

# Higher-rank Polymorphism: Type Inference and Extensions

*by*

**Ningning Xie**  
(谢宁宁)



A thesis submitted in partial fulfillment of the requirements for  
the degree of Doctor of Philosophy  
at The University of Hong Kong

February 2021



Abstract of thesis entitled  
**“Higher-rank Polymorphism: Type Inference and Extensions”**

Submitted by  
**Ningning Xie**

for the degree of Doctor of Philosophy  
at The University of Hong Kong  
in February 2021

---



# DECLARATION

I declare that this thesis represents my own work, except where due acknowledgment is made, and that it has not been previously included in a thesis, dissertation or report submitted to this University or to any other institution for a degree, diploma or other qualifications.

.....

**Ningning Xie**

February 2021



## ACKNOWLEDGMENTS





# CONTENTS

DECLARATION	I
ACKNOWLEDGMENTS	III
LIST OF FIGURES	IX
LIST OF TABLES	XI
I PROLOGUE	1
1 INTRODUCTION	3
1.1 Contributions . . . . .	3
1.2 Organization . . . . .	5
2 BACKGROUND	7
2.1 The Hindley-Milner Type System . . . . .	7
2.1.1 Declarative System . . . . .	7
2.1.2 Principal Type Scheme . . . . .	9
2.1.3 Algorithmic Type System . . . . .	10
2.2 The Odersky-Läufer Type System . . . . .	10
2.2.1 Higher-rank Types . . . . .	11
2.2.2 Declarative System . . . . .	12
2.2.3 Relating to HM . . . . .	14
2.3 The Dunfield-Krishnaswami Type System . . . . .	15
2.3.1 Bidirectional Type Checking . . . . .	15
2.3.2 Declarative System . . . . .	16
2.3.3 Algorithmic Type System . . . . .	19

II	BIDIRECTIONAL TYPE CHECKING WITH APPLICATION MODE	21
3	HIGHER-RANK POLYMORPHISM WITH APPLICATION MODE	23
3.1	Introduction and Motivation . . . . .	23
3.1.1	Revisiting Bidirectional Type Checking . . . . .	23
3.1.2	Type Checking with The Application Mode . . . . .	24
3.1.3	Benefits of Information Flowing from Arguments to Functions . . .	27
3.1.4	Type Inference of Higher-rank Types . . . . .	28
3.2	Declarative System . . . . .	30
3.2.1	Syntax . . . . .	31
3.2.2	Type System . . . . .	31
3.2.3	Subtyping . . . . .	35
3.3	Type-directed Translation . . . . .	37
3.3.1	Target Language . . . . .	37
3.3.2	Subtyping Coercions . . . . .	38
3.3.3	Type-Directed Translation of Typing . . . . .	41
3.3.4	Type Safety . . . . .	41
3.3.5	Coherence . . . . .	41
3.4	Type Inference Algorithm . . . . .	43
3.5	Discussion . . . . .	44
3.5.1	Combining Application and Checking Modes . . . . .	44
3.5.2	Additional Constructs . . . . .	45
3.5.3	More Expressive Type Applications . . . . .	46
3.5.4	Dependent Type Systems . . . . .	48
III	HIGHER-RANK POLYMORPHISM AND GRADUAL TYPING	51
4	GRADUALLY TYPED HIGHER-RANK POLYMORPHISM	53
4.1	Introduction and Motivation . . . . .	53
4.1.1	Background: Gradual Typing . . . . .	53
4.1.2	Motivation: Gradually Typed Higher-Rank Polymorphism . . . . .	55
4.1.3	Application: Efficient (Partly) Typed Encodings of ADTs . . . . .	56
4.2	Revisiting Consistent Subtyping . . . . .	59
4.2.1	Consistency and Subtyping . . . . .	59
4.2.2	Towards Consistent Subtyping . . . . .	62
4.2.3	Abstracting Gradual Typing . . . . .	65

4.2.4	Directed Consistency . . . . .	66
4.2.5	Consistent Subtyping Without Existentials . . . . .	67
4.3	Gradually Typed Implicit Polymorphism . . . . .	68
4.3.1	Typing in Detail . . . . .	69
4.3.2	Type-directed Translation . . . . .	69
4.3.3	Correctness Criteria . . . . .	73
4.4	Algorithmic Type System . . . . .	76
4.4.1	Algorithmic Consistent Subtyping . . . . .	78
4.4.2	Instantiation . . . . .	80
4.4.3	Algorithmic Typing . . . . .	82
4.4.4	Decidability . . . . .	84
4.4.5	Context Extension . . . . .	86
4.4.6	Soundness . . . . .	87
4.4.7	Completeness . . . . .	88
4.5	Simple Extensions and Variants . . . . .	89
4.5.1	Top Types . . . . .	89
4.5.2	A More Declarative Type System . . . . .	91
5	RESTORING THE DYNAMIC GRADUAL GUARANTEE WITH TYPE PARAMETERS . . . . .	93
5.1	Declarative Type System . . . . .	93
5.2	Substitutions and Representative Translations . . . . .	95
5.3	Dynamic Gradual Guarantee, Reloaded . . . . .	98
5.4	Extended Algorithmic Type System . . . . .	98
5.4.1	Extended Algorithmic Consistent Subtyping . . . . .	101
5.4.2	Extended Instantiation . . . . .	102
5.4.3	Algorithmic Typing and Metatheory . . . . .	104
5.4.4	Restricted Generalization . . . . .	104
IV	TYPE INFERENCE WITH PROMOTION . . . . .	107
6	HIGHER-RANK TYPE INFERENCE WITH TYPE PROMOTION . . . . .	109
6.1	Introduction and Motivation . . . . .	109
6.1.1	Background: Type Inference in Contexts . . . . .	109
6.1.2	Our Approach: Type Promotion . . . . .	111
6.1.3	Polymorphic Promotion . . . . .	112

## Contents

6.2	Unification for Simply Typed Lambda Calculus . . . . .	113
6.2.1	Declarative System . . . . .	114
6.2.2	Algorithmic System . . . . .	114
6.2.3	Soundness and Completeness . . . . .	115
6.3	Subtyping for Higher-Rank Polymorphism . . . . .	116
6.3.1	Declarative System . . . . .	116
6.3.2	Algorithmic System . . . . .	117
6.3.3	Soundness and Completeness . . . . .	119
6.4	Discussion . . . . .	120
6.4.1	Promoting Dependent Types . . . . .	120
6.4.2	Promoting Gradual Types . . . . .	121
7	KIND INFERENCE FOR DATATYPES . . . . .	123
7.1	Introduction and Motivation . . . . .	123
7.1.1	Kind Inference for Datatypes . . . . .	123
7.1.2	Kind Inference in Haskell98 . . . . .	125
7.1.3	Kind Inference in Modern GHC Haskell . . . . .	125
7.1.4	Desirable Properties for Kind Inference . . . . .	129
V	RELATED AND FUTURE WORK . . . . .	131
8	RELATED WORK . . . . .	133
9	FUTURE WORK . . . . .	135
VI	EPILOGUE . . . . .	137
10	CONCLUSION . . . . .	139
	BIBLIOGRAPHY . . . . .	141
VII	TECHNICAL APPENDIX . . . . .	147

# LIST OF FIGURES

2.1	Syntax and static semantics of the Hindley-Milner type system. . . . .	8
2.2	Subtyping in the Hindley-Milner type system. . . . .	9
2.3	Syntax of the Odersky-Läufer type system. . . . .	12
2.4	Well-formedness of types in the Odersky-Läufer type system. . . . .	12
2.5	Static semantics of the Odersky-Läufer type system. . . . .	13
2.6	Syntax of the Dunfield-Krishnaswami Type System . . . . .	16
2.7	Static semantics of the Dunfield-Krishnaswami type system. . . . .	17
3.1	Syntax of System AP. . . . .	31
3.2	Typing rules of System AP. . . . .	32
3.3	Syntax and typing rules of System F. . . . .	38
3.4	Subtyping translation rules of System AP. . . . .	39
3.5	Typing translation rules of System AP. . . . .	40
3.6	Type erasure and eta-id equality of System F. . . . .	42
4.1	Subtyping and type consistency in $\mathbf{FOb}_{<}^?$ . . . . .	54
4.2	Syntax of types, consistency, subtyping and well-formedness of types in declarative GPC. . . . .	60
4.3	Examples that break the original definition of consistent subtyping. . . . .	63
4.4	Observations of consistent subtyping . . . . .	63
4.5	Example that is fixed by the new definition of consistent subtyping. . . . .	64
4.6	Consistent Subtyping for implicit polymorphism. . . . .	68
4.7	Syntax of expressions and declarative typing of declarative GPC . . . . .	70
4.8	Less Precision . . . . .	75
4.9	Syntax and well-formedness of the algorithmic GPC . . . . .	77
4.10	Algorithmic consistent subtyping . . . . .	79
4.11	Algorithmic instantiation . . . . .	81
4.12	Algorithmic typing . . . . .	83
4.13	Context extension . . . . .	87

## *List of Figures*

5.1	Syntax of types, and consistent subtyping in the extended declarative system.	94
5.2	Syntax of types, contexts and consistent subtyping in the extended algorithmic system. . . . .	99
5.3	Extended algorithmic consistent subtyping . . . . .	100
5.4	Instantiation in the extended algorithmic system . . . . .	103
6.1	Types, contexts, unification and promotion of algorithmic STLC . . . . .	114
6.2	Types, contexts, subtyping and (polymorphic) promotion of the algorithmic system . . . . .	118

## LIST OF TABLES





# PART I

## PROLOGUE



# 1 INTRODUCTION

mention that in this thesis when we say “higher-rank polymorphism” we mean “predicative implicit higher-rank polymorphism”.

## 1.1 CONTRIBUTIONS

In summary the contributions of this thesis are:

- Part II**
- Chapter 3 proposes a new design for type inference of higher-rank polymorphism.
    - We design a variant of bi-directional type checking where the inference mode is combined with a new, so-called, application mode. The application mode naturally propagates type information from arguments to the functions.
    - With the application mode, we give a new design for type inference of higher-rank polymorphism, which generalizes the HM type system, supports a polymorphic let as syntactic sugar, and infers higher rank types. We present a syntax-directed specification, an elaboration semantics to System F, and an algorithmic type system with completeness and soundness proofs.
  - Chapter 6 presents a new approach for implementing unification.
    - We propose a process named *promotion*, which, given a unification variable and a type, promotes the type so that all unification variables in the type are well-typed with regard to the unification variable.
    - We apply promotion in a new implementation of the unification procedure in higher-rank polymorphism, and show that the new implementation is sound and complete.
- ??
- Chapter 4 extends higher-rank polymorphism with gradual types.
    - We define a framework for consistent subtyping with

- - ★ a new definition of consistent subtyping that subsumes and generalizes that of Siek and Taha [2007a] and can deal with polymorphism and top types;
  - ★ and a syntax-directed version of consistent subtyping that is sound and complete with respect to our definition of consistent subtyping, but still guesses instantiations.
- Based on consistent subtyping, we present the calculus GPC. We prove that our calculus satisfies the static aspects of the refined criteria for gradual typing [Siek et al. 2015], and is type-safe by a type-directed translation to  $\lambda B$  [Ahmed et al. 2009].
- We present a sound and complete bidirectional algorithm for implementing the declarative system based on the design principle of Garcia and Cimini [2015].
- Chapter 7 further explores the design of promotion in the context of kind inference for datatypes.
  - We formalize Haskell98’s datatype declarations, providing both a declarative specification and syntax-driven algorithm for kind inference. We prove that the algorithm is sound and observe how Haskell98’s technique of defaulting unconstrained kinds to  $\star$  leads to incompleteness. We believe that ours is the first formalization of this aspect of Haskell98.
  - We then present a type and kind language that is unified and dependently typed, modeling the challenging features for kind inference in modern Haskell. We include both a declarative specification and a syntax-driven algorithm. The algorithm is proved sound, and we observe where and why completeness fails. In the design of our algorithm, we must choose between completeness and termination; we favor termination but conjecture that an alternative design would regain completeness. Unlike other dependently typed languages, we retain the ability to infer top-level kinds instead of relying on compulsory annotations.

Many metatheory in the paper comes with Coq proofs, including type safety, coherence, etc.<sup>1</sup>

---

<sup>1</sup>For convenience, whenever possible, definitions, lemmas and theorems have hyperlinks (click ) to their Coq counterparts.

## 1.2 ORGANIZATION

This thesis is largely based on the publications by the author [Xie et al. 2018, 2019a,b; Xie and Oliveira 2017, 2018], as indicated below.

**Chapter 3:** Ningning Xie and Bruno C. d. S. Oliveira. 2018. “Let Arguments Go First”. In *European Symposium on Programming (ESOP)*.

**Chapter 6:** Ningning Xie and Bruno C. d. S. Oliveira. 2017. “Towards Unification for Dependent Types” (Extended abstract), In *Draft Proceedings of Trends in Functional Programming (TFP)*.

**Chapter 4:** Ningning Xie, Xuan Bi, and Bruno C. d. S. Oliveira. 2018. “Consistent Subtyping for All”. In *European Symposium on Programming (ESOP)*.

Ningning Xie, Xuan Bi, Bruno C. d. S. Oliveira, and Tom Schrijvers. 2019. “Consistent Subtyping for All”. In *ACM Transactions on Programming Languages and Systems (TOPLAS)*.

**Chapter 7:** Ningning Xie, Richard Eisenberg and Bruno C. d. S. Oliveira. 2020. “Kind Inference for Datatypes”. In *Symposium on Principles of Programming Languages (POPL)*.

---



## 2 BACKGROUND

This chapter sets the stage for type systems in later chapters. Section 2.1 reviews the Hindley-Milner type system [Damas and Milner 1982; Hindley 1969; Milner 1978], a classical type system for the lambda calculus with parametric polymorphism. Section 2.2 presents the Odersky-Läufer type system [Odersky and Läufer 1996], which extends upon the Hindley-Milner type system by putting higher-rank type annotations to work. Finally in Section 2.3 we introduce the Dunfield-Krishnaswami type system, a bidirectional higher-rank type system.

### 2.1 THE HINDLEY-MILNER TYPE SYSTEM

The global type-inference algorithms employed in modern functional languages such as ML, Haskell and OCaml, are derived from the Hindley-Milner type system. The Hindley-Milner type system, hereafter referred to as HM, is a polymorphic type discipline first discovered in Hindley [1969], later rediscovered by Milner [1978], and also closely formalized by Damas and Milner [1982]. In what follows, we first review its declarative specification, then discuss the property of principality, and finally talk briefly about its algorithmic system.

#### 2.1.1 DECLARATIVE SYSTEM

The declarative system of HM is given in Figure 2.1.

**SYNTAX.** The expressions  $e$  include variables  $x$ , literals  $n$ , lambda abstractions  $\lambda x. e$ , applications  $e_1 e_2$  and **let**  $x = e_1$  **in**  $e_2$ . Note here lambda abstractions have no type annotations, and the type information is to be reconstructed by the type system.

Types consist of polymorphic types  $\sigma$  and monomorphic types (monotypes)  $\tau$ . A polymorphic type is a sequence of universal quantifications (which can be empty) followed by a monotype  $\tau$ , which can be the integer type  $\text{Int}$ , type variables  $a$  and function types  $\tau_1 \rightarrow \tau_2$ .

A context  $\Psi$  tracks the type information for variables. We implicitly assume items in a context are distinct throughout the thesis.

## 2 Background

Expressions	$e ::= x \mid n \mid \lambda x. e \mid e_1 e_2 \mid \mathbf{let} x = e_1 \mathbf{in} e_2$
Types	$\sigma ::= \forall \bar{a}_i^i. \tau$
Monotypes	$\tau ::= \mathbf{Int} \mid a \mid \tau_1 \rightarrow \tau_2$
Contexts	$\Psi ::= \bullet \mid \Psi, x : \sigma$

  

$\Psi \vdash^{HM} e : \sigma$

(Typing)

  

$\frac{\text{HM-VAR} \quad (x : \sigma) \in \Psi}{\Psi \vdash^{HM} x : \sigma}$	$\frac{\text{HM-INT}}{\Psi \vdash^{HM} n : \mathbf{Int}}$	$\frac{\text{HM-LAM} \quad \Psi, x : \tau_1 \vdash^{HM} e : \tau_2}{\Psi \vdash^{HM} \lambda x. e : \tau_1 \rightarrow \tau_2}$
$\frac{\text{HM-APP} \quad \Psi \vdash^{HM} e_1 : \tau_1 \rightarrow \tau_2 \quad \Psi \vdash^{HM} e_2 : \tau_1}{\Psi \vdash^{HM} e_1 e_2 : \tau_2}$	$\frac{\text{HM-LET} \quad \Psi \vdash^{HM} e_1 : \sigma \quad \Psi, x : \sigma \vdash^{HM} e_2 : \tau}{\Psi \vdash^{HM} \mathbf{let} x = e_1 \mathbf{in} e_2 : \tau}$	
$\frac{\text{HM-GEN} \quad \bar{a}_i^i \notin \text{FV}(\Psi) \quad \Psi \vdash^{HM} e : \tau}{\Psi \vdash^{HM} e : \forall \bar{a}_i^i. \tau}$	$\frac{\text{HM-INST} \quad \Psi \vdash^{HM} e : \forall \bar{a}_i^i. \tau}{\Psi \vdash^{HM} e : \tau[\bar{a}_i \mapsto \bar{\tau}_i^i]}$	

Figure 2.1: Syntax and static semantics of the Hindley-Milner type system.

**TYPING.** The declarative typing judgment  $\Psi \vdash^{HM} e : \sigma$  derives the type  $\sigma$  of the expression  $e$  under the context  $\Psi$ . Rule [HM-VAR](#) fetches a polymorphic type  $x : \sigma$  from the context. Literals always have the integer type (rule [HM-INT](#)). For lambdas (rule [HM-LAM](#)), since there is no type given for the binder, the system *guesses* a *monotype*  $\tau_1$  as the type of  $x$ , and derives the type  $\tau_2$  for the body  $e$ , returning a function  $\tau_1 \rightarrow \tau_2$ . Function types are eliminated by applications. In rule [HM-APP](#), the type of the argument must match the parameter's type  $\tau_1$ , and the whole application returns type  $\tau_2$ .

Rule [HM-LET](#) is the key rule for flexibility in HM, where a *polymorphic* expression can be defined, and later instantiated with different types in the call sites. In this rule, the expression  $e_1$  has a polymorphic type  $\sigma$ , and the rule adds  $x : \sigma$  into the context to type-check  $e_2$ .

Rule [HM-GEN](#) and rule [HM-INST](#) correspond to *generalization* and *instantiation* respectively. In rule [HM-GEN](#), we can generalize over type variables  $\bar{a}_i^i$  which are not bound in the type context  $\Psi$ . In rule [HM-INST](#), we can instantiate the type variables with arbitrary *monotypes*.



$$\boxed{\vdash^{HM} \sigma_1 <: \sigma_2} \quad (Subtyping)$$

$$\begin{array}{c}
\text{HM-S-REFL} \\
\hline
\vdash^{HM} \tau <: \tau
\end{array}
\quad
\begin{array}{c}
\text{HM-S-FORALLR} \\
a \notin \text{FV}(\sigma_1) \quad \vdash^{HM} \sigma_1 <: \sigma_2 \\
\hline
\vdash^{HM} \sigma_1 <: \forall a. \sigma_2
\end{array}
\quad
\begin{array}{c}
\text{HM-S-FORALLL} \\
\vdash^{HM} \sigma_1[a \mapsto \tau] <: \sigma_2 \\
\hline
\vdash^{HM} \forall a. \sigma_1 <: \sigma_2
\end{array}$$

Figure 2.2: Subtyping in the Hindley-Milner type system.

## 2.1.2 PRINCIPAL TYPE SCHEME

One salient feature of HM is that the system enjoys the existence of *principal types*, without requiring any type annotations. Before we present the definition of principal types, let's first define the *subtyping* relation among types.

The judgment  $\vdash^{HM} \sigma_1 <: \sigma_2$ , given in Figure 2.2, reads that  $\sigma_1$  is a subtype of  $\sigma_2$ . The subtyping relation indicates that  $\sigma_1$  is more *general* than  $\sigma_2$ : for any instantiation of  $\sigma_2$ , we can find an instantiation of  $\sigma_1$  to make two types match. Rule **HM-S-REFL** is simply reflexive for monotypes. Rule **HM-S-FORALLR** has a polymorphic type  $\forall a. \sigma_2$  on the right hand side. In order to prove the subtyping relation for *all* possible instantiations of  $a$ , we *skolemize*  $a$ , by making sure  $a$  does not appear in  $\sigma_1$  (up to  $\alpha$ -renaming). In this case, if  $\sigma_1$  is still a subtype of  $\sigma_2$ , we are sure then whatever  $a$  can be instantiated to,  $\sigma_1$  can be instantiated to match  $\sigma_2$ . In rule **HM-S-FORALLL**, by contrast, the  $a$  in  $\forall a. \sigma_1$  can be instantiated to any monotype to match the right hand side. Here are some examples of the subtyping relation:

$$\begin{array}{l}
\vdash^{HM} \text{Int} \rightarrow \text{Int} <: \text{Int} \rightarrow \text{Int} \\
\vdash^{HM} \forall a. a \rightarrow a <: \text{Int} \rightarrow \text{Int}
\end{array}$$

Given the subtyping relation, now we can formally state that HM enjoys *principality*. That is, for every well-typed expression in HM, there exists one type for the expression, which is more general than any other types the expression can derive. Formally,

**Theorem 2.1** (Principality for HM). *If  $\Psi \vdash^{HM} e : \sigma$ , then there exists  $\sigma'$  such that  $\Psi \vdash^{HM} e : \sigma'$ , and for all  $\sigma$  such that  $\Psi \vdash^{HM} e : \sigma$ , we have  $\vdash^{HM} \sigma' <: \sigma$ .*

Consider the expression  $\lambda x. x$ . It has a principal type  $\forall a. a \rightarrow a$ , which is more general than other options, e.g.,  $\text{Int} \rightarrow \text{Int}$ ,  $(\text{Int} \rightarrow \text{Int}) \rightarrow (\text{Int} \rightarrow \text{Int})$ , etc.

