \downarrow\uparrow \beta^+ \rightarrow \downarrow\uparrow \gamma^+ \rightarrow N$, which, however is not the greatest one.

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. Since M_1 and M_2 are not equivalent, it 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 a.k.a. 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,

- (ii) one can introduce redundant quantifiers. For example, $\forall \alpha^+, \beta^+. \uparrow \alpha^+$ is equivalent but not equal to $\forall \alpha^+. \uparrow \alpha^+$;
- (i) one can reorder quantifiers. For example, $\forall \alpha^+, \beta^+. \alpha^+ \rightarrow \beta^+ \rightarrow \gamma^-$ is equivalent but not equal to $\forall \alpha^+, \beta^+. \beta^+ \rightarrow \alpha^+ \rightarrow \gamma^-$;
- (iii) the transformations (i) and (ii) can happen 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

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

The normalization algorithm is shown in fig. 4. The steps (i-ii) are implemented by the ordering function, which takes a set of variables $vars$ and a type and returns a list of variables from $vars$ that occur in the type in the order of their first occurrence. Its formal definition can be found in ??.

$\boxed{\text{nf}(N) = M}$ $\frac{}{\text{nf}(\alpha^-) = \alpha^-} \quad (\text{VAR}_{-}^{\text{NF}})$ $\frac{\text{nf}(P) = Q}{\text{nf}(\uparrow P) = \uparrow Q} \quad (\uparrow^{\text{NF}})$ $\frac{\text{nf}(P) = Q \quad \text{nf}(N) = M}{\text{nf}(P \rightarrow N) = Q \rightarrow M} \quad (\rightarrow^{\text{NF}})$ $\frac{\text{nf}(N) = N' \quad \text{ord } \vec{\alpha}^+ \text{ in } N' = \vec{\alpha}^+}{\text{nf}(\forall \vec{\alpha}^+. N) = \forall \vec{\alpha}^+. N'} \quad (\forall^{\text{NF}})$ <p style="margin-top: 10px;">ord $vars$ in N returns a list of variables $vars \cap \text{fv}(N)$ in the order of their first occurrence in N</p>	$\boxed{\text{nf}(P) = Q}$ $\frac{}{\text{nf}(\alpha^+) = \alpha^+} \quad (\text{VAR}_{+}^{\text{NF}})$ $\frac{\text{nf}(N) = M}{\text{nf}(\downarrow N) = \downarrow M} \quad (\downarrow^{\text{NF}})$ $\frac{\text{nf}(P) = P' \quad \text{ord } \vec{\alpha}^- \text{ in } P' = \vec{\alpha}^-}{\text{nf}(\exists \vec{\alpha}^-. P) = \exists \vec{\alpha}^-. P'} \quad (\exists^{\text{NF}})$ <p style="margin-top: 10px;">ord $vars$ in P returns a list of variables $vars \cap \text{fv}(P)$ in the order of their first occurrence in P</p>
--	---

Fig. 4. Type Normalization Procedure

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

Property 5 (Correctness of normalization).

- For N and M well-formed in T , $T \vdash N \simeq^{\leq} M$ is equivalent to $\text{nf}(N) = \text{nf}(M)$;
- + analogously, for P and Q well-formed in T , $T \vdash P \simeq^{\leq} Q$ is equivalent to $\text{nf}(P) = \text{nf}(Q)$.

3.3 Type Inference

The declarative specification of the type inference is shown in fig. 5. The positive typing judgment $T; \Gamma \vdash v : P$ is read as “under the type context T and variable context Γ , the term v is allowed to have the 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 of $T; \Gamma \vdash N \bullet \vec{v} \Rightarrow M$, which reads “under the type context T and variable context Γ , the application of a function of type N to the list of arguments \vec{v} is allowed to have the type M ”.

$\boxed{T; \Gamma \vdash c : N}$ Negative typing

$$\begin{array}{c}
 \frac{T \vdash P \quad T; \Gamma, x : P \vdash c : N}{T; \Gamma \vdash \lambda x : P. c : P \rightarrow N} \quad (\lambda^{\text{INF}}) \\
 \\
 \frac{T, \alpha^+; \Gamma \vdash c : N}{T; \Gamma \vdash \Lambda \alpha^+. c : \forall \alpha^+. N} \quad (\Lambda^{\text{INF}}) \\
 \\
 \frac{T; \Gamma \vdash v : P}{T; \Gamma \vdash \text{return } v : \uparrow P} \quad (\text{RET}^{\text{INF}}) \\
 \\
 \frac{T; \Gamma \vdash v : P \quad T; \Gamma, x : P \vdash c : N}{T; \Gamma \vdash \text{let } x = v; c : N} \quad (\text{LET}^{\text{INF}}) \\
 \\
 \frac{T; \Gamma \vdash v : \downarrow M \quad T; \Gamma \vdash M \bullet \vec{v} \Rightarrow \uparrow Q \text{ unique} \quad T; \Gamma, x : Q \vdash c : N}{T; \Gamma \vdash \text{let } x = v(\vec{v}); c : N} \quad (\text{LET}_{@}^{\text{INF}})
 \end{array}$$

$$\begin{array}{c}
 \frac{T \vdash P \quad T; \Gamma \vdash v : \downarrow M \quad T; \Gamma \vdash M \bullet \vec{v} \Rightarrow \uparrow Q \quad T \vdash \uparrow Q \leq \uparrow P \quad T; \Gamma, x : P \vdash c : N}{T; \Gamma \vdash \text{let } x : P = v(\vec{v}); c : N} \quad (\text{LET}_{@}^{\text{INF}}) \\
 \\
 \frac{T; \Gamma \vdash v : \exists \vec{\alpha}. P \quad \text{nf}(\exists \vec{\alpha}. P) = \exists \vec{\alpha}. P \quad T, \vec{\alpha}; \Gamma, x : P \vdash c : N \quad T \vdash N}{T; \Gamma \vdash \text{let}^{\exists}(\vec{\alpha}, x) = v; c : N} \quad (\text{LET}_{\exists}^{\text{INF}}) \\
 \\
 \frac{T \vdash M \quad T; \Gamma \vdash c : N \quad T \vdash N \leq M}{T; \Gamma \vdash (c : M) : M} \quad (\text{ANN}_{-}^{\text{INF}}) \\
 \\
 \frac{T; \Gamma \vdash c : N \quad T \vdash N \simeq^{\leq} N'}{T; \Gamma \vdash c : N'} \quad (\simeq_{-}^{\text{INF}})
 \end{array}$$

$\boxed{T; \Gamma \vdash v : P}$ Positive typing

$$\begin{array}{c}
 \frac{x : P \in \Gamma}{T; \Gamma \vdash x : P} \quad (\text{VAR}^{\text{INF}}) \\
 \\
 \frac{T; \Gamma \vdash c : N}{T; \Gamma \vdash \{c\} : \downarrow N} \quad (\{\}^{\text{INF}}) \\
 \\
 \frac{T \vdash Q \quad T; \Gamma \vdash v : P \quad T \vdash Q \geq P}{T; \Gamma \vdash (v : Q) : Q} \quad (\text{ANN}_{+}^{\text{INF}}) \\
 \\
 \frac{T; \Gamma \vdash v : P \quad T \vdash P \simeq^{\leq} P'}{T; \Gamma \vdash v : P'} \quad (\simeq_{+}^{\text{INF}})
 \end{array}$$

