

# Consistent Subtyping for All

Ningning Xie<sup>✉</sup>, Xuan Bi, and Bruno C. d. S. Oliveira

The University of Hong Kong  
{nnxie, xbi, bruno}@cs.hku.hk

**Abstract.** Consistent subtyping is employed in some gradual type systems to validate type conversions. The original definition by Siek and Taha serves as a guideline for designing gradual type systems with subtyping. Polymorphic types à la System F also induce a subtyping relation that relates polymorphic types to their instantiations. However Siek and Taha’s definition is not adequate for polymorphic subtyping. The first goal of this paper is to propose a generalization of consistent subtyping that is adequate for polymorphic subtyping, and subsumes the original definition by Siek and Taha. The new definition of consistent subtyping provides novel insights with respect to previous polymorphic gradual type systems, which did not employ consistent subtyping. The second goal of this paper is to present a gradually typed calculus for implicit (higher-rank) polymorphism that uses our new notion of consistent subtyping. We develop both declarative and (bidirectional) algorithmic versions for the type system. We prove that the new calculus satisfies all static aspects of the refined criteria for gradual typing, which are mechanically formalized using the Coq proof assistant.

## 1 Introduction

Gradual typing [21] is an increasingly popular topic in both programming language practice and theory. On the practical side there is a growing number of programming languages adopting gradual typing. Those languages include Clojure [6], Python [27], TypeScript [5], Hack [26], and the addition of Dynamic to C# [4], to cite a few. On the theoretical side, recent years have seen a large body of research that defines the foundations of gradual typing [13, 8, 9], explores their use for both functional and object-oriented programming [21, 22], as well as its applications to many other areas [24, 3].

A key concept in gradual type systems is *consistency* [21]. Consistency weakens type equality to allow for the presence of *unknown* types. In some gradual type systems with subtyping, consistency is combined with subtyping to give rise to the notion of *consistent subtyping* [22]. Consistent subtyping is employed by gradual type systems to validate type conversions arising from conventional subtyping. One nice feature of consistent subtyping is that it is derivable from the more primitive notions of *consistency* and *subtyping*. As Siek and Taha [22] put it this shows that “*gradual typing and subtyping are orthogonal and can be combined in a principled fashion*”. Thus consistent subtyping is often used as a guideline for designing gradual type systems with subtyping.

Unfortunately, as noted by Garcia et al. [13], notions of consistency and/or consistent subtyping “become more difficult to adapt as type systems get more complex”. In particular, for the case of type systems with subtyping, certain kinds of subtyping do not fit well with the original definition of consistent subtyping by Siek and Taha [22]. One important case where such mismatch happens is in type systems supporting implicit (higher-rank) polymorphism [18, 11]. It is well-known that polymorphic types à la System F induce a subtyping relation that relates polymorphic types to their instantiations [17, 16]. However Siek and Taha’s definition is not adequate for this kind of subtyping. Moreover the current framework for *Abstracting Gradual Typing* (AGT) [13] also does not account for polymorphism, with the authors acknowledging that this is one of the interesting avenues for future work.

Existing work on gradual type systems with polymorphism does not use consistent subtyping. The Polymorphic Blame Calculus ( $\lambda\mathbf{B}$ ) [1] is an *explicitly* polymorphic calculus with explicit casts, which is often used as a target language for gradual type systems with polymorphism. In  $\lambda\mathbf{B}$  a notion of *compatibility* is employed to validate conversions allowed by casts. Interestingly  $\lambda\mathbf{B}$  allows conversions from polymorphic types to their instantiations. For example, it is possible to cast a value with type  $\forall a.a \rightarrow a$  into  $\text{Int} \rightarrow \text{Int}$ . Thus an important remark here is that while  $\lambda\mathbf{B}$  is explicitly polymorphic, casting and conversions are closer to *implicit* polymorphism. That is, in a conventional explicitly polymorphic calculus (such as System F), the primary notion is type equality, where instantiation is not taken into account. Thus the types  $\forall a.a \rightarrow a$  and  $\text{Int} \rightarrow \text{Int}$  are deemed *incompatible*. However in *implicitly* polymorphic calculi [18, 11]  $\forall a.a \rightarrow a$  and  $\text{Int} \rightarrow \text{Int}$  are deemed *compatible*, since the latter type is an instantiation of the former. Therefore  $\lambda\mathbf{B}$  is in a sense a hybrid between implicit and explicit polymorphism, utilizing type equality (à la System F) for validating applications, and *compatibility* for validating casts.

An alternative approach to polymorphism has recently been proposed by Igarashi et al. [14]. Like  $\lambda\mathbf{B}$  their calculus is explicitly polymorphic. However, in that work they employ type consistency to validate cast conversions, and forbid conversions from  $\forall a.a \rightarrow a$  to  $\text{Int} \rightarrow \text{Int}$ . This makes their casts closer to explicit polymorphism, in contrast to  $\lambda\mathbf{B}$ . Nonetheless, there is still same flavour of implicit polymorphism in their calculus when it comes to interactions between dynamically typed and polymorphically typed code. For example, in their calculus type consistency allows types such as  $\forall a.a \rightarrow \text{Int}$  to be related to  $\star \rightarrow \text{Int}$ , where some sort of (implicit) polymorphic subtyping is involved.

The first goal of this paper is to study the gradually typed subtyping and consistent subtyping relations for *predicative implicit polymorphism*. To accomplish this, we first show how to reconcile consistent subtyping with polymorphism by generalizing the original consistent subtyping definition by Siek and Taha. The new definition of consistent subtyping can deal with polymorphism, and preserves the orthogonality between consistency and subtyping. To slightly rephrase Siek and Taha, the motto of our paper is that:

*Gradual typing and **polymorphism** are orthogonal and can be combined in a principled fashion.*<sup>1</sup>

With the insights gained from our work, we argue that, for implicit polymorphism, Ahmed et al.’s notion of compatibility is too permissive (i.e. too many programs are allowed to type-check), and that Igarashi et al.’s notion of type consistency is too conservative. As a step towards an algorithmic version of consistent subtyping, we present a syntax-directed version of consistent subtyping that is sound and complete with respect to our formal definition of consistent subtyping. The syntax-directed version of consistent subtyping is remarkably simple and well-behaved, without the ad-hoc *restriction* operator [22]. Moreover, to further illustrate the generality of our consistent subtyping definition, we show that it can also account for *top types*, which cannot be dealt with by Siek and Taha’s definition either.

The second goal of this paper is to present a (source-level) gradually typed calculus for (predicative) implicit higher-rank polymorphism that uses our new notion of consistent subtyping. As far as we are aware, there is no work on bridging the gap between implicit higher-rank polymorphism and gradual typing, which is interesting for two reasons. On one hand, modern functional languages (such as Haskell) employ sophisticated type-inference algorithms that, aided by type annotations, can deal with implicit higher-rank polymorphism. So a natural question is how gradual typing can be integrated in such languages. On the other hand, there is several existing work on integrating *explicit* polymorphism into gradual typing [1, 14]. Yet no work investigates how to move such expressive power into a source language with implicit polymorphism. Therefore as a step towards gradualizing such type systems, this paper develops both declarative and algorithmic versions for a gradual type system with implicit higher-rank polymorphism. The new calculus brings the expressive power of full implicit higher-rank polymorphic into a gradually typed source language. We prove that our calculus satisfies all of the *static* aspects of the refined criteria for gradual typing [25], while discussing some issues related with the *dynamic guarantee*.

In summary, the contributions of this paper are:

- We define a framework for consistent subtyping with:
  - a new definition of consistent subtyping that subsumes and generalizes that of Siek and Taha, and can deal with polymorphism and top types.
  - a syntax-directed version of consistent subtyping that is sound and complete with respect to our definition of consistent subtyping, but still guesses polymorphic instantiations.
- Based on consistent subtyping, we present a declarative gradual type system with predicative implicit higher-rank polymorphism. We prove that our calculus satisfies the static aspects of the refined criteria for gradual typing [25], and is type-safe by a type-directed translation to  $\lambda\mathbf{B}$ , and thus hereditarily preserves parametricity [2].

---

<sup>1</sup> Note here that we borrow Siek and Taha’s motto mostly to talk about the static semantics. As Ahmed et al. [1] show there are several non-trivial interactions between polymorphism and casts at the level of the dynamic semantics.

$$\begin{array}{c}
\boxed{A <: B} \\
\text{Int } <: \text{Int} \quad \text{Bool } <: \text{Bool} \quad \text{Float } <: \text{Float} \quad \text{Int } <: \text{Float} \\
\frac{B_1 <: A_1 \quad A_2 <: B_2}{A_1 \rightarrow A_2 <: B_1 \rightarrow B_2} \quad [l_i : A_i^{i \in 1 \dots n+m}] <: [l_i : A_i^{i \in 1 \dots n}] \quad \star <: \star \\
\boxed{A \sim B} \\
A \sim A \quad A \sim \star \quad \star \sim A \quad \frac{A_1 \sim B_1 \quad A_2 \sim B_2}{A_1 \rightarrow A_2 \sim B_1 \rightarrow B_2} \quad \frac{A_i \sim B_i}{[l_i : A_i] \sim [l_i : B_i]}
\end{array}$$

Fig. 1: Subtyping and type consistency in  $\mathbf{FOb}_{<}^?$ .

- We present a complete and sound bidirectional algorithm for implementing the declarative system based on the design principle of Garcia and Cimini [12] and the approach of Dunfield and Krishnaswami [11].
- All of the metatheory of this paper, except some manual proofs for the algorithmic type system, has been mechanically formalized in Coq<sup>2</sup>.

## 2 Background and Motivation

In this section we review a simple gradually typed language with objects [22], to introduce the concept of consistency subtyping. We also briefly talk about the Odersky-Läufer type system for higher-rank types [17], which serves as the original language on which our gradually typed calculus with implicit higher-rank polymorphism is based.

### 2.1 Gradual Subtyping

Siek and Taha [22] developed a gradual typed system for object-oriented languages that they call  $\mathbf{FOb}_{<}^?$ . Central to gradual typing is the concept of *consistency* (written  $\sim$ ) between gradual types, which are types that may involve the unknown type  $\star$ . The intuition is that consistency relaxes the structure of a type system to tolerate unknown positions in a gradual type. They also defined the subtyping relation in a way that static type safety is preserved. Their key insight is that the unknown type  $\star$  is neutral to subtyping, with only  $\star <: \star$ . Both relations are found in Fig. 1.