## 2 Background

### 2.1.3 ALGORITHMIC TYPE SYSTEM

The declarative specification of HM given in Figure 2.1 does not directly lead to an algorithm. In particular, the system is not *syntax-directed*, and there are still many guesses in the system, such as in rule [HM-LAM](#).

**SYNTAX-DIRECTED SYSTEM.** A type system is *syntax-directed*, if the typing rules are completely driven by the syntax of expressions; in other words, there is exactly one typing rule for each syntactic form of expressions. However, in Figure 2.1, the rule for generalization (rule [HM-GEN](#)) and instantiation (rule [HM-INST](#)) can be applied anywhere.

A syntax-directed presentation of HM can be easily derived. In particular, from the typing rules we observe that, except for fetching a variable from the context (rule [HM-VAR](#)), the only place where a polymorphic type can be generated is for the let expressions (rule [HM-LET](#)). Thus, a syntax-directed system of HM can be presented as the original system, with instantiation applied to only variables, and generalization applied to only let expressions. Specifically,

$$\begin{array}{c}
 \text{HM-VAR-INST} \\
 \frac{(x : \forall \bar{a}_i^i. \tau) \in \Psi}{\Psi \vdash^{HM} x : \tau[\bar{a}_i \mapsto \bar{\tau}_i^i]}
 \end{array}
 \qquad
 \begin{array}{c}
 \text{HM-LET-GEN} \\
 \frac{\Psi \vdash^{HM} e_1 : \tau \quad \bar{a}_i^i = \text{FV}(\tau) - \text{FV}(\Psi) \quad \Psi, x : \forall \bar{a}_i^i. \tau \vdash^{HM} e_2 : \tau}{\Psi \vdash^{HM} \text{let } x = e_1 \text{ in } e_2 : \tau}
 \end{array}$$

**TYPE INFERENCE.** The guessing part of the system can be deterministically solved by the technique of *type inference*. There exists a sound and complete type inference algorithm for HM [Damas and Milner 1982], which has served as the basis for the type inference algorithm for many other systems [Odersky and Läufer 1996; Peyton Jones et al. 2007], including the system presented in Chapter 3. We will discuss more about it in Chapter 3.

## 2.2 THE ODERSKY-LÄUFER TYPE SYSTEM

The HM system is simple, flexible and powerful. Yet, since the type annotations in lambda abstractions are always missing, HM only derives polymorphic types of *rank 1*. That is, universal quantifiers only appear at the top level. Polymorphic types are of *higher-rank*, if universal quantifiers can appear anywhere in a type.

Essentially higher-rank types enable much of the expressive power of System F, with the advantage of implicit polymorphism. Complete type inference for System F is known to be undecidable [Wells 1999]. Odersky and Läufer [1996] proposed a type system, hereafter

referred to as OL, which extends HM by allowing lambda abstractions to have explicit *higher-rank* types as type annotations. As a motivation, consider the following program<sup>1</sup>:

```
(\f. (f 1, f 'a')) (\x. x)
```

which is not typeable under HM because it fails to infer the type of  $f$ :  $F$  is supposed to be polymorphic as it is applied to two arguments of different types. With OL we can add the type annotation for  $f$ :

```
(\f :  $\forall a. a \rightarrow a$ . (f 1, f 'a')) (\x. x)
```

Note that the first function now has a rank-2 type, as the polymorphic type  $\forall a. a \rightarrow a$  appears in the argument position of a function:

```
(\f :  $\forall a. a \rightarrow a$ . (f 1, f 'a')) : ( $\forall a. a \rightarrow a$ )  $\rightarrow$  (Int, Char)
```

In the rest of this section, we first give the definition of the rank of a type, and then present the declarative specification of OL, and show that OL is a conservative extension of HM.

### 2.2.1 HIGHER-RANK TYPES

We define the rank of types as follows.

**Definition 1** (Type rank). The *rank* of a type is the depth at which universal quantifiers appear contravariantly [Kfoury and Tiuryn 1992]. Formally,

---

$\text{rank}(\tau)$	$=$	0
$\text{rank}(\sigma_1 \rightarrow \sigma_2)$	$=$	$\max(\text{rank}(\sigma_1) + 1, \text{rank}(\sigma_2))$
$\text{rank}(\forall a. \sigma)$	$=$	$\max(1, \text{rank}(\sigma))$

---

Below we give some examples:

$\text{rank}(\text{Int} \rightarrow \text{Int})$	$=$	0
$\text{rank}(\forall a. a \rightarrow a)$	$=$	1
$\text{rank}(\text{Int} \rightarrow (\forall a. a \rightarrow a))$	$=$	1
$\text{rank}((\forall a. a \rightarrow a) \rightarrow \text{Int})$	$=$	2

From the definition, we can see that monotypes always have rank 0, and the polymorphic types in HM ( $\sigma$  in Figure 2.1) has at most rank 1.

<sup>1</sup>For the purpose of illustration, we assume basic constructs like booleans and pairs in examples.

## 2 Background

Expressions	$e ::=$	$x \mid n \mid \lambda x : \sigma. e \mid \lambda x. e \mid e_1 e_2 \mid \text{let } x = e_1 \text{ in } e_2$
Types	$\sigma ::=$	$\text{Int} \mid a \mid \sigma_1 \rightarrow \sigma_2 \mid \forall a. \sigma$
Monotypes	$\tau ::=$	$\text{Int} \mid a \mid \tau_1 \rightarrow \tau_2$
Contexts	$\Psi ::=$	$\bullet \mid \Psi, x : \sigma \mid \Psi, a$

Figure 2.3: Syntax of the Odersky-Läufer type system.

$\boxed{\Psi \vdash^{OL} \sigma}$	(Type Well-formedness)		
OL-WF-INT	OL-WF-TVAR	OL-WF-ARROW	OL-WF-FORALL
$\frac{}{\Psi \vdash^{OL} \text{Int}}$	$\frac{a \in \Psi}{\Psi \vdash^{OL} a}$	$\frac{\Psi \vdash^{OL} \sigma_1 \quad \Psi \vdash^{OL} \sigma_2}{\Psi \vdash^{OL} \sigma_1 \rightarrow \sigma_2}$	$\frac{\Psi, a \vdash^{OL} \sigma}{\Psi \vdash^{OL} \forall a. \sigma}$

Figure 2.4: Well-formedness of types in the Odersky-Läufer type system.

### 2.2.2 DECLARATIVE SYSTEM

**SYNTAX.** The syntax of OL is given in Figure 2.3. Comparing to HM, we observe the following differences.

First, expressions  $e$  include not only unannotated lambda abstractions  $\lambda x. e$ , but also annotated lambda abstractions  $\lambda x : \sigma. e$ , where the type annotation  $\sigma$  can be a polymorphic type. Thus unlike HM, the argument type for a function is not limited to a monotype.

Second, the polymorphic types  $\sigma$  now include the integer type  $\text{Int}$ , type variables  $a$ , functions  $\sigma_1 \rightarrow \sigma_2$  and universal quantifications  $\forall a. \sigma$ . Since the argument type in a function can be polymorphic, we see that OL supports *arbitrary* rank of types. The definition of monotypes remains the same, with polymorphic types still subsuming monotypes.

Finally, in addition to variable types, the contexts  $\Psi$  now also keep track of type variables. Note that in the original work in Odersky and Läufer [1996], the system, much like HM, does not track type variables; instead, it explicitly checks that type variables are fresh with respect to a context or a type when needed. Here we include type variables in contexts, as it sets us well for the Dunfield-Krishnaswami type system to be introduced in the next section. Moreover, it provides a complete view of possible formalisms of contexts in a type system with generalization.

Now since the context tracks type variables, we define the notion of *well-formedness* of types, given in Figure 2.4. For a type to be well-formedness, it must have all its free variable bound in the context. All rules are straightforward.

**TYPE SYSTEM.** The typing rules for OL are given in Figure 2.5.

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

(Typing)

$$\begin{array}{c}
 \text{OL-VAR} \\
 \frac{(x : \sigma) \in \Psi}{\Psi \vdash^{OL} x : \sigma}
 \end{array}
 \quad
 \begin{array}{c}
 \text{OL-INT} \\
 \frac{}{\Psi \vdash^{OL} n : \text{Int}}
 \end{array}
 \quad
 \begin{array}{c}
 \text{OL-LAMANN} \\
 \frac{\Psi, x : \sigma_1 \vdash^{OL} e : \sigma_2}{\Psi \vdash^{OL} \lambda x : \sigma_1. e : \sigma_1 \rightarrow \sigma_2}
 \end{array}$$
  

$$\begin{array}{c}
 \text{OL-LAM} \\
 \frac{\Psi \vdash^{OL} \tau \quad \Psi, x : \tau \vdash^{OL} e : \sigma}{\Psi \vdash^{OL} \lambda x. e : \tau \rightarrow \sigma}
 \end{array}
 \quad
 \begin{array}{c}
 \text{OL-APP} \\
 \frac{\Psi \vdash^{OL} e_1 : \sigma_1 \rightarrow \sigma_2 \quad \Psi \vdash^{OL} e_2 : \sigma_1}{\Psi \vdash^{OL} e_1 e_2 : \sigma_2}
 \end{array}$$
  

$$\begin{array}{c}
 \text{OL-LET} \\
 \frac{\Psi \vdash^{OL} e_1 : \sigma_1 \quad \Psi, x : \sigma_1 \vdash^{OL} e_2 : \sigma_2}{\Psi \vdash^{OL} \text{let } x = e_1 \text{ in } e_2 : \sigma_2}
 \end{array}
 \quad
 \begin{array}{c}
 \text{OL-GEN} \\
 \frac{\Psi, a \vdash^{OL} e : \sigma}{\Psi \vdash^{OL} e : \forall a. \sigma}
 \end{array}$$
  

$$\begin{array}{c}
 \text{OL-SUB} \\
 \frac{\Psi \vdash^{OL} e : \sigma_1 \quad \Psi \vdash^{OL} \sigma_1 <: \sigma_2}{\Psi \vdash^{OL} e : \sigma_2}
 \end{array}$$

$\Psi \vdash^{OL} \sigma_1 <: \sigma_2$

(Subtyping)

$$\begin{array}{c}
 \text{OL-S-TVAR} \\
 \frac{a \in \Psi}{\Psi \vdash^{OL} a <: a}
 \end{array}
 \quad
 \begin{array}{c}
 \text{OL-S-INT} \\
 \frac{}{\Psi \vdash^{OL} \text{Int} <: \text{Int}}
 \end{array}
 \quad
 \begin{array}{c}
 \text{OL-S-ARROW} \\
 \frac{\Psi \vdash^{OL} \sigma_3 <: \sigma_1 \quad \Psi \vdash^{OL} \sigma_2 <: \sigma_4}{\Psi \vdash^{OL} \sigma_1 \rightarrow \sigma_2 <: \sigma_3 \rightarrow \sigma_4}
 \end{array}$$
  

$$\begin{array}{c}
 \text{OL-S-FORALLL} \\
 \frac{\Psi \vdash^{OL} \tau \quad \Psi \vdash^{OL} \sigma[a \mapsto \tau] <: \sigma_2}{\Psi \vdash^{OL} \forall a. \sigma_1 <: \sigma_2}
 \end{array}
 \quad
 \begin{array}{c}
 \text{OL-S-FORALLR} \\
 \frac{\Psi, a \vdash^{OL} \sigma_1 <: \sigma_2}{\Psi \vdash^{OL} \sigma_1 <: \forall a. \sigma_2}
 \end{array}$$

Figure 2.5: Static semantics of the Odersky-Läufer type system.

## 2 Background

Rule **OL-VAR** and rule **OL-INT** are the same as that of HM. Rule **OL-LAMANN** type-checks annotated lambda abstractions, by simply putting  $x : \sigma$  into the context to type the body. For unannotated lambda abstractions in rule **OL-LAM**, the system still guesses a mere monotype. That is, the system never guesses a polymorphic type for lambdas; instead, an explicit polymorphic type annotation is required. Rule **OL-APP** and rule **OL-LET** are similar as HM, except that polymorphic types may appear in return types. In the generalization rule **OL-GEN**, we put a fresh type variable  $a$  into the context, and the return type  $\sigma$  is then generalized over  $a$ , returning  $\forall a. \sigma$ .

The subsumption rule **OL-SUB** is crucial for OL, which allows an expression of type  $\sigma_1$  to have type  $\sigma_2$  with  $\sigma_1$  being a subtype of  $\sigma_2$  ( $\Psi \vdash^{OL} \sigma_1 <: \sigma_2$ ). Note that the instantiation rule **HM-INST** in HM is a special case of rule **OL-SUB**, as we have  $\forall \bar{a}_i^i. \tau <: \tau[\bar{a}_i \mapsto \bar{\tau}_i^i]$  by applying rule **HM-S-FORALL** repeatedly.

The subtyping relation of OL  $\Psi \vdash^{OL} \sigma_1 <: \sigma_2$  also generalizes the subtyping relation of HM. In particular, in rule **OL-S-ARROW**, functions are *contravariant* on arguments, and *covariant* on return types. This rule allows us to compare higher-rank polymorphic types, rather than just polymorphic types with universal quantifiers only at the top level. For example,

$$\begin{array}{ll} \Psi \vdash^{OL} \forall a. a \rightarrow a & <: \text{Int} \rightarrow \text{Int} \\ \Psi \vdash^{OL} \text{Int} \rightarrow (\forall a. a \rightarrow a) & <: \text{Int} \rightarrow (\text{Int} \rightarrow \text{Int}) \\ \Psi \vdash^{OL} (\text{Int} \rightarrow \text{Int}) \rightarrow \text{Int} & <: (\forall a. a \rightarrow a) \rightarrow \text{Int} \end{array}$$

**PREDICATIVITY.** In a system with high-ranker types, one important design decision to make is whether the system is *predicative* or *impredicative*. A system is predicative, if the type variable bound by a universal quantifier is only allowed to be substituted by a monotype; otherwise it is impredicative. It is well-known that general type inference for impredicativity is undecidable [Wells 1999]. OL is predicative, which can be seen from rule **OL-S-FORALL**. We focus only on predicative type systems throughout the thesis.

### 2.2.3 RELATING TO HM

It can be proved that OL is a conservative extension of HM. That is, every well-typed expression in HM is well-typed in OL, modulo the different representation of contexts.

**Theorem 2.2** (Odersky-Läufer type system conservative over Hindley-Milner type system). *If  $\Psi \vdash^{HM} e : \sigma$ , suppose  $\Psi'$  is  $\Psi$  extended with type variables in  $\Psi$  and  $\sigma$ , then  $\Psi' \vdash^{OL} e : \sigma$ .*

Moreover, since OL is predicative and only guesses monotypes for unannotated lambda abstractions, its algorithmic system can be implemented as a direct extension of the one for HM.

## 2.3 THE DUNFIELD-KRISHNASWAMI TYPE SYSTEM

Both HM and OL derive only monotypes for unannotated lambda abstractions. OL improves on HM by allowing polymorphic lambda abstractions but requires the polymorphic type annotations to be given explicitly. The Dunfield-Krishnaswami type system [Dunfield and Krishnaswami 2013], hereafter referred to as DK, give a *bidirectional* account of higher-rank polymorphism, where type information can be propagated through the syntax tree. Therefore, it is possible for a variable bound in a lambda abstraction without explicit type annotations to get a polymorphic type. In this section, we first review the idea of bidirectional type checking, and then present the declarative DK and discuss its algorithm.

### 2.3.1 BIDIRECTIONAL TYPE CHECKING

Bidirectional type checking has been known in the folklore of type systems for a long time. It was popularized by Pierce and Turner’s work on local type inference [Pierce and Turner 2000]. Local type inference was introduced as an alternative to HM type systems, which could easily deal with polymorphic languages with subtyping. The key idea in local type inference is simple. The “local” in local type inference comes from the fact that:

*“... missing annotations are recovered using only information from adjacent nodes in the syntax tree, without long-distance constraints such as unification variables.”*

Bidirectional type checking is one component of local type inference that, aided by some type annotations, enables type inference in an expressive language with polymorphism and subtyping. In its basic form typing is split into *inference* and *checking* modes. The most salient feature of a bidirectional type-checker is when information deduced from inference mode is used to guide checking of an expression in checked mode.

Since Pierce and Turner’s work, various other authors have proved the effectiveness of bidirectional type checking in several other settings, including many different systems with subtyping [Davies and Pfenning 2000; Dunfield and Pfenning 2004], systems with dependent types [Asperti et al. 2012; Coquand 1996; Löh et al. 2010; Xi and Pfenning 1999], etc.

In particular, bidirectional type checking has also been combined with HM-style techniques for providing type inference in the presence of higher-rank type, including DK and Peyton Jones et al. [2007]. Let’s revisit the example in Section 2.2:

## 2 Background

Expressions	$e$	$::=$	$x \mid n \mid \lambda x : \sigma. e \mid \lambda x. e \mid e_1 e_2 \mid e : \sigma$
Types	$\sigma$	$::=$	$\text{Int} \mid a \mid \sigma_1 \rightarrow \sigma_2 \mid \forall a. \sigma$
Monotypes	$\tau$	$::=$	$\text{Int} \mid a \mid \tau_1 \rightarrow \tau_2$
Contexts	$\Psi$	$::=$	$\bullet \mid \Psi, x : \sigma \mid \Psi, a$

Figure 2.6: Syntax of the Dunfield-Krishnaswami Type System

```
(\f. (f 1, f 'a')) (\x. x)
```

which is not typeable in HM as it they fail to infer the type of  $f$ . In OL, it can be type-checked by adding a polymorphic type annotation on  $f$ . In DK, we can also add a polymorphic type annotation on  $f$ . But with bi-directional type checking, the type annotation can be propagated from somewhere else. For example, we can rewrite this program as:

```
((\f. (f 1, f 'c')) : (\forall a. a \rightarrow a) \rightarrow (\text{Int}, \text{Char})) (\x . x)
```

Here the type of  $f$  can be easily derived from the type signature using checking mode in bi-directional type checking.

### 2.3.2 DECLARATIVE SYSTEM

**SYNTAX.** The syntax of the DK is given in Figure 2.6. Comparing to OL, only the definition of expressions slightly differs. First, the expressions  $e$  in DK have no let expressions. Dunfield and Krishnaswami [2013] omitted let-bindings from the formal development, but argued that restoring let-bindings is easy, as long as they get no special treatment incompatible with substitution (e.g., a syntax-directed HM does polymorphic generalization only at let-bindings). Second, DK has annotated expressions  $e : \sigma$ , in which the type annotation can be propagated into the expression, as we will see shortly.

The definitions of types and contexts are the same as in OL. Thus, DK also shares the same well-formedness definition as in OL (Figure 2.4). We thus omit the definitions, but use  $\Psi \vdash^{DK} \sigma$  to denote the corresponding judgment in DK.

**TYPE SYSTEM.** Figure 2.7 presents the typing rules for DK. The system uses bidirectional type checking to accommodate polymorphism. Traditionally, two modes are employed in bidirectional systems: the inference mode  $\Psi \vdash^{DK} e \Rightarrow \sigma$ , which takes a term  $e$  and produces a type  $\sigma$ , similar to the judgment  $\Psi \vdash^{HM} e : \sigma$  or  $\Psi \vdash^{OL} e : \sigma$  in previous systems; the checking mode  $\Psi \vdash^{DK} e \Leftarrow \sigma$ , which takes a term  $e$  and a type  $\sigma$  as input, and ensures that the term  $e$  checks against  $\sigma$ . We first discuss rules in the inference mode.



$\Psi \vdash^{DK} e \Rightarrow \sigma$	<i>(Type Inference)</i>
<div style="display: flex; justify-content: space-between;"> <div style="width: 30%;"> <math display="block">\frac{\text{DK-INF-VAR} \quad (x : \sigma) \in \Psi}{\Psi \vdash^{DK} x \Rightarrow \sigma}</math> </div> <div style="width: 30%;"> <math display="block">\frac{\text{DK-INF-INT}}{\Psi \vdash^{DK} n \Rightarrow \text{Int}}</math> </div> <div style="width: 30%;"> <math display="block">\frac{\text{DK-INF-LAM} \quad \Psi \vdash^{DK} \tau_1 \rightarrow \tau_2 \quad \Psi, x : \tau_1 \vdash^{DK} e \Rightarrow \tau_2}{\Psi \vdash^{DK} \lambda x. e \Rightarrow \tau_1 \rightarrow \tau_2}</math> </div> </div>	
<div style="display: flex; justify-content: space-between;"> <div style="width: 45%;"> <math display="block">\frac{\text{DK-INF-APP} \quad \Psi \vdash^{DK} e_1 \Rightarrow \sigma \quad \Psi \vdash^{DK} \sigma \triangleright \sigma_1 \rightarrow \sigma_2 \quad \Psi \vdash^{DK} e_2 \Leftarrow \sigma_1}{\Psi \vdash^{DK} e_1 e_2 \Rightarrow \sigma_2}</math> </div> <div style="width: 45%;"> <math display="block">\frac{\text{DK-INF-ANNO} \quad \Psi \vdash^{DK} e \Leftarrow \sigma}{\Psi \vdash^{DK} e : \sigma \Rightarrow \sigma}</math> </div> </div>	
$\Psi \vdash^{DK} e \Leftarrow \sigma$	<i>(Type Checking)</i>
<div style="display: flex; justify-content: space-between;"> <div style="width: 30%;"> <math display="block">\frac{\text{DK-CHK-INT}}{\Psi \vdash^{DK} n \Leftarrow \text{Int}}</math> </div> <div style="width: 30%;"> <math display="block">\frac{\text{DK-CHK-LAM} \quad \Psi, x : \sigma_1 \vdash^{DK} e \Leftarrow \sigma_2}{\Psi \vdash^{DK} \lambda x. e \Leftarrow \sigma_1 \rightarrow \sigma_2}</math> </div> <div style="width: 30%;"> <math display="block">\frac{\text{DK-CHK-GEN} \quad \Psi, a \vdash^{DK} e \Leftarrow \sigma}{\Psi \vdash^{DK} e \Leftarrow \forall a. \sigma}</math> </div> </div>	
$\frac{\text{DK-CHK-SUB} \quad \Psi \vdash^{DK} e \Rightarrow \sigma_1 \quad \Psi \vdash^{DK} \sigma_1 <: \sigma_2}{\Psi \vdash^{DK} e \Leftarrow \sigma_2}$	
$\Psi \vdash^{DK} \sigma_1 \triangleright \sigma_2$	<i>(Matching)</i>
<div style="display: flex; justify-content: space-between;"> <div style="width: 45%;"> <math display="block">\frac{\text{DK-M-FORALL} \quad \Psi \vdash^{DK} \tau \quad \Psi \vdash^{DK} \sigma[a \mapsto \tau] \triangleright \sigma_1 \rightarrow \sigma_2}{\Psi \vdash^{DK} \forall a. \sigma \triangleright \sigma_1 \rightarrow \sigma_2}</math> </div> <div style="width: 45%;"> <math display="block">\frac{\text{DK-M-ARR}}{\Psi \vdash^{DK} \sigma_1 \rightarrow \sigma_2 \triangleright \sigma_1 \rightarrow \sigma_2}</math> </div> </div>	

Figure 2.7: Static semantics of the Dunfield-Krishnaswami type system.