$\boxed{T; \Gamma \vdash N \bullet \vec{v} \Rightarrow M}$ Application typing

$$\begin{array}{c}
 \frac{T \vdash N \simeq^{\leq} N'}{T; \Gamma \vdash N \bullet \cdot \Rightarrow N'} \quad (\emptyset_{\bullet \Rightarrow}^{\text{INF}}) \\
 \\
 \frac{T; \Gamma \vdash v : P \quad T \vdash Q \geq P \quad T; \Gamma \vdash N \bullet \vec{v} \Rightarrow M}{T; \Gamma \vdash Q \rightarrow N \bullet v, \vec{v} \Rightarrow M} \quad (\rightarrow_{\bullet \Rightarrow}^{\text{INF}}) \\
 \\
 \frac{\vec{v} \neq \cdot \quad \vec{\alpha}^+ \neq \cdot \quad T \vdash \sigma : \vec{\alpha}^+ \quad T; \Gamma \vdash [\sigma] N \bullet \vec{v} \Rightarrow M}{T; \Gamma \vdash \forall \vec{\alpha}^+. N \bullet \vec{v} \Rightarrow M} \quad (\forall_{\bullet \Rightarrow}^{\text{INF}})
 \end{array}$$

Fig. 5. Declarative Inference

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

Variables. Rule (VAR^{INF}) allows to infer the type of a variable from the context. In literature can be found 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 (\simeq_+^{INF}) .

Annotations. Subtyping is also used by the annotation rules ($\text{ANN}_-^{\text{INF}}$) and ($\text{ANN}_+^{\text{INF}}$). 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 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 (RET^{INF}) and ($\{\}^{\text{INF}}$) add the corresponding shifts 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. 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 difference is that the former rule is for the *unannotated* let binder, and thus it requires the resulting type of application to be unique (up to equivalence), so that the type of the bound variable x is known before it is put into the context. 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 $T \vdash \uparrow Q \leq \uparrow P$. This judgment is more restrictive than checking bare $T \vdash P \geq Q$, however, it is necessary to make the algorithm complete as it allows us to preserve certain invariants (see ??). In ?? we discuss how this restriction can be relaxed together with invariant shift subtyping.

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_+^{INF}) and (\simeq_-^{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 $(\emptyset_{\bullet}^{\text{INF}})$. 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_+^{INF}) , however, the local relaxation is sufficient to prove the global property.

Application of a Polymorphic Type \forall . The complexity of the system is hidden in the rules, whose output type is not immediately defined by their input and the output of their premises (a.k.a. not mode-correct [Dunfield et al. 2020]). In our typing system, such rule is $(\forall_{\bullet}^{\text{INF}})$: the instantiation of the quantifying variables is not known a priori. The algorithm we present in ?? 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 important 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 (??), this subtyping check will provide the constraints we need to resolve the delayed instantiations of the quantifying variables.

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

Property 6 (Declarative Typing is Defined up to Equivalence). *Let us assume that $T \vdash \Gamma_1 \simeq^< \Gamma_2$, i.e., the corresponding types assigned by Γ_1 and Γ_2 are equivalent in T . Also, let us assume that $T \vdash N_1 \simeq^< N_2$, $T \vdash P_1 \simeq^< P_2$, and $T \vdash M_1 \simeq^< M_2$. Then*

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

Ilya: Other properties?

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, must be mode-correct ([Dunfield et al. 2020]), i.e., the 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.

Negative Algorithmic Variables

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

Positive Algorithmic Variables

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

Negative Algorithmic Types

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

Positive Algorithmic Types

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

Algorithmic Type Context

$\Upsilon = \{\hat{\alpha}_1^\pm, \dots, \hat{\alpha}_n^\pm\}$ where $\hat{\alpha}_1^\pm, \dots, \hat{\alpha}_n^\pm$ are pairwise distinct

Constraint Type Context

$\Sigma ::= \{\hat{\alpha}_1^\pm \{T_1\}, \dots, \hat{\alpha}_n^\pm \{T_n\}\}$ where $\hat{\alpha}_1^\pm, \dots, \hat{\alpha}_n^\pm$ are pairwise distinct

Fig. 6. Algorithmic Syntax

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 (a.k.a. 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 $\vec{\widehat{\alpha}}^+ / \vec{\alpha}^+$ and $\vec{\widehat{\alpha}}^- / \vec{\alpha}^-$.

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

Algorithmic Contexts and Well-formedness. To specify when algorithmic types are well-formed, we define algorithmic contexts Υ as sets of algorithmic variables. Then $T; \Upsilon \vdash P$ and $T; \Upsilon \vdash N$ represent the well-formedness judgment of algorithmic terms defined as expected. Informally, they check that all free declarative variables are in T , and all free algorithmic variables are in Υ . In addition to the rules repeating the declarative definition ??, we have two base cases for the algorithmic variables: $(\text{UVar}_+^{\text{WF}})$ and $(\text{UVar}_-^{\text{WF}})$ (see section 4.1)

$$\begin{array}{c} \vdots \\ \frac{\widehat{\alpha}^+ \in \Upsilon}{T; \Upsilon \vdash \widehat{\alpha}^+} \quad (\text{UVar}_+^{\text{WF}}) \end{array} \qquad \begin{array}{c} \vdots \\ \frac{\widehat{\alpha}^- \in \Upsilon}{T; \Upsilon \vdash \widehat{\alpha}^-} \quad (\text{UVar}_-^{\text{WF}}) \end{array}$$

Fig. 7. Well-formedness of Algorithmic Types

Algorithmic Normalization. Similarly to well-formedness, the normalization of algorithmic types is defined by extending the declarative definition with the algorithmic variables. To the rules repeating the declarative normalization, we add rules saying that normalization is trivial on algorithmic variables (see section 4.1).

$$\begin{array}{c} \vdots \\ \frac{}{\text{nf}(\widehat{\alpha}^+) = \widehat{\alpha}^+} \quad (\text{UVar}_+^{\text{NF}}) \end{array} \qquad \begin{array}{c} \vdots \\ \frac{}{\text{nf}(\widehat{\alpha}^-) = \widehat{\alpha}^-} \quad (\text{UVar}_-^{\text{NF}}) \end{array}$$

Fig. 8. Normalization of Algorithmic Types

4.2 Type Constraints

As the algorithm proceeds, it accumulates the information about the algorithmic type variables in the form of *constraints*. In our system, the constraints can be of two kinds: *subtyping constraints* and *unification constraints*. The subtyping constraint can only have a positive shape $\widehat{\alpha}^+ \geq P$, i.e., it restricts a positive algorithmic variable to be a supertype of a certain declarative type—this is one of the invariants that we preserve in the algorithm. The unification constraint can have either a

positive form $\widehat{\alpha}^+ \approx P$ or a negative form $\widehat{\alpha}^- \approx N$, however, the right-hand side of the constraint cannot contain algorithmic type variables. The set of constraints is denoted as C . We assume that each algorithmic variable can be restricted by at most one constraint.

We separately define UC as a set consisting of unification constraints only. This is done to simplify the representation of the algorithm. The unification algorithm, which we use as a subroutine of the subtyping algorithm, can only produce unification constraints. A set of unification constraints can be resolved in a simpler way than a general constraint set. This way, the separation of the unification constraint resolution into a separate procedure allows us to better decompose the structure of the algorithm, and thus, simplify the inductive proofs.

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

Fig. 9. Constraint Entries and Sets

Constraint Contexts. When one instantiates an algorithmic variable, they 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 a *constraint context* Σ —a set of pairs associating algorithmic variables and declarative contexts.

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

4.3 Subtyping Algorithm

For convenience and scalability, we decompose the subtyping algorithm into several procedures. Figure 10 shows these procedures and the dependencies between them: arrows denote the invocation of one procedure from another. The label «nf» annotating arrows means that the calling procedure normalizes the input before passing it to the callee.