A primary contribution of their work is to show that consistency and subtyping are orthogonal. To compose subtyping and consistency, Siek and Taha defined *consistent subtyping* (written  $\lesssim$ ) in two equivalent ways:

<sup>2</sup> All supplementary materials are available at <https://bitbucket.org/xieningning/consistent-subtyping>

**Definition 1 (Consistent Subtyping à la Siek and Taha [22]).**

- $A \lesssim B$  if and only if  $A \sim C$  and  $C <: B$  for some  $C$ .
- $A \lesssim B$  if and only if  $A <: C$  and  $C \sim B$  for some  $C$ .

Both definitions are non-deterministic because of the intermediate type  $C$ . To remove non-determinism, they proposed a so-called *restriction operator*, written  $A|_B$  that masks off the parts of a type  $A$  that are unknown in a type  $B$ .

$$\begin{aligned}
 A|_B = & \mathbf{case} \ A, B \ \mathbf{of} \ | \ (-, \star) \Rightarrow \star \\
 & | \ A_1 \rightarrow A_2, B_1 \rightarrow B_2 = A_1|_{B_1} \rightarrow A_2|_{B_2} \\
 & | \ [l_1 : A_1, \dots, l_n : A_n], [l_1 : B_1, \dots, l_m : B_m] \ \mathbf{if} \ n \leq m \Rightarrow [l_1 : A_1|_{B_1}, \dots, l_n : A_n|_{B_n}] \\
 & | \ [l_1 : A_1, \dots, l_n : A_n], [l_1 : B_1, \dots, l_m : B_m] \ \mathbf{if} \ n > m \Rightarrow \\
 & \quad [l_1 : A_1|_{B_1}, \dots, l_m : A_m|_{B_m}, \dots, l_n : A_n] \\
 & | \ \mathbf{otherwise} \Rightarrow A
 \end{aligned}$$

With the restriction operator, consistent subtyping is simply defined as  $A \lesssim B \equiv A|_B <: B|_A$ . Then they proved that this definition is equivalent to Definition 1.

## 2.2 The Odersky-Läufer Type System

The calculus we are combining gradual typing with is the well-established predicative type system for higher-rank types proposed by Odersky and Läufer [17]. One difference is that, for simplicity, we do not account for a let expression, as there is already existing work about gradual type systems with let expressions and let generalization (for example, see Garcia and Cimini [12]). Similar techniques can be applied to our calculus to enable let generalization.

The syntax of the type system, along with the typing and subtyping judgments is given in Fig. 2. An implicit assumption throughout the paper is that variables in contexts are distinct. We save the explanations for the static semantics to Section 4, where we present our gradually typed version of the calculus.

## 2.3 Motivation: Gradually Typed Higher-Rank Polymorphism

Our work combines implicit (higher-rank) polymorphism with gradual typing. As is well known, a gradually typed language supports both fully static and fully dynamic checking of program properties, as well as the continuum between these two extremes. It also offers programmers fine-grained control over the static-to-dynamic spectrum, i.e., a program can be evolved by introducing more or less precise types as needed [13].

Haskell is a language that supports implicit higher-rank polymorphism, but no gradual typing. Therefore some programs that are safe at run-time may be rejected due to the conservativity of the type system. For example, consider the following Haskell program adapted from Peyton Jones et al. [18]:

```
foo :: ([Int], [Char])
foo = let f x = (x [1, 2], x ['a', 'b']) in f reverse
```

Expressions	$e ::= x \mid n \mid \lambda x : A. e \mid \lambda x. e \mid e e$
Types	$A, B ::= \text{Int} \mid a \mid A \rightarrow B \mid \forall a. A$
Monotypes	$\tau, \sigma ::= \text{Int} \mid a \mid \tau \rightarrow \sigma$
Contexts	$\Psi ::= \emptyset \mid \Psi, x : A \mid \Psi, a$

$\Psi \vdash^{OL} e : A$

$\frac{x : A \in \Psi}{\Psi \vdash^{OL} x : A} \text{VAR}$	$\frac{}{\Psi \vdash^{OL} n : \text{Int}} \text{NAT}$	$\frac{\Psi, x : A \vdash^{OL} e : B}{\Psi \vdash^{OL} \lambda x : A. e : A \rightarrow B} \text{LAMANN}$
$\frac{\Psi \vdash^{OL} e_1 : A_1 \rightarrow A_2 \quad \Psi \vdash^{OL} e_2 : A_1}{\Psi \vdash^{OL} e_1 e_2 : A_2} \text{APP}$	$\frac{\Psi \vdash^{OL} e : A_1 \quad \Psi \vdash A_1 <: A_2}{\Psi \vdash^{OL} e : A_2} \text{SUB}$	
$\frac{\Psi, x : \tau \vdash^{OL} e : B}{\Psi \vdash^{OL} \lambda x. e : \tau \rightarrow B} \text{LAM}$	$\frac{\Psi, a \vdash^{OL} e : A}{\Psi \vdash^{OL} e : \forall a. A} \text{GEN}$	

$\Psi \vdash A <: B$

$\frac{a \in \Psi}{\Psi \vdash a <: a} \text{CS-TVAR}$	$\frac{}{\Psi \vdash \text{Int} <: \text{Int}} \text{CS-INT}$	$\frac{\Psi \vdash \tau \quad \Psi \vdash A[a \mapsto \tau] <: B}{\Psi \vdash \forall a. A <: B} \text{FORALLL}$
$\frac{\Psi, a \vdash A <: B}{\Psi \vdash A <: \forall a. B} \text{FORALLR}$	$\frac{\Psi \vdash B_1 <: A_1 \quad \Psi \vdash A_2 <: B_2}{\Psi \vdash A_1 \rightarrow A_2 <: B_1 \rightarrow B_2} \text{CS-FUN}$	

Fig. 2: Syntax and static semantics of the Odersky-Läufer type system.

This program is rejected by Haskell’s type checker because Haskell implements the Damas-Milner rule that a lambda-bound argument (such as  $x$ ) can only have a monotype, i.e., the type checker can only assign  $x$  the type  $\mathbf{Int} \rightarrow \mathbf{Int}$ , or  $\mathbf{Char} \rightarrow \mathbf{Char}$ , but not  $\forall a.[a] \rightarrow [a]$ . Finding such manual polymorphic annotations can be non-trivial. Instead of rejecting the program outright, due to missing type annotations, gradual typing provides a simple alternative by giving  $x$  the unknown type (denoted  $\star$ ). With such typing the same program type-checks and produces  $([2, 1], [b', a'])$ . By running the program, programmers can gain some additional insight about the run-time behaviour. Then, with such insight, they can also give  $x$  a more precise type  $(\forall a.[a] \rightarrow [a])$  a posteriori so that the program continues to type-check via implicit polymorphism and also grants more static safety. In this paper, we envision such a language that combines the benefits of both implicit higher-rank polymorphism and gradual typing.

### 3 Revisiting Consistent Subtyping

In this section we explore the design space of consistent subtyping. We start with the definitions of consistency and subtyping for polymorphic types, and

Types	$A, B ::= \text{Int} \mid a \mid A \rightarrow B \mid \forall a.A \mid \star$			
Monotypes	$\tau, \sigma ::= \text{Int} \mid a \mid \tau \rightarrow \sigma$			
Contexts	$\Psi ::= \emptyset \mid \Psi, x : A \mid \Psi, a$			
$A \sim B$				
$A \sim A$	$A \sim \star$	$\star \sim A$	$\frac{A_1 \sim B_1 \quad A_2 \sim B_2}{A_1 \rightarrow A_2 \sim B_1 \rightarrow B_2}$	$\frac{A \sim B}{\forall a.A \sim \forall a.B}$
$\Psi \vdash A <: B$				
$\frac{\Psi, a \vdash A <: B}{\Psi \vdash A <: \forall a.B}$	$\text{S-FORALLR}$	$\frac{\Psi \vdash \tau \quad \Psi \vdash A[a \mapsto \tau] <: B}{\Psi \vdash \forall a.A <: B}$	$\text{S-FORALLL}$	$\frac{a \in \Psi}{\Psi \vdash a <: a}$
$\frac{}{\Psi \vdash \text{Int} <: \text{Int}}$	$\text{S-INT}$	$\frac{\Psi \vdash B_1 <: A_1 \quad \Psi \vdash A_2 <: B_2}{\Psi \vdash A_1 \rightarrow A_2 <: B_1 \rightarrow B_2}$	$\text{S-FUN}$	$\frac{}{\Psi \vdash \star <: \star}$
		$\text{S-UNKNOWN}$		

Fig. 3: Syntax of types, consistency, and subtyping in the declarative system.

compare with some relevant work. We then discuss the design decisions involved towards our new definition of consistent subtyping, and justify the new definition by demonstrating its equivalence with that of Siek and Taha [22] and the AGT approach [13] on simple types.

The syntax of types is given at the top of Fig. 3. We write  $A, B$  for types. Types are either the integer type  $\text{Int}$ , type variables  $a$ , functions types  $A \rightarrow B$ , universal quantification  $\forall a.A$ , or the unknown type  $\star$ . Though we only have one base type  $\text{Int}$ , we also use  $\text{Bool}$  for the purpose of illustration. Note that monotypes  $\tau$  contain all types other than the universal quantifier and the unknown type  $\star$ . We will discuss this restriction when we present the subtyping rules. Contexts  $\Psi$  are *ordered* lists of type variable declarations and term variables.

### 3.1 Consistency and Subtyping

We start by giving the definitions of consistency and subtyping for polymorphic types, and comparing our definitions with the compatibility relation by Ahmed et al. [1] and type consistency by Igarashi et al. [14].

*Consistency.* The key observation here is that consistency is mostly a structural relation, except that the unknown type  $\star$  can be regarded as any type. Following this observation, we naturally extend the definition from Fig. 1 with polymorphic types, as shown at the middle of Fig. 3. In particular a polymorphic type  $\forall a.A$  is consistent with another polymorphic type  $\forall a.B$  if  $A$  is consistent with  $B$ .

*Subtyping.* We express the fact that one type is a polymorphic generalization of another by means of the subtyping judgment  $\Psi \vdash A <: B$ . Compared with the subtyping rules of Odersky and Läufer [17] in Fig. 2, the only addition is

the neutral subtyping of  $\star$ . Notice that, in the rule S-FORALL, the universal quantifier is only allowed to be instantiated with a *monotype*. The judgment  $\Psi \vdash \tau$  checks all the type variables in  $\tau$  are bound in the context  $\Psi$ . For space reasons, we omit the definition. According to the syntax in Fig. 3, monotypes do not include the unknown type  $\star$ . This is because if we were to allow the unknown type to be used for instantiation, we could have  $\forall a.a \rightarrow a <: \star \rightarrow \star$  by instantiating  $a$  with  $\star$ . Since  $\star \rightarrow \star$  is consistent with any functions  $A \rightarrow B$ , for instance,  $\text{Int} \rightarrow \text{Bool}$ , this means that we could provide an expression of type  $\forall a.a \rightarrow a$  to a function where the input type is supposed to be  $\text{Int} \rightarrow \text{Bool}$ . However, as we might expect,  $\forall a.a \rightarrow a$  is definitely not compatible with  $\text{Int} \rightarrow \text{Bool}$ . This does not hold in any polymorphic type systems without gradual typing. So the gradual type system should not accept it either. (This is the so-called *conservative extension* property that will be made precise in Section 4.3.)