## 2 Background

**TYPE INFERENCE.** Rule **DK-INF-VAR** and rule **DK-INF-INT** are straightforward. To infer unannotated lambdas, rule **DK-INF-LAM** guesses a monotype. For an application  $e_1 e_2$ , rule **DK-INF-APP** first infers the type  $\sigma$  of the expression  $e_1$ . Then, because  $e_1$  is applied to an argument, the type  $\sigma$  is decomposed into a function type  $\sigma_1 \rightarrow \sigma_2$ , using the matching judgment (discussed shortly). Now since the function expects an argument of type  $\sigma_1$ , the rule proceeds by checking  $e_2$  against  $\sigma_1$ . Similarly, for an annotated expression  $e : \sigma$ , rule **DK-INF-ANNO** simply checks  $e$  against  $\sigma$ . Both rules (rule **DK-INF-APP** and rule **DK-INF-ANNO**) have mode switched from inference to checking.

**TYPE CHECKING.** Now we turn to the checking mode. When an expression is checked against a type, the expression is expected to have that type. More importantly, the checking mode allows us to push the type information into the expressions.

Rule **DK-CHK-INT** checks literals again the integer type  $\text{Int}$ . Rule **DK-CHK-LAM** is where the system benefits from bidirectional type checking: the type information gets pushed inside an lambda. For an unannotated lambda abstraction  $\lambda x. e$ , recall that in the inference mode, we can only guess a monotype for  $x$ . With the checking mode, when  $\lambda x. e$  is checked against  $\sigma_1 \rightarrow \sigma_2$ , we do not need to guess any type. Instead,  $x$  gets directly the (possibly polymorphic) argument type  $\sigma_1$ . Then the rule proceeds by checking  $e$  with  $\sigma_2$ , allowing the type information to be pushed further inside. Note how rule **DK-CHK-LAM** improves over HM and OL, by allowing lambda abstractions to have a polymorphic argument type without requiring type annotations.

Rule **DK-CHK-GEN** deals with a polymorphic type  $\forall a. \sigma$ , by putting the (fresh) type variable  $a$  into the context to check  $e$  against  $\sigma$ . Rule **DK-CHK-SUB** switches the mode from checking to inference: an expression  $e$  can be checked against  $\sigma_2$ , if  $e$  infers the type  $\sigma_1$  and  $\sigma_1$  is a subtype of  $\sigma_2$ .

**MATCHING.** In rule **DK-INF-APP** where we type-check an application  $e_1 e_2$ , we derive that  $e_1$  has type  $\sigma$ , but  $e_1$  must have a function type so that it can be applied to an argument. The *matching* judgment instantiates  $\sigma$  into a function.

Matching has two straightforward rules: rule **DK-M-FORALL** instantiates a polymorphic type, by substituting  $a$  with a well-formed monotype  $\tau$ , and continues matching on  $\sigma[a \mapsto \tau]$ ; rule **DK-M-ARR** returns the function type directly.

In Dunfield and Krishnaswami [2013], they use an *application judgment* instead of matching. The application judgment  $\Psi \vdash^{DK} \sigma_1 \cdot e \Rightarrow \sigma_2$ , whose definition is given below, is interpreted as, when we apply an expression of type  $\sigma_1$  to the expression  $e$ , we get a return type  $\sigma_2$ .

$$\boxed{\Psi \vdash^{DK} \sigma_1 \cdot e \Rightarrow \sigma_2} \quad (\text{Application})$$

$$\begin{array}{c}
 \text{DK-APP-FORALL} \\
 \frac{\Psi \vdash^{DK} \tau \quad \Psi \vdash^{DK} \sigma[a \mapsto \tau] \cdot e \Rightarrow \sigma_1 \rightarrow \sigma_2}{\Psi \vdash^{DK} \forall a. \sigma \cdot e \Rightarrow \sigma_1 \rightarrow \sigma_2}
 \end{array}
 \quad
 \begin{array}{c}
 \text{DK-APP-ARR} \\
 \frac{\Psi \vdash^{DK} e \Leftarrow \sigma_1}{\Psi \vdash^{DK} \sigma_1 \rightarrow \sigma_2 \cdot e \Rightarrow \sigma_2}
 \end{array}$$

With the application judgment, rule **DK-INF-APP** is replaced by rule **DK-INF-APP2**.

$$\begin{array}{c}
 \text{DK-INF-APP2} \\
 \frac{\Psi \vdash^{DK} e_1 \Rightarrow \sigma \quad \Psi \vdash^{DK} \sigma \cdot e_2 \Rightarrow \sigma_2}{\Psi \vdash^{DK} e_1 e_2 \Rightarrow \sigma_2}
 \end{array}$$

It can be easily shown that the presentation of rule **DK-INF-APP** with matching is equivalent to that of rule **DK-INF-APP2** with the application judgment. Essentially, they both make sure that the expression being applied has an arrow type  $\sigma_1 \rightarrow \sigma_2$ , and then check the argument against  $\sigma_1$ .

We prefer the presentation of rule **DK-INF-APP** with matching, as matching is a simple and independent process whose purpose is clear. In contrast, it is relatively less comprehensible with rule **DK-INF-APP2** and the application judgment, where all three forms of the judgment (inference, checking, application) are mutually dependent.

**SUBTYPING.** DK shares the same subtyping relation as of OL. We thus omit the definition and use  $\Psi \vdash^{DK} \sigma_1 <: \sigma_2$  to denote the subtyping relation in DK.

### 2.3.3 ALGORITHMIC TYPE SYSTEM

Dunfield and Krishnaswami [2013] also presented a sound and complete bidirectional algorithmic type system. The key idea of the algorithm is using *ordered* algorithmic contexts for storing existential variables and their solutions. Comparing to the algorithm for HM, they argued that their algorithm is remarkably simple. The algorithm is later discussed and used in Part III and Part IV. We will discuss more about it there.



## PART II

# BIDIRECTIONAL TYPE CHECKING WITH APPLICATION MODE



# 3 HIGHER-RANK POLYMORPHISM WITH APPLICATION MODE

We have seen in Section 2.3 that bi-directional type checking is an useful and versatile tool for type checking and type inference. In traditional bi-directional type-checking, type information flows from functions to arguments (e.g., rule [DK-IN-APP](#) in Section 2.3.2). In this section, we present a novel variant of bi-directional type checking where the type information flows from arguments to functions. This variant retains the inference mode, but adds a so-called *application* mode. Such design can remove annotations that basic bi-directional type checking cannot, and is useful when type information from arguments is required to type-check the functions being applied.

We illustrate our novel design of bi-directional type-checking using System AP, a lambda calculus with implicit higher-rank polymorphism. This section first presents the declarative, syntax-directed type system of System AP in Section 3.2. The interesting aspects about the new type system are: 1) the typing rules, which employ a combination of the inference mode and the *application* mode; 2) the novel subtyping relation under an application context. Later, we prove our type system is type-safe by a type directed translation to System F in Section 3.3. An algorithmic type system is discussed in Section 3.4.

## 3.1 INTRODUCTION AND MOTIVATION

### 3.1.1 REVISITING BIDIRECTIONAL TYPE CHECKING

Traditional type checking rules can be heavyweight on annotations, in the sense that lambda-bound variables always need explicit annotations. As we have seen in Section 2.3, bidirectional type checking [Pierce and Turner 2000] provides an alternative, which allows types to propagate downward the syntax tree. For example, in the expression  $(\lambda f : \text{Int} \rightarrow \text{Int}. f) (\lambda y. y)$ , the type of  $y$  is provided by the type annotation on  $f$ . This is supported by the bidirectional typing rule [DK-INF-APP](#) for applications:

$$\frac{\text{DK-INF-APP} \quad \Psi \vdash^{DK} e_1 \Rightarrow \sigma \quad \Psi \vdash^{DK} \sigma \triangleright \sigma_1 \rightarrow \sigma_2 \quad \Psi \vdash^{DK} e_2 \Leftarrow \sigma_1}{\Psi \vdash^{DK} e_1 e_2 \Rightarrow \sigma_2}$$

Specifically, if we know that the type of  $e_1$  is a function from  $\sigma_1 \rightarrow \sigma_2$ , we can check that  $e_2$  has type  $\sigma_1$ . Notice that here the type information flows from functions to arguments.

One guideline for designing bidirectional type checking rules [Dunfield and Pfenning 2004] is to distinguish introduction rules from elimination rules. Constructs which correspond to introduction forms are *checked* against a given type, while constructs corresponding to elimination forms *infer* (or synthesize) their types. For instance, under this design principle, the introduction rule for literals is supposed to be in checking mode, as in the rule [DK-CHK-INT](#):

$$\frac{\text{DK-CHK-INT}}{\Psi \vdash^{DK} n \Leftarrow \text{Int}}$$

Unfortunately, this means that the trivial program 1 cannot type-check, which in this case has to be rewritten to  $1 : \text{Int}$ .

In this particular case, bidirectional type checking goes against its original intention of removing burden from programmers, since a seemingly unnecessary annotation is needed. Therefore, in practice, bidirectional type systems do not strictly follow the guideline, and usually have additional inference rules for the introduction form of constructs. For literals, the corresponding rule is rule [DK-INF-INT](#).

$$\frac{\text{DK-INF-INT}}{\Psi \vdash^{DK} n \Rightarrow \text{Int}}$$

Now we can type check 1, but the price to pay is that two typing rules for literals are needed. Worse still, the same criticism applies to other constructs (e.g., pairs). This shows one drawback of bidirectional type checking: often to minimize annotations, many rules are duplicated for having both inference and checking mode, which scales up with the typing rules in a type system.

#### 3.1.2 TYPE CHECKING WITH THE APPLICATION MODE

We propose a variant of bidirectional type checking with a new *application mode* (unrelated to the application judgment in DK). The application mode preserves the advantage of bidirectional type checking, namely many redundant annotations are removed, while certain



programs can type check with even fewer annotations. Also, with our proposal, the inference mode is a special case of the application mode, so it does not produce duplications of rules in the type system. Additionally, the checking mode can still be *easily* combined into the system. The essential idea of the application mode is to enable the type information flow in applications to propagate from arguments to functions (instead of from functions to arguments as in traditional bidirectional type checking).

To motivate the design of bidirectional type checking with an application mode, consider the simple expression

$(\lambda x. x) 1$

This expression cannot type check in traditional bidirectional type checking, because unannotated abstractions, as a construct which correspond to introduction forms, only have a checking mode, so annotations are required<sup>1</sup>. For example,  $((\lambda x. x) : \mathbf{Int} \rightarrow \mathbf{Int}) 1$ .

In this example we can observe that if the type of the argument is accounted for in inferring the type of  $\lambda x. x$ , then it is actually possible to deduce that the lambda expression has type  $\mathbf{Int} \rightarrow \mathbf{Int}$ , from the argument 1.

**THE APPLICATION MODE.** If types flow from the arguments to the function, an alternative idea is to push the type of the arguments into the typing of the function, as follows,

$$\frac{\text{APP} \quad \Psi \vdash e_2 \Rightarrow \sigma_1 \quad \Psi; \Sigma, \sigma_1 \vdash e_1 \Rightarrow \sigma \rightarrow \sigma_2}{\Psi; \Sigma \vdash e_1 e_2 \Rightarrow \sigma_2}$$

In this rule, there are two kinds of judgments. The first judgment is just the usual inference mode, which is used to infer the type of the argument  $e_2$ . The second judgment, the application mode, is similar to the inference mode, but it has an additional context  $\Sigma$ . The context  $\Sigma$  is a stack that tracks the types of the arguments of outer applications. In the rule for application, the type of the argument  $e_2$  synthesizes its type  $\sigma_1$ , which then is pushed into the application context  $\Sigma$  for inferring the type of  $e_1$ . Applications are themselves in the application mode, since they can be in the context of an outer application.

---

<sup>1</sup>It type-checks in DK, because in DK rules for lambdas are duplicated for having both inference (integrated with type inference techniques) and checking mode.

### 3 Higher-Rank Polymorphism with Application Mode

Lambda expressions can now make use of the application context, leading to the following rule:

$$\text{LAM} \quad \frac{\Psi, x : \sigma; \Sigma \vdash e \Rightarrow \sigma_2}{\Psi; \Sigma, \sigma \vdash \lambda x. e \Rightarrow \sigma \rightarrow \sigma_2}$$

The type  $\sigma$  that appears last in the application context serves as the type for  $x$ , and type checking continues with a smaller application context and  $x : \sigma$  in the typing context. Therefore, using the rule [APP](#) and rule [LAM](#), the expression  $(\lambda x. x) 1$  can type-check without annotations, since the type  $\text{Int}$  of the argument  $1$  is used as the type of the binding  $x$ .

Note that, since the examples so far are based on simple types, obviously they can be solved by integrating type inference and relying on techniques like unification or constraint solving (as in DK). However, here the point is that the application mode helps to reduce the number of annotations *without requiring such sophisticated techniques*. Also, the application mode helps with situations where those techniques cannot be easily applied, such as type systems with subtyping.

**INTERPRETATION OF THE APPLICATION MODE.** As we have seen, the guideline for designing bi-directional type checking [Dunfield and Pfenning 2004], based on introduction and elimination rules, is often not enough in practice. This leads to extra introduction rules in the inference mode. The application mode does not distinguish between introduction rules and elimination rules. Instead, to decide whether a rule should be in inference or application mode, we need to think whether the expression can be applied or not. Variables, lambda expressions and applications are all examples of expressions that can be applied, and they should have application mode rules. However literals or pairs cannot be applied and should have inference rules. For example, type checking pairs would simply have the inference mode. Nevertheless elimination rules of pairs could have non-empty application contexts (see Section 3.5.2 for details). In the application mode, arguments are always inferred first in applications and propagated through application contexts. An empty application context means that an expression is not being applied to anything, which allows us to model the inference mode as a particular case<sup>2</sup>.

**PARTIAL TYPE CHECKING.** The inference mode synthesizes the type of an expression, and the checked mode checks an expression against some type. A natural question is how do these

<sup>2</sup>Although the application mode generalizes the inference mode, we refer to them as two different modes. Thus the variant of bi-directional type checking in this paper is interpreted as a type system with both *inference* and *application* modes.

modes compare to application mode. An answer is that, in some sense: the application mode is stronger than inference mode, but weaker than checked mode. Specifically, the inference mode means that we know nothing about the type an expression before hand. The checked mode means that the whole type of the expression is already known before hand. With the application mode we know some partial type information about the type of an expression: we know some of its argument types (since it must be a function type when the application context is non-empty), but not the return type.

Instead of nothing or all, this partialness gives us a finer grain notion on how much we know about the type of an expression. For example, assume  $e : \sigma_1 \rightarrow \sigma_2 \rightarrow \sigma_3$ . In the inference mode, we only have  $e$ . In the checked mode, we have both  $e$  and  $\sigma_1 \rightarrow \sigma_2 \rightarrow \sigma_3$ . In the application mode, we have  $e$ , and maybe an empty context (which degenerates into inference mode), or an application context  $\sigma_1$  (we know the type of first argument), or an application context  $\sigma_1, \sigma_2$  (we know the types of both arguments).

**TRADE-OFFS.** Note that the application mode is *not* conservative over traditional bidirectional type checking due to the different information flow. However, it provides a new design choice for type inference/checking algorithms, especially for those where the information about arguments is useful. Therefore we next discuss some benefits of the application mode for two interesting cases where functions are either variables; or lambda (or type) abstractions.

### 3.1.3 BENEFITS OF INFORMATION FLOWING FROM ARGUMENTS TO FUNCTIONS

**LOCAL CONSTRAINT SOLVER FOR FUNCTION VARIABLES.** Many type systems, including type systems with *implicit polymorphism* and/or *static overloading*, need information about the types of the arguments when type checking function variables. For example, in conventional functional languages with implicit polymorphism, function calls such as (id 1) where  $\text{id} : \forall a. (a \rightarrow a)$ , are *pervasive*. In such a function call the type system must instantiate  $a$  to  $\text{Int}$ . Dealing with such implicit instantiation gets trickier in systems with *higher-rank types*. For example, Peyton Jones et al. [2007] require additional syntactic forms and relations, whereas DK add a special purpose matching or the application judgment.

With the application mode, all the type information about the arguments being applied is available in application contexts and can be used to solve instantiation constraints. To exploit such information, the type system employs a special subtyping judgment called *application subtyping*, with the form  $\Sigma \vdash \sigma_1 <: \sigma_2$ . Unlike conventional subtyping, computationally  $\Psi$  and  $\sigma_1$  are interpreted as inputs and  $\sigma_2$  as output. In above example, we have that  $\text{Int} \vdash \forall a. a \rightarrow a <: \sigma$  and we can determine that  $a = \text{Int}$  and  $\sigma = \text{Int} \rightarrow \text{Int}$ . In this way, type

system is able to solve the constraints *locally* according to the application contexts since we no longer need to propagate the instantiation constraints to the typing process.

DECLARATION DESUGARING FOR LAMBDA ABSTRACTIONS. An interesting consequence of the usage of an application mode is that it enables the following **let** sugar:

$$\mathbf{let} \ x = e_1 \ \mathbf{in} \ e_2 \rightsquigarrow (\lambda x. e_2) e_1$$

Such syntactic sugar for **let** is, of course, standard. However, in the context of implementations of typed languages it normally requires extra type annotations or a more sophisticated type-directed translation. Type checking  $(\lambda x. e_2) e_1$  would normally require annotations (for example a higher-rank type annotation for  $x$  as in OL and DK), or otherwise such annotation should be inferred first. Nevertheless, with the application mode no extra annotations/inference is required, since from the type of the argument  $e_1$  it is possible to deduce the type of  $x$ . Generally speaking, with the application mode *annotations are never needed for applied lambdas*. Thus **let** can be the usual sugar from the untyped lambda calculus, including HM-style **let** expression and even type declarations.

#### 3.1.4 TYPE INFERENCE OF HIGHER-RANK TYPES

We believe the application mode can be integrated into many traditional bidirectional type systems. In this chapter, we focus on integrating the application mode into a bidirectional type system with higher-rank types. Our paper [Xie and Oliveira 2018] includes another application to System F.

Consider again the motivation example used in Section 2.2:

$(\backslash f. (f \ 1, f \ \text{'a'})) (\backslash x. x)$

which is not typeable in HM, but can be rewritten to include type annotations in OL and DK. For example, both in OL and DK we can write:

$(\backslash f: (\forall a. a \rightarrow a). (f \ 1, f \ \text{'c'})) (\backslash x. x)$

However, although some redundant annotations are removed by bidirectional type checking, the burden of inferring higher-rank types is still carried by programmers: they are forced to add polymorphic annotations to help with the type derivation of higher-rank types. For the above example, the type annotation is still *provided by programmers*, even though the necessary type information can be derived intuitively without any annotations:  $f$  is applied to  $\backslash x. x$ , which is of type  $\forall a. a \rightarrow a$ .

TYPE INFERENCE FOR HIGHER-RANK TYPES WITH THE APPLICATION MODE. Using our bi-directional type system with an application mode, the original expression can type check without annotations or rewrites:  $(\backslash f. (f\ 1, f\ 'c')) (\backslash x. x)$ .

This result comes naturally if we allow type information flow from arguments to functions. For inferring polymorphic types for arguments, we use *generalization*. In the above example, we first infer the type  $\forall a. a \rightarrow a$  for the argument, then pass the type to the function. A nice consequence of such an approach is that, as mentioned before, HM-style polymorphic **let** expressions are simply regarded as syntactic sugar to a combination of lambda/application:

$$\text{let } x = e_1 \text{ in } e_2 \rightsquigarrow (\lambda x. e_2) e_1$$

With this approach, nested lets can lead to types which are *more general* than HM. For example, consider the following expression:

**let**  $s = \backslash x. x$  **in** **let**  $t = \backslash y. s$  **in**  $e$

The type of  $s$  is  $\forall a. a \rightarrow a$  after generalization. Because  $t$  returns  $s$  as a result, we might expect  $t: \forall b. b \rightarrow (\forall a. a \rightarrow a)$ , which is what our system will return. However, HM will return type  $t: \forall b. \forall a. b \rightarrow (a \rightarrow a)$ , as it can only return rank 1 types, which is less general than the previous one according to the subtyping relation for polymorphic types in OL (Figure 2.5).

CONSERVATIVITY OVER THE HINDLEY-MILNER TYPE SYSTEM. Our type system is a conservative extension over HM, in the sense that every program that can type-check in HM is accepted in our type system. This result is not surprising: after desugaring **let** into a lambda and an application, programs remain typeable.