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 T , a constraint context Σ , 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 T , but the first type (N and P) may contain algorithmic variables, whose instantiation contexts are specified by Σ .

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.

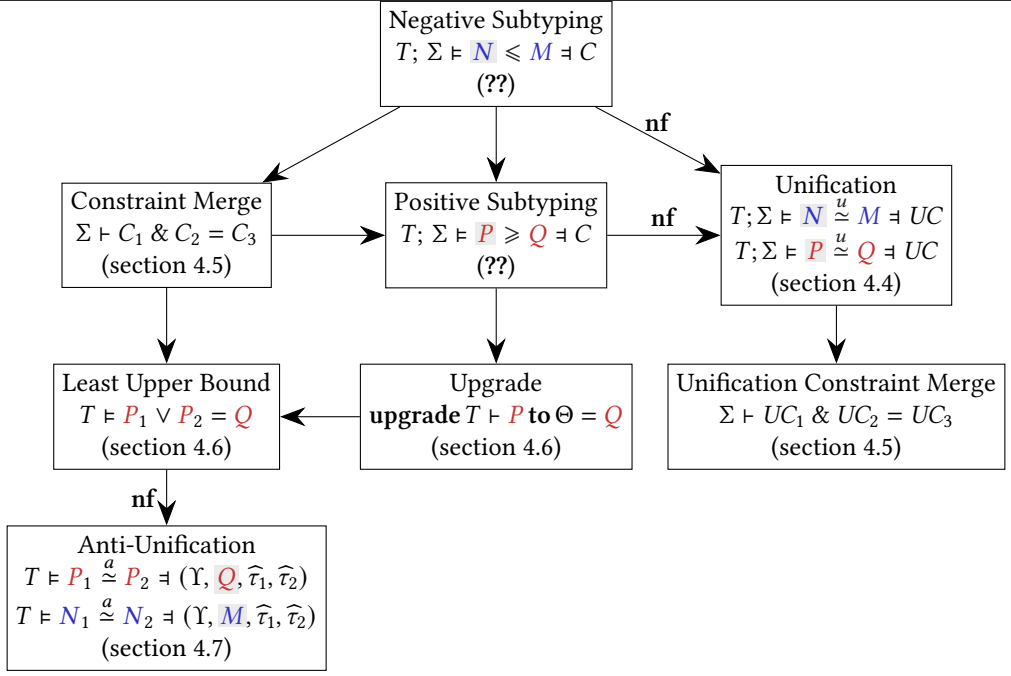


Fig. 10. Dependency graph of the subtyping algorithm

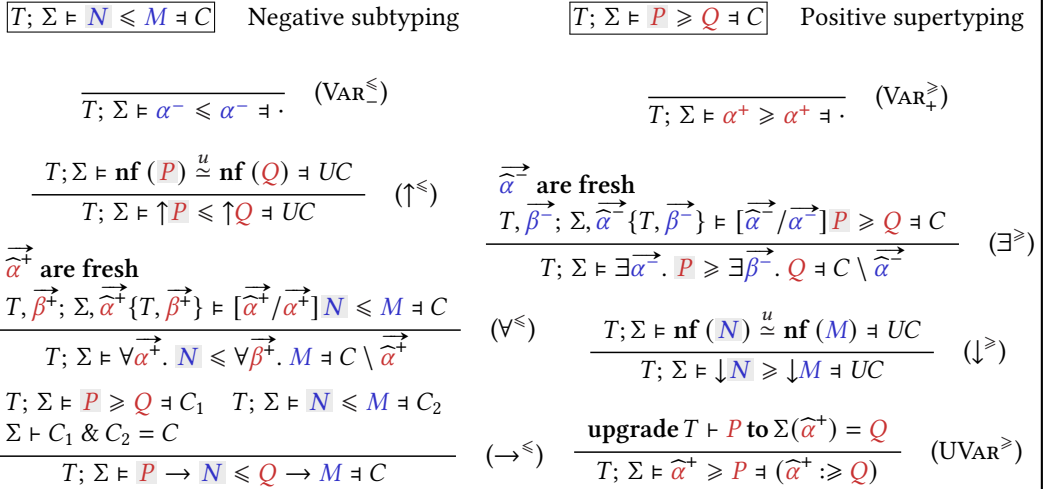


Fig. 11. Subtyping Algorithm

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, which make it equivalent to the right-hand side. For this purpose, the algorithm invokes the unification procedure $??$ preceded by 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. The rule introduces fresh algorithmic variables $\vec{\alpha}^+$ or $\vec{\alpha}^-$, puts them into the constraint context Σ (specifying that they must be instantiated in the extended context $T, \vec{\beta}^+$ or $T, \vec{\beta}^-$) 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 constraint 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 later in $??$) and returned as the output.

Algorithmic Variable. If one of the sides of the subtyping is a unification variable, the algorithm must create 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 $\widehat{\alpha}^+ \geq P$, which is covered by (UVar^{\geq}) .

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

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^{\geq}) *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 T , a constraint context Σ , 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 T and Σ , whereas the right-hand side is well-formed declaratively in T .

Since only the left-hand side may contain algorithmic variables, that the unification instantiates, we could have called this procedure *matching*. However, in $??$ we will discuss several modifications

of the type system, where this invariant is not preserved, and thus, this procedure becomes a genuine 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.

$T; \Sigma \models \mathbf{N} \stackrel{u}{\simeq} \mathbf{M} \dashv UC$ Negative unification	$T; \Sigma \models \mathbf{P} \stackrel{u}{\simeq} \mathbf{Q} \dashv UC$ Positive unification
$\frac{}{T; \Sigma \models \alpha^- \stackrel{u}{\simeq} \alpha^- \dashv \cdot} \quad (\text{VAR}^-_u)$	$\frac{}{T; \Sigma \models \alpha^+ \stackrel{u}{\simeq} \alpha^+ \dashv \cdot} \quad (\text{VAR}^+_u)$
$\frac{T; \Sigma \models \mathbf{P} \stackrel{u}{\simeq} \mathbf{Q} \dashv UC}{T; \Sigma \models \uparrow \mathbf{P} \stackrel{u}{\simeq} \uparrow \mathbf{Q} \dashv UC} \quad (\uparrow^u_u)$	$\frac{T; \Sigma \models \mathbf{N} \stackrel{u}{\simeq} \mathbf{M} \dashv UC}{T; \Sigma \models \downarrow \mathbf{N} \stackrel{u}{\simeq} \downarrow \mathbf{M} \dashv UC} \quad (\downarrow^u_u)$
$\frac{T; \Sigma \models \mathbf{P} \stackrel{u}{\simeq} \mathbf{Q} \dashv UC_1 \quad T; \Sigma \models \mathbf{N} \stackrel{u}{\simeq} \mathbf{M} \dashv UC_2}{T; \Sigma \models \mathbf{P} \rightarrow \mathbf{N} \stackrel{u}{\simeq} \mathbf{Q} \rightarrow \mathbf{M} \dashv UC_1 \& UC_2} \quad (\rightarrow^u_u)$	$\frac{T, \vec{\alpha}^-; \Sigma \models \mathbf{P} \stackrel{u}{\simeq} \mathbf{Q} \dashv UC}{T; \Sigma \models \exists \vec{\alpha}^-. \mathbf{P} \stackrel{u}{\simeq} \exists \vec{\alpha}^-. \mathbf{Q} \dashv UC} \quad (\exists^u_u)$
$\frac{T, \vec{\alpha}^+; \Sigma \models \mathbf{N} \stackrel{u}{\simeq} \mathbf{M} \dashv UC}{T; \Sigma \models \forall \vec{\alpha}^+. \mathbf{N} \stackrel{u}{\simeq} \forall \vec{\alpha}^+. \mathbf{M} \dashv UC} \quad (\forall^u_u)$	$\frac{\Sigma(\widehat{\alpha}^+) \vdash P}{T; \Sigma \models \widehat{\alpha}^+ \stackrel{u}{\simeq} P \dashv (\widehat{\alpha}^+ := P)} \quad (\text{UVar}^+_u)$
$\frac{\Sigma(\widehat{\alpha}^-) \vdash N}{T; \Sigma \models \widehat{\alpha}^- \stackrel{u}{\simeq} N \dashv (\widehat{\alpha}^- := N)} \quad (\text{UVar}^-_u)$	

Fig. 12. Unification Algorithm

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 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_u) and (\uparrow^u_u) require the 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_u) and (\exists^u_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_u) unifies two functional types. First, it unifies the argument types and their result types recursively. Then it merges the resulting constraints using the constraint merge procedure (merge) .