Importantly there is a subtle but crucial distinction between a type variable and the unknown type, although they all represent a kind of “arbitrary” type. The unknown type stands for the absence of type information: it could be *any type* at *any instance*. Therefore, the unknown type is consistent with any type, and additional type-checks have to be performed at runtime. On the other hand, a type variable indicates *parametricity*. In other words, a type variable can only be instantiated to a single type. For example, in the type  $\forall a.a \rightarrow a$ , the two occurrences of  $a$  represent an arbitrary but single type (e.g.,  $\text{Int} \rightarrow \text{Int}$ ,  $\text{Bool} \rightarrow \text{Bool}$ ), while  $\star \rightarrow \star$  could be an arbitrary function (e.g.,  $\text{Int} \rightarrow \text{Bool}$ ) at runtime.

*Comparison with Other Relations.* In other polymorphic gradual calculi, consistency and subtyping are often mixed up to some extent. In  $\lambda\mathbf{B}$  [1], the compatibility relation for polymorphic types is defined as follows:

$$\frac{A < B}{A < \forall X.B} \text{COMP-ALLR} \qquad \frac{A[X \mapsto \star] < B}{\forall X.A < B} \text{COMP-ALLL}$$

Notice that, in rule COMP-ALLL, the universal quantifier is *always* instantiated to  $\star$ . However, this way,  $\lambda\mathbf{B}$  allows  $\forall a.a \rightarrow a < \text{Int} \rightarrow \text{Bool}$ , which as we discussed before might not be what we expect. Indeed  $\lambda\mathbf{B}$  relies on sophisticated runtime checks to rule out such instances of the compatibility relation a posteriori.

Igarashi et al. [14] introduced the so-called *quasi-polymorphic* types for types that may be used where a  $\forall$ -type is expected, which is important for their purpose of conservativity over System F. Their type consistency relation, involving polymorphism, is defined as follows<sup>3</sup>:

$$\frac{A \sim B}{\forall a.A \sim \forall a.B} \qquad \frac{A \sim B \quad B \neq \forall a.B' \quad \star \in \text{Types}(B)}{\forall a.A \sim B}$$

Compared with our consistency definition in Fig. 3, their first rule is the same as ours. The second rule says that a non  $\forall$ -type can be consistent with a  $\forall$ -type only if it contains  $\star$ . In this way, their type system is able to reject  $\forall a.a \rightarrow a \sim$

<sup>3</sup> This is a simplified version.



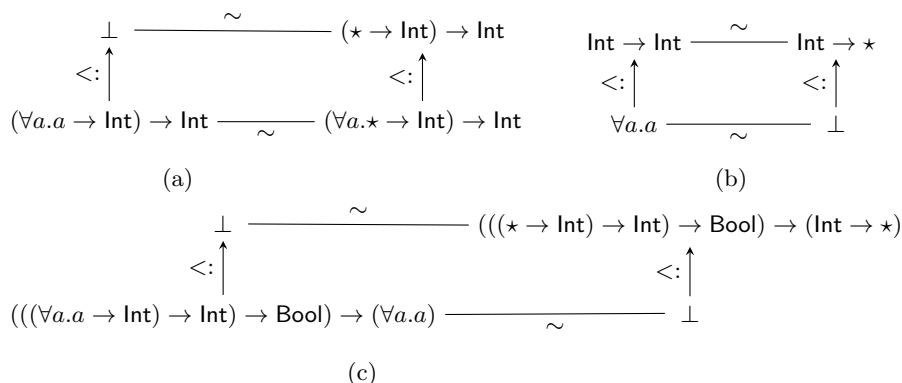


Fig. 4: Examples that break the original definition of consistent subtyping.

$\text{Int} \rightarrow \text{Bool}$ . However, in order to keep conservativity, they also reject  $\forall a.a \rightarrow a \sim \text{Int} \rightarrow \text{Int}$ , which is perfectly sensible in their setting (i.e., explicit polymorphism). However with implicit polymorphism, we would expect  $\forall a.a \rightarrow a$  to be related with  $\text{Int} \rightarrow \text{Int}$ , since  $a$  can be instantiated to  $\text{Int}$ .

Nonetheless, when it comes to interactions between dynamically typed and polymorphically typed terms, both relations allow  $\forall a.a \rightarrow \text{Int}$  to be related with  $\star \rightarrow \text{Int}$  for example, which in our view, is some sort of (implicit) polymorphic subtyping combined with type consistency, and that should be derivable by the more primitive notions in the type system (instead of inventing new relations). One of our design principles is that subtyping and consistency is *orthogonal*, and can be naturally superimposed, echoing the same opinion of Siek and Taha [22].

### 3.2 Towards Consistent Subtyping

With the definitions of consistency and subtyping, the question now is how to compose these two relations so that two types can be compared in a way that takes these two relations into account.

Unfortunately, the original definition of Siek and Taha (Definition 1) does not work well with our definitions of consistency and subtyping for polymorphic types. Consider two types:  $(\forall a.a \rightarrow \text{Int}) \rightarrow \text{Int}$ , and  $(\star \rightarrow \text{Int}) \rightarrow \text{Int}$ . The first type can only reach the second type in one way (first by applying consistency, then subtyping), but not the other way, as shown in Fig. 4a. We use  $\perp$  to mean that we cannot find such a type. Similarly, there are situations where the first type can only reach the second type by the other way (first applying subtyping, and then consistency), as shown in Fig. 4b.

What is worse, if those two examples are composed in a way that those types all appear co-variantly, then the resulting types cannot reach each other in either way. For example, Fig. 4c shows such two types by putting a  $\text{Bool}$  type in the middle, and neither definition of consistent subtyping works.

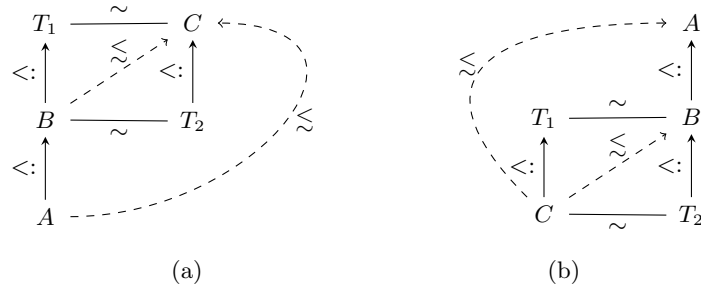


Fig. 5: Observations of consistent subtyping

*Observations on Consistent Subtyping Based on Information Propagation.* In order to develop the correct definition of consistent subtyping for polymorphic types, we need to understand how consistent subtyping works. We first review two important properties of subtyping: (1) subtyping induces the subsumption rule: if  $A <: B$ , then an expression of type  $A$  can be used where  $B$  is expected; (2) subtyping is transitive: if  $A <: B$ , and  $B <: C$ , then  $A <: C$ . Though consistent subtyping takes the unknown type into consideration, the subsumption rule should also apply: if  $A \lesssim B$ , then an expression of type  $A$  can also be used where  $B$  is expected, given that there might be some information lost by consistency. A crucial difference from subtyping is that consistent subtyping is *not* transitive because information can only be lost once (otherwise, any two types are a consistent subtype of each other). Now consider a situation where we have both  $A <: B$ , and  $B \lesssim C$ , this means that  $A$  can be used where  $B$  is expected, and  $B$  can be used where  $C$  is expected, with possibly some loss of information. In other words, we should expect that  $A$  can be used where  $C$  is expected, since there is at most one-time loss of information.

**Observation 1** *If  $A <: B$ , and  $B \lesssim C$ , then  $A \lesssim C$ .*

This is reflected in Fig. 5a. A symmetrical observation is given in Fig. 5b:

**Observation 2** *If  $C \lesssim B$ , and  $B <: A$ , then  $C \lesssim A$ .*

From the above observations, we see what the problem is with the original definition. In Fig. 5a, if  $B$  can reach  $C$  by  $T_1$ , then by subtyping transitivity,  $A$  can reach  $C$  by  $T_1$ . However, if  $B$  can only reach  $C$  by  $T_2$ , then  $A$  cannot reach  $C$  through the original definition. A similar problem is shown in Fig. 5b.

However, it turns out that those two problems can be fixed using the same strategy: instead of taking one-step subtyping and one-step consistency, our definition of consistent subtyping allows types to take *one-step subtyping*, *one-step consistency*, and *one more step subtyping*. Specifically,  $A <: B \sim T_2 <: C$  (in Fig. 5a) and  $C <: T_1 \sim B <: A$  (in Fig. 5b) have the same relation chain: subtyping, consistency, and subtyping.

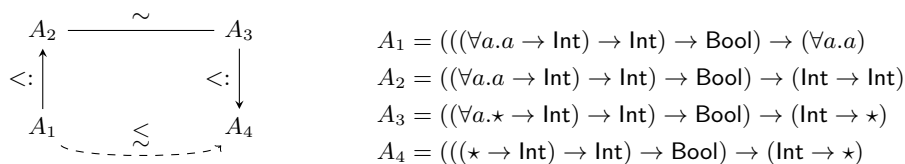


Fig. 6: Example that is fixed by the new definition of consistent subtyping.

*Definition of Consistent subtyping.* From the above discussion, we are ready to modify Definition 1, and adapt it to our notation:

**Definition 2 (Consistent Subtyping).**

$$\frac{\Psi \vdash A <: C \quad C \sim D \quad \Psi \vdash D <: B}{\Psi \vdash A \lesssim B}$$

With Definition 2, Figure 6 illustrates the correct relation chain for the broken example shown in Fig. 4c. At first sight, Definition 2 seems worse than the original: we need to guess *two* types! It turns out that Definition 2 is a generalization of Definition 1, and they are equivalent in the system of Siek and Taha [22]. However, more generally, Definition 2 is compatible with polymorphic types.

**Proposition 1 (Generalization of Consistent Subtyping).**

- Definition 2 subsumes Definition 1.
- Definition 1 is equivalent to Definition 2 in the system of Siek and Taha.