COMPARING PREDICATIVE HIGHER-RANK TYPE INFERENCE SYSTEMS. We will give a full discussion and comparison of related work in Section 8. Among those works, we believe DK and the work by Peyton Jones et al. [2007] are the most closely related work to our system. Both their systems and ours are based on a *predicative* type system: universal quantifiers can only be instantiated by monotypes. So we would like to emphasize our system's properties in relation to those works. In particular, here we discuss two interesting differences, and also briefly (and informally) discuss how the works compare in terms of expressiveness.

1) Inference of higher-rank types. In both works, every polymorphic type inferred by the system must correspond to one annotation provided by the programmer. However, in our system, some higher-rank types can be inferred from the expression itself without any annotation. The motivating expression above provides an example of this.

2) Where are annotations needed? Since type annotations are useful for inferring higher rank types, a clear answer to the question where annotations are needed is necessary so that programmers know when they are required to write annotations. To this question, previous systems give a concrete answer: only on the binding of polymorphic types. Our answer is slightly different: only on the bindings of polymorphic types in abstractions *that are not applied to arguments*. Roughly speaking this means that our system ends up with fewer or smaller annotations.

3) Expressiveness. Based on these two answers, it may seem that our system should accept all expressions that are typeable in their system. However, this is not true because the application mode is *not* conservative over traditional bi-directional type checking. Consider the expression:

$$(\backslash f : (\forall a. a \rightarrow a) \rightarrow (\text{nat}, \text{char}). f) (\backslash g. (g\ 1, g\ 'a'))$$

which is typeable in their system. In this case, even if  $g$  is a polymorphic binding without a type annotation the expression can still type-check. This is because the original application rule propagates the information from the outer binding into the inner expressions. Note that the fact that such expression type-checks does not contradict their guideline of providing type annotations for every polymorphic binder. Programmers that strictly follow their guideline can still add a polymorphic type annotation for  $g$ . However it does mean that it is a little harder to understand where annotations for polymorphic binders can be *omitted* in their system. This requires understanding how the applications in checked mode operate.

In our system the above expression is not typeable, as a consequence of the information flow in the application mode. However, following our guideline for annotations leads to a program that can be type-checked with a smaller annotation:

$$(\backslash f. f) (\backslash g : (\forall a. a \rightarrow a). (g\ 1, g\ 'a')).$$

This means that our work is not conservative over their work, which is due to the design choice of the application typing rule. Nevertheless, we can always rewrite programs using our guideline, which often leads to fewer/smaller annotations.

## 3.2 DECLARATIVE SYSTEM

This section first presents the declarative, *syntax-directed* specification of AP. The interesting aspects about the new type system are: 1) the typing rules, which employ a combination of inference and application modes; 2) the novel subtyping relation under an application context.

Expressions	$e$	$::=$	$x \mid n \mid \lambda x : \sigma. e \mid \lambda x. e \mid e_1 e_2$
Types	$\sigma$	$::=$	$\text{Int} \mid a \mid \sigma_1 \rightarrow \sigma_2 \mid \forall a. \sigma$
Monotypes	$\tau$	$::=$	$\text{Int} \mid a \mid \tau_1 \rightarrow \tau_2$
Contexts	$\Psi$	$::=$	$\bullet \mid \Psi, x : \sigma$
Application Contexts	$\Sigma$	$::=$	$\bullet \mid \Sigma, \sigma$

Figure 3.1: Syntax of System AP.

### 3.2.1 SYNTAX

The syntax of the language is given in Figure 3.1.

**EXPRESSIONS.** The definition of expressions  $e$  include variables ( $x$ ), integers ( $n$ ), annotated lambda abstractions ( $\lambda x : \sigma. e$ ), lambda abstractions ( $\lambda x. e$ ), and applications ( $e_1 e_2$ ). Notably, the syntax does not include a **let** expression (**let**  $x = e_1$  **in**  $e_2$ ). Let expressions can be regarded as the standard syntax sugar  $(\lambda x. e_2) e_1$ , as illustrated in more detail later.

**TYPES.** Types include the integer type  $\text{Int}$ , type variables ( $a$ ), functions ( $\sigma_1 \rightarrow \sigma_2$ ) and polymorphic types ( $\forall a. \sigma$ ). Monotypes are types without universal quantifiers.

**CONTEXTS.** Typing contexts  $\Psi$  are standard: they map a term variable  $x$  to its type  $\sigma$ . In this system, the context is modeled in a HM-style context (Section 2.1), which does not track type variables. Again, we implicitly assume that all the variables in  $\Psi$  are distinct.

The main novelty lies in the *application contexts*  $\Sigma$ , which are the main data structure needed to allow types to flow from arguments to functions. Application contexts are modeled as a stack. The stack collects the types of arguments in applications. The context is a stack because if a type is pushed last then it will be popped first. For example, inferring expression  $e$  under application context  $(a, \text{Int})$ , means  $e$  is now being applied to two arguments  $e_1, e_2$ , with  $e_1 : \text{Int}$ ,  $e_2 : a$ , so  $e$  should be of type  $\text{Int} \rightarrow a \rightarrow \sigma$  for some  $\sigma$ .

### 3.2.2 TYPE SYSTEM

The top part of Figure 3.2 gives the typing rules for our language. The judgment  $\Psi; \Sigma \vdash^{AP} e \Rightarrow \sigma$  is read as: under typing context  $\Psi$ , and application context  $\Sigma$ ,  $e$  has type  $\sigma$ . The standard inference mode  $\Psi \vdash^{AP} e \Rightarrow \sigma$  can be regarded as a special case when the application context is empty. Note that the variable names are assumed to be fresh enough when new variables are added into the typing context, or when generating new type variables.

### 3 Higher-Rank Polymorphism with Application Mode

$\boxed{\Psi \vdash^{AP} e \Rightarrow \sigma}$			(Typing Inference)
$\frac{\text{AP-INF-INT}}{\Psi \vdash^{AP} n \Rightarrow \text{Int}}$	$\frac{\text{AP-INF-LAM}}{\Psi, x : \tau \vdash^{AP} e \Rightarrow \sigma} \quad \Psi \vdash^{AP} \lambda x. e \Rightarrow \tau \rightarrow \sigma$	$\frac{\text{AP-INF-LAMANN}}{\Psi \vdash^{AP} \lambda x : \sigma_1. e \Rightarrow \sigma_1 \rightarrow \sigma_2} \quad \Psi, x : \sigma_1 \vdash^{AP} e \Rightarrow \sigma_2$	
$\boxed{\Psi; \Sigma \vdash^{AP} e \Rightarrow \sigma}$			(Typing Application Mode)
$\frac{\text{AP-APP-VAR}}{(x : \sigma_1) \in \Psi \quad \Sigma \vdash^{AP} \sigma_1 <: \sigma_2} \quad \Psi; \Sigma \vdash^{AP} x \Rightarrow \sigma_2$	$\frac{\text{AP-APP-LAM}}{\Psi; \Sigma, \sigma_1 \vdash^{AP} \lambda x. e \Rightarrow \sigma_1 \rightarrow \sigma_2} \quad \Psi, x : \sigma_1 \vdash^{AP} e \Rightarrow \sigma_2$		
$\frac{\text{AP-APP-LAMANN}}{\Psi; \Sigma, \sigma_2 \vdash^{AP} \lambda x : \sigma_1. e \Rightarrow \sigma_2 \rightarrow \sigma_3} \quad \vdash^{AP} \sigma_2 <: \sigma_1 \quad \Psi, x : \sigma_1 \vdash^{AP} e \Rightarrow \sigma_3$			
$\frac{\text{AP-APP-APP}}{\Psi; \Sigma \vdash^{AP} e_1 e_2 \Rightarrow \sigma_3} \quad \begin{array}{l} \Psi \vdash^{AP} e_2 \Rightarrow \sigma_1 \quad \bar{a}_i^i = \text{FV}(\sigma_1) - \text{FV}(\Psi) \\ \sigma_2 = \forall \bar{a}_i^i. \sigma_1 \quad \Psi; \Sigma, \sigma_2 \vdash^{AP} e_1 \Rightarrow \sigma_2 \rightarrow \sigma_3 \end{array}$			
$\boxed{\vdash^{AP} \sigma_1 <: \sigma_2}$			(Subtyping)
$\frac{\text{AP-S-TVAR}}{\vdash^{AP} a <: a}$	$\frac{\text{AP-S-INT}}{\vdash^{AP} \text{Int} <: \text{Int}}$	$\frac{\text{AP-S-ARROW}}{\vdash^{AP} \sigma_1 \rightarrow \sigma_2 <: \sigma_3 \rightarrow \sigma_4} \quad \vdash^{AP} \sigma_3 <: \sigma_1 \quad \vdash^{AP} \sigma_2 <: \sigma_4$	
$\frac{\text{AP-S-FORALLL}}{\vdash^{AP} \forall a. \sigma_1 <: \sigma_2} \quad \vdash^{AP} \sigma[a \mapsto \tau] <: \sigma_2$	$\frac{\text{AP-S-FORALLR}}{\vdash^{AP} \sigma_1 <: \forall a. \sigma_2} \quad a \notin \text{FV}(\sigma_1) \quad \vdash^{AP} \sigma_1 <: \sigma_2$		
$\boxed{\Sigma \vdash^{AP} \sigma_1 <: \sigma_2}$			(Application Subtyping)
$\frac{\text{AP-AS-EMPTY}}{\bullet \vdash^{AP} \sigma <: \sigma}$	$\frac{\text{AP-AS-FORALL}}{\Sigma, \sigma_3 \vdash^{AP} \forall a. \sigma_1 <: \sigma_2} \quad \Sigma, \sigma_3 \vdash^{AP} \sigma_1[a \mapsto \tau] <: \sigma_2$	$\frac{\text{AP-AS-ARROW}}{\Sigma, \sigma_3 \vdash^{AP} \sigma_1 \rightarrow \sigma_2 <: \sigma_3 \rightarrow \sigma_4} \quad \vdash^{AP} \sigma_3 <: \sigma_1 \quad \Sigma \vdash^{AP} \sigma_2 <: \sigma_4$	

Figure 3.2: Typing rules of System AP.



We discuss the rules when the application context is empty first. Those rules are unsurprising. Rule **AP-INF-INT** shows that integer literals are only inferred to have type `Int` under an empty application context. This is obvious since an integer cannot accept any arguments. Rule **AP-INF-LAM** deals with lambda abstractions when the application context is empty. In this situation, a monotype  $\tau$  is *guessed* for the argument, just like previous calculi. Rule **AP-INF-LAMANN** also works as expected: a new variable  $x$  is put with its type  $\sigma$  into the typing context, and inference continues on the abstraction body.

Now we turn to the cases when the application context is not empty. Rule **AP-APP-VAR** says that if  $x : \sigma_1$  is in the typing context, and  $\sigma_1$  is a subtype of  $\sigma_2$  under application context  $\Sigma$ , then  $x$  has type  $\sigma_2$ . It depends on the subtyping rules that are explained in Section 3.2.3.

Rule **AP-APP-LAM** shows the strength of application contexts. It states that, without annotations, if the application context is non-empty, a type can be popped from the application context to serve as the type for  $x$ . Inference of the body then continues with the rest of the application context. This is possible, because the expression  $\lambda x. e$  is being applied to an argument of type  $\sigma_1$ , which is the type at the top of the application context stack.

For lambda abstraction with annotations  $\lambda x : \sigma_1. e$ , if the application context has type  $\sigma_2$ , then rule **AP-APP-LAMANN** first checks that  $\sigma_2$  is a subtype of  $\sigma_1$  before putting  $x : \sigma_1$  in the typing context. However, note that it is always possible to remove annotations in an abstraction if it has been applied to some arguments.

Rule **AP-APP-APP** pushes types into the application context. The application rule first infers the type of the argument  $e_2$  with type  $\sigma_1$ . Then the type  $\sigma_1$  is generalized in the same way as the HM type system. The resulting generalized type is  $\sigma_2$ . Thus the type of  $e_1$  is now inferred under an application context extended with type  $\sigma_2$ . The generalization step is important to infer higher ranked types: since  $\sigma_2$  is a possibly polymorphic type, which is the argument type of  $e_1$ , then  $e_1$  is of possibly a higher rank type.

**LET EXPRESSIONS.** The language does not have built-in let expressions, but instead supports **let** as syntactic sugar. Recall the syntactic-directed typing rule (rule **HM-LET-GEN**) for let expressions with generalization in the HM system. With slight reformat to match AP, we get (without the gray-shaded part):

$$\frac{\Psi \vdash e_1 \Rightarrow \sigma_1 \quad \overline{a_i}^i = \text{FV}(\tau) - \text{FV}(\Psi) \quad \sigma_2 = \forall \overline{a_i}^i. \sigma_1 \quad \Psi, x : \sigma_2; \Sigma \vdash e_2 \Rightarrow \sigma_3}{\Psi; \Sigma \vdash \text{let } x = e_1 \text{ in } e_2 \Rightarrow \sigma_3}$$

where we do generalization on the type of  $e_1$ , which is then assigned as the type of  $x$  while inferring  $e_2$ . Adapting this rule to our system with application contexts would result in the

gray-shaded part, where the application context is only used for  $e_2$ , because  $e_2$  is the expression being applied. If we desugar the let expression (**let**  $x = e_1$  **in**  $e_2$ ) to  $(\lambda x. e_2) e_1$ , we have the following derivation:

$$\frac{\Psi \vdash e_1 \Rightarrow \sigma_1 \quad \bar{a}_i^i = \text{FV}(\sigma_1) - \text{FV}(\Psi) \quad \sigma_2 = \forall \bar{a}_i^i. \sigma_1 \quad \frac{\Psi, x : \sigma_2; \Sigma \vdash e_2 \Rightarrow \sigma_3}{\Psi; \Sigma, \sigma_2 \vdash \lambda x. e_2 \Rightarrow \sigma_2 \rightarrow \sigma_3}}{\Psi; \Sigma \vdash (\lambda x. e_2) e_1 \Rightarrow \sigma_3}$$

The type  $\sigma_2$  is now pushed into application context in rule **AP-APP-APP**, and then assigned to  $x$  in rule **AP-APP-LAM**. Comparing this with the typing derivations for let expressions, we now have the same preconditions. Thus we can see that the rules in Figure 3.2 are sufficient to express an HM-style polymorphic let construct.

**METATHEORY.** The type system enjoys several interesting properties, especially lemmas about application contexts. Before we present those lemmas, we need a helper definition of what it means to use arrows on application contexts.

**Definition 2** ( $\Sigma \rightarrow \sigma$ ). If  $\Sigma = \sigma_1, \sigma_2, \dots, \sigma_n$ , then  $\Sigma \rightarrow \sigma$  means the function type  $\sigma_n \rightarrow \dots \rightarrow \sigma_2 \rightarrow \sigma_1 \rightarrow \sigma$ .

Such definition is useful to reason about the typing result with application contexts. One specific property is that the application context determines the form of the typing result.

**Lemma 3.1** ( $\Sigma$  Coincides with Typing Results). *If  $\Psi; \Sigma \vdash^{AP} e \Rightarrow \sigma$ , then for some  $\sigma'$ , we have  $\sigma = \Sigma \rightarrow \sigma'$ .*

Having this lemma, we can always use the judgment  $\Psi; \Sigma \vdash^{AP} e \Rightarrow \Sigma \rightarrow \sigma'$  instead of  $\Psi; \Sigma \vdash^{AP} e \Rightarrow \sigma$ .

In traditional bi-directional type checking, we often have one subsumption rule that transfers between inference and checked mode, which states that if an expression can be inferred to some type, then it can be checked with this type (e.g., rule **DK-CHK-SUB** in DK). In our system, we regard the normal inference mode  $\Psi \vdash^{AP} e \Rightarrow \sigma$  as a special case, when the application context is empty. We can also turn from normal inference mode into application mode with an application context.

**Lemma 3.2** (Subsumption). *If  $\Psi \vdash^{AP} e \Rightarrow \Sigma \rightarrow \sigma$ , then  $\Psi; \Sigma \vdash^{AP} e \Rightarrow \Sigma \rightarrow \sigma$ .*

This lemma is actually a special case for the following one:

**Lemma 3.3** (Generalized Subsumption). *If  $\Psi; \Sigma_1 \vdash^{AP} e \Rightarrow \Sigma_1 \rightarrow \Sigma_2 \rightarrow \sigma$ , then  $\Psi; \Sigma_2, \Sigma_1 \vdash^{AP} e \Rightarrow \Sigma_1 \rightarrow \Sigma_2 \rightarrow \sigma$ .*

The relationship between our system and standard Hindley Milner type system (HM) can be established through the desugaring of let expressions. Namely, if  $e$  is typeable in HM, then the desugared expression  $e'$  is typeable in our system, with a more general typing result.

**Lemma 3.4** (AP Conservative over HM). *If  $\Psi \vdash^{HM} e : \sigma$ , and desugaring let expression in  $e$  gives back  $e'$ , then for some  $\sigma'$ , we have  $\Psi \vdash^{AP} e' \Rightarrow \sigma'$ , and  $A' <: A$ .*

### 3.2.3 SUBTYPING

We present our subtyping rules at the bottom of Figure 3.2. Interestingly, our subtyping has two different forms.

**SUBTYPING.** The first subtyping judgment  $\vdash^{AP} \sigma_1 <: \sigma_2$  follows OL, but in HM-style; that is, without tracking type variables. Rule **AP-S-FORALLR** states  $\sigma_1$  is subtype of  $\forall a. \sigma_2$  only if  $\sigma_1$  is a subtype of  $\sigma_2$ , with the assumption  $a$  is a fresh variable. Rule **AP-S-FORALLL** says  $\forall a. \sigma_1$  is a subtype of  $\sigma_2$  if we can instantiate it with some  $\tau$  and show the result is a subtype of  $\sigma_2$ .

**APPLICATION SUBTYPING.** The typing rule **AP-APP-VAR** uses the second subtyping judgment  $\Sigma \vdash^{AP} \sigma_1 <: \sigma_2$ . To motivate this new kind of judgment, consider the expression `id 1` for example, whose derivation is stuck at rule **AP-APP-VAR** (here we assume  $\text{id} : \forall a. a \rightarrow a \in \Psi$ ):

$$\frac{\Psi \vdash^{AP} 1 \Rightarrow \text{Int} \quad \emptyset = \text{FV}(\text{Int}) - \text{FV}(\Psi) \quad \frac{\text{id} : \forall a. a \rightarrow a \in \Psi \quad ???}{\Psi; \text{Int} \vdash^{AP} \text{id} \Rightarrow ?} \text{AP-APP-VAR}}{\Psi \vdash^{AP} \text{id } 1 \Rightarrow ?} \text{AP-APP-APP}$$

Here we know that  $\text{id} : \forall a. a \rightarrow a$  and also, from the application context, that `id` is applied to an argument of type `Int`. Thus we need a mechanism for solving the instantiation  $a = \text{Int}$  and return a supertype  $\text{Int} \rightarrow \text{Int}$  as the type of `id`. This is precisely what the application subtyping achieves: resolve instantiation constraints according to the application context. Notice that unlike existing works (Peyton Jones et al. [2007] or DK), application subtyping provides a way to solve instantiation more *locally*, since it does not mutually depend on typing.

Back to the rules in Figure 3.2, one way to understand the judgment  $\Sigma \vdash^{AP} \sigma_1 <: \sigma_2$  from a computational point-of-view is that the type  $\sigma_2$  is a *computed* output, rather than an input. In other words  $\sigma_2$  is determined from  $\Sigma$  and  $\sigma_1$ . This is unlike the judgment  $\vdash^{AP} \sigma_1 <: \sigma_2$ , where both  $\sigma_1$  and  $\sigma_2$  would be computationally interpreted as inputs. Therefore it is not possible to view  $\vdash^{AP} \sigma_1 <: \sigma_2$  as a special case of  $\Sigma \vdash^{AP} \sigma_1 <: \sigma_2$  where  $\Sigma$  is empty.

There are three rules dealing with application contexts. Rule **AP-AS-EMPTY** is for case when the application context is empty. Because it is empty, we have no constraints on the type, so we return it back unchanged. Note that this is where HM-style systems (also Peyton Jones et al. [2007]) would normally use an instantiation rule (e.g. rule **HM-INST** in HM) to remove top-level quantifiers. Our system does not need the instantiation rule, because in applications, type information flows from arguments to the function, instead of function to arguments. In the latter case, the instantiation rule is needed because a function type is wanted instead of a polymorphic type. In our approach, instantiation of type variables is avoided unless necessary.

The two remaining rules apply when the application context is non-empty, for polymorphic and function types respectively. Note that we only need to deal with these two cases because  $\text{Int}$  or type variables  $a$  cannot have a non-empty application context. In rule **AP-AS-FORALL**, we instantiate the polymorphic type with some  $\tau$ , and continue. This instantiation is forced by the application context. In rule **AP-AS-ARROW**, one function of type  $\sigma_1 \rightarrow \sigma_2$  is now being applied to an argument of type  $\sigma_3$ . So we check  $\vdash^{AP} \sigma_3 <: \sigma_1$ . Then we continue with  $\sigma_2$  and the rest application context, and return  $\sigma_3 \rightarrow \sigma_4$  as the result type of the function.

**METATHEORY.** Application subtyping is novel in our system, and it enjoys some interesting properties. For example, similarly to typing, the application context decides the form of the supertype.

**Lemma 3.5** ( $\Sigma$  Coincides with Subtyping Results). *If  $\Sigma \vdash^{AP} \sigma_1 <: \sigma_2$ , then for some  $\sigma_3$ ,  $\sigma_2 = \Sigma \rightarrow \sigma_3$ .*

Therefore we can always use the judgment  $\Sigma \vdash^{AP} \sigma_1 <: \Sigma \rightarrow \sigma_2$ , instead of  $\Sigma \vdash^{AP} \sigma_1 <: \sigma_2$ . Application subtyping is also reflexive and transitive. Interestingly, in those lemmas, if we remove all applications contexts, they are exactly the reflexivity and transitivity of traditional subtyping.