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 subroutines discussed here, 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 constraint 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, except that it makes sure that the resulting instantiations agree with the input constraint context Σ . 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 T can be merged as shown in fig. 13. Suppose that e_1 and e_2 are input entries. The result of the merge $e_1 \& e_2$ must be the weakest entry which implies both e_1 and e_2 .

$[T \vdash e_1 \& e_2 = e_3]$ Subtyping Constraint Entry Merge

$$\begin{array}{c}
 \frac{T \models P_1 \vee P_2 = Q}{T \vdash (\hat{\alpha}^+ : \geq P_1) \& (\hat{\alpha}^+ : \geq P_2) = (\hat{\alpha}^+ : \geq Q)} \quad (\geq \&^+ \geq) \\
 \\
 \frac{T; \cdot \models P \geq Q \dashv \cdot}{T \vdash (\hat{\alpha}^+ : \simeq P) \& (\hat{\alpha}^+ : \geq Q) = (\hat{\alpha}^+ : \simeq P)} \quad (\simeq \&^+ \geq) \\
 \\
 \frac{T; \cdot \models Q \geq P \dashv \cdot}{T \vdash (\hat{\alpha}^+ : \geq P) \& (\hat{\alpha}^+ : \simeq Q) = (\hat{\alpha}^+ : \simeq Q)} \quad (\geq \&^+ \simeq) \\
 \\
 \frac{\text{nf}(P) = \text{nf}(P')}{T \vdash (\hat{\alpha}^+ : \simeq P) \& (\hat{\alpha}^+ : \simeq P') = (\hat{\alpha}^+ : \simeq P)} \quad (\simeq \&^+ \simeq) \\
 \\
 \frac{\text{nf}(N) = \text{nf}(N')}{T \vdash (\hat{\alpha}^- : \simeq N) \& (\hat{\alpha}^- : \simeq N') = (\hat{\alpha}^- : \simeq N)} \quad (\simeq \&^- \simeq)
 \end{array}$$

Fig. 13. 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 .

- If e_2 is also a restricting 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 $(\simeq \&^- \simeq)$.

- If e_2 is a supertype constraint $\widehat{\alpha}^+ \triangleright P$, the algorithm must check that the type assigned by e_1 is a supertype of P . The corresponding symmetric rules are $(\triangleright \&^+ \simeq)$ and $(\simeq \&^+ \triangleright)$.

If both input entries are supertype constraints: $\widehat{\alpha}^+ \triangleright P$ and $\widehat{\alpha}^+ \triangleright Q$, then their conjunction is $\widehat{\alpha}^+ \triangleright P \vee Q$, as given by $(\triangleright \&^+ \triangleright)$. The least upper bound $\neg P \vee Q$ is the least supertype of both P and Q , and this way, $\widehat{\alpha}^+ \triangleright P \vee Q$ is the weakest constraint entry that implies $\widehat{\alpha}^+ \triangleright P$ and $\widehat{\alpha}^+ \triangleright Q$. The algorithm for finding the least upper bound is discussed in ??.

Merging Constraint Sets. The algorithm for merging constraint sets is shown in fig. 14. 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. 13) of matching entries.

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

Then $\Sigma \vdash C_1 \& C_2 = C$ 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
- $\Sigma(\widehat{\alpha}^\pm) \vdash e_1 \& e_2 = e$ for some $e_1 \in C_1$ and $e_2 \in C_2$ such that e_1 and e_2 both restrict variable $\widehat{\alpha}^\pm$.

Fig. 14. Constraint Merge

As shown in fig. 13, 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 ell-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 T , the *upgrade* operation finds the least among those supertypes of P that are well-formed in a smaller context $\Theta \subseteq T$. For given two types P_1 and P_2 well-formed in T , the *least upper bound* operation finds the least among common supertypes of P_1 and P_2 well-formed in T . These algorithms are shown in fig. 15.

$\text{upgrade } T \vdash P \text{ to } \Theta = Q$ Type Upgrade

$T \models P_1 \vee P_2 = Q$ Least Upper Bound

$$\begin{array}{c}
 \begin{array}{l}
 T = \Theta, \vec{\alpha}^\pm \\
 \vec{\beta}^\pm \text{ are fresh } \vec{\gamma}^\pm \text{ are fresh} \\
 \Theta, \vec{\beta}^\pm, \vec{\gamma}^\pm \models [\vec{\beta}^\pm / \vec{\alpha}^\pm] P \vee [\vec{\gamma}^\pm / \vec{\alpha}^\pm] P = Q
 \end{array} \\
 \hline
 \text{upgrade } T \vdash P \text{ to } \Theta = Q \quad (\text{UPG})
 \end{array}
 \qquad
 \begin{array}{c}
 \frac{T, \vec{\alpha}^\pm, \vec{\beta}^\pm \models P_1 \vee P_2 = Q}{T \models \exists \vec{\alpha}^\pm. P_1 \vee \exists \vec{\beta}^\pm. P_2 = Q} \quad (\exists^\vee) \\
 \frac{}{T \models \alpha^+ \vee \alpha^+ = \alpha^+} \quad (\text{VAR}^\vee) \\
 \frac{T \models \text{nf}(\downarrow N) \stackrel{a}{\simeq} \text{nf}(\downarrow M) \ni (\Upsilon, \vec{P}, \widehat{\tau}_1, \widehat{\tau}_2)}{T \models \downarrow N \vee \downarrow M = \exists \vec{\alpha}^\pm. [\vec{\alpha}^\pm / \Upsilon] \vec{P}} \quad (\downarrow^\vee)
 \end{array}$$

Fig. 15. Type Upgrade and Least Upper Bound Algorithms

The Type Upgarde. 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 $\Theta = T \setminus \vec{\alpha}^\pm$ (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 $\vec{\alpha}^-$ and $\vec{\beta}^-$ 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 $(Y, P, \widehat{\tau}_1, \widehat{\tau}_2)$, where important for us is P —the pattern and Y —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 [**todo**], 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 \widehat{Y}^- , with a pair of substitutions $\widehat{Y}^- \mapsto N_1$ and $\widehat{Y}^- \mapsto N_2$. But it cannot be more specific such as $\forall \beta^+. \widehat{Y}^+ \rightarrow \uparrow \beta^+$ (since the hole cannot be positive) or $\forall \beta^+. \widehat{Y}^-$ (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. 16. 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 $T \models P_1 \stackrel{a}{\simeq} P_2 \dashv (\Upsilon, \mathcal{Q}, \widehat{\tau}_1, \widehat{\tau}_2)$. As an input, it takes a context T , 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: Υ —a set of ‘holes’ represented by negative algorithmic variables, \mathcal{Q} —a pattern represented as a positive algorithmic type, whose algorithmic variables are in Υ , and two substitutions $\widehat{\tau}_1$ and $\widehat{\tau}_2$ instantiating the variables from Υ such that $[\widehat{\tau}_1] \mathcal{Q} = P_1$ and $[\widehat{\tau}_2] \mathcal{Q} = P_2$.