### 3.3 Abstracting Gradual Typing

Garcia et al. [13] presented a new foundation for gradual typing that they call the *Abstracting Gradual Typing* (AGT) approach. In the AGT approach, gradual types are interpreted as sets of static types, where static types refer to types containing no unknown types. In this interpretation, predicates and functions on static types can then be lifted to apply to gradual types. Central to their approach is the so-called *concretization* function. For simple types, a concretization  $\gamma$  from gradual types to a set of static types<sup>4</sup> is defined as follows:

**Definition 3 (Concretization).**

$$\gamma(\text{Int}) = \{\text{Int}\} \quad \gamma(A \rightarrow B) = \gamma(A) \rightarrow \gamma(B) \quad \gamma(\star) = \{\text{All static types}\}$$

Based on the concretization function, subtyping between static types can be lifted to gradual types, resulting in the consistent subtyping relation:

**Definition 4 (Consistent Subtyping in AGT).**  $A \lesssim B$  if and only if  $A_1 <: B_1$  for some  $A_1 \in \gamma(A)$ ,  $B_1 \in \gamma(B)$ .

<sup>4</sup> For simplification, we directly regard type constructor  $\rightarrow$  as a set-level operator.

Later they proved that this definition of consistent subtyping coincides with that of Siek and Taha [22] (Definition 1). By Proposition 1, we can directly conclude that our definition coincides with AGT:

**Proposition 2 (Equivalence to AGT on Simple Types).**  $A \lesssim B$  iff  $A \widetilde{<} B$ .

However, AGT does not show how to deal with polymorphism (e.g. the interpretation of type variables) yet. Still, as noted by Garcia et al. [13], it is a promising line of future work for AGT, and the question remains whether our definition would coincide with it.

Another note related to AGT is that the definition is later adopted by Castagna and Lanvin [7], where the static types  $A_1, B_1$  in Definition 4 can be algorithmically computed by also accounting for top and bottom types.

### 3.4 Directed Consistency

*Directed consistency* [15] is defined in terms of precision and static subtyping:

$$\frac{A' \sqsubseteq A \quad A <: B \quad B' \sqsubseteq B}{A' \lesssim B'}$$

The judgment  $A \sqsubseteq B$  is read “ $A$  is less precise than  $B$ ”. In their setting, precision is defined for type constructors and subtyping for static types. If we interpret this definition from AGT’s point of view, finding a more precise static type<sup>5</sup> has the same effect as concretization. Namely,  $A' \sqsubseteq A$  implies  $A \in \gamma(A')$  and  $B' \sqsubseteq B$  implies  $B \in \gamma(B')$ . Therefore we consider this definition as AGT-style. From this perspective, this definition naturally coincides with Definition 2.

The value of their definition is that consistent subtyping is derived compositionally from *static subtyping* and *precision*. These are two more atomic relations. At first sight, their definition looks very similar to Definition 2 (replacing  $\sqsubseteq$  by  $<:$  and  $<:$  by  $\sim$ ). Then a question arises as to *which one is more fundamental*. To answer this, we need to discuss the relation between consistency and precision.

*Relating Consistency and Precision.* Precision is a partial order (anti-symmetric and transitive), while consistency is symmetric but not transitive. Nonetheless, precision and consistency are related by the following proposition:

**Proposition 3 (Consistency and Precision).**

- If  $A \sim B$ , then there exists (static)  $C$ , such that  $A \sqsubseteq C$ , and  $B \sqsubseteq C$ .
- If for some (static)  $C$ , we have  $A \sqsubseteq C$ , and  $B \sqsubseteq C$ , then we have  $A \sim B$ .

It may seem that precision is a more atomic relation, since consistency can be derived from precision. However, recall that consistency is in fact an equivalence relation lifted from static types to gradual types. Therefore defining consistency independently is straightforward, and it is theoretically viable to validate the

<sup>5</sup> The definition of precision of types is given in appendix.

definition of consistency directly. On the other hand, precision is usually connected with the gradual criteria [25], and finding a correct partial order that adheres to the criteria is not always an easy task. For example, Igarashi et al. [14] argued that term precision for System  $F_G$  is actually nontrivial, leaving the gradual guarantee of the semantics as a conjecture. Thus precision can be difficult to extend to more sophisticated type systems, e.g. dependent types.

Still, it is interesting that those two definitions illustrate the correspondence of different foundations (on simple types): one is defined directly on gradual types, and the other stems from AGT, which is based on static subtyping.

### 3.5 Consistent Subtyping Without Existentials

Definition 2 serves as a fine specification of how consistent subtyping should behave in general. But it is inherently non-deterministic because of the two intermediate types  $C$  and  $D$ . As with Definition 1, we need a combined relation to directly compare two types. A natural attempt is to try to extend the restriction operator for polymorphic types. Unfortunately, as we show below, this does not work. However it is possible to devise an equivalent inductive definition instead.

*Attempt to Extend the Restriction Operator.* Suppose that we try to extend the restriction operator to account for polymorphic types. The original restriction operator is structural, meaning that it works for types of similar structures. But for polymorphic types, two input types could have different structures due to universal quantifiers, e.g.  $\forall a.a \rightarrow \text{Int}$  and  $(\text{Int} \rightarrow \star) \rightarrow \text{Int}$ . If we try to mask the first type using the second, it seems hard to maintain the information that  $a$  should be instantiated to a function while ensuring that the return type is masked. There seems to be no satisfactory way to extend the restriction operator in order to support this kind of non-structural masking.

*Interpretation of the Restriction Operator and Consistent Subtyping.* If the restriction operator cannot be extended naturally, it is useful to take a step back and revisit what the restriction operator actually does. For consistent subtyping, two input types could have unknown types in different positions, but we only care about the known parts. What the restriction operator does is (1) erase the type information in one type if the corresponding position in the other type is the unknown type; and (2) compare the resulting types using the normal subtyping relation. The example below shows the masking-off procedure for the types  $\text{Int} \rightarrow \star \rightarrow \text{Bool}$  and  $\text{Int} \rightarrow \text{Int} \rightarrow \star$ . Since the known parts have the relation that  $\text{Int} \rightarrow \star \rightarrow \star <: \text{Int} \rightarrow \star \rightarrow \star$ , we conclude that  $\text{Int} \rightarrow \star \rightarrow \text{Bool} \lesssim \text{Int} \rightarrow \text{Int} \rightarrow \star$ .

$$\begin{array}{l} \text{Int} \rightarrow \boxed{\star} \rightarrow \boxed{\text{Bool}} \\ \text{Int} \rightarrow \boxed{\text{Int}} \rightarrow \boxed{\star} \end{array} \left. \begin{array}{l} | \text{Int} \rightarrow \text{Int} \rightarrow \star = \text{Int} \rightarrow \star \rightarrow \star \\ | \text{Int} \rightarrow \star \rightarrow \text{Bool} = \text{Int} \rightarrow \star \rightarrow \star \end{array} \right) <:$$

Here differences of the types in boxes are erased because of the restriction operator. Now if we compare the types in boxes directly instead of through the lens of the restriction operator, we can observe that the *consistent subtyping*

$$\boxed{\Psi \vdash A \lesssim B}$$

$$\frac{\Psi, a \vdash A \lesssim B}{\Psi \vdash A \lesssim \forall a. B} \text{CS-FORALLR} \qquad \frac{\Psi \vdash \tau \quad \Psi \vdash A[a \mapsto \tau] \lesssim B}{\Psi \vdash \forall a. A \lesssim B} \text{CS-FORALLL}$$

$$\frac{\Psi \vdash B_1 \lesssim A_1 \quad \Psi \vdash A_2 \lesssim B_2}{\Psi \vdash A_1 \rightarrow A_2 \lesssim B_1 \rightarrow B_2} \text{CS-FUN} \qquad \frac{a \in \Psi}{\Psi \vdash a \lesssim a} \text{CS-TVAR} \qquad \frac{}{\Psi \vdash \text{Int} \lesssim \text{Int}} \text{CS-INT}$$

$$\frac{}{\Psi \vdash \star \lesssim A} \text{CS-UNKNOWNL} \qquad \frac{}{\Psi \vdash A \lesssim \star} \text{CS-UNKNOWNR}$$

Fig. 7: Consistent Subtyping for implicit polymorphism.

relation always holds between the unknown type and an arbitrary type. We can interpret this observation directly from Definition 2: the unknown type is neutral to subtyping ( $\star <: \star$ ), the unknown type is consistent with any type ( $\star \sim A$ ), and subtyping is reflexive ( $A <: A$ ). Therefore, *the unknown type is a consistent subtype of any type* ( $\star \lesssim A$ ), and *vice versa* ( $A \lesssim \star$ ). Note that this interpretation provides a general recipe on how to lift a (static) subtyping relation to a (gradual) consistent subtyping relation, as discussed below.

*Defining Consistent Subtyping Directly.* From the above discussion, we can define the consistent subtyping relation directly, *without* resorting to subtyping or consistency at all. The key idea is that we replace  $<:$  with  $\lesssim$  in Fig. 3, get rid of rule S-UNKNOWN and add two extra rules concerning  $\star$ , resulting in the rules of consistent subtyping in Fig. 7. Of particular interest are the rules CS-UNKNOWNL and CS-UNKNOWNR, both of which correspond to what we just said: the unknown type is a consistent subtype of any type, and vice versa. From now on, we use the symbol  $\lesssim$  to refer to the consistent subtyping relation in Fig. 7. What is more, we can prove that those two are equivalent<sup>6</sup>:

**Theorem 1**  $\Psi \vdash A \lesssim B \Leftrightarrow \Psi \vdash A <: C, C \sim D, \Psi \vdash D <: B$  for some  $C, D$ .