**Lemma 3.6** (Reflexivity of Application Subtyping).  $\Sigma \vdash^{AP} \Sigma \rightarrow \sigma <: \Sigma \rightarrow \sigma$ .

**Lemma 3.7** (Transitivity of Application Subtyping). *If  $\Sigma_1 \vdash^{AP} \sigma_1 <: \Sigma_1 \rightarrow \sigma_2$ , and  $\Sigma_2 \vdash^{AP} \sigma_2 <: \Sigma_2 \rightarrow \sigma_3$ , then  $\Sigma_2, \Sigma_1 \vdash^{AP} \sigma_1 <: \Sigma_1 \rightarrow \Sigma_2 \rightarrow \sigma_3$ .*

Finally, we can convert between subtyping and application subtyping. We can remove the application context and still get a subtyping relation:

**Lemma 3.8** ( $\Sigma \vdash^{AP}$  to  $\vdash^{AP}$ ). *If  $\Sigma \vdash^{AP} \sigma_1 <: \sigma_2$ , then  $\vdash^{AP} \sigma_1 <: \sigma_2$ .*

Transferring from subtyping to application subtyping will result in a more general type.

**Lemma 3.9** ( $\vdash^{AP}$  to  $\vdash^{AP}$ ). *If  $\vdash^{AP} \sigma_1 <: \Sigma \rightarrow \sigma_2$ , then for some  $\sigma_3$ , we have  $\Sigma \vdash^{AP} \sigma_1 <: \Sigma \rightarrow \sigma_3$ , and  $\vdash^{AP} \sigma_3 <: \sigma_2$ .*

This lemma may not seem intuitive at first glance. Consider a concrete example. Consider the derivation for  $\vdash^{AP} \text{Int} \rightarrow \forall a. a <: \text{Int} \rightarrow \text{Int}$ :

$$\frac{\frac{\vdash^{AP} \text{Int} <: \text{Int} \quad \text{AP-S-INT} \quad \frac{\vdash^{AP} \text{Int} <: \text{Int} \quad \text{AP-S-INT}}{\vdash^{AP} \forall a. a <: \text{Int}} \quad \text{AP-S-FORALLL}}{\vdash^{AP} \text{Int} \rightarrow \forall a. a <: \text{Int} \rightarrow \text{Int}} \quad \text{AP-S-ARROW}$$

and for  $\text{Int} \vdash^{AP} \text{Int} \rightarrow \forall a. a <: \text{Int} \rightarrow \forall a. a$ :

$$\frac{\frac{\vdash^{AP} \text{Int} <: \text{Int} \quad \text{AP-S-INT}}{\text{Int} \vdash^{AP} \text{Int} \rightarrow \forall a. a <: \text{Int} \rightarrow \forall a. a} \quad \frac{\vdash^{AP} \forall a. a <: \forall a. a \quad \text{AP-AS-EMPTY}}{\text{Int} \vdash^{AP} \text{Int} \rightarrow \forall a. a <: \text{Int} \rightarrow \forall a. a} \quad \text{AP-AS-ARROW}$$

The former one, holds because we have  $\vdash^{AP} \forall a. a <: \text{Int}$  in the return type. But in the latter one, after  $\text{Int}$  is consumed from application context, we eventually reach rule [AP-AS-EMPTY](#), which always returns the original type back.

### 3.3 TYPE-DIRECTED TRANSLATION

This section discusses the type-directed translation of System AP into a variant of System F that is also used in Peyton Jones et al. [2007]. The translation is shown to be coherent and type safe. The later result implies the type-safety of the source language. To prove coherency, we need to decide when two translated terms are the same using  $\eta$ -id equality, and show that the translation is unique up to this definition.

#### 3.3.1 TARGET LANGUAGE

The syntax and typing rules of our target language are given in Figure 3.3.

Expressions include variables  $x$ , integers  $n$ , annotated abstractions  $\lambda x : \sigma. s$ , type-level abstractions  $\Lambda a. s$ , and  $s_1 s_2$  for term application, and  $s \sigma$  for type application. The types and the typing contexts are the same as our system, where typing contexts does not track type

Expressions	$s, f ::= x \mid n \mid \lambda x : \sigma. s \mid \Lambda a. s \mid s_1 s_2 \mid s \sigma$
Types	$\sigma ::= \text{Int} \mid a \mid \sigma_1 \rightarrow \sigma_2 \mid \forall a. \sigma$
Contexts	$\Psi ::= \bullet \mid \Psi, x : \sigma$

  

$\Psi \vdash^F s : \sigma$

*(Typing)*

  

$\frac{\text{F-VAR} \quad (x : \sigma) \in \Psi}{\Psi \vdash^F x : \sigma}$	$\frac{\text{F-INT}}{\Psi \vdash^F n : \text{Int}}$	$\frac{\text{F-LAMANN} \quad \Psi, x : \sigma_1 \vdash^F s : \sigma_2}{\Psi \vdash^F \lambda x : \sigma_1. s : \sigma_1 \rightarrow \sigma_2}$
$\frac{\text{F-APP} \quad \Psi \vdash^F s_1 : \sigma_1 \rightarrow \sigma_2 \quad \Psi \vdash^F s_2 : \sigma_1}{\Psi \vdash^F s_1 s_2 : \sigma_2}$	$\frac{\text{F-TABS} \quad \Psi \vdash^F s : \sigma \quad a \notin \text{FV}(\Psi)}{\Psi \vdash^F \Lambda a. s : \forall a. \sigma}$	
$\frac{\text{F-TAPP} \quad \Psi \vdash^F s : \forall a. \sigma_1}{\Psi \vdash^F s \sigma_2 : \sigma_1[a \mapsto \sigma_2]}$		

Figure 3.3: Syntax and typing rules of System F.

variables. In translation, we use  $f$  to refer to the coercion function produced by subtyping translation, and  $s$  to refer to the translated term in System F.

Most typing rules are straightforward. Rule **F-TABS** types a type abstraction with explicit generalization, while rule **F-TAPP** types a type application with explicit instantiation.

### 3.3.2 SUBTYPING COERCIONS

The type-directed translation of subtyping is shown in Figure 3.4. The translation follows the subtyping relations from Figure 3.2, but adds a new component. The judgment  $\vdash^{AP} \sigma_1 <: \sigma_2 \rightsquigarrow f$  is read as: if  $\vdash^{AP} \sigma_1 <: \sigma_2$  holds, it can be translated to a coercion function  $f$  in System F. The coercion function produced by subtyping is used to transform values from one type to another, so we should have  $\bullet \vdash^F f : \sigma_1 \rightarrow \sigma_2$ .

Rule **AP-S-INT** and rule **AP-S-TVAR** produce identity functions, since the source type and target type are the same. In rule **AP-S-ARROW**, the coercion function  $f_1$  of type  $\sigma_3 \rightarrow \sigma_1$  is applied to  $y$  to get a value of type  $\sigma_1$ . Then the resulting value becomes an argument to  $x$ , and a value of type  $\sigma_2$  is obtained. Finally we apply  $f_2$  to the value of type  $\sigma_2$ , so that a value of type  $\sigma_4$  is computed. In rule **A-PS-FORALLL**, the input argument is a polymorphic type, so we feed the type  $\tau$  to it and apply the coercion function  $f$  from the precondition. Rule **AP-S-**

$\boxed{\vdash^{AP} \sigma_1 <: \sigma_2 \rightsquigarrow f}$  (Subtyping Translation)

AP-S-TVAR

$$\frac{}{\vdash^{AP} a <: a \rightsquigarrow \lambda x : a. x}$$

AP-S-INT

$$\frac{}{\vdash^{AP} \text{Int} <: \text{Int} \rightsquigarrow \lambda x : \text{Int}. x}$$

AP-S-ARROW

$$\frac{\vdash^{AP} \sigma_3 <: \sigma_1 \rightsquigarrow f_1 \quad \vdash^{AP} \sigma_2 <: \sigma_4 \rightsquigarrow f_2}{\vdash^{AP} \sigma_1 \rightarrow \sigma_2 <: \sigma_3 \rightarrow \sigma_4 \rightsquigarrow \lambda x : \sigma_1 \rightarrow \sigma_2. \lambda y : \sigma_3. f_2(x(f_1 y))}$$

AP-S-FORALLL

$$\frac{\vdash^{AP} \sigma[a \mapsto \tau] <: \sigma_2 \rightsquigarrow f}{\vdash^{AP} \forall a. \sigma_1 <: \sigma_2 \rightsquigarrow \lambda x : \forall a. \sigma_1. f(x\tau)}$$

AP-S-FORALLR

$$\frac{a \notin \text{FV}(\sigma_1) \quad \vdash^{AP} \sigma_1 <: \sigma_2 \rightsquigarrow f}{\vdash^{AP} \sigma_1 <: \forall a. \sigma_2 \rightsquigarrow \lambda x : \sigma_1. \Lambda a. f x}$$

$\boxed{\Sigma \vdash^{AP} \sigma_1 <: \sigma_2 \rightsquigarrow f}$  (Application Subtyping)

AP-AS-EMPTY

$$\frac{}{\bullet \vdash^{AP} \sigma <: \sigma \rightsquigarrow \lambda x : \sigma. x}$$

AP-AS-FORALL

$$\frac{\Sigma, \sigma_3 \vdash^{AP} \sigma_1[a \mapsto \tau] <: \sigma_2 \rightsquigarrow f}{\Sigma, \sigma_3 \vdash^{AP} \forall a. \sigma_1 <: \sigma_2 \rightsquigarrow \lambda x : \forall a. \sigma_1. f(x\tau)}$$

AP-AS-ARROW

$$\frac{\vdash^{AP} \sigma_3 <: \sigma_1 \rightsquigarrow f_1 \quad \Sigma \vdash^{AP} \sigma_2 <: \sigma_4 \rightsquigarrow f_2}{\Sigma, \sigma_3 \vdash^{AP} \sigma_1 \rightarrow \sigma_2 <: \sigma_3 \rightarrow \sigma_4 \rightsquigarrow \lambda x : \sigma_1 \rightarrow \sigma_2. \lambda y : \sigma_3. f_2(x(f_1 y))}$$

Figure 3.4: Subtyping translation rules of System AP.

$\Psi \vdash^{AP} e \Rightarrow \sigma \rightsquigarrow s$

*(Typing Inference)*

$$\frac{\text{AP-INF-INT}}{\Psi \vdash^{AP} n \Rightarrow \text{Int} \rightsquigarrow n}$$

$$\frac{\text{AP-INF-LAM} \quad \Psi, x : \tau \vdash^{AP} e \Rightarrow \sigma \rightsquigarrow s}{\Psi \vdash^{AP} \lambda x. e \Rightarrow \tau \rightarrow \sigma \rightsquigarrow \lambda x : \tau. s}$$

$$\frac{\text{AP-INF-LAMANN} \quad \Psi, x : \sigma_1 \vdash^{AP} e \Rightarrow \sigma_2 \rightsquigarrow s}{\Psi \vdash^{AP} \lambda x : \sigma_1. e \Rightarrow \sigma_1 \rightarrow \sigma_2 \rightsquigarrow \lambda x : \sigma_1. s}$$

$\Psi; \Sigma \vdash^{AP} e \Rightarrow \sigma \rightsquigarrow s$

*(Typing Application Mode)*

$$\frac{\text{AP-APP-VAR} \quad (x : \sigma_1) \in \Psi \quad \Sigma \vdash^{AP} \sigma_1 <: \sigma_2 \rightsquigarrow f}{\Psi; \Sigma \vdash^{AP} x \Rightarrow \sigma_2 \rightsquigarrow f x}$$

$$\frac{\text{AP-APP-LAM} \quad \Psi, x : \sigma_1 \vdash^{AP} e \Rightarrow \sigma_2 \rightsquigarrow s}{\Psi; \Sigma, \sigma_1 \vdash^{AP} \lambda x. e \Rightarrow \sigma_1 \rightarrow \sigma_2 \rightsquigarrow \lambda x : \sigma_1. s}$$

$$\frac{\text{AP-APP-LAMANN} \quad \vdash^{AP} \sigma_2 <: \sigma_1 \rightsquigarrow f \quad \Psi, x : \sigma_1 \vdash^{AP} e \Rightarrow \sigma_3 \rightsquigarrow s}{\Psi; \Sigma, \sigma_2 \vdash^{AP} \lambda x : \sigma_1. e \Rightarrow \sigma_2 \rightarrow \sigma_3 \rightsquigarrow \lambda y : \sigma_2. (\lambda x : \sigma_1. s) (f y)}$$

$$\frac{\text{AP-APP-APP} \quad \Psi \vdash^{AP} e_2 \Rightarrow \sigma_1 \rightsquigarrow s_2 \quad \overline{a_i}^i = \text{FV}(\sigma_1) - \text{FV}(\Psi) \quad \sigma_2 = \forall \overline{a_i}^i. \sigma_1 \quad \Psi; \Sigma, \sigma_2 \vdash^{AP} e_1 \Rightarrow \sigma_2 \rightarrow \sigma_3 \rightsquigarrow s_1}{\Psi; \Sigma \vdash^{AP} e_1 e_2 \Rightarrow \sigma_3 \rightsquigarrow s_1 (\Lambda \overline{a_i}^i. s_2)}$$

Figure 3.5: Typing translation rules of System AP.



**FORALLR** uses the coercion  $f$  and, in order to produce a polymorphic type, we add one type abstraction to turn it into a coercion of type  $\sigma_1 \rightarrow \forall a. \sigma_2$ .

The second part of the subtyping translation deals with coercions generated by application subtyping. The judgment  $\Sigma \vdash^{AP} \sigma <: \sigma_2 \rightsquigarrow f$  is read as: if  $\Sigma \vdash^{AP} \sigma <: \sigma_2$  holds, it can be translated to a coercion function  $f$  in System F. If we compare two parts, we can see application contexts play no role in the generation of the coercion. Notice the similarity between rule **AP-S-TVAR** and rule **AP-AS-EMPTY**, between rule **AP-S-FORALLR** and rule **AP-AS-FORALL**, and between rule **AP-S-ARROW** and rule **AP-AS-ARROW**. We therefore omit more explanations.

### 3.3.3 TYPE-DIRECTED TRANSLATION OF TYPING

The type directed translation of typing is shown in the Figure 3.5, which extends the rules in Figure 3.1. The judgment  $\Psi; \Sigma \vdash^{AP} e \Rightarrow \sigma \rightsquigarrow s$  is read as: if  $\Psi; \Sigma \vdash^{AP} e \Rightarrow \sigma$  holds, the expression can be translated to term  $s$  in System F. The judgment  $\Psi \vdash^{AP} e \Rightarrow \sigma \rightsquigarrow s$  is the special case when the application context is empty.

Most translation rules are straightforward. In rule **AP-APP-VAR**,  $x$  is translated to  $f x$ , where  $f$  is the coercion function generated from subtyping. Rule **AP-APP-LAMANN** applies the coercion function  $f$  to  $y$ , then feeds  $y$  to the function generated from the abstraction. When generalizing over a type, rule **AP-APP-APP** generate type-level abstractions.

### 3.3.4 TYPE SAFETY

We prove that our system is type safe by proving that the translation produces well-typed terms.

**Lemma 3.10** (Soundness of Typing Translation). *If  $\Psi; \Sigma \vdash^{AP} e \Rightarrow \sigma \rightsquigarrow s$ , then  $\Psi \vdash^F s : \sigma$ .*

The lemma relies on the translation of subtyping to produce type-correct coercions.

**Lemma 3.11** (Soundness of Subtyping Translation).

1. *If  $\vdash^{AP} \sigma <: \sigma_2 \rightsquigarrow f$ , then  $\bullet \vdash^F f : \sigma \rightarrow \sigma_2$ .*
2. *If  $\Sigma \vdash^{AP} \sigma <: \sigma_2 \rightsquigarrow f$ , then  $\bullet \vdash^F f : \sigma \rightarrow \sigma_2$ .*

### 3.3.5 COHERENCE

One problem with the translation is that there are multiple targets corresponding to one expression. This is because in original system there are multiple choices when instantiating a

$ x $	$=$	$ x $	$ \Lambda a. s $	$=$	$ s $
$ n $	$=$	$ n $	$ s_1 s_2 $	$=$	$ s_1   s_2 $
$ \lambda x : \sigma. s $	$=$	$\lambda x.  s $	$ s \sigma $	$=$	$ s $

  

$f_1 =_{\eta id} f_2$

(Eta-id Equality)

  

$\frac{\text{ETA-REDUCE} \quad x \notin \text{FV}(f)}{\lambda x. f x =_{\eta id} f}$	$\frac{\text{ETA-ID}}{(\lambda x. x) f =_{\eta id} f}$	$\frac{\text{ETA-APP} \quad f_1 =_{\eta id} f'_1 \quad f_2 =_{\eta id} f'_2}{f_1 f_2 =_{\eta id} f'_1 f'_2}$
$\frac{\text{ETA-LAM} \quad f =_{\eta id} f'}{\lambda x. f =_{\eta id} \lambda x. f'}$	$\frac{\text{ETA-REFL}}{f =_{\eta id} f}$	$\frac{\text{ETA-SYMM} \quad f =_{\eta id} f'}{f' =_{\eta id} f}$
$\frac{\text{ETA-TRANS} \quad f_1 =_{\eta id} f_2 \quad f_2 =_{\eta id} f_3}{f_1 =_{\eta id} f_3}$		

Figure 3.6: Type erasure and eta-id equality of System F.

polymorphic type, or guessing the annotation for unannotated lambda abstraction (rule [AP-INF-LAM](#)). For each choice, the corresponding target will be different. For example, expression  $\lambda x. x$  can be type checked with  $\text{Int} \rightarrow \text{Int}$ , or  $a \rightarrow a$ , and the corresponding targets are  $\lambda x : \text{Int}. x$ , and  $\lambda x : a. x$ .

Therefore, in order to prove the translation is coherent, we turn to prove that all the translations have the same operational semantics. There are two steps towards the goal: type erasure, and considering  $\eta$  expansion and identity functions.

**TYPE ERASURE.** Since type information is useless after type-checking, we erase the type information of the targets for comparison. The erasure process is given at the top of Figure 3.6.

The erasure process is standard, where we erase the type annotation in abstractions, and remove type abstractions and type applications. The calculus after erasure is the untyped lambda calculus.

**ETA-ID EQUALITY.** However, even if we have type erasure, multiple targets for one expression can still be syntactically different. The problem is that we can insert more coercion functions in one translation than another, since an expression can have a more polymorphic type in one derivation than another one. So we need a more refined definition of equality instead of syntactic equality.

We use a similar definition of eta-id equality as in Chen [2003], given in Figure 3.6. In  $=_{\eta id}$  equality, two expressions are regarded as equivalent if they can turn into the same expression

through  $\eta$ -reduction or removal of redundant identity functions. The  $=_{\eta id}$  relation is reflexive, symmetrical, and transitive. As a small example, we can show that  $\lambda x. (\lambda y. y) f x =_{\eta id} f$ .

$$\frac{\frac{\frac{}{f =_{\eta id} f} \text{ETA-REFL}}{(\lambda y. y) f =_{\eta id} f} \text{ETA-ID}}{\lambda x. (\lambda y. y) f x =_{\eta id} f} \text{ETA-REDUCE}$$

Now we first prove that the erasure of the translation result of subtyping is always  $=_{\eta id}$  to an identity function.

**Lemma 3.12** (Subtyping Coercions eta-id equal to id).

- if  $\vdash^{AP} \sigma_1 <: \sigma_2 \rightsquigarrow f$ , then  $|f| =_{\eta id} \lambda x. x$ .
- if  $\Sigma \vdash^{AP} \sigma_1 <: \sigma_2 \rightsquigarrow f$ , then  $|f| =_{\eta id} \lambda x. x$ .

We then prove that our translation actually generates a *unique* target:

**Lemma 3.13** (Coherence). If  $\Psi_1; \Sigma_1 \vdash^{AP} e \Rightarrow \sigma \rightsquigarrow s_1$ , and  $\Psi_2; \Sigma_2 \vdash^{AP} e \Rightarrow \sigma_2 \rightsquigarrow s_2$ , then  $|s_1| =_{\eta id} |s_2|$ .

### 3.4 TYPE INFERENCE ALGORITHM

Even though our specification is syntax-directed, it does not directly lead to an algorithm, because there are still many guesses in the system, such as in rule **AP-INF-LAM**. This subsection presents a brief introduction of the algorithm, which closely follows the approach by Peyton Jones et al. [2007].

Instead of guessing, the algorithm creates *meta* type variables  $\hat{\alpha}, \hat{\beta}$  which are waiting to be solved. The judgment for the algorithmic type system is

$$(S_0, N_0); \Psi; \Sigma \vdash^{AP} e \Rightarrow \sigma \hookrightarrow (S_1, N_1)$$

Here we use  $N$  as name supply, from which we can always extract new names. Also, every time a meta type variable is solved, we need to record its solution. We use  $S$  as a notation for the substitution that maps meta type variables to their solutions. For example, rule **AP-INF-LAM** becomes

$$\frac{(S_0, N_0); \Psi, x : \hat{\beta} \vdash^{AP} \lambda x. e \Rightarrow \sigma \hookrightarrow (S_1, N_1)}{(S_0, N_0 \hat{\beta}); \Psi \vdash^{AP} \lambda x. e \Rightarrow \hat{\beta} \rightarrow \sigma \hookrightarrow (S_1, N_1)} \text{AP-INF-LAM-ALGO}$$

Comparing it to rule [AP-INF-LAM](#),  $\tau$  is replaced by a new meta type variable  $\hat{\beta}$  from name supply  $N_0\hat{\beta}$ . But despite of the name supply and substitution, the rule retains the structure of rule [AP-INF-LAM](#).

Having the name supply and substitutions, the algorithmic system is a direct extension of the specification in Figure 3.2, with a process to do unifications that solve meta type variables. Such unification process is quite standard and similar to the one used in the Hindley-Milner system. We proved our algorithm is sound and complete with respect to the specification. Here *fmv* means free meta type variables.