$$\begin{array}{c}
 \boxed{T \models P_1 \stackrel{a}{\simeq} P_2 \dashv (\Upsilon, \mathcal{Q}, \widehat{\tau}_1, \widehat{\tau}_2)} \\
 \\
 \frac{}{T \models \alpha^+ \stackrel{a}{\simeq} \alpha^+ \dashv (\cdot, \alpha^+, \cdot, \cdot)} \quad (\text{VAR}_+^a) \\
 \\
 \frac{T \models N_1 \stackrel{a}{\simeq} N_2 \dashv (\Upsilon, \mathcal{M}, \widehat{\tau}_1, \widehat{\tau}_2)}{T \models \downarrow N_1 \stackrel{a}{\simeq} \downarrow N_2 \dashv (\Upsilon, \downarrow \mathcal{M}, \widehat{\tau}_1, \widehat{\tau}_2)} \quad (\downarrow \stackrel{a}{\simeq}) \\
 \\
 \frac{\vec{\alpha} \cap T = \emptyset \quad T \models P_1 \stackrel{a}{\simeq} P_2 \dashv (\Upsilon, \mathcal{Q}, \widehat{\tau}_1, \widehat{\tau}_2)}{T \models \exists \vec{\alpha}. P_1 \stackrel{a}{\simeq} \exists \vec{\alpha}. P_2 \dashv (\Upsilon, \exists \vec{\alpha}. \mathcal{Q}, \widehat{\tau}_1, \widehat{\tau}_2)} \quad (\exists \stackrel{a}{\simeq}) \\
 \\
 \boxed{T \models N_1 \stackrel{a}{\simeq} N_2 \dashv (\Upsilon, \mathcal{M}, \widehat{\tau}_1, \widehat{\tau}_2)} \\
 \\
 \frac{}{T \models \alpha^- \stackrel{a}{\simeq} \alpha^- \dashv (\cdot, \alpha^-, \cdot, \cdot)} \quad (\text{VAR}_-^a) \\
 \\
 \frac{T \models P_1 \stackrel{a}{\simeq} P_2 \dashv (\Upsilon, \mathcal{Q}, \widehat{\tau}_1, \widehat{\tau}_2)}{T \models \uparrow P_1 \stackrel{a}{\simeq} \uparrow P_2 \dashv (\Upsilon, \uparrow \mathcal{Q}, \widehat{\tau}_1, \widehat{\tau}_2)} \quad (\uparrow \stackrel{a}{\simeq}) \\
 \\
 \frac{\vec{\alpha}^+ \cap T = \emptyset \quad T \models N_1 \stackrel{a}{\simeq} N_2 \dashv (\Upsilon, \mathcal{M}, \widehat{\tau}_1, \widehat{\tau}_2)}{T \models \forall \vec{\alpha}^+. N_1 \stackrel{a}{\simeq} \forall \vec{\alpha}^+. N_2 \dashv (\Upsilon, \forall \vec{\alpha}^+. \mathcal{M}, \widehat{\tau}_1, \widehat{\tau}_2)} \quad (\forall \stackrel{a}{\simeq}) \\
 \\
 \frac{T \models P_1 \stackrel{a}{\simeq} P_2 \dashv (\Upsilon_1, \mathcal{Q}, \widehat{\tau}_1, \widehat{\tau}_2) \quad T \models N_1 \stackrel{a}{\simeq} N_2 \dashv (\Upsilon_2, \mathcal{M}, \widehat{\tau}_1', \widehat{\tau}_2')}{T \models P_1 \rightarrow N_1 \stackrel{a}{\simeq} P_2 \rightarrow N_2 \dashv (\Upsilon_1 \cup \Upsilon_2, \mathcal{Q} \rightarrow \mathcal{M}, \widehat{\tau}_1 \cup \widehat{\tau}_1', \widehat{\tau}_2 \cup \widehat{\tau}_2')} \quad (\rightarrow \stackrel{a}{\simeq}) \\
 \\
 \frac{\text{if other rules are not applicable} \quad T \vdash N \quad T \vdash M}{T \models N \stackrel{a}{\simeq} M \dashv (\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. 16. Anti-Unification Algorithm

At the high level, the algorithm scheme follows the standard approach [todo] consisting of 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 . The name of 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 T is *not* extended with the quantified variables. In this algorithm, T 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:

- $Y_1 \cup Y_2$ —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 T . 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 $\hat{\gamma}^-$ (with mappings $\hat{\gamma}^- \mapsto \uparrow\alpha^+$ and $\hat{\gamma}^- \mapsto \uparrow\beta^+$), 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. 17.

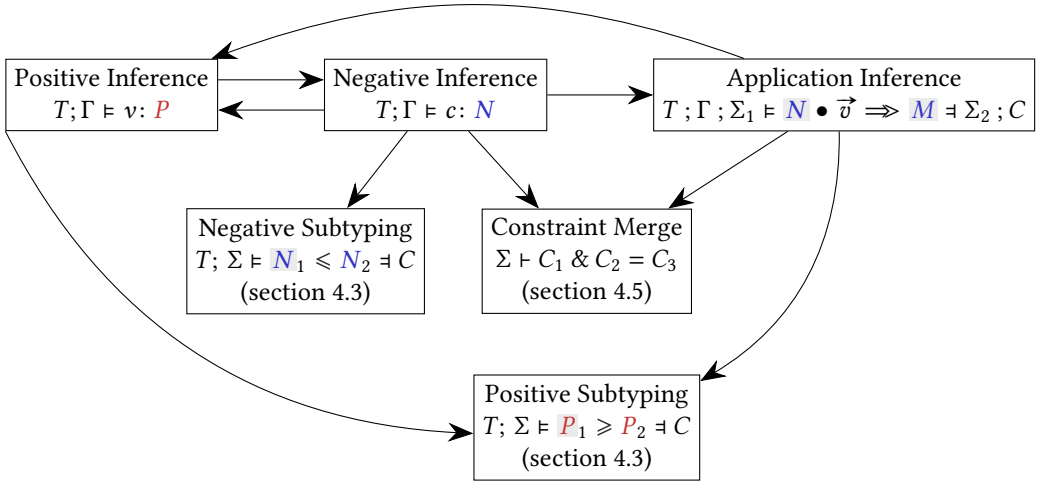


Fig. 17. Dependency graph of the typing algorithm

The positive and the negative type inference judgments have symmetric forms: $T; \Gamma \models v : P$ and $T; \Gamma \models c : N$. Both of these algorithms take as an input typing context T , a variable context Γ , and a term (a value or a computation) taking its type variables from T , 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 $T; \Gamma; \Sigma_1 \vdash N \bullet \vec{v} \Rightarrow M \models \Sigma_2; C$. As an input, it takes three contexts: typing context T , a variable context Γ , and a constraint context Σ_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 Σ_1 , in other words, $T; \text{dom}(\Sigma_1) \vdash N$. 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 Σ_2 —an updated constraint context and C —a set of subtyping constraints. Together Σ_2 and C specify how the algorithmic variables must be instantiated.

The inference rules are shown in fig. 18. 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 T . 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 $T \vdash N$.