## 4 Gradually Typed Implicit Polymorphism

In Section 3 we introduced the consistent subtyping relation that accommodates polymorphic types. In this section we continue with the development by giving a declarative type system for predicative implicit polymorphism that employs the consistent subtyping relation. The declarative system itself is already quite interesting as it is equipped with both higher-rank polymorphism and the unknown type. The syntax of expressions in the declarative system is given below:

$$\text{Expressions } e ::= x \mid n \mid \lambda x : A. e \mid \lambda x. e \mid e e$$

<sup>6</sup> Theorems with  $\mathcal{T}$  are those proved in Coq. The same applies to Lemmas.

$$\boxed{\Psi \vdash e : A \rightsquigarrow s}$$

$$\frac{x : A \in \Psi}{\Psi \vdash x : A \rightsquigarrow x} \text{VAR} \quad \frac{}{\Psi \vdash n : \text{Int} \rightsquigarrow n} \text{NAT} \quad \frac{\Psi, a \vdash e : A \rightsquigarrow s}{\Psi \vdash e : \forall a. A \rightsquigarrow \Lambda a. s} \text{GEN}$$

$$\frac{\Psi, x : A \vdash e : B \rightsquigarrow s}{\Psi \vdash \lambda x : A. e : A \rightarrow B \rightsquigarrow \lambda x : A. s} \text{LAMANN} \quad \frac{\Psi, x : \tau \vdash e : B \rightsquigarrow s}{\Psi \vdash \lambda x. e : \tau \rightarrow B \rightsquigarrow \lambda x : \tau. s} \text{LAM}$$

$$\frac{\Psi \vdash e_1 : A \rightsquigarrow s_1 \quad \Psi \vdash A \triangleright A_1 \rightarrow A_2 \quad \Psi \vdash e_2 : A_3 \rightsquigarrow s_2 \quad \Psi \vdash A_3 \lesssim A_1}{\Psi \vdash e_1 e_2 : A_2 \rightsquigarrow ((A \hookrightarrow A_1 \rightarrow A_2) s_1) ((A_3 \hookrightarrow A_1) s_2)} \text{APP}$$

$$\boxed{\Psi \vdash A \triangleright A_1 \rightarrow A_2}$$

$$\frac{\Psi \vdash \tau \quad \Psi \vdash A[a \mapsto \tau] \triangleright A_1 \rightarrow A_2}{\Psi \vdash \forall a. A \triangleright A_1 \rightarrow A_2} \text{M-FORALL}$$

$$\frac{}{\Psi \vdash (A_1 \rightarrow A_2) \triangleright (A_1 \rightarrow A_2)} \text{M-ARR} \quad \frac{}{\Psi \vdash \star \triangleright \star \rightarrow \star} \text{M-UNKNOWN}$$

Fig. 8: Declarative typing

#### 4.1 Typing in Detail

Figure 8 gives the typing rules for our declarative system (the reader is advised to ignore the gray-shaded parts for now). Rule VAR extracts the type of the variable from the typing context. Rule NAT always infers integer types. Rule LAMANN puts  $x$  with type annotation  $A$  into the context, and continues type checking the body  $e$ . Rule LAM assigns a monotype  $\tau$  to  $x$ , and continues type checking the body  $e$ . Gradual types and polymorphic types are introduced via annotations explicitly. Rule GEN puts a fresh type variable  $a$  into the type context and generalizes the typing result  $A$  to  $\forall a. A$ . Rule APP first infers the type of  $e_1$ , then the matching judgment  $\Psi \vdash A \triangleright A_1 \rightarrow A_2$  extracts the domain type  $A_1$  and the codomain type  $A_2$  from type  $A$ . The type  $A_3$  of the argument  $e_2$  is then compared with  $A_1$  using the consistent subtyping judgment.

*Matching.* The matching judgment of Siek et al. [25] can be extended to polymorphic types naturally, resulting in  $\Psi \vdash A \triangleright A_1 \rightarrow A_2$ . In M-FORALL, a monotype  $\tau$  is guessed to instantiate the universal quantifier  $a$ . This rule is inspired by the *application judgment*  $\Phi \vdash A \bullet e \Rightarrow C$  [11], which says that if we apply a term of type  $A$  to an argument  $e$ , we get something of type  $C$ . If  $A$  is a polymorphic type, the judgment works by guessing instantiations until it reaches an arrow type. Matching further simplifies the application judgment, since it is independent of typing. Rule M-ARR and M-UNKNOWN are the same as Siek et al. [25]. M-ARR returns the domain type  $A_1$  and range type  $A_2$  as expected. If the input is  $\star$ , then M-UNKNOWN returns  $\star$  as both the type for the domain and the range.

Note that matching saves us from having a subsumption rule (SUB in Fig. 2). the subsumption rule is incompatible with consistent subtyping, since the latter is not transitive. A discussion of a subsumption rule based on normal subtyping can be found in the appendix.

## 4.2 Type-directed Translation

We give the dynamic semantics of our language by translating it to  $\lambda\mathbf{B}$ . Below we show a subset of the terms in  $\lambda\mathbf{B}$  that are used in the translation:

$$\text{Terms } s ::= x \mid n \mid \lambda x : A. s \mid \Lambda a. s \mid s_1 s_2 \mid \langle A \hookrightarrow B \rangle s$$

A cast  $\langle A \hookrightarrow B \rangle s$  converts the value of term  $s$  from type  $A$  to type  $B$ . A cast from  $A$  to  $B$  is permitted only if the types are *compatible*, written  $A \prec B$ , as briefly mentioned in Section 3.1. The syntax of types in  $\lambda\mathbf{B}$  is the same as ours.

The translation is given in the gray-shaded parts in Fig. 8. The only interesting case here is to insert explicit casts in the application rule. Note that there is no need to translate matching or consistent subtyping, instead we insert the source and target types of a cast directly in the translated expressions, thanks to the following two lemmas:

**Lemma 1** ( $\triangleright$  to  $\prec$ ) *If  $\Psi \vdash A \triangleright A_1 \rightarrow A_2$ , then  $A \prec A_1 \rightarrow A_2$ .*

**Lemma 2** ( $\lesssim$  to  $\prec$ ) *If  $\Psi \vdash A \lesssim B$ , then  $A \prec B$ .*

In order to show the correctness of the translation, we prove that our translation always produces well-typed expressions in  $\lambda\mathbf{B}$ . By Lemmas 1 and 2, we have the following theorem:

**Theorem 2 (Type Safety)** *If  $\Psi \vdash e : A \rightsquigarrow s$ , then  $\Psi \vdash^B s : A$ .*

*Parametricity.* An important semantic property of polymorphic types is *relational parametricity* [19]. The parametricity property says that all instances of a polymorphic function should behave *uniformly*. A classic example is a function with the type  $\forall a. a \rightarrow a$ . The parametricity property guarantees that a value of this type must be either the identity function (i.e.,  $\lambda x. x$ ) or the undefined function (one which never returns a value). However, with the addition of the unknown type  $\star$ , careful measures are to be taken to ensure parametricity. This is exactly the circumstance that  $\lambda\mathbf{B}$  was designed to address. Ahmed et al. [2] proved that  $\lambda\mathbf{B}$  satisfies relational parametricity. Based on their result, and by Theorem 2, parametricity is preserved in our system.

*Ambiguity from Casts.* The translation does not always produce a unique target expression. This is because when we guess a monotype  $\tau$  in rule M-FORALL and CS-FORALL, we could have different choices, which inevitably leads to different types. Unlike (non-gradual) polymorphic type systems [18, 11], the choice of monotypes could affect runtime behaviour of the translated programs, since they could appear inside the explicit casts. For example, the following shows



two possible translations for the same source expression  $\lambda x : \star. f x$ , where the type of  $f$  is instantiated to  $\text{Int} \rightarrow \text{Int}$  and  $\text{Bool} \rightarrow \text{Bool}$ , respectively:

$$\begin{aligned}
& f : \forall a. a \rightarrow a \vdash (\lambda x : \star. f x) : \star \rightarrow \text{Int} \\
& \quad \rightsquigarrow (\lambda x : \star. (\langle \forall a. a \rightarrow a \hookrightarrow \text{Int} \rightarrow \text{Int} \rangle f) (\langle \star \hookrightarrow \text{Int} \rangle x)) \\
& f : \forall a. a \rightarrow a \vdash (\lambda x : \star. f x) : \star \rightarrow \text{Bool} \\
& \quad \rightsquigarrow (\lambda x : \star. (\langle \forall a. a \rightarrow a \hookrightarrow \text{Bool} \rightarrow \text{Bool} \rangle f) (\langle \star \hookrightarrow \text{Bool} \rangle x))
\end{aligned}$$

If we apply  $\lambda x : \star. f x$  to 3, which is fine since the function can take any input, the first translation runs smoothly in  $\lambda\text{B}$ , while the second one will raise a cast error ( $\text{Int}$  cannot be cast to  $\text{Bool}$ ). Similarly, if we apply it to  $\text{true}$ , then the second succeeds while the first fails. The culprit lies in the highlighted parts where any instantiation of  $a$  would be put inside the explicit cast. More generally, any choice introduces an explicit cast to that type in the translation, which causes a runtime cast error if the function is applied to a value whose type does not match the guessed type. Note that this does not compromise the type safety of the translated expressions, since cast errors are part of the type safety guarantees.

*Coherence.* The ambiguity of translation seems to imply that the declarative system is *incoherent*. A semantics is coherent if distinct typing derivations of the same typing judgment possess the same meaning [20]. We argue that the declarative system is “coherent up to cast errors” in the sense that a well-typed program produces a unique value, or results in a cast error. In the above example, whatever the translation might be, applying  $\lambda x : \star. f x$  to 3 either results in a cast error, or produces 3, nothing else.

This discrepancy is due to the guessing nature of the *declarative* system. As far as the declarative system is concerned, both  $\text{Int} \rightarrow \text{Int}$  and  $\text{Bool} \rightarrow \text{Bool}$  are equally acceptable. But this is not the case at runtime. The acute reader may have found that the *only* appropriate choice is to instantiate  $f$  to  $\star \rightarrow \star$ . However, as specified by rule M-FORALL in Fig. 8, we can only instantiate type variables to monotypes, but  $\star$  is *not* a monotype! We will get back to this issue in Section 6.2 after we present the corresponding algorithmic system in Section 5.

### 4.3 Correctness Criteria

Siek et al. [25] present a set of properties that a well-designed gradual typing calculus must have, which they call the refined criteria. Among all the criteria, those related to the static aspects of gradual typing are well summarized by Cimini and Siek [8]. Here we review those criteria and adapt them to our notation. We have proved in Coq that our type system satisfies all these criteria.

#### Lemma 3 (Correctness Criteria)

- *Conservative extension:* for all static  $\Psi$ ,  $e$ , and  $A$ ,
  - if  $\Psi \vdash^{OL} e : A$ , then there exists  $B$ , such that  $\Psi \vdash e : B$ , and  $\Psi \vdash B <: A$ .
  - if  $\Psi \vdash e : A$ , then  $\Psi \vdash^{OL} e : A$

- **Monotonicity w.r.t. precision:** for all  $\Psi, e, e', A$ , if  $\Psi \vdash e : A$ , and  $e' \sqsubseteq e$ , then  $\Psi \vdash e' : B$ , and  $B \sqsubseteq A$  for some  $B$ .
- **Type Preservation of cast insertion:** for all  $\Psi, e, A$ , if  $\Psi \vdash e : A$ , then  $\Psi \vdash e : A \rightsquigarrow s$ , and  $\Psi \vdash^B s : A$  for some  $s$ .
- **Monotonicity of cast insertion:** for all  $\Psi, e_1, e_2, e'_1, e'_2, A$ , if  $\Psi \vdash e_1 : A \rightsquigarrow e'_1$ , and  $\Psi \vdash e_2 : A \rightsquigarrow e'_2$ , and  $e_1 \sqsubseteq e_2$ , then  $\Psi \vdash e'_1 \sqsubseteq^B e'_2$ .