**Theorem 3.14** (Soundness). *If  $(\square, N_0); \Psi \vdash^{AP} e \Rightarrow \sigma \hookrightarrow (S_1, N_1)$ , then for any substitution  $V$  with  $\text{dom}(V) = \text{fmv}(S_1\Psi, S_1\sigma)$ , we have  $VS_1\Psi \vdash^{AP} e \Rightarrow VS_1\sigma$ .*

**Theorem 3.15** (Completeness). *If  $\Psi \vdash^{AP} e \Rightarrow \sigma$ , then for a fresh  $N_0$ , we have  $(\square, N_0); \Psi \vdash^{AP} e \Rightarrow \sigma_2 \hookrightarrow (S_1, N_1)$ , and for some  $S_2$ , if  $\bar{a}_i^i = \text{FV}(\Psi) - \text{FV}(S_2S_1\sigma_2)$ , and  $\bar{b}_i^i = \text{FV}(\Psi) - \text{FV}(\sigma)$ , we have  $\vdash^{AP} \forall \bar{a}_i^i. S_2S_1\sigma_2 <: \forall \bar{b}_i^i. \sigma$ .*

## 3.5 DISCUSSION

### 3.5.1 COMBINING APPLICATION AND CHECKING MODES

Although the application mode provides us with alternative design choices in a bi-directional type system, a checking mode can still be *easily* added. One motivation for the checking mode would be annotated expressions  $e : \sigma$ , where the type of expressions is known and is therefore used to check expressions, as in DK.

Consider adding  $e : \sigma$  for introducing the checking mode for the language. Notice that, since the checking mode is stronger than application mode, when entering checking mode the application context is no longer useful. Instead we use application subtyping to satisfy the application context requirements. A possible typing rule for annotation expressions is:

$$\frac{\text{AP-APP-ANNO} \quad \Psi \vdash^{AP} e \Leftarrow \sigma_1 \quad \Sigma \vdash^{AP} \sigma_1 <: \sigma_2}{\Psi; \Sigma \vdash^{AP} e : \sigma_1 \Rightarrow \sigma_2}$$

Here,  $e$  is checking using its annotation  $\sigma_1$ , and then we instantiate  $\sigma_1$  to  $\sigma_2$  using subtyping with application context  $\Sigma$ .

Now we can have a rule set of the checking mode for all expressions, much like those checking rules in DK. For example, one useful rule for abstractions in checking mode could be rule [AP-CHK-LAM](#), where the parameter type  $\sigma_1$  serves as the type of  $x$ , and typing checks

the body with  $\sigma_2$ . Also, combined with the information flow, the checking rule for application checks the function with the full type.

$$\frac{\text{AP-CHK-LAM} \quad \Psi, x : \sigma_1 \vdash^{AP} e \Leftarrow \sigma_2}{\Psi \vdash^{AP} \lambda x. e \Leftarrow \sigma_1 \rightarrow \sigma_2}$$

Moreover, combined with the information flow, the checked rule for application checks the function with the full type.

$$\frac{\text{AP-CHK-APP} \quad \Psi \vdash^{AP} e_2 \Rightarrow \sigma_1 \quad \Psi \vdash^{AP} e_1 \Leftarrow \sigma_1 \rightarrow \sigma_2}{\Psi \vdash^{AP} e_1 e_2 \Leftarrow \sigma_2}$$

Note that adding expression annotations might bring convenience for programmers, since annotations can be more freely placed in a program. For example,  $(\lambda x. x \ 1) : (\text{Int} \rightarrow \text{Int}) \rightarrow \text{Int}$  becomes valid. However this does not add expressive power, since programs that are typeable under expression annotations, would remain typeable after moving the annotations to bindings. For example the previous program is equivalent to  $(\lambda x : \text{Int} \rightarrow \text{Int}. x \ 1)$ .

This discussion is a sketch. We have not defined the corresponding declarative system nor algorithm. However we believe that the addition of a checked mode will *not* bring surprises to the meta-theory.

### 3.5.2 ADDITIONAL CONSTRUCTS

In this section, we show that the application mode is compatible with other constructs, by discussing how to add support for pairs in the language. A similar methodology would apply to other constructs like sum types, data types, if-then-else expressions and so on.

The introduction rule for pairs must be in the inference mode with an empty application context. Also, the subtyping rule for pairs is as expected.

$$\frac{\text{AP-INF-PAIR} \quad \Psi \vdash^{AP} e_1 \Rightarrow \sigma_1 \quad \Psi \vdash^{AP} e_2 \Rightarrow \sigma_2}{\Psi \vdash^{AP} (e_1, e_2) \Rightarrow (\sigma_1, \sigma_2)} \quad \frac{\text{AP-S-PAIR} \quad \vdash^{AP} \sigma_1 <: \sigma_3 \quad \vdash^{AP} \sigma_2 <: \sigma_4}{\vdash^{AP} (\sigma_1, \sigma_2) <: (\sigma_3, \sigma_4)}$$

The application mode can apply to the elimination constructs of pairs. If one component of the pair is a function, for example,  $\text{fst } (\lambda x. x, 1) \ 2$ , then it is possible to have a judgment

### 3 Higher-Rank Polymorphism with Application Mode

with a non-empty application context. Therefore, we can use the application subtyping to account for the application contexts:

$$\begin{array}{c}
 \text{AP-APP-FST} \\
 \frac{\Psi \vdash^{AP} e \Rightarrow (\sigma_1, \sigma_2) \quad \Sigma \vdash^{AP} \sigma_1 <: \sigma_3}{\Psi; \Sigma \vdash^{AP} \mathbf{fst} e \Rightarrow \sigma_3} \\
 \\
 \text{AP-APP-SND} \\
 \frac{\Psi \vdash^{AP} e \Rightarrow (\sigma_1, \sigma_2) \quad \Sigma \vdash^{AP} \sigma_2 <: \sigma_3}{\Psi; \Sigma \vdash^{AP} \mathbf{snd} e \Rightarrow \sigma_3}
 \end{array}$$

However, in polymorphic type systems, we need to take the subsumption rule into consideration. For example, in the expression  $(\lambda x : \forall a. (a, b). \mathbf{fst} x)$ ,  $\mathbf{fst}$  is applied to a polymorphic type. Interestingly, instead of a non-deterministic subsumption rule, having polymorphic types actually leads to a simpler solution. According to the philosophy of the application mode, the types of the arguments always flow into the functions. Therefore, instead of regarding  $\mathbf{fst} e$  as an expression form, where  $e$  is itself an argument, we could regard  $\mathbf{fst}$  as a function on its own, whose type is  $\forall a. \forall b. (a, b) \rightarrow a$ . Then as in the variable case, we use the subtyping rule to deal with application contexts. Thus the typing rules for  $\mathbf{fst}$  and  $\mathbf{snd}$  can be modeled as:

$$\begin{array}{c}
 \text{AP-APP-FST-VAR} \\
 \frac{\Sigma \vdash^{AP} \forall a. \forall b. (a, b) \rightarrow a <: \sigma}{\Psi; \Sigma \vdash^{AP} \mathbf{fst} \Rightarrow \sigma} \\
 \\
 \text{AP-APP-SND-VAR} \\
 \frac{\Sigma \vdash^{AP} \forall a. \forall b. (a, b) \rightarrow b <: \sigma}{\Psi; \Sigma \vdash^{AP} \mathbf{snd} \Rightarrow \sigma}
 \end{array}$$

Note that another way to model those two rules would be to simply have an initial typing environment  $\Psi_{init} \equiv \mathbf{fst} : \forall a. \forall b. (a, b) \rightarrow a, \mathbf{snd} : \forall a. \forall b. (a, b) \rightarrow b$ . In this case the elimination of pairs be dealt directly by the rule for variables.

An extended version of the calculus extended with rules for pairs (rule [AP-INF-PAIR](#), rule [AP-S-PAIR](#), rule [AP-APP-FST-VAR](#) and rule [AP-APP-SND-VAR](#)), has been formally studied. All the theorems presented before hold with the extension of pairs.

#### 3.5.3 MORE EXPRESSIVE TYPE APPLICATIONS

The design choice of propagating arguments to functions was subject to consideration in the original work on local type inference [Pierce and Turner 2000], but was rejected due to possible non-determinism introduced by explicit type applications:

*“It is possible, of course, to come up with examples where it would be beneficial to synthesize the argument types first and then use the resulting information to avoid*

*type annotations in the function part of an application expression....Unfortunately this refinement does not help infer the type of polymorphic functions. For example, we cannot uniquely determine the type of  $x$  in the expression  $(\text{fun}[X](x) e) [\text{Int}] 3$ ."*

As a response to this challenge, we also present an application of the application mode to a variant of System F [Xie and Oliveira 2018]. The development of the calculus shows that the application mode can actually work well with calculi with explicit type applications. Here we explain the key ideas of the design of the system, but refer to Xie and Oliveira [2018] for more details.

To explain the new design, consider the expression:

$$\Lambda a. (\lambda x : a. x + 1) \text{Int}$$

which is not typeable in the traditional type system for System F. In System F the lambda abstractions do not account for the context of possible function applications. Therefore when type checking the inner body of the lambda abstraction, the expression  $x + 1$  is ill-typed, because all that is known is that  $x$  has the (abstract) type  $a$ .

If we are allowed to propagate type information from arguments to functions, then we can verify that  $a = \text{Int}$  and  $x + 1$  is well-typed. The key insight in the new type system is to use contexts to track type equalities induced by type applications. This enables us to type check expressions such as the body of the lambda above ( $x + 1$ ). The key rules for type abstractions and type applications are:

$$\frac{\Psi; \Sigma, [[\Psi]\sigma_1] \vdash^{AP} e \Rightarrow \sigma_2}{\Psi; \Sigma \vdash^{AP} e \sigma_1 \Rightarrow \sigma_2} \text{AP-APP-TAPP} \qquad \frac{\Psi, a = \sigma_1; \Sigma \vdash^{AP} e \Rightarrow \sigma_2}{\Psi; \Sigma, [\sigma_1] \vdash^{AP} \Lambda a. e \Rightarrow \sigma_2} \text{AP-APP-TLAM}$$

For type applications, rule **AP-APP-TAPP** stores the type argument  $\sigma_1$  into the application context. Since  $\Psi$  tracks type equalities, we apply  $\Psi$  as a type substitution to  $\sigma_1$  (i.e.,  $[\Psi]\sigma_1$ ). Moreover, to distinguish between type arguments and types of term arguments, we put type arguments in brackets (i.e.,  $[[\Psi]\sigma_1]$ ). For type abstractions (rule **AP-APP-TLAM**), if the application context is non-empty, we put a new type equality between the type variable  $a$  and the type argument  $\sigma_1$  into the context.

Now, back to the problematic expression  $(\text{fun}[X](x) e) [\text{Int}] 3$ , the type of  $x$  can be inferred as either  $X$  or  $\text{Int}$  since they are actually equivalent.

**SUGAR FOR TYPE SYNONYMS.** In the same way that we can regard **let** expressions as syntactic sugar, in the new type system we further *gain built-in type synonyms for free*. A *type synonym*

is a new name for an existing type. Type synonyms are common in languages such as Haskell. In our calculus a simple form of type synonyms can be desugared as follows:

$$\mathbf{type} \ a = \sigma \ \mathbf{in} \ e \rightsquigarrow (\Lambda a. e) \sigma$$

One practical benefit of such syntactic sugar is that it enables a direct encoding of a System F-like language with declarations (including type-synonyms). Although declarations are often viewed as a routine extension to a calculus, and are not formally studied, they are highly relevant in practice. Therefore, a more realistic formalization of a programming language should directly account for declarations. By providing a way to encode declarations, our new calculus enables a simple way to formalize declarations.

**TYPE ABSTRACTION.** The type equalities introduced by type applications may seem like we are breaking System F type abstraction. However, we argue that *type abstraction* is still supported by our System F variant. For example:

$$\mathbf{let} \ inc = \Lambda a. \lambda x : a. x + 1 \ \mathbf{in} \ inc \ \mathbf{Int} \ 1$$

(after desugaring) does *not* type-check, as in a System-F like language. In our type system lambda abstractions that are immediately applied to an argument, and unapplied lambda abstractions behave differently. Unapplied lambda abstractions are just like System F abstractions and retain type abstraction. The example above illustrates this. In contrast the typeable example  $(\Lambda a. \lambda x : a. x + 1) \ \mathbf{Int}$ , which uses a lambda abstraction directly applied to an argument, can be regarded as the desugared expression for  $\mathbf{type} \ a = \mathbf{Int} \ \mathbf{in} \ \lambda x : a. x + 1$ .

#### 3.5.4 DEPENDENT TYPE SYSTEMS

**Ningning:** Move to future work.

One remark about the application mode is that the same idea is possibly applicable to systems with advanced features, where type inference is sophisticated or even undecidable. One promising application is, for instance, dependent type systems [Xi and Pfenning 1999]. Type systems with dependent types usually unify the syntax for terms and types, with a single lambda abstraction generalizing both type and lambda abstractions. Unfortunately, this means that the **let** desugar is not valid in those systems. As a concrete example, consider desugaring the expression  $\mathbf{let} \ a = \mathbf{Int} \ \mathbf{in} \ \lambda x : a. x + 1$  into  $(\Lambda a. \lambda x : a. x + 1) \ \mathbf{Int}$ , which is ill-typed because the type of  $x$  in the abstraction body is  $a$  and not  $\mathbf{Int}$ .



Because **let** cannot be encoded, declarations cannot be encoded either. Modeling declarations in dependently typed languages is a subtle matter, and normally requires some additional complexity [Severi and Poll 1994].

We believe that the same technique presented in Section 3.5.3 can be adapted into a dependently typed language to enable a **let** encoding. In a dependent type system with unified syntax for terms and types, we can combine the two forms in the typing context, i.e.,  $x : \sigma$  and  $a = \sigma$ , into a unified form  $x = e : \sigma$ . Then we can combine two application rules rule **AP-APP-APP** and rule **AP-APP-TAPP** into rule **AP-APP-DAPP**, and also two abstraction rules rule **AP-APP-LAM** and rule **AP-APP-TLAM** into rule **AP-APP-DLAM**.

$$\frac{\Psi \vdash^{AP} e_2 \Rightarrow \sigma_1 \quad \Psi; \Sigma, e_2 : \sigma_1 \vdash^{AP} e_1 \Rightarrow \sigma_2}{\Psi; \Sigma \vdash^{AP} e_1 e_2 \Rightarrow \sigma_2} \text{AP-APP-DAPP}$$

$$\frac{\Psi, x = e_1 : \sigma_1; \Sigma \vdash^{AP} e \Rightarrow \sigma_2}{\Psi; \Sigma, e_1 : \sigma_1 \vdash^{AP} \lambda x. e \Rightarrow \sigma_2} \text{AP-APP-DLAM}$$

With such rules it would be possible to handle declarations easily in dependent type systems. Note this is still a rough idea and we have not fully worked out the typing rules for this type system yet.



## PART III

# HIGHER-RANK POLYMORPHISM AND GRADUAL TYPING



# 4 GRADUALLY TYPED HIGHER-RANK POLYMORPHISM

*Consistent subtyping* is employed in some gradual type systems to validate type conversions. The original definition by Siek and Taha [2007b] 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. This section first proposes a generalization of consistent subtyping (Section 4.2) 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.

We then present GPC, a gradually typed calculus for implicit higher-rank polymorphism that uses our new notion of consistent subtyping. We develop both declarative (Section 4.3) and bi-directional algorithmic versions (Section 4.4) for the type system. The algorithmic version employs techniques developed by DK [Dunfield and Krishnaswami 2013] for higher-rank polymorphism to deal with instantiation.

## 4.1 INTRODUCTION AND MOTIVATION

### 4.1.1 BACKGROUND: GRADUAL TYPING

Siek and Taha [2007b] developed a gradual type system for object-oriented languages that they call  $\text{FOb}_{<}^?$ . Central to gradual typing is the concept of *consistency* (written  $\sim$ ) between gradual types, which are types that may involve the unknown type  $?$ . 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  $?$  is neutral to subtyping, with only  $? <: ?$ . Both relations are defined in Figure 4.1.

A primary contribution of their work is to show that consistency and subtyping are orthogonal. However, the orthogonality of consistency and subtyping does not lead to a de-

$\boxed{\sigma_1 <: \sigma_2}$				(Subtyping)
$\overline{\text{Int} <: \text{Int}}$	$\overline{\text{Bool} <: \text{Bool}}$	$\overline{\text{Float} <: \text{Float}}$	$\overline{\text{Int} <: \text{Float}}$	
$\frac{\sigma_3 <: \sigma_1 \quad \sigma_2 <: \sigma_4}{\sigma_1 \rightarrow \sigma_2 <: \sigma_3 \rightarrow \sigma_4}$		$\frac{\overline{\sigma_i <: \sigma'_i}^i}{[\overline{l_i : \sigma_i}^i] <: [\overline{l_i : \sigma'_i}^i]}$	$\overline{? <: ?}$	
$\boxed{\sigma_1 \sim \sigma_2}$				(Type Consistency)
$\overline{\sigma \sim \sigma}$	$\overline{\sigma \sim ?}$	$\overline{? \sim \sigma}$	$\frac{\sigma_1 \sim \sigma_3 \quad \sigma_2 \sim \sigma_4}{\sigma_1 \rightarrow \sigma_2 \sim \sigma_3 \rightarrow \sigma_4}$	$\frac{\overline{\sigma_i \sim \sigma'_i}^i}{[\overline{l_i : \sigma_i}^i] \sim [\overline{l_i : \sigma'_i}^i]}$

 Figure 4.1: Subtyping and type consistency in  $\text{FOb}_{<}^?$ .

terministic relation. Thus Siek and Taha defined *consistent subtyping* (written  $\lesssim$ ) based on a *restriction operator*, written  $\sigma_1|_{\sigma_2}$  that “masks off” the parts of type  $\sigma_1$  that are unknown in type  $\sigma_2$ . For example,

$$\begin{aligned} \text{Int} \rightarrow \text{Int}|_{\text{Bool} \rightarrow \text{Bool}} &= \text{Int} \rightarrow ? \\ \text{Bool} \rightarrow ?|_{\text{Int} \rightarrow \text{Int}} &= \text{Bool} \rightarrow ? \end{aligned}$$

The definition of the restriction operator is given below:

$$\begin{aligned} \sigma|_{\sigma'} &= \text{case } (\sigma, \sigma') \text{ of} \\ &| (\_, ?) \Rightarrow ? \\ &| (\sigma_1 \rightarrow \sigma_2, \sigma'_1 \rightarrow \sigma'_2) \Rightarrow \sigma_1|_{\sigma'_1} \rightarrow \sigma_2|_{\sigma'_2} \\ &| ([l_1 : \sigma_1, \dots, l_n : \sigma_n], [l_1 : \sigma'_1, \dots, l_m : \sigma'_m]) \text{ if } n \leq m \Rightarrow [l_1 : \sigma_1|_{\sigma'_1}, \dots, l_n : \sigma_n|_{\sigma'_n}] \\ &| ([l_1 : \sigma_1, \dots, l_n : \sigma_n], [l_1 : \sigma'_1, \dots, l_m : \sigma'_m]) \text{ if } n > m \Rightarrow [l_1 : \sigma_1|_{\sigma'_1}, \dots, l_m : \sigma_m|_{\sigma'_m}, \dots, l_n : \sigma_n] \\ &| (\_, \_) \Rightarrow \sigma \end{aligned}$$

With the restriction operator, consistent subtyping is simply defined as:

**Definition 3** (Algorithmic Consistent Subtyping of Siek and Taha [2007b]).  $\sigma_1 \lesssim \sigma_2 \equiv \sigma_1|_{\sigma_2} <: \sigma_2|_{\sigma_1}$ .

Later they show a proposition that consistent subtyping is equivalent to two declarative definitions, which we refer to as the strawman for *declarative* consistent subtyping because it

servers as a good guideline on superimposing consistency and subtyping. Both definitions are non-deterministic because of the intermediate type  $\sigma_3$ .

**Definition 4** (Strawman for Declarative Consistent Subtyping). The following two are equivalent:

1.  $\sigma_1 \lesssim \sigma_2$  if and only if  $\sigma_1 \sim \sigma_3$  and  $\sigma_3 <: \sigma_2$  for some  $\sigma_3$ .
2.  $\sigma_1 \lesssim \sigma_2$  if and only if  $\sigma_1 <: \sigma_3$  and  $\sigma_3 \sim \sigma_2$  for some  $\sigma_3$ .

In our later discussion, it will always be clear which definition we are referring to. For example, we focus more on Definition 4 in Section 4.2.2, and more on Definition 3 in Section 4.2.5.

#### 4.1.2 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 [Garcia et al. 2016].

Haskell is a language renowned for its advanced type system, but it does not feature gradual typing. Of particular interest to us is its support for implicit higher-rank polymorphism, which is supported via explicit type annotations. In Haskell some programs that are safe at run-time may be rejected due to the conservativity of the type system. For example, consider again the example from Section 2.2:

```
(\f. (f 1, f 'a')) (\x. x)
```

This program is rejected by Haskell’s type checker because Haskell implements the HM rule that a lambda-bound argument (such as  $f$ ) can only have a monotype, i.e., the type checker can only assign  $f$  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, especially when the program scales up and the annotation is long and complicated.

Instead of rejecting the program outright, due to missing type annotations, gradual typing provides a simple alternative by giving  $f$  the unknown type  $?$ . With this type the same program type-checks and produces  $(1, 'a')$ . By running the program, programmers can gain more insight about its run-time behaviour. Then, with this insight, they can also give  $f$  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.

### 4.1.3 APPLICATION: EFFICIENT (PARTLY) TYPED ENCODINGS OF ADTs

We illustrate two concrete applications of gradually typed higher-rank polymorphism related to algebraic datatypes. The first application shows how gradual typing helps in defining Scott encodings of algebraic datatypes [Curry et al. 1958; Parigot 1992], which are impossible to encode in plain System F. The second application shows how gradual typing makes it easy to model and use heterogeneous containers.