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 ($\text{LET}_{\text{@}}^{\text{INF}}$) 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 $\uparrow Q$, that must be instantiated to a subtype of the annotation $\uparrow P$. Then premise $T; \Sigma \vdash \uparrow Q \leq \uparrow P \models C_2$ together with $\Sigma \vdash C_1 \ \& \ C_2 = C$ check whether the instantiation to the annotated type $\uparrow P$ is possible, and if it is, the algorithm infers the type of the body in the extended context, and returns it as the result.

Rule ($\text{LET}_{\text{@}}^{\text{INF}}$) works similarly to (LET^{INF}). However, since there is no annotation, instead of checking that the inferred type $\uparrow Q$ can be a subtype of the annotation, the algorithm checks that $\uparrow Q$ is unique. As we prove, uniqueness means that all the algorithmic variables of $\uparrow Q$ are sufficiently restricted by C . This way, uniqueness is checked by the combination of $\text{fav } Q = \text{dom}(C)$

$\boxed{T; \Gamma \vdash v : P}$ Positive typing

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

 $\boxed{T; \Gamma \vdash c : N}$ Negative typing

$$\frac{T \vdash M \quad T; \Gamma \vdash c : N}{T; \cdot \vdash N \leq M \dashv} \quad (\text{ANN}_-^{\text{INF}}) \quad \frac{T; \Gamma \vdash v : P}{T; \Gamma \vdash \mathbf{return} v : \uparrow P} \quad (\text{RET}^{\text{INF}})$$

$$\frac{T \vdash P \quad T; \Gamma, x : P \vdash c : N}{T; \Gamma \vdash \lambda x : P. c : \mathbf{nf}(P \rightarrow N)} \quad (\lambda^{\text{INF}}) \quad \frac{T; \Gamma \vdash v : P \quad T; \Gamma, x : P \vdash c : N}{T; \Gamma \vdash \mathbf{let} x = v; c : N} \quad (\text{LET}^{\text{INF}})$$

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

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

$$\frac{T; \Gamma \vdash v : \downarrow M \quad T; \Gamma; \cdot \vdash M \bullet \vec{v} \Rightarrow \uparrow Q \dashv \Sigma; C \quad \mathbf{fav} Q = \mathbf{dom}(C) \quad C \text{ singular with } \hat{\sigma} \quad T; \Gamma, x : [\hat{\sigma}] Q \vdash c : N}{T; \Gamma \vdash \mathbf{let} x = v(\vec{v}); c : N} \quad (\text{LET}_{@}^{\text{INF}})$$

 $\boxed{T; \Gamma; \Sigma_1 \vdash N \bullet \vec{v} \Rightarrow M \dashv \Sigma_2; C}$ Application typing

$$\frac{}{T; \Gamma; \Sigma \vdash N \bullet \cdot \Rightarrow \mathbf{nf}(N) \dashv \Sigma; \cdot} \quad (\emptyset^{\text{INF}}) \quad \frac{T; \Gamma \vdash v : P \quad T; \Sigma \vdash Q \geq P \dashv C_1 \quad T; \Gamma; \Sigma \vdash N \bullet \vec{v} \Rightarrow M \dashv \Sigma'; C_2 \quad \Sigma \vdash C_1 \ \& \ C_2 = C}{T; \Gamma; \Sigma \vdash Q \rightarrow N \bullet v, \vec{v} \Rightarrow M \dashv \Sigma'; C} \quad (\rightarrow^{\text{INF}})$$

$$\frac{T; \Gamma; \Sigma, \vec{\alpha}^+ \{T\} \vdash [\vec{\alpha}^+ / \alpha^+] N \bullet \vec{v} \Rightarrow M \dashv \Sigma'; C \quad \vec{\alpha}^+ \text{ are fresh} \quad \vec{v} \neq \cdot \quad \vec{\alpha}^+ \neq \cdot}{T; \Gamma; \Sigma \vdash \forall \vec{\alpha}^+. N \bullet \vec{v} \Rightarrow M \dashv \Sigma'; C|_{\mathbf{fav}(N) \cup \mathbf{fav}(M)}} \quad (\forall^{\text{INF}})$$

Fig. 18. Algorithmic Type Inference

and C singular with $\hat{\sigma}$. Together, these two premises guarantee that the only possible instantiation of Q is $[\hat{\sigma}] Q$.

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 $\vec{v} \neq \cdot$.

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 $\mathcal{Q} \rightarrow \mathcal{N}$. First, it infers the type of the first argument v , and then, calling the algorithmic subtyping, finds C_1 —the minimal constraint ensuring that \mathcal{Q} is a supertype of the type of v . Then it makes a recursive call applying \mathcal{N} to the rest of the arguments and merges the resulting constraint with C_1 .

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 \mathcal{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— \mathcal{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 \mathcal{M} , that is $T ; \Gamma \vdash [\hat{\sigma}] \mathcal{N} \bullet \vec{v} \Rightarrow [\hat{\sigma}] \mathcal{M}$.

The application inference completeness means that if there exists a substitution $\hat{\sigma}$ and the resulting type \mathcal{M} ensuring the declarative inference $T ; \Gamma \vdash [\hat{\sigma}] \mathcal{N} \bullet \vec{v} \Rightarrow \mathcal{M}$ then the algorithm succeeds and gives the most general result \mathcal{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 \mathcal{M} , if it ensures the declarative inference, then $\hat{\sigma}$ can be extended in a C_0 -complying way to equate \mathcal{M}_0 with \mathcal{M} .

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

- + $T ; \Gamma \models v : \mathcal{P}$ implies $T ; \Gamma \vdash v : \mathcal{P}$,
- $T ; \Gamma \models c : \mathcal{N}$ implies $T ; \Gamma \vdash c : \mathcal{N}$,
- $T ; \Gamma ; \Sigma \models \mathcal{N} \bullet \vec{v} \Rightarrow \mathcal{M} \models \Sigma' ; C$ implies $T ; \Gamma \vdash [\hat{\sigma}] \mathcal{N} \bullet \vec{v} \Rightarrow [\hat{\sigma}] \mathcal{M}$, for any instantiation of $\hat{\sigma}$ satisfying constraints C , assuming $T \vdash^2 \Sigma$ and $T ; \text{dom}(\Sigma) \vdash \mathcal{N}$ and that \mathcal{N} is free from negative algorithmic variables.

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

- + $T ; \Gamma \vdash v : \mathcal{P}$ implies $T ; \Gamma \models v : \text{nf}(\mathcal{P})$,
- $T ; \Gamma \vdash c : \mathcal{N}$ implies $T ; \Gamma \models c : \text{nf}(\mathcal{N})$,

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.

- If $T ; \Gamma \vdash [\hat{\sigma}] \mathbf{N} \bullet \vec{\sigma} \Rightarrow M$ where (1) $T \vdash^{\neg} \Sigma$, (2) $T \vdash M$, (3) $T ; \mathbf{dom}(\Sigma) \vdash \mathbf{N}$ (free from negative algorithmic variables), and (4) $\Sigma \vdash \hat{\sigma} : \mathbf{fav}(\mathbf{N})$, then there exist M_0 , Σ_0 , and C_0 such that
 - (1) $T ; \Gamma ; \Sigma \vdash \mathbf{N} \bullet \vec{\sigma} \Rightarrow M_0 \equiv \Sigma_0 ; C_0$ and
 - (2) for any other $\hat{\sigma}$ and M (where $\Sigma \vdash \hat{\sigma} : \mathbf{fav}(\mathbf{N})$ and $T \vdash M$) such that $T ; \Gamma \vdash [\hat{\sigma}] \mathbf{N} \bullet \vec{\sigma} \Rightarrow M$, there exists $\hat{\sigma}'$ such that (a) $\Sigma_0 \vdash \hat{\sigma}' : \mathbf{fav} \mathbf{N} \cup \mathbf{fav} M_0$ and $\Sigma_0 \vdash \hat{\sigma}' : C_0$, (b) $\Sigma \vdash \hat{\sigma}' \simeq^{\leq} \hat{\sigma} : \mathbf{fav} \mathbf{N}$, and (c) $T \vdash [\hat{\sigma}'] M_0 \simeq^{\leq} M$.