The first criterion states that the gradual type system should be a conservative extension of the original system. In other words, a *static* program that is typeable in the Odersky-Läufer type system if and only if it is typeable in the gradual type system. A static program is one that does not contain any type  $\star$ <sup>7</sup>. However since our gradual type system does not have the subsumption rule, it produces more general types.

The second criterion states that if a typeable expression loses some type information, it remains typeable. This criterion depends on the definition of the precision relation, written  $A \sqsubseteq B$ , which is given in the appendix. The relation intuitively captures a notion of types containing more or less unknown types ( $\star$ ). The precision relation over types lifts to programs, i.e.,  $e_1 \sqsubseteq e_2$  means that  $e_1$  and  $e_2$  are the same program except that  $e_2$  has more unknown types.

The first two criteria are fundamental to gradual typing. They explain for example why these two programs  $(\lambda x : \text{Int}. x + 1)$  and  $(\lambda x : \star. x + 1)$  are typeable, as the former is typeable in the Odersky-Läufer type system and the latter is a less-precise version of it.

The last two criteria relate the compilation to the cast calculus. The third criterion is essentially the same as Theorem 2, given that a target expression should always exist, which can be easily seen from Fig. 8. The last criterion ensures that the translation must be monotonic over the precision relation  $\sqsubseteq$ .

As for the dynamic guarantee, things become a bit murky for two reasons: (1) as we discussed before, our declarative system is incoherent in that the runtime behaviour of the same source program can vary depending on the particular translation; (2) it is still unknown whether dynamic guarantee holds in  $\lambda B$ . We will have more discussion on the dynamic guarantee in Section 6.3.

## 5 Algorithmic Type System

In this section we give a bidirectional account of the algorithmic type system that implements the declarative specification. The algorithm is largely inspired by the algorithmic bidirectional system of Dunfield and Krishnaswami [11] (henceforth DK system). However our algorithmic system differs from theirs in three aspects: 1) the addition of the unknown type  $\star$ ; 2) the use of the matching judgment; and 3) the approach of *gradual inference only producing static types* [12]. We then prove that our algorithm is both sound and complete with respect to the declarative type system. Full proofs can be found in the appendix.

<sup>7</sup> Note that the term *static* has appeared several times with different meanings.

Expressions	$e ::= x \mid n \mid \lambda x : A. e \mid \lambda x. e \mid e e \mid e : A$
Types	$A, B ::= \text{Int} \mid a \mid \hat{a} \mid A \rightarrow B \mid \forall a. A \mid \star$
Monotypes	$\tau, \sigma ::= \text{Int} \mid a \mid \hat{a} \mid \tau \rightarrow \sigma$
Contexts	$\Gamma, \Delta, \Theta ::= \emptyset \mid \Gamma, x : A \mid \Gamma, a \mid \Gamma, \hat{a} \mid \Gamma, \hat{a} = \tau$
Complete Contexts	$\Omega ::= \emptyset \mid \Omega, x : A \mid \Omega, a \mid \Omega, \hat{a} = \tau$

Fig. 9: Syntax of the algorithmic system

$$\boxed{\Gamma \vdash A \lesssim B \vdash \Delta}$$

$$\frac{}{\Gamma[a] \vdash a \lesssim a \vdash \Gamma[a]} \text{ACS-TVAR} \qquad \frac{}{\Gamma[\hat{a}] \vdash \hat{a} \lesssim \hat{a} \vdash \Gamma[\hat{a}]} \text{ACS-EXVAR}$$

$$\frac{}{\Gamma \vdash \text{Int} \lesssim \text{Int} \vdash \Gamma} \text{ACS-INT} \qquad \frac{}{\Gamma \vdash \star \lesssim A \vdash \Gamma} \text{ACS-UNKNOWNL} \qquad \frac{}{\Gamma \vdash A \lesssim \star \vdash \Gamma} \text{ACS-UNKNOWNR}$$

$$\frac{\Gamma \vdash B_1 \lesssim A_1 \vdash \Theta \quad \Theta \vdash [\Theta]A_2 \lesssim [\Theta]B_2 \vdash \Delta}{\Gamma \vdash A_1 \rightarrow A_2 \lesssim B_1 \rightarrow B_2 \vdash \Delta} \text{ACS-FUN}$$

$$\frac{\Gamma, a \vdash A \lesssim B \vdash \Delta, a, \Theta}{\Gamma \vdash A \lesssim \forall a. B \vdash \Delta} \text{ACS-FORALLR} \qquad \frac{\Gamma, \hat{a} \vdash A[a \mapsto \hat{a}] \lesssim B \vdash \Delta}{\Gamma \vdash \forall a. A \lesssim B \vdash \Delta} \text{ACS-FORALLL}$$

$$\frac{\hat{a} \notin \text{fv}(A) \quad \Gamma[\hat{a}] \vdash \hat{a} \lesssim A \vdash \Delta}{\Gamma[\hat{a}] \vdash \hat{a} \lesssim A \vdash \Delta} \text{ACS-INSTL} \qquad \frac{\hat{a} \notin \text{fv}(A) \quad \Gamma[\hat{a}] \vdash A \lesssim \hat{a} \vdash \Delta}{\Gamma[\hat{a}] \vdash A \lesssim \hat{a} \vdash \Delta} \text{ACS-INSTR}$$

Fig. 10: Algorithmic consistent subtyping

*Algorithmic Contexts.* The algorithmic context  $\Gamma$  is an *ordered* list containing declarations of type variables  $a$  and term variables  $x : A$ . Unlike declarative contexts, algorithmic contexts also contain declarations of existential type variables  $\hat{a}$ , which can be either unsolved (written  $\hat{a}$ ) or solved to some monotype (written  $\hat{a} = \tau$ ). Complete contexts  $\Omega$  are those that contain no unsolved existential type variables. Figure 9 shows the syntax of the algorithmic system. Apart from expressions in the declarative system, we have annotated expressions  $e : A$ .

### 5.1 Algorithmic Consistent Subtyping and Instantiation

Figure 10 shows the algorithmic consistent subtyping rules. The first five rules do not manipulate contexts. Rule ACS-FUN is a natural extension of its declarative counterpart. The output context of the first premise is used by the second premise, and the output context of the second premise is the output context of the conclusion. Note that we do not simply check  $A_2 \lesssim B_2$ , but apply  $\Theta$  to both types (e.g.,  $[\Theta]A_2$ ). This is to maintain an important invariant that types are fully applied under input context  $\Gamma$  (they contain no existential variables already solved in  $\Gamma$ ). The same invariant applies to every algorithmic judgment.

$$\boxed{\Gamma \vdash \hat{a} \lesssim A \dashv \Delta}$$

$$\frac{\Gamma \vdash \tau}{\Gamma, \hat{a}, \Gamma' \vdash \hat{a} \lesssim \tau \dashv \Gamma, \hat{a} = \tau, \Gamma'} \text{INSTLSOLVE} \quad \frac{}{\Gamma[\hat{a}][\hat{b}] \vdash \hat{a} \lesssim \hat{b} \dashv \Gamma[\hat{a}][\hat{b} = \hat{a}]} \text{INSTLREACH}$$

$$\frac{}{\Gamma[\hat{a}] \vdash \hat{a} \lesssim \star \dashv \Gamma[\hat{a}]} \text{INSTLSOLVEU} \quad \frac{\Gamma[\hat{a}], b \vdash \hat{a} \lesssim B \dashv \Delta, b, \Delta'}{\Gamma[\hat{a}] \vdash \hat{a} \lesssim \forall b. B \dashv \Delta} \text{INSTLALLR}$$

$$\frac{\Gamma[\hat{a}_2, \hat{a}_1, \hat{a} = \hat{a}_1 \rightarrow \hat{a}_2] \vdash A_1 \lesssim \hat{a}_1 \dashv \Theta \quad \Theta \vdash \hat{a}_2 \lesssim [\Theta]A_2 \dashv \Delta}{\Gamma[\hat{a}] \vdash \hat{a} \lesssim A_1 \rightarrow A_2 \dashv \Delta} \text{INSTLARR}$$

Fig. 11: Algorithmic instantiation

Rule ACS-FORALLR looks similar to its declarative counterpart, except that we need to drop the trailing context  $a, \Theta$  from the concluding output context since they become out of scope. Rule ACS-FORALLL generates a fresh existential variable  $\hat{a}$ , and replaces  $a$  with  $\hat{a}$  in the body  $A$ . The new existential variable  $\hat{a}$  is then added to the premise's input context. As a side note, when both types are quantifiers, then either ACS-FORALLR or ACS-FORALLL could be tried. In practice, one can apply ACS-FORALLR eagerly. The last two rules together check consistent subtyping with an unsolved existential variable on one side and an arbitrary type on the other side by the help of the instantiation judgment.

The judgment  $\Gamma \vdash \hat{a} \lesssim A \dashv \Delta$  defined in Fig. 11 instantiates unsolved existential variables. Judgment  $\hat{a} \lesssim A$  reads “instantiate  $\hat{a}$  to a consistent subtype of  $A$ ”. For space reasons, we omit its symmetric judgement  $\Gamma \vdash A \lesssim \hat{a} \dashv \Delta$ . Rule INSTLSOLVE and rule INSTLREACH set  $\hat{a}$  to  $\tau$  and  $\hat{b}$  in the output context, respectively. Rule INSTLSOLVEU is similar to ACS-UNKNOWNR in that we put no constraint on  $\hat{a}$  when it meets the unknown type  $\star$ . This design decision reflects the point that type inference only produces static types [12]. We will get back to this point in Section 6.2. Rule INSTLALLR is the instantiation version of rule ACS-FORALLR. The last rule INSTLARR applies when  $\hat{a}$  meets a function type. It follows that the solution must also be a function type. That is why, in the first premise, we generate two fresh existential variables  $\hat{a}_1$  and  $\hat{a}_2$ , and insert them just before  $\hat{a}$  in the input context, so that the solution of  $\hat{a}$  can mention them. Note that  $A_1 \lesssim \hat{a}_1$  switches to the other instantiation judgment.