Our calculus does not provide built-in support for algebraic datatypes (ADTs). Nevertheless, the calculus is expressive enough to support efficient function-based encodings of (optionally polymorphic) ADTs<sup>1</sup>. This offers an immediate way to model algebraic datatypes in our calculus without requiring extensions to our calculus or, more importantly, to its target—the polymorphic blame calculus. While we believe that such extensions are possible, they would likely require non-trivial extensions to the polymorphic blame calculus. Thus the alternative of being able to model algebraic datatypes without extending  $\lambda B$  is appealing. The encoding also paves the way to provide built-in support for algebraic datatypes in the source language, while elaborating them via the encoding into  $\lambda B$ .

**CHURCH AND SCOTT ENCODINGS.** It is well-known that polymorphic calculi such as System F can encode datatypes via Church encodings. However these encodings have well-known drawbacks. In particular, some operations are hard to define, and they can have a time complexity that is greater than that of the corresponding functions for built-in algebraic datatypes. A good example is the definition of the predecessor function for Church numerals [Church 1941]. Its definition requires significant ingenuity (while it is trivial with built-in algebraic datatypes), and it has *linear* time complexity (versus the *constant* time complexity for a definition using built-in algebraic datatypes).

An alternative to Church encodings are the so-called Scott encodings [Curry et al. 1958]. They address the two drawbacks of Church encodings: they allow simple definitions that directly correspond to programs implemented with built-in algebraic datatypes, and those definitions have the same time complexity to programs using algebraic datatypes.

Unfortunately, Scott encodings, or more precisely, their typed variant [Parigot 1992], cannot be expressed in System F: in the general case they require recursive types, which System

---

<sup>1</sup>In a type system with impure features, such as non-termination or exceptions, the encoded types can have valid inhabitants with side-effects, which means we only get the *lazy* version of those datatypes.



F does not support. However, with gradual typing, we can remove the need for recursive types, thus enabling Scott encodings in our calculus.

A SCOTT ENCODING OF PARAMETRIC LISTS. Consider for instance the typed Scott encoding of parametric lists in a system with implicit polymorphism:

$$\begin{aligned}\text{List } a &\triangleq \mu L. \forall b. b \rightarrow (a \rightarrow L \rightarrow b) \rightarrow b \\ \text{nil} &\triangleq \mathbf{fold}_{\text{List } a} (\lambda m. \lambda c. m) : \forall a. \text{List } a \\ \text{cons} &\triangleq \lambda x. \lambda xs. \mathbf{fold}_{\text{List } a} (\lambda m. \lambda c. c \ x \ xs) : \forall a. a \rightarrow \text{List } a \rightarrow \text{List } a\end{aligned}$$

This encoding requires both polymorphic and recursive types<sup>2</sup>. Like System F, our calculus only supports the former, but not the latter. Nevertheless, gradual types still allow us to use the Scott encoding in a partially typed fashion. The trick is to omit the recursive type binder  $\mu L$  and replace the recursive occurrence of  $L$  by the unknown type  $?$ :

$$\text{List}_? a \triangleq \forall b. b \rightarrow (a \rightarrow ? \rightarrow b) \rightarrow b$$

As a consequence, we need to replace the term-level witnesses of the iso-recursion by explicit type annotations to respectively forget or recover the type structure of the recursive occurrences:

$$\begin{aligned}\mathbf{fold}_{\text{List}_? a} &\triangleq \lambda x. x : (\forall b. b \rightarrow (a \rightarrow \text{List}_? a \rightarrow b) \rightarrow b) \rightarrow \text{List}_? a \\ \mathbf{unfold}_{\text{List}_? a} &\triangleq \lambda x. x : \text{List}_? a \rightarrow (\forall b. b \rightarrow (a \rightarrow \text{List}_? a \rightarrow b) \rightarrow b)\end{aligned}$$

With the reinterpretation of **fold** and **unfold** as functions instead of built-in primitives, we have exactly the same definitions of  $\text{nil}_?$  and  $\text{cons}_?$ .

Note that when we elaborate our calculus into the polymorphic blame calculus, the above type annotations give rise to explicit casts. For instance, after elaboration  $\mathbf{fold}_{\text{List}_? a} e$  results in the cast  $\langle (\forall b. b \rightarrow (a \rightarrow \text{List}_? a \rightarrow b) \rightarrow b) \hookrightarrow \text{List}_? a \rangle s$  where  $s$  is the elaboration of  $e$ .

In order to perform recursive traversals on lists, e.g., to compute their length, we need a fixpoint combinator like the Y combinator. Unfortunately, this combinator cannot be assigned a type in the simply typed lambda calculus or System F. Yet, we can still provide a gradual typing for it in our system.

$$\text{fix} \triangleq \lambda f. (\lambda x : ?. f(x \ x)) (\lambda x : ?. f(x \ x)) : \forall a. (a \rightarrow a) \rightarrow a$$

<sup>2</sup>Here we use iso-recursive types, but equi-recursive types can be used too.

This allows us for instance to compute the length of a list.

$$\text{length} \triangleq \text{fix } (\lambda \text{len}. \lambda l. \text{zero}_? (\lambda xs. \text{succ}_? (\text{len } xs)))$$

Here  $\text{zero}_? : \text{Nat}_?$  and  $\text{succ}_? : \text{Nat}_? \rightarrow \text{Nat}_?$  are the encodings of the constructors for natural numbers  $\text{Nat}_?$ . In practice, for performance reasons, we could extend our language with a **letrec** construct in a standard way to support general recursion, instead of defining a fixpoint combinator.

Observe that the gradual typing of lists still enforces that all elements in the list are of the same type. For instance, a heterogeneous list like  $\text{cons}_? \text{zero}_? (\text{cons}_? \text{true}_? \text{nil}_?)$ , is rejected because  $\text{zero}_? : \text{Nat}_?$  and  $\text{true}_? : \text{Bool}_?$  have different types.

**HETEROGENEOUS CONTAINERS.** Heterogeneous containers are datatypes that can store data of different types, which is very useful in various scenarios. One typical application is that an XML element is heterogeneously typed. Moreover, the result of a SQL query contains heterogeneous rows.

In statically typed languages, there are several ways to obtain heterogeneous lists. For example, in Haskell, one option is to use *dynamic types*. Haskell's library **Data.Dynamic** provides the special type **Dynamic** along with its injection **toDyn** and projection **fromDyn**. The drawback is that the code is littered with **toDyn** and **fromDyn**, which obscures the program logic. One can also use the **HList** library [Kiselyov et al. 2004], which provides strongly typed data structures for heterogeneous collections. The library requires several Haskell extensions, such as multi-parameter classes [Peyton Jones et al. 1997] and functional dependencies [Jones 2000]. With fake dependent types [McBride 2002], heterogeneous vectors are also possible with type-level constructors.

In our type system, with explicit type annotations that set the element types to the unknown type we can disable the homogeneous typing discipline for the elements and get gradually typed heterogeneous lists<sup>3</sup>. Such gradually typed heterogeneous lists are akin to Haskell's approach with Dynamic types, but much more convenient to use since no injections and projections are needed, and the  $?$  type is built-in and natural to use.

An example of such gradually typed heterogeneous collections is:

$$l \triangleq \text{cons}_? (\text{zero}_? : ?) (\text{cons}_? (\text{true}_? : ?) \text{nil}_?)$$

Here we annotate each element with type annotation  $?$  and the type system is happy to type-check that  $l : \text{List}_? ?$ . Note that we are being meticulous about the syntax, but with proper

<sup>3</sup>This argument is based on the extended type system in Chapter 5.

implementation of the source language, we could write more succinct programs akin to Haskell's syntax, such as `[0, True]`.

## 4.2 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 compare with some relevant work. We then discuss the design decisions involved in our new definition of consistent subtyping, and justify the new definition by demonstrating its equivalence with that of Siek and Taha [2007b] and the AGT approach [Garcia et al. 2016] on simple types.

The syntax of types is given at the top of Figure 4.2. Types  $\sigma$  are either the integer type `Int`, type variables  $a$ , function types  $\sigma_1 \rightarrow \sigma_2$ , universal quantification  $\forall a. \sigma$ , or the unknown type `?`. Note that monotypes  $\tau$  contain all types other than the universal quantifier and the unknown type `?`. We will discuss this restriction when we present the subtyping rules. Contexts  $\Psi$  are *ordered* lists of type variable declarations and term variables.

### 4.2.1 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. [2009] and type consistency by Igarashi et al. [2017].

**CONSISTENCY.** The key observation here is that consistency is mostly a structural relation, except that the unknown type `?` can be regarded as any type. In other words, consistency is an equivalence relation lifted from static types to gradual types [Garcia et al. 2016]. Following this observation, we naturally extend the definition from Figure 4.1 with polymorphic types, as shown in the middle of Figure 4.2. In particular a polymorphic type  $\forall a. \sigma$  is consistent with another polymorphic type  $\forall a. \sigma_2$  if  $\sigma$  is consistent with  $\sigma_2$ .

**SUBTYPING.** We express the fact that one type is a polymorphic generalization of another by means of the subtyping judgment  $\Psi \vdash^G \sigma <: \sigma_2$ . Compared with the subtyping rules of Odersky and Läufer [1996] in Figure 2.5, the only addition is the neutral subtyping of `?`. Notice that, in rule **GPC-S-FORALLL**, the universal quantifier is only allowed to be instantiated with a *monotype*. The judgment  $\Psi \vdash^G \sigma$  checks whether all the type variables in  $\sigma$  are bound in the context  $\Psi$ . According to the syntax in Figure 4.2, monotypes do not include the unknown type `?`. This is because if we were to allow the unknown type to be used for instantiation, we could have  $\forall a. a \rightarrow a <: ? \rightarrow ?$  by instantiating  $a$  with `?`. Since  $? \rightarrow ?$  is

Types	$\sigma ::=$	$\text{Int} \mid a \mid \sigma_1 \rightarrow \sigma_2 \mid \forall a. \sigma \mid ?$
Monotypes	$\tau ::=$	$\text{Int} \mid a \mid \tau_1 \rightarrow \tau_2$
Contexts	$\Psi ::=$	$\bullet \mid \Psi, x : \sigma \mid \Psi, a$

  

$\sigma \sim \sigma_2$

(Type Consistency)

  

$\frac{}{\sigma \sim \sigma}$

$\frac{}{\sigma \sim ?}$

$\frac{}{? \sim \sigma}$

$\frac{\sigma_1 \sim \sigma_3 \quad \sigma_2 \sim \sigma_4}{\sigma_1 \rightarrow \sigma_2 \sim \sigma_3 \rightarrow \sigma_4}$

$\frac{\sigma \sim \sigma_2}{\forall a. \sigma \sim \forall a. \sigma_2}$

  

$\Psi \vdash^G \sigma <: \sigma_2$

(Subtyping)

  

$\frac{a \in \Psi}{\Psi \vdash^G a <: a} \text{ GPC-S-TVAR}$

$\frac{}{\Psi \vdash^G \text{Int} <: \text{Int}} \text{ GPC-S-INT}$

  

$\frac{\Psi \vdash^G \sigma_3 <: \sigma_1 \quad \Psi \vdash^G \sigma_2 <: \sigma_4}{\Psi \vdash^G \sigma_1 \rightarrow \sigma_2 <: \sigma_3 \rightarrow \sigma_4} \text{ GPC-S-ARROW}$

  

$\frac{\Psi \vdash^G \tau \quad \Psi \vdash^G \sigma[a \mapsto \tau] <: \sigma_2}{\Psi \vdash^G \forall a. \sigma <: \sigma_2} \text{ GPC-S-FORALLL}$

$\frac{\Psi, a \vdash^G \sigma <: \sigma_2}{\Psi \vdash^G \sigma <: \forall a. \sigma_2} \text{ GPC-S-FORALLR}$

  

$\frac{}{\Psi \vdash^G ? <: ?} \text{ GPC-S-UNKNOWN}$

  

$\Psi \vdash^G \sigma$

(Well-formedness of types)

  

$\frac{}{\Psi \vdash^G \text{Int}}$

$\frac{}{\Psi \vdash^G ?}$

$\frac{a \in \Psi}{\Psi \vdash^G a}$

$\frac{\Psi \vdash^G \sigma \quad \Psi \vdash^G \sigma_2}{\Psi \vdash^G \sigma \rightarrow \sigma_2}$

$\frac{\Psi, a \vdash^G \sigma}{\Psi \vdash^G \forall a. \sigma}$

Figure 4.2: Syntax of types, consistency, subtyping and well-formedness of types in declarative GPC.

consistent with any functions  $\sigma_1 \rightarrow \sigma_2$ , 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 know,  $\forall a. a \rightarrow a$  is definitely not compatible with  $\text{Int} \rightarrow \text{Bool}$ . Indeed, 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 *conservative extension* property that will be made precise in Section 4.3.3.)

Importantly there is a subtle distinction between a type variable and the unknown type, although they both 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  $? \rightarrow ?$  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  $\lambda\text{B}$  [Ahmed et al. 2009], the compatibility relation for polymorphic types is defined as follows:

$$\frac{\sigma_1 \prec \sigma_2}{\sigma_1 \prec \forall a. \sigma_2} \text{ COMP-ALLR} \qquad \frac{\sigma_1[a \mapsto ?] \prec \sigma_2}{\forall a. \sigma_1 \prec \sigma_2} \text{ COMP-ALLL}$$

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

Igarashi et al. [2017] 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>4</sup>:

$$\frac{\sigma \sim \sigma_2}{\forall a. \sigma \sim \forall a. \sigma_2} \qquad \frac{\sigma \sim \sigma_2 \quad \sigma_2 \neq \forall a. \sigma'_2 \quad ? \in \text{Types}(\sigma_2)}{\forall a. \sigma \sim \sigma_2}$$

<sup>4</sup>This is a simplified version. These two rules are presented in Section 3.1 in their paper as one of the key ideas of the design of type consistency, which are later amended with *labels*.

Compared with our consistency definition in Figure 4.2, 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  $?$ . In this way, their type system is able to reject  $\forall a. a \rightarrow a \sim \text{Int} \rightarrow \text{Bool}$ . 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 of 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  $? \rightarrow \text{Int}$  for example, which in our view, is a kind 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 are *orthogonal*, and can be naturally superimposed, echoing the opinion of Siek and Taha [2007b].

#### 4.2.2 TOWARDS CONSISTENT SUBTYPING

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

Unfortunately, the strawman version of consistent subtyping (Definition 4) 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  $(? \rightarrow \text{Int}) \rightarrow \text{Int}$ . 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 Figure 4.3a. We use  $\emptyset$  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 Figure 4.3b.

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, Figure 4.3c shows two such types by putting a  $\text{Bool}$  type in the middle, and neither definition of consistent subtyping works.

**OBSERVATIONS ON CONSISTENT SUBTYPING BASED ON INFORMATION PROPAGATION.** In order to develop a 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  $\sigma_1 <: \sigma_2$ , then an expression of type  $\sigma_1$  can be used where  $\sigma_2$  is expected; (2) subtyping is transitive: if  $\sigma_1 <: \sigma_2$ , and  $\sigma_2 <: \sigma_3$ , then  $\sigma_1 <: \sigma_3$ . Though consistent subtyping takes the unknown type into consid-



Figure 4.3: Examples that break the original definition of consistent subtyping.



Figure 4.4: Observations of consistent subtyping

eration, the subsumption rule should also apply: if  $\sigma_1 \lesssim \sigma_2$ , then an expression of type  $\sigma_1$  can also be used where  $\sigma_2$  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  $\sigma_1 <: \sigma_2$ , and  $\sigma_2 \lesssim \sigma_3$ , this means that  $\sigma_1$  can be used where  $\sigma_2$  is expected, and  $\sigma_2$  can be used where  $\sigma_3$  is expected, with possibly some loss of information. In other words, we should expect that  $\sigma_1$  can be used where  $\sigma_3$  is expected, since there is at most one-time loss of information.

**Observation 1.** If  $\sigma_3 \lesssim \sigma_2$ , and  $\sigma_2 <: \sigma_1$ , then  $\sigma_3 \lesssim \sigma_1$ .

This is reflected in Figure 4.4a. A symmetrical observation is given in Figure 4.4b:

**Observation 2.** If  $\sigma_3 \lesssim \sigma_2$ , and  $\sigma_2 <: \sigma$ , then  $\sigma_3 \lesssim \sigma$ .



$$\begin{aligned}
\sigma_1 &= (((\forall a. a \rightarrow \text{Int}) \rightarrow \text{Int}) \rightarrow \text{Bool}) \rightarrow (\forall a. a) \\
\sigma_2 &= (((\forall a. a \rightarrow \text{Int}) \rightarrow \text{Int}) \rightarrow \text{Bool}) \rightarrow (\text{Int} \rightarrow \text{Int}) \\
\sigma_3 &= (((\forall a. ? \rightarrow \text{Int}) \rightarrow \text{Int}) \rightarrow \text{Bool}) \rightarrow (\text{Int} \rightarrow ?) \\
\sigma_4 &= (((? \rightarrow \text{Int}) \rightarrow \text{Int}) \rightarrow \text{Bool}) \rightarrow (\text{Int} \rightarrow ?)
\end{aligned}$$

Figure 4.5: Example that is fixed by the new definition of consistent subtyping.

From the above observations, we see what the problem is with the original definition. In Fig. 4.4a, if  $\sigma_2$  can reach  $\sigma_3$  by  $\sigma'$ , then by subtyping transitivity,  $\sigma_1$  can reach  $\sigma_3$  by  $\sigma'$ . However, if  $\sigma_2$  can only reach  $\sigma_3$  by  $\sigma''$ , then  $\sigma$  cannot reach  $\sigma_3$  through the original definition. A similar problem is shown in Fig. 4.4b.

It turns out that these 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,  $\sigma_1 <: \sigma_2 \sim \sigma'' <: \sigma_3$  (in Figure 4.4a) and  $\sigma_3 <: \sigma' \sim \sigma_2 <: \sigma$  (in Figure 4.4b) have the same relation chain: subtyping, consistency, and subtyping.

**DEFINITION OF CONSISTENT SUBTYPING.** From the above discussion, we are ready to modify Definition 4, and adapt it to our notation:

**Definition 5** (Consistent Subtyping).  $\Psi \vdash^G \sigma_1 \lesssim \sigma_2$  if and only if  $\Psi \vdash^G \sigma_1 <: \sigma'$ ,  $\sigma' \sim \sigma''$  and  $\Psi \vdash^G \sigma'' <: \sigma_2$  for some  $\sigma'$  and  $\sigma''$ .

With Definition 5, Figure 4.5 illustrates the correct relation chain for the broken example shown in Figure 4.3c.

At first sight, Definition 5 seems worse than the original: we need to guess *two* types! It turns out that Definition 5 is a generalization of Definition 4, and they are equivalent in the system of Siek and Taha [2007b]. However, more generally, Definition 5 is compatible with polymorphic types.

**Proposition 4.1** (Generalization of Declarative Consistent Subtyping).

- *Definition 5 subsumes Definition 4.*  
In Definition 5, by choosing  $\sigma'' = \sigma_2$ , we have  $\sigma_1 <: \sigma'$  and  $\sigma' \sim \sigma_2$ ; by choosing  $\sigma' = \sigma_1$ , we have  $\sigma_1 \sim \sigma''$ , and  $\sigma'' <: \sigma_2$ .



- Definition 4 is equivalent to Definition 5 in the system of Siek and Taha.  
 If  $\sigma_1 <: \sigma'$ ,  $\sigma' \sim \sigma''$ , and  $\sigma'' <: \sigma_2$ , by Definition 4,  $\sigma_1 \sim \sigma_3$ ,  $\sigma_3 <: \sigma''$  for some  $\sigma_3$ .  
 By subtyping transitivity,  $\sigma_3 <: \sigma_2$ . So  $\sigma_1 \lesssim \sigma_2$  by  $\sigma_1 \sim \sigma_3$  and  $\sigma_3 <: \sigma_2$ .

#### 4.2.3 ABSTRACTING GRADUAL TYPING

Garcia et al. [2016] 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 is defined as follows:

**Definition 6** (Concretization).

$$\begin{aligned} \gamma(\text{Int}) &= \{\text{Int}\} \\ \gamma(\sigma_1 \rightarrow \sigma_2) &= \{\sigma'_1 \rightarrow \sigma'_2 \mid \sigma'_1 \in \gamma(\sigma_1), \sigma'_2 \in \gamma(\sigma_2)\} \\ \gamma(?) &= \{\text{All static types}\} \end{aligned}$$

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

**Definition 7** (Consistent Subtyping in AGT).  $\sigma_1 \widetilde{<} \sigma_2$  if and only if  $\sigma'_1 <: \sigma'_2$  for some *static* types  $\sigma'_1$  and  $\sigma'_2$  such that  $\sigma'_1 \in \gamma(\sigma_1)$  and  $\sigma'_2 \in \gamma(\sigma_2)$ .

Later they proved that this definition of consistent subtyping coincides with that of Definition 4. By Proposition 4.1, we can directly conclude that our definition coincides with AGT:

**Corollary 4.2** (Equivalence to AGT on Simple Types).  $\sigma_1 \lesssim \sigma_2$  if and only if  $\sigma_1 \widetilde{<} \sigma_2$ .

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. [2016], 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 [2017] in a gradual type system with union and intersection types, where the static types  $\sigma'_1, \sigma'_2$  in Definition 7 can be algorithmically computed by also accounting for top and bottom types.