The proof of soundness and completeness result is done gradually for all the subroutines, following the structure of the algorithm (figs. 10 and 17) 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 T , we have $T \vdash \mathbf{N} \simeq^{\leq} M \iff \mathbf{nf}(\mathbf{N}) = \mathbf{nf}(M)$ and $T \vdash P \simeq^{\leq} Q \iff \mathbf{nf}(P) = \mathbf{nf}(Q)$.*

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

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 T , 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 $T \models P_1 \stackrel{a}{=} P_2 \equiv (\mathbf{Y}, \mathbf{Q}, \widehat{\tau}_1, \widehat{\tau}_2)$ then (1) $T ; \mathbf{Y} \vdash \mathbf{Q}$, (2) $T ; \cdot \vdash \widehat{\tau}_i : \mathbf{Y}$ for $i \in \{1, 2\}$ are anti-unification substitutions (in particular, \mathbf{Y} contains only negative algorithmic variables), and (3) $[\widehat{\tau}_i] \mathbf{Q} = P_i$ for $i \in \{1, 2\}$.*

Lemma (Completeness of (Positive) Anti-Unification). *Assuming that P_1 and P_2 are normalized terms well-formed in T and there exist $(\mathbf{Y}', \mathbf{Q}', \widehat{\tau}'_1, \widehat{\tau}'_2)$ such that (1) $T ; \mathbf{Y}' \vdash \mathbf{Q}'$, (2) $T ; \cdot \vdash \widehat{\tau}'_i : \mathbf{Y}'$ for $i \in \{1, 2\}$ are anti-unification substitutions, and (3) $[\widehat{\tau}'_i] \mathbf{Q}' = P_i$ for $i \in \{1, 2\}$.*

Then the anti-unification algorithm terminates, that is there exists $(\mathbf{Y}, \mathbf{Q}, \widehat{\tau}_1, \widehat{\tau}_2)$ such that $T \models P_1 \stackrel{a}{=} P_2 \equiv (\mathbf{Y}, \mathbf{Q}, \widehat{\tau}_1, \widehat{\tau}_2)$.

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 $T ; \Upsilon_2 \vdash \widehat{\rho} : \Upsilon_1$ replacing the placeholders Υ_1 with types that themselves can contain placeholders from Υ_2 .

Lemma (Initiality of Anti-Unification). *Assume that P_1 and P_2 are normalized types well-formed in T , and the anti-unification algorithm succeeds: $T \models P_1 \stackrel{a}{\approx} P_2 \dashv (\Upsilon, \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2)$. Then $(\Upsilon, \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2)$ is more specific than any other sound anti-unifier $(\Upsilon', \underline{Q}', \widehat{\tau}'_1, \widehat{\tau}'_2)$, i.e., if (1) $T ; \Upsilon' \vdash \underline{Q}'$, (2) $T ; \cdot \vdash \widehat{\tau}'_i : \Upsilon'$ for $i \in \{1, 2\}$, and (3) $[\widehat{\tau}'_i] \underline{Q}' = P_i$ for $i \in \{1, 2\}$ then there exists a ‘refining’ substitution $\widehat{\rho}$ such that $T ; \Upsilon \vdash \widehat{\rho} : (\Upsilon')|_{\text{fav } \underline{Q}'}$ and $[\widehat{\rho}] \underline{Q}' = \underline{Q}$.*

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 $(\Upsilon, \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2)$, there are no two different placeholders $\widehat{\beta}^-$ mapped to the same pair of types by $\widehat{\tau}_1$ and $\widehat{\tau}_2$. 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 $\downarrow \widehat{\alpha}^- \rightarrow \widehat{\alpha}^-$, but not (less specific) $\downarrow \widehat{\alpha}^- \rightarrow \widehat{\beta}^-$. 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 $T \models P_1 \stackrel{a}{\approx} P_2 \dashv (\Upsilon, \underline{Q}, \widehat{\tau}_1, \widehat{\tau}_2)$ then for any $\widehat{\beta}^- \in \Upsilon$, $\widehat{\beta}^- = \widehat{\alpha}^-_{\{[\widehat{\tau}_1] \widehat{\beta}^-, [\widehat{\tau}_2] \widehat{\beta}^-\}}$.*

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 $T \vdash P_1$, and $T \vdash P_2$, if $T \models P_1 \vee P_2 = \underline{Q}$ then*

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

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

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 T , the set of normalized T -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 T -formed supertypes are the $(T, \overrightarrow{\beta^-})$ -formed supertypes of P' not using β^- ;

- if P is a downshift type $\downarrow M$, its supertypes have form $\exists \vec{\alpha}^{\rightarrow}. \downarrow M'$ such that there exists a T -formed instantiation of $\vec{\alpha}^{\rightarrow}$ in $\downarrow M'$ that makes $\downarrow M'$ equal to $\downarrow M$, i.e., $[\vec{N}/\vec{\alpha}^{\rightarrow}] \downarrow M' = \downarrow M$.

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 $T = \Theta, \vec{\alpha}^{\rightarrow}$, if $\text{upgrade } T \vdash P$ to $\Theta = Q$ then (1) $\Theta \vdash Q$ (2) $T \vdash Q \geq P$*

Lemma (Upgrade Completeness). *Assuming P is well-formed in $T = \Theta, \vec{\alpha}^{\rightarrow}$, for any Q' such that Q' is a Θ -formed upper bound of P , i.e., (1) $\Theta \vdash Q'$ and (2) $T \vdash Q' \geq P$, there exists Q such that $(\text{upgrade } T \vdash P \text{ to } \Theta = Q)$ and $\Theta \vdash Q' \geq Q$.*

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 ??, 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: $T \vdash [\hat{\sigma}] P \geq Q$. 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 $T \vdash^{\geq} \Sigma$ and $T \vdash Q$ and $T; \text{dom}(\Sigma) \vdash P$ and $T; \Sigma \models P \geq Q \models C$, then $\Sigma \vdash C : \text{fav } P$ and for any $\hat{\sigma}$ such that $\Sigma \vdash \hat{\sigma} : C$, we have $T \vdash [\hat{\sigma}] P \geq Q$.*

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 $T \vdash^{\geq} \Sigma$, $T \vdash Q$ and $T; \text{dom}(\Sigma) \vdash P$. Then if there exists $\hat{\sigma}$ such that $\Sigma \vdash \hat{\sigma} : \text{fav}(P)$ and $T \vdash [\hat{\sigma}] P \geq Q$, then the subtyping algorithm succeeds: $T; \Sigma \models P \geq Q \models C$.*

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 $T \vdash^{\geq} \Sigma$ and $T \vdash M$ and $T; \text{dom}(\Sigma) \vdash N$ and N does not contain negative unification variables ($\vec{\alpha}^- \notin \text{fav } N$). Then for any $\Sigma \vdash \hat{\sigma} : \text{fav}(N)$ such that $T \vdash [\hat{\sigma}] N \leq M$, the subtyping algorithm succeeds: $T; \Sigma \models N \leq M \models C$.*

5.5 Singularity

The rule ($\text{LET}_{@}^{\text{INF}}$) of the declarative inference says that one can omit the type annotation of a let-binding if the inferred type of the bound application is unique. Accordingly, the algorithm, that returns the constraint set, must be able to check that the set of constraints has a single solution.

and also to find this solution. This subroutine is called singularity, and its correctness is formulated as follows.