## 5.2 Algorithmic Typing

We now turn to the algorithmic typing rules in Fig. 12. The algorithmic system uses bidirectional type checking to accommodate polymorphism. Most of them are quite standard. Perhaps rule AAPP (which differs significantly from that in the DK system) deserves attention. It relies on the algorithmic matching judgment  $\Gamma \vdash A \triangleright A_1 \rightarrow A_2 \dashv \Delta$ . Rule AM-FORALLL replaces  $a$  with a

$$\boxed{\Gamma \vdash e \Rightarrow A \dashv \Delta}$$

$$\frac{(x : A) \in \Gamma}{\Gamma \vdash x \Rightarrow A \dashv \Gamma} \text{AVAR} \qquad \frac{}{\Gamma \vdash n \Rightarrow \text{Int} \dashv \Gamma} \text{ANAT}$$

$$\frac{\Gamma, \widehat{a}, \widehat{b}, x : \widehat{a} \vdash e \Leftarrow \widehat{b} \dashv \Delta, x : \widehat{a}, \Theta}{\Gamma \vdash \lambda x. e \Rightarrow \widehat{a} \rightarrow \widehat{b} \dashv \Delta} \text{ALAMU} \qquad \frac{\Gamma, x : A \vdash e \Rightarrow B \dashv \Delta, x : A, \Theta}{\Gamma \vdash \lambda x : A. e \Rightarrow A \rightarrow B \dashv \Delta} \text{ALAMANNA}$$

$$\frac{\Gamma \vdash A \quad \Gamma \vdash e \Leftarrow A \dashv \Delta}{\Gamma \vdash e : A \Rightarrow A \dashv \Delta} \text{AANNO}$$

$$\frac{\Gamma \vdash e_1 \Rightarrow A \dashv \Theta_1 \quad \Theta_1 \vdash [\Theta_1]A \triangleright A_1 \rightarrow A_2 \dashv \Theta_2 \quad \Theta_2 \vdash e_2 \Leftarrow [\Theta_2]A_1 \dashv \Delta}{\Gamma \vdash e_1 e_2 \Rightarrow A_2 \dashv \Delta} \text{AAPP}$$

$$\boxed{\Gamma \vdash e \Leftarrow A \dashv \Delta}$$

$$\frac{\Gamma, x : A \vdash e \Leftarrow B \dashv \Delta, x : A, \Theta}{\Gamma \vdash \lambda x. e \Leftarrow A \rightarrow B \dashv \Delta} \text{ALAM} \qquad \frac{\Gamma, a \vdash e \Leftarrow A \dashv \Delta, a, \Theta}{\Gamma \vdash e \Leftarrow \forall a. A \dashv \Delta} \text{AGEN}$$

$$\frac{\Gamma \vdash e \Rightarrow A \dashv \Theta \quad \Theta \vdash [\Theta]A \lesssim [\Theta]B \dashv \Delta}{\Gamma \vdash e \Leftarrow B \dashv \Delta} \text{ASUB}$$

$$\boxed{\Gamma \vdash A \triangleright A_1 \rightarrow A_2 \dashv \Delta}$$

$$\frac{\Gamma, \widehat{a} \vdash A[a \mapsto \widehat{a}] \triangleright A_1 \rightarrow A_2 \dashv \Delta}{\Gamma \vdash \forall a. A \triangleright A_1 \rightarrow A_2 \dashv \Delta} \text{AM-FORALL} \qquad \frac{}{\Gamma \vdash (A_1 \rightarrow A_2) \triangleright (A_1 \rightarrow A_2) \dashv \Gamma} \text{AM-ARR}$$

$$\frac{}{\Gamma \vdash \star \triangleright \star \rightarrow \star \dashv \Gamma} \text{AM-UNKNOWN} \qquad \frac{}{\Gamma[\widehat{c}] \vdash \widehat{c} \triangleright \widehat{a} \rightarrow \widehat{b} \dashv \Gamma[\widehat{a}, \widehat{b}, \widehat{c} = \widehat{a} \rightarrow \widehat{b}]} \text{AM-VAR}$$

Fig. 12: Algorithmic typing

fresh existential variable  $\widehat{a}$ , thus eliminating guessing. Rule AM-ARR and AM-UNKNOWN correspond directly to the declarative rules. Rule AM-VAR, which has no corresponding declarative version, is similar to INSTRARR/INSTLARR: we create  $\widehat{a}$  and  $\widehat{b}$  and add  $\widehat{c} = \widehat{a} \rightarrow \widehat{b}$  to the context.

### 5.3 Completeness and Soundness

We prove that the algorithmic rules are sound and complete with respect to the declarative specifications. We need an auxiliary judgment  $\Gamma \longrightarrow \Delta$  that captures a notion of information increase from input contexts  $\Gamma$  to output contexts  $\Delta$  [11].

*Soundness.* Roughly speaking, soundness of the algorithmic system says that given an expression  $e$  that type checks in the algorithmic system, there exists a

corresponding expression  $e'$  that type checks in the declarative system. However there is one complication:  $e$  does not necessarily have more annotations than  $e'$ . For example, by ALAM we have  $\lambda x. x \Leftarrow (\forall a.a) \rightarrow (\forall a.a)$ , but  $\lambda x. x$  itself cannot have type  $(\forall a.a) \rightarrow (\forall a.a)$  in the declarative system. To circumvent that, we add an annotation to the lambda abstraction, resulting in  $\lambda x : (\forall a.a). x$ , which is typeable in the declarative system with the same type. To relate  $\lambda x. x$  and  $\lambda x : (\forall a.a). x$ , we erase all annotations on both expressions. The definition of erasure  $[\cdot]$  is standard and thus omitted.

**Theorem 1 (Soundness of Algorithmic Typing)** *Given  $\Delta \longrightarrow \Omega$ ,*

1. *If  $\Gamma \vdash e \Rightarrow A \dashv \Delta$  then  $\exists e'$  such that  $[\Omega]\Delta \vdash e' : [\Omega]A$  and  $[e] = [e']$ .*
2. *If  $\Gamma \vdash e \Leftarrow A \dashv \Delta$  then  $\exists e'$  such that  $[\Omega]\Delta \vdash e' : [\Omega]A$  and  $[e] = [e']$ .*

*Completeness.* Completeness of the algorithmic system is the reverse of soundness: given a declarative judgment of the form  $[\Omega]\Gamma \vdash [\Omega]\dots$ , we want to get an algorithmic derivation of  $\Gamma \vdash \dots \dashv \Delta$ . It turns out that completeness is a bit trickier to state in that the algorithmic rules generate existential variables on the fly, so  $\Delta$  could contain unsolved existential variables that are not found in  $\Gamma$ , nor in  $\Omega$ . Therefore the completeness proof must produce another complete context  $\Omega'$  that extends both the output context  $\Delta$ , and the given complete context  $\Omega$ . As with soundness, we need erasure to relate both expressions.

**Theorem 2 (Completeness of Algorithmic Typing)** *Given  $\Gamma \longrightarrow \Omega$  and  $\Gamma \vdash A$ , if  $[\Omega]\Gamma \vdash e : A$  then there exist  $\Delta, \Omega', A'$  and  $e'$  such that  $\Delta \longrightarrow \Omega'$  and  $\Omega \longrightarrow \Omega'$  and  $\Gamma \vdash e' \Rightarrow A' \dashv \Delta$  and  $A = [\Omega']A'$  and  $[e] = [e']$ .*

## 6 Discussion

### 6.1 Top Types

To demonstrate that our definition of consistent subtyping (Definition 2) is applicable to other features, we show how to extend our approach to **Top** types with all the desired properties preserved.

In order to preserve the orthogonality between subtyping and consistency, we require  $\top$  to be a common supertype of all static types, as shown in rule S-TOP. This rule might seem strange at first glance, since even if we remove the requirement  $A$  *static*, the rule seems reasonable. However, an important point is that because of the orthogonality between subtyping and consistency, subtyping itself should not contain a potential information loss! Therefore, subtyping instances such as  $\star <: \top$  are not allowed. For consistency, we add the rule that  $\top$  is consistent with  $\top$ , which is actually included in the original reflexive rule  $A \sim A$ . For consistent subtyping, every type is a consistent subtype of  $\top$ , for example,  $\text{Int} \rightarrow \star \lesssim \top$ .

$$\frac{A \text{ static}}{\Psi \vdash A <: \top} \text{S-TOP} \quad \top \sim \top \quad \frac{}{\Psi \vdash A \lesssim \top} \text{CS-TOP}$$

It is easy to verify that Definition 2 is still equivalent to that in Fig. 7 extended with rule CS-TOP. That is, Theorem 1 holds:

**Proposition 4 (Extension with  $\top$ ).**  $\Psi \vdash A \lesssim B \Leftrightarrow \Psi \vdash A <: C, C \sim D, \Psi \vdash D <: B$ , for some  $C, D$ .

We extend the definition of concretization (Definition 3) with  $\top$  by adding another equation  $\gamma(\top) = \{\top\}$ . Note that Castagna and Lanvin [7] also have this equation in their calculus. It is easy to verify that Proposition 2 still holds:

**Proposition 5 (Equivalent to AGT on  $\top$ ).**  $A \lesssim B$  if only if  $A \widetilde{<} B$ .

*Siek and Taha's Definition of Consistent Subtyping Does Not Work for  $\top$ .* As the analysis in Section 3.2,  $\text{Int} \rightarrow \star \lesssim \top$  only holds when we first apply consistency, then subtyping. However we cannot find a type  $A$  such that  $\text{Int} \rightarrow \star <: A$  and  $A \sim \top$ . Also we have a similar problem in extending the restriction operator: *non-structural* masking between  $\text{Int} \rightarrow \star$  and  $\top$  cannot be easily achieved.

## 6.2 Interpretation of the Dynamic Semantics

In Section 4.2 we have seen an example where a source expression could produce two different target expressions with different runtime behaviour. As we explained, this is due to the guessing nature of the declarative system, and from the typing point of view, no type is particularly better than others. However, in practice, this is not desirable. Let us revisit the same example, now from the algorithmic point of view (we omit the translation for space reasons):

$$f : \forall a. a \rightarrow a \vdash (\lambda x : \star. f x) \Rightarrow \star \rightarrow \hat{a} \dashv f : \forall a. a \rightarrow a, \hat{a}$$

Compared with declarative typing, which produces many types ( $\star \rightarrow \text{Int}$ ,  $\star \rightarrow \text{Bool}$ , and so on), the algorithm computes the type  $\star \rightarrow \hat{a}$  with  $\hat{a}$  unsolved in the output context. What can we know from the output context? The only thing we know is that  $\hat{a}$  is not constrained at all! However, it is possible to make a more refined distinction between different kinds of existential variables. The first kind of existential variables are those that indeed have no constraints at all, as they do not affect the dynamic semantics. The second kind of existential variables (as in this example) are those where the only constraint is that *the variable was once compared with an unknown type* [12].