The soundness 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 $\Sigma \vdash C : \Upsilon$, and C singular with $\widehat{\sigma}$. Then $\Sigma \vdash \widehat{\sigma} : \Upsilon$, $\Sigma \vdash \widehat{\sigma} : C$ and for any $\widehat{\sigma}'$ such that $\Sigma \vdash \widehat{\sigma}' : C$, $\Sigma \vdash \widehat{\sigma}' \simeq^{\leq} \widehat{\sigma} : \Upsilon$.*

Lemma (Completeness of Singularity). *For a given $\Sigma \vdash C$, suppose that all the substitutions satisfying C are equivalent on $\mathbf{dom}(C)$. In other words, suppose that there exists $\Sigma \vdash \widehat{\sigma}_1 : \mathbf{dom}(C)$ such that for any $\Sigma \vdash \widehat{\sigma} : \mathbf{dom}(C)$, $\Sigma \vdash \widehat{\sigma} : C$ implies $\Sigma \vdash \widehat{\sigma} \simeq^{\leq} \widehat{\sigma}_1 : \mathbf{dom}(C)$. 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 $T \vdash \Gamma$ and $T \vdash^{\geq} \Sigma$. Then*

- + If $T; \Gamma \models v : P$ and $T; \Gamma \models v : P'$ then $P = P'$.
- If $T; \Gamma \models c : N$ and $T; \Gamma \models c : N'$ then $N = N'$.
- If $T; \Gamma; \Sigma \models N \bullet \vec{v} \Rightarrow M \models \Sigma'; C$ and $T; \Gamma; \Sigma \models N \bullet \vec{v} \Rightarrow M' \models \Sigma'; C'$ then $M = M'$, $\Sigma = \Sigma'$, 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 $(\text{LET}_{\odot}^{\text{INF}})$. To prove the soundness, in other words, to transform this rule into its declarative counterpart $(\text{LET}_{\odot}^{\text{INF}})$, 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 (??), 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 $T; \Gamma \vdash v : P$ or $T; \Gamma \vDash c : N$), it is always zero. For application inference judgments ($T; \Gamma \vdash N \bullet \vec{v} \Rightarrow M$ or $T; \Gamma; \Sigma \vDash N \bullet \vec{v} \Rightarrow M \vDash \Sigma'; C$), 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 ($\forall_{\bullet \Rightarrow}^{\text{INF}}$) and ($\forall_{\bullet \Rightarrow}^{\text{INF}}$) 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:

$$\begin{aligned} &+ \text{size}(T; \Gamma \vdash v : P) = (\text{size}(v), 0); &+ \text{size}(T; \Gamma \vDash v : P) = (\text{size}(v), 0); \\ &- \text{size}(T; \Gamma \vdash c : N) = (\text{size}(c), 0); &- \text{size}(T; \Gamma \vDash c : N) = (\text{size}(c), 0); \\ &\bullet \text{size}(T; \Gamma \vdash N \bullet \vec{v} \Rightarrow M) = &\bullet \text{size}(T; \Gamma; \Sigma \vDash N \bullet \vec{v} \Rightarrow M \vDash \Sigma'; C) = \\ &(\text{size}(\vec{v}), \text{npq}(N)); &(\text{size}(\vec{v}), \text{npq}(N)). \end{aligned}$$

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.,

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

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_{-}^{INF}) and (\simeq_{+}^{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 $T; \Gamma \vDash \text{let } x = v(\vec{v}); c' : N$, we can assume, that *completeness* holds for a term of shape $T; \Gamma \vdash M \bullet \vec{v} \Rightarrow K$, 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 $T; \Gamma \vdash M \bullet \vec{v} \Rightarrow \uparrow[\widehat{\sigma}] Q$ constructed by the induction hypothesis) is unique.

6 EXTENSIONS

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:

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

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:

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

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:

$$\begin{array}{c} \frac{T, \vec{\alpha}^+ \vdash N \leq M}{T \vdash \forall \vec{\alpha}^+. N \leq \forall \vec{\alpha}^+. M} \quad (E\forall_R^\leq) \qquad \frac{M \neq \forall \vec{\beta}^+. M' \quad T \vdash \sigma : \vec{\alpha}^+ \quad T \vdash [\sigma]N \leq M}{T \vdash \forall \vec{\alpha}^+. N \leq M} \quad (E\forall_L^\leq) \\[10pt] \frac{T, \vec{\alpha}^+ \vdash P \geq Q}{T \vdash \exists \vec{\alpha}^+. P \geq \exists \vec{\alpha}^+. Q} \quad (E\exists_R^\geq) \qquad \frac{Q \neq \exists \vec{\beta}^+. Q' \quad T \vdash \sigma : \vec{\alpha}^+ \quad T \vdash [\sigma]P \geq Q}{T \vdash \exists \vec{\alpha}^+. P \geq Q} \quad (E\exists_L^\geq) \end{array}$$

From the perspective of mutual subtyping, these changes fix the order of the quantifiers for an equivalence class. Moreover, the equivalence degenerates to equality (alpha-equivalence).

To accommodate these changes in the *algorithm*, it suffices to (i) replace the normalization procedure with identity: $\text{nf}(N) \stackrel{\text{def}}{=} N$, $\text{nf}(P) \stackrel{\text{def}}{=} P$, (ii) modify the least upper bound polymorphic rule (\exists^\vee) in the following way:

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

and (iii) replace the subtyping polymorphic rule (\forall^\leq) by the following pair of rules:

$$\frac{T; \vec{\alpha}^+; \Sigma \models N \leq M \dashv C}{T; \Sigma \models \forall \vec{\alpha}^+. N \leq \forall \vec{\alpha}^+. M \dashv C} \quad (\forall_R^{\leq}) \qquad \frac{\vec{\alpha}^+ \text{ are fresh } M \neq \forall \vec{\beta}^+. M' \quad T; \Sigma, \vec{\alpha}^+ \{T\} \models [\vec{\alpha}^+ / \vec{\alpha}^+] N \leq M \dashv C}{T; \Sigma \models \forall \vec{\alpha}^+. N \leq M \dashv C \setminus \vec{\alpha}^+} \quad (\forall_L^{\leq})$$

6.2 Bounded Quantification

It is possible to smoothly extend the type system with bounded quantifiers. In particular, we can add lower bounds to polymorphic \forall -quantifiers: $\forall(\vec{\alpha}^+ \geq \vec{P}). N$ with the expected subtyping specification:

???

6.3 Bidirectionalization and Invariants Relaxation

The algorithm we provide requires that all lambda functions are annotated. To augment the system's expressiveness and lessen this requirement, we can employ bidirectionalization.

However, this approach will disrupt the crucial invariants of the system and necessitate a more sophisticated constraint solver. Specifically, the generated constraints may contain algorithmic variables on both sides of the subtyping relation. If the constraint set only imposes an equivalence restriction, this task aligns with pattern-unification. Nevertheless, the subtyping restrictions render the problem undecidable. For instance, ...

7 CONCLUSION

[Botlan et al. 2003] [dunfieldBidirectionalTyping2020]

REFERENCES

- Didier Le Botlan and Didier Rémy (Aug. 2003). "MLF Raising ML to the Power of System F." In: *ICFP '03*. Uppsala, Sweden: ACM Press, pp. 52–63.
- 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.
- Jana Dunfield and Neel Krishnaswami (Nov. 2020). "Bidirectional Typing." In: arXiv: 1908.05839.
- 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.
- 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.
- 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 Pawel 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

1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567

1568