To emphasize the difference and have better support for dynamic semantics, we could have *gradual variables* in addition to existential variables, with the difference that only unsolved gradual variables are allowed to be unified with the unknown type. An irreversible transition from existential variables to gradual variables occurs when an existential variable is compared with  $\star$ . After the algorithm terminates, we can set all unsolved existential variables to be any (static) type (or more precisely, as Garcia and Cimini [12], with *static type parameters*), and all unsolved gradual variables to be  $\star$  (or *gradual type parameters*). However, this approach requires a more sophisticated declarative/algorithmic type system than the ones presented in this paper, where we only produce static monotypes

in type inference. We believe this is a typical trade-off in existing gradual type systems with inference [23, 12]. Here we suppress the complexity of dynamic semantics in favour of the conciseness of static typing.

### 6.3 The Dynamic Guarantee

In Section 4.3 we mentioned that the dynamic guarantee is closely related to the coherence issue. To aid discussion, we first give the definition of dynamic guarantee as follows:

**Definition 5 (Dynamic guarantee).** *Suppose  $e' \sqsubseteq e$ ,  $\emptyset \vdash e : A \rightsquigarrow s$  and  $\emptyset \vdash e' : A' \rightsquigarrow s'$ , if  $s \Downarrow v$ , then  $s' \Downarrow v'$  and  $v' \sqsubseteq v$ .*

The dynamic guarantee says that if a gradually typed program evaluates to a value, then removing type annotations always produces a program that evaluates to an equivalent value (modulo type annotations). Now apparently the coherence issue of the declarative system breaks the dynamic guarantee. For instance:

$$(\lambda f : \forall a.a \rightarrow a. \lambda x : \text{Int}. f x) (\lambda x.x) 3 \quad (\lambda f : \forall a.a \rightarrow a. \lambda x : \star. f x) (\lambda x.x) 3$$

The left one evaluates to 3, whereas its less precise version (right) will give a cast error if  $a$  is instantiated to **Bool** for example.

As discussed in Section 6.2, we could design a more sophisticated declarative/algorithmic type system where coherence is retained. However, even with a coherent source language, the dynamic guarantee is still a question. Currently, the dynamic guarantee for our target language  $\lambda\mathbf{B}$  is still an open question. According to Igarashi et al. [14], the difficulty lies in the definition of term precision that preserves the semantics.

## 7 Related Work

Along the way we discussed some of the most relevant work to motivate, compare and promote our gradual typing design. In what follows, we briefly discuss related work on gradual typing and polymorphism.

*Gradual Typing* The seminal paper by Siek and Taha [21] is the first to propose gradual typing. The original proposal extends the simply typed lambda calculus by introducing the unknown type  $\star$  and replacing type equality with type consistency. Later Siek and Taha [22] incorporated gradual typing into a simple object oriented language, and showed that subtyping and consistency are orthogonal – an insight that partly inspired our work. We show that subtyping and consistency are orthogonal in a much richer type system with higher-rank polymorphism. Siek et al. [25] proposed a set of criteria that provides important guidelines for designers of gradually typed languages. Cimini and Siek [8] introduced the *Gradualizer*, a general methodology for generating gradual type systems from static type systems. Later they also develop an algorithm to generate dynamic semantics [9]. Garcia et al. [13] introduced the AGT approach based on abstract interpretation.



*Gradual Type Systems with Explicit Polymorphism* Ahmed et al. [1] proposed  $\lambda\mathbf{B}$  that extends the blame calculus [29] to incorporate polymorphism. The key novelty of their work is to use dynamic sealing to enforce parametricity. Devriese et al. [10] proved that embedding of System F terms into  $\lambda\mathbf{B}$  is not fully abstract. Igarashi et al. [14] also studied integrating gradual typing with parametric polymorphism. They proposed System  $F_G$ , a gradually typed extension of System F, and System  $F_C$ , a new polymorphic blame calculus. As has been discussed extensively, their definition of type consistency does not apply to our setting (implicit polymorphism). All of these approaches mix consistency with subtyping to some extent, which we argue should be orthogonal.

*Gradual Type Inference* Siek and Vachharajani [23] studied unification-based type inference for gradual typing, where they show why three straightforward approaches fail to meet their design goals. Their type system infers gradual types, which results in a complicated type system and inference algorithm. Garcia and Cimini [12] presented a new approach where gradual type inference only produces static types, which is adopted in our type system. They also deal with let-polymorphism (rank 1 types). However none of these works deals with higher-ranked implicit polymorphism.

*Higher-rank Implicit Polymorphism* Odersky and Läufer [17] introduced a type system for higher-rank types. Based on that, Peyton Jones et al. [18] developed an approach for type checking higher-rank predicative polymorphism. Dunfield and Krishnaswami [11] proposed a bidirectional account of higher-rank polymorphism, and an algorithm for implementing the declarative system, which serves as a sole inspiration for our algorithmic system. The key difference, however, is the integration of gradual typing. Vytiniotis et al. [28] defers static type errors to runtime, which is fundamentally different from gradual typing, where programmers can control over static or runtime checks by precision of the annotations.

## 8 Conclusion

In this paper, we present a generalized definition of consistent subtyping, which is proved to be applicable to both polymorphic and top types. Based on the new definition of consistent subtyping, we have developed a gradually typed calculus with predicative implicit higher-rank polymorphism, and an algorithm to implement it. As future work, we are interested to investigate if our results can scale to real world languages and other programming language features.

## Acknowledgements

We thank Ronald Garcia and the anonymous reviewers for their helpful comments. This work has been sponsored by the Hong Kong Research Grant Council projects number 17210617 and 17258816.

## Bibliography

- [1] Amal Ahmed, Robert Bruce Findler, Jeremy G. Siek, and Philip Wadler. Blame for all. In *Proceedings of the 38th Symposium on Principles of Programming Languages*, 2011.
- [2] Amal Ahmed, Dustin Jamner, Jeremy G. Siek, and Philip Wadler. Theorems for free for free: Parametricity, with and without types. In *Proceedings of the 22nd International Conference on Functional Programming*, 2017.
- [3] Felipe Bañados Schwerter, Ronald Garcia, and Éric Tanter. A theory of gradual effect systems. In *Proceedings of the 19th International Conference on Functional Programming*, 2014.
- [4] Gavin Bierman, Erik Meijer, and Mads Torgersen. Adding dynamic types to c#. In *Proceedings of the European Conference on Object-Oriented Programming*, 2010.
- [5] Gavin Bierman, Martín Abadi, and Mads Torgersen. Understanding type-script. In *Proceedings of the 28th European Conference on Object-Oriented Programming*, 2014.
- [6] Ambrose Bonnaire-Sergeant, Rowan Davies, and Sam Tobin-Hochstadt. Practical optional types for clojure. In *Programming Languages and Systems*. 2016.
- [7] Giuseppe Castagna and Victor Lanvin. Gradual typing with union and intersection types. *Proc. ACM Program. Lang.*, 1(ICFP):41:1–41:28, August 2017.
- [8] Matteo Cimini and Jeremy G. Siek. The gradualizer: A methodology and algorithm for generating gradual type systems. In *Proceedings of the 43rd Symposium on Principles of Programming Languages*, 2016.
- [9] Matteo Cimini and Jeremy G. Siek. Automatically generating the dynamic semantics of gradually typed languages. In *Proceedings of the 44th Symposium on Principles of Programming Languages*, 2017.
- [10] Dominique Devriese, Marco Patrignani, and Frank Piessens. Parametricity versus the universal type. *Proceedings of the ACM on Programming Languages*, 2(POPL):38, 2017.
- [11] Joshua Dunfield and Neelakantan R Krishnaswami. Complete and easy bidirectional typechecking for higher-rank polymorphism. In *International Conference on Functional Programming*, 2013.
- [12] Ronald Garcia and Matteo Cimini. Principal type schemes for gradual programs. In *Proceedings of the 42nd Symposium on Principles of Programming Languages*, 2015.
- [13] Ronald Garcia, Alison M Clark, and Éric Tanter. Abstracting gradual typing. In *Proceedings of the 43rd Symposium on Principles of Programming Languages*, 2016.
- [14] Yuu Igarashi, Taro Sekiyama, and Atsushi Igarashi. On polymorphic gradual typing. In *Proceedings of the 22nd International Conference on Functional Programming*, 2017.

- [15] Khurram A. Jafery and Joshua Dunfield. Sums of uncertainty: Refinements go gradual. In *Proceedings of the 44th Symposium on Principles of Programming Languages*, 2017.
- [16] John C Mitchell. Polymorphic type inference and containment. In *Logical foundations of functional programming*, 1990.
- [17] Martin Odersky and Konstantin Läufer. Putting type annotations to work. In *Proceedings of the 23rd Symposium on Principles of Programming Languages*, 1996.
- [18] Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Mark Shields. Practical type inference for arbitrary-rank types. *Journal of Functional Programming*, 17(1):1–82, 2007.
- [19] John C. Reynolds. Types, abstraction and parametric polymorphism. In *Proceedings of the IFIP 9th World Computer Congress*, 1983.
- [20] John C. Reynolds. The coherence of languages with intersection types. In *Proceedings of the International Conference on Theoretical Aspects of Computer Software*, 1991.
- [21] Jeremy G. Siek and Walid Taha. Gradual typing for functional languages. In *Proceedings of the 2006 Scheme and Functional Programming Workshop*, 2006.
- [22] Jeremy G. Siek and Walid Taha. Gradual typing for objects. In *European Conference on Object-Oriented Programming*, 2007.
- [23] Jeremy G. Siek and Manish Vachharajani. Gradual typing with unification-based inference. In *Proceedings of the 2008 Symposium on Dynamic Languages*, 2008.
- [24] Jeremy G. Siek and Philip Wadler. The key to blame: Gradual typing meets cryptography (draft), 2016.
- [25] Jeremy G. Siek, Michael M Vitousek, Matteo Cimini, and John Tang Boyland. Refined criteria for gradual typing. In *LIPICs-Leibniz International Proceedings in Informatics*, 2015.
- [26] Julien Verlaguet. Facebook: Analyzing php statically. In *Proceedings of Commercial Users of Functional Programming*, 2013.
- [27] Michael M. Vitousek, Andrew M. Kent, Jeremy G. Siek, and Jim Baker. Design and evaluation of gradual typing for python. In *Proceedings of the 10th Symposium on Dynamic languages*, 2014.
- [28] Dimitrios Vytiniotis, Simon Peyton Jones, and José Pedro Magalhães. Equality proofs and deferred type errors: A compiler pearl. In *Proceedings of the 17th International Conference on Functional Programming, ICFP '12*, New York, NY, USA, 2012.
- [29] Philip Wadler and Robert Bruce Findler. Well-typed programs can't be blamed. In *Proceedings of the 18th European Symposium on Programming Languages and Systems*, 2009.