## 4.2.4 DIRECTED CONSISTENCY

*Directed consistency* [Jafery and Dunfield 2017] is defined in terms of precision and subtyping:

$$\frac{\sigma'_1 \sqsubseteq \sigma_1 \quad \sigma_1 <: \sigma_2 \quad \sigma'_2 \sqsubseteq \sigma_2}{\sigma'_1 \lesssim \sigma'_2}$$

The judgment  $\sigma_1 \sqsubseteq \sigma_2$  is read “ $\sigma_1$  is less precise than  $\sigma_2$ ”.<sup>5</sup> In their setting, precision is first defined for type constructors and then lifted to gradual types, and subtyping is defined for gradual types. If we interpret this definition from the AGT point of view, finding a more precise static type has the same effect as concretization. Namely,  $\sigma'_1 \sqsubseteq \sigma_1$  implies  $\sigma_1 \in \gamma(\sigma'_1)$  and  $\sigma'_2 \sqsubseteq \sigma_2$  implies  $\sigma_2 \in \gamma(\sigma'_2)$  if  $\sigma_1$  and  $\sigma_2$  are static types. Therefore we consider this definition as AGT-style. From this perspective, this definition naturally coincides with Definition 7, and by Corollary 4.2, it coincides with Definition 5.

The value of their definition is that consistent subtyping is derived compositionally from *gradual subtyping* and *precision*. Arguably, gradual types play a role in both definitions, which is different from Definition 5 where subtyping is neutral to unknown types. Still, the definition is interesting as it takes precision into consideration, rather than consistency. Then a question arises as to *how are consistency and precision related*.

**CONSISTENCY AND PRECISION.** Precision is a partial order (anti-symmetric and transitive), while consistency is symmetric but not transitive. Recall that consistency is in fact an equivalence relation lifted from static types to gradual types [Garcia et al. 2016], which embodies the key role of gradual types in typing. Therefore defining consistency independently is straightforward, and it is theoretically viable to validate the definition of consistency directly. On the other hand, precision is usually connected with the gradual criteria [Siek et al. 2015], and finding a correct partial order that adheres to the criteria is not always an easy task. For example, Igarashi et al. [2017] argued that term precision for gradual System F 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.

Nonetheless, in our system, precision and consistency can be related by the following lemma:

**Lemma 4.3** (Consistency and Precision).

- If  $\sigma_1 \sim \sigma_2$ , then there exists (static)  $\sigma_3$ , such that  $\sigma_1 \sqsubseteq \sigma_3$ , and  $\sigma_2 \sqsubseteq \sigma_3$ .

<sup>5</sup>Jafery and Dunfield actually read  $\sigma_1 \sqsubseteq \sigma_2$  as “ $\sigma_1$  is more precise than  $\sigma_2$ ”. We, however, use the “less precise” notation (which is also adopted by Cimini and Siek [2016]) throughout the paper. The full rules can be found in Figure 4.8.

- If for some (static)  $\sigma_3$ , we have  $\sigma_1 \sqsubseteq \sigma_3$ , and  $\sigma_2 \sqsubseteq \sigma_3$ , then we have  $\sigma_1 \sigma_2$ .

#### 4.2.5 CONSISTENT SUBTYPING WITHOUT EXISTENTIALS

Definition 5 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  $\sigma'$  and  $\sigma''$ . As Definition 3, 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 Definition 3 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  $(\text{Int} \rightarrow ?) \rightarrow \text{Int}$ . 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  $\text{Int} \rightarrow ? \rightarrow \text{Bool}$  and  $\text{Int} \rightarrow \text{Int} \rightarrow ?$ . Since the known parts have the relation that  $\text{Int} \rightarrow ? \rightarrow ? <: \text{Int} \rightarrow ? \rightarrow ?$ , we conclude that  $\text{Int} \rightarrow ? \rightarrow \text{Bool} \lesssim \text{Int} \rightarrow \text{Int} \rightarrow ?$ .

$$\begin{array}{l} \text{Int} \rightarrow \boxed{?} \rightarrow \boxed{\text{Bool}} \\ \text{Int} \rightarrow \boxed{\text{Int}} \rightarrow \boxed{?} \end{array} \quad \begin{array}{l} | \text{Int} \rightarrow \text{Int} \rightarrow ? = \text{Int} \rightarrow ? \rightarrow ? \\ | \text{Int} \rightarrow ? \rightarrow \text{Bool} = \text{Int} \rightarrow ? \rightarrow ? \end{array} \Bigg) <:$$

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 relation always holds between the unknown type and an arbitrary type*. We can interpret this observation directly from Definition 5: the unknown

$\boxed{\Psi \vdash^G \sigma_1 \lesssim \sigma_2}$				(Consistent Subtyping)
$\frac{\text{GPC-CS-TVAR} \quad a \in \Psi}{\Psi \vdash^G a \lesssim a}$	$\frac{\text{GPC-CS-INT}}{\Psi \vdash^G \text{Int} \lesssim \text{Int}}$	$\frac{\text{GPC-CS-ARROW} \quad \Psi \vdash^G \sigma_3 \lesssim \sigma_1 \quad \Psi \vdash^G \sigma_2 \lesssim \sigma_4}{\Psi \vdash^G \sigma_1 \rightarrow \sigma_2 \lesssim \sigma_3 \rightarrow \sigma_4}$	$\frac{\text{GPC-CS-FORALLR} \quad \Psi, a \vdash^G \sigma_1 \lesssim \sigma_2}{\Psi \vdash^G \sigma_1 \lesssim \forall a. \sigma_2}$	
$\frac{\text{GPC-CS-FORALLL} \quad \Psi \vdash^G \tau \quad \Psi \vdash^G \sigma_1[a \mapsto \tau] \lesssim \sigma_2}{\Psi \vdash^G \forall a. \sigma_1 \lesssim \sigma_2}$	$\frac{\text{GPC-CS-UNKNOWNL}}{\Psi \vdash^G ? \lesssim \sigma}$	$\frac{\text{GPC-CS-UNKNOWNR}}{\Psi \vdash^G \sigma \lesssim ?}$		

Figure 4.6: Consistent Subtyping for implicit polymorphism.

type is neutral to subtyping ( $? <: ?$ ), the unknown type is consistent with any type ( $? \sim \sigma$ ), and subtyping is reflexive ( $\sigma <: \sigma$ ). Therefore, *the unknown type is a consistent subtype of any type* ( $? \lesssim \sigma$ ), *and vice versa* ( $\sigma \lesssim ?$ ). Note that this interpretation provides a general recipe for lifting 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 Figure 4.2, get rid of rule **GPC-S-UNKNOWN** and add two extra rules concerning  $?$ , resulting in the rules of consistent subtyping in Figure 4.6. Of particular interest are the rules **GPC-CS-UNKNOWNL** and **GPC-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 Figure 4.6. What is more, we can prove that the two definitions are equivalent.

**Theorem 4.4.**  $\Psi \vdash^G \sigma_1 \lesssim \sigma_2 \Leftrightarrow \Psi \vdash^G \sigma_1 <: \sigma', \sigma' \sim \sigma'', \Psi \vdash^G \sigma'' <: \sigma_2$  for some  $\sigma', \sigma''$ .

### 4.3 GRADUALLY TYPED IMPLICIT POLYMORPHISM

In Section 4.2 we introduced our 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, GPC, 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 at the top of Figure 4.7. The definition of expressions are the same as of OL in Figure 2.3. Meta-variable  $e$  ranges over expressions. Expressions include variables  $x$ , integers  $n$ , annotated lambda abstractions  $\lambda x : \sigma. e$ , un-annotated lambda abstractions  $\lambda x. e$ , applications  $e_1 e_2$ , and let expressions **let**  $x = e_1$  **in**  $e_2$ .

#### 4.3.1 TYPING IN DETAIL

Figure 4.7 gives the typing rules for our declarative system (the reader is advised to ignore the gray-shaded parts for now). Rule **GPC-VAR** extracts the type of the variable from the typing context. Rule **GPC-INT** always infers integer types. Rule **GPC-LAMANN** puts  $x$  with type annotation  $\sigma$  into the context, and continues type checking the body  $e$ . Rule **GPC-LAM** assigns a monotype  $\tau$  to  $x$ , and continues type checking the body  $e$ . Gradual types and polymorphic types are introduced via explicit annotations. Rule **GPC-GEN** puts a fresh type variable  $a$  into the type context and generalizes the typing result  $\sigma$  to  $\forall a. \sigma$ . Rule **GPC-LET** infers the type  $\sigma$  of  $e_1$ , then puts  $x : \sigma$  in the context to infer the type of  $e_2$ . Rule **GPC-APP** first infers the type of  $e_1$ , then the matching judgment  $\Psi \vdash^G \sigma \triangleright \sigma_1 \rightarrow \sigma_2$  extracts the domain type  $\sigma_1$  and the codomain type  $\sigma_2$  from type  $\sigma$ . The type  $\sigma_3$  of the argument  $e_2$  is then compared with  $\sigma_1$  using the consistent subtyping judgment.

**MATCHING.** The matching judgment of Siek et al. [2015] is extended to polymorphic types naturally, resulting in  $\Psi \vdash^G \sigma \triangleright \sigma_1 \rightarrow \sigma_2$ . Note that the matching rules generalize that of DK in Figure 2.7 with the unknown type. In rule **GPC-M-FORALL**, a monotype  $\tau$  is guessed to instantiate the universal quantifier  $a$ . If  $\sigma$  is a polymorphic type, the judgment works by guessing instantiations until it reaches an arrow type. Rule **GPC-M-ARR** returns the domain type  $\sigma_1$  and range type  $\sigma_2$  as expected. If the input is  $?$ , then rule **GPC-M-UNKNOWN** returns  $?$  as both the type for the domain and the range.

Note that in GPC, matching saves us from having a subsumption rule (rule **OL-SUB** in Fig. 2.5). 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 Section 4.5.2.

#### 4.3.2 TYPE-DIRECTED TRANSLATION

We give the dynamic semantics of our language by translating it to  $\lambda B$  [Ahmed et al. 2009]. Below we show a subset of the terms in  $\lambda B$  that are used in the translation:

$$\lambda B \text{ Terms } \quad s ::= x \mid n \mid \lambda x : \sigma. s \mid \Lambda a. s \mid s_1 s_2 \mid \langle \sigma \hookrightarrow \sigma_2 \rangle s$$



Figure 4.7: Syntax of expressions and declarative typing of declarative GPC

A cast  $\langle \sigma_1 \hookrightarrow \sigma_2 \rangle s$  converts the value of term  $s$  from type  $\sigma_1$  to type  $\sigma_2$ . A cast from  $\sigma_1$  to  $\sigma_2$  is permitted only if the types are *compatible*, written  $\sigma_1 \hookrightarrow \sigma_2$ , as briefly mentioned in Section 4.2.1. The syntax of types in  $\lambda B$  is the same as ours.

The translation is given in the gray-shaded parts in Figure 4.7. 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 4.5** ( $\triangleright$  to  $\hookrightarrow$ ). *If  $\Psi \vdash^G \sigma \triangleright \sigma_1 \rightarrow \sigma_2$ , then  $\sigma \hookrightarrow \sigma_1 \rightarrow \sigma_2$ .*

**Lemma 4.6** ( $\lesssim$  to  $\hookrightarrow$ ). *If  $\Psi \vdash^G \sigma \lesssim \sigma_2$ , then  $\sigma \hookrightarrow \sigma_2$ .*

In order to show the correctness of the translation, we prove that our translation always produces well-typed expressions in  $\lambda B$ . By Lemmas 4.5 and 4.6, we have the following theorem:

**Theorem 4.7** (Type Safety). *If  $\Psi \vdash^G e : \sigma \rightsquigarrow s$ , then  $\Psi \vdash^B s : \sigma$ .*

**PARAMETRICITY.** An important semantic property of polymorphic types is *relational parametricity* [Reynolds 1983]. 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  $?$ , careful measures are to be taken to ensure parametricity. Our translation target  $\lambda B$  is taken from Ahmed et al. [2009], where relational parametricity is enforced by dynamic sealing [Matthews and Ahmed 2008; Neis et al. 2009], but there is no rigorous proof. Later, Ahmed et al. [2009] imposed a syntactic restriction on terms of  $\lambda B$ , where all type abstractions must have *values* as their body. With this invariant, they proved that the restricted  $\lambda B$  satisfies relational parametricity. It remains to see if our translation process can be adjusted to target restricted  $\lambda B$ . One possibility is to impose similar restriction to the rule [GPC-GEN](#):

$$\frac{\Psi, a \vdash^G e : \sigma \rightsquigarrow v}{\Psi \vdash^G e : \forall a. \sigma \rightsquigarrow \Lambda a. v} \text{GPC-GEN2}$$

where we only generate type abstractions if the inner body is a value. However, the type system with this rule is a weaker calculus, which is not a conservative extension of the OL type system.

**AMBIGUITY FROM CASTS.** The translation does not always produce a unique target expression. This is because when guessing some monotype  $\tau$  in rule **GPC-M-FORALL** and rule **GPC-CS-FORALL**, we could have many choices, which inevitably leads to different types. This is usually not a problem for (non-gradual) System F-like systems [Dunfield and Krishnaswami 2013; Peyton Jones et al. 2007] because they adopt a type-erasure semantics [Pierce 2002]. However, in our case, the choice of monotypes may affect the runtime behaviour of translated programs, since they could appear inside the explicit casts. For instance, the following example shows two possible translations for the same source expression  $(\lambda x : ?.fx) : ? \rightarrow \text{Int}$ , where the type of  $f$  is instantiated to  $\text{Int} \rightarrow \text{Int}$  and  $\text{Bool} \rightarrow \text{Int}$ , respectively:

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

If we apply  $\lambda x : ?.fx$  to 3, which is fine since the function can take any input, the first translation runs smoothly in  $\lambda B$ , while the second one will raise a cast error ( $\text{Int}$  cannot be cast to  $\text{Bool}$ ). Similarly, if we apply it to `true`, then the second succeeds while the first fails. The culprit lies in the highlighted parts where the instantiation of  $a$  appears in 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.

The semantic discrepancy is due to the guessing nature of the *declarative* system. As far as the static semantics is concerned, both  $\text{Int} \rightarrow \text{Int}$  and  $\text{Bool} \rightarrow \text{Int}$  are equally acceptable. But this is not the case at runtime. The astute reader may have found that the *only* appropriate choice is to instantiate the type of  $f$  to  $? \rightarrow \text{Int}$  in the matching judgment. However, as specified by rule **GPC-M-FORALL** in Figure 4.7, we can only instantiate type variables to monotypes, but  $?$  is *not* a monotype! We will get back to this issue in Chapter 5.

**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 [Reynolds 1991]. 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, suppose  $f$  is defined as  $(\lambda x. 1)$ , then whatever the



translation might be, applying  $(\lambda x : ?.f x)$  to 3 either results in a cast error, or produces 1, nothing else.

We defined contextual equivalence [Morris Jr 1969] to formally characterize that two open expressions have the same behavior. The definition of contextual equivalence requires a notion of well-typed expression contexts  $\sigma_3$ , written  $\mathcal{C} : (\Psi \vdash^B \sigma) \rightsquigarrow (\Psi' \vdash^B \sigma')$ . The definitions of contexts and context typing are standard and thus omitted. As is common, we first define contextual approximation. In our setting, we need to relax the notion of contextual approximation of  $\lambda B$  [Ahmed et al. 2009] to also take into consideration of cast errors. We write  $\Psi \vdash s_1 \preceq_{ctx} s_2 : \sigma$  to say that  $s_2$  mimics the behaviour of  $s_1$  at type  $\sigma$  in the sense that whenever a program containing  $s_1$  reduces to an integer, replacing it with  $s_2$  either reduces to the same integer, or emits a cast error. We restrict the program results to integers to eliminate the role of types in values. If it is not an integer, it is always possible to embed it into another context that reduces to an integer. Then we write  $\Psi \vdash s_1 \simeq_{ctx} s_2 : \sigma$  to say  $s_1$  and  $s_2$  are contextually equivalent, that is, they approximate each other.

**Definition 8** (Contextual Approximation and Equivalence up to Cast Errors).

$$\begin{aligned} \Psi \vdash s_1 \preceq_{ctx} s_2 : \sigma &\triangleq \Psi \vdash^B s_1 : \sigma \wedge \Psi \vdash^B s_2 : \sigma \wedge \\ &\text{for all } \mathcal{C} : (\Psi \vdash^B \sigma) \rightsquigarrow (\bullet \vdash^B \text{Int}) \implies \\ &\quad \mathcal{C}\{s_1\} \Downarrow n \implies (\mathcal{C}\{s_2\} \Downarrow n \vee \mathcal{C}\{s_2\} \Downarrow \text{blame}) \\ \Psi \vdash s_1 \simeq_{ctx} s_2 : \sigma &\triangleq \Psi \vdash s_1 \preceq_{ctx} s_2 : \sigma \wedge \Psi \vdash s_2 \preceq_{ctx} s_1 : \sigma \end{aligned}$$

Before presenting the formal definition of coherence, first we observe that after erasing types and casts, all translations of the same expression are exactly the same. This is easy to see by examining each elaboration rule. We use  $\lfloor s \rfloor$  to denote an expression in  $\lambda B$  after erasure.

**Lemma 4.8.** *If  $\Psi \vdash^G e : \sigma \rightsquigarrow s_1$ , and  $\Psi \vdash^G e : \sigma \rightsquigarrow s_2$ , then  $\lfloor s_1 \rfloor \equiv_\alpha \lfloor s_2 \rfloor$ .*

Second, at runtime, the only role of types and casts is to emit cast errors caused by type mismatch. Therefore, By Lemma 4.8 coherence follows as a corollary:

**Lemma 4.9** (Coherence up to cast errors). *For any expression  $e$  such that  $\Psi \vdash^G e : \sigma \rightsquigarrow s_1$  and  $\Psi \vdash^G e : \sigma \rightsquigarrow s_2$ , we have  $\Psi \vdash s_1 \simeq_{ctx} s_2 : \sigma$ .*

#### 4.3.3 CORRECTNESS CRITERIA

Siek et al. [2015] present a set of properties, the *refined criteria*, that a well-designed gradual typing calculus must have. Among all the criteria, those related to the static aspects of gradual typing are well summarized by Cimini and Siek [2016]. 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 4.10** (Correctness Criteria).

- **Conservative extension:** for all static  $\Psi$ ,  $e$ , and  $\sigma$ ,
  - if  $\Psi \vdash^{OL} e : \sigma$ , then there exists  $\sigma_2$ , such that  $\Psi \vdash^G e : \sigma_2$ , and  $\Psi \vdash^G \sigma_2 <: \sigma$ .
  - if  $\Psi \vdash^G e : \sigma$ , then  $\Psi \vdash^{OL} e : \sigma$
- **Monotonicity w.r.t. precision:** for all  $\Psi, e, e', \sigma$ , if  $\Psi \vdash^G e : \sigma$ , and  $e' \sqsubseteq e$ , then  $\Psi \vdash^G e' : \sigma_2$ , and  $\sigma_2 \sqsubseteq \sigma$  for some  $B$ .
- **Type Preservation of cast insertion:** for all  $\Psi, e, \sigma$ , if  $\Psi \vdash^G e : \sigma$ , then  $\Psi \vdash^G e : \sigma \rightsquigarrow s$ , and  $\Psi \vdash^B s : \sigma$  for some  $s$ .
- **Monotonicity of cast insertion:** for all  $\Psi, e_1, e_2, s_1, s_2, \sigma$ , if  $\Psi \vdash^G e_1 : \sigma \rightsquigarrow s_1$ , and  $\Psi \vdash^G e_2 : \sigma \rightsquigarrow s_2$ , and  $e_1 \sqsubseteq e_2$ , then  $\Psi \vdash s_1 \sqsubseteq^B s_2$ .

The first criterion states that the gradual type system should be a conservative extension of the original system. In other words, a *static* program is typeable in the OL 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  $?$ <sup>6</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  $\sigma_1 \sqsubseteq \sigma_2$ , which is given in Figure 4.8. The relation intuitively captures a notion of types containing more or less unknown types ( $?$ ). 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  $e_1$  has more unknown types.

The first two criteria are fundamental to gradual typing. They explain for example why these two programs  $\lambda x : \text{Nat}. x + 1$  and  $\lambda x : ?. x + 1$  are typeable, as the former is typeable in the OL 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 4.7, given that a target expression should always exist, which can be easily seen from Fig. 4.7. The last criterion ensures that the translation must be monotonic over the precision relation  $\sqsubseteq$ . Ahmed et al. [2009] does not include a formal definition of precision, but an *approximation* definition and a *simulation relation*. Here we adapt the

<sup>6</sup>Note that the term *static* has appeared several times with different meanings.

$\sigma_1 \sqsubseteq \sigma_2$		(Type Precision)	
GPC-L-UNKNOWN	GPC-L-INT	GPC-L-ARROW	GPC-L-TVAR
$\frac{}{? \sqsubseteq \sigma}$	$\frac{}{\text{Int} \sqsubseteq \text{Int}}$	$\frac{\sigma_1 \sqsubseteq \sigma_3 \quad \sigma_2 \sqsubseteq \sigma_4}{\sigma_1 \rightarrow \sigma_2 \sqsubseteq \sigma_3 \rightarrow \sigma_4}$	$\frac{}{a \sqsubseteq a}$
GPC-L-FORALL			
$\frac{\sigma_1 \sqsubseteq \sigma_2}{\forall a. \sigma_1 \sqsubseteq \forall a. \sigma_2}$			
$e_1 \sqsubseteq e_2$		(Term Precision)	
GPC-LE-REFL	GPC-LE-LAMANN	GPC-LE-APP	
$\frac{}{e \sqsubseteq e}$	$\frac{\sigma_1 \sqsubseteq \sigma_2 \quad e_1 \sqsubseteq e_2}{\lambda x : \sigma_1. e_1 \sqsubseteq \lambda x : \sigma_2. e_2}$	$\frac{e_1 \sqsubseteq e_3 \quad e_2 \sqsubseteq e_4}{e_1 e_2 \sqsubseteq e_3 e_4}$	
$s_1 \sqsubseteq s_2$		(Term Precision in $\lambda B$ )	
B-LE-VAR	B-LE-NAT	B-LE-LAMANN	B-LE-TABS
$\frac{}{x \sqsubseteq x}$	$\frac{}{n \sqsubseteq n}$	$\frac{\sigma_1 \sqsubseteq \sigma_2 \quad s_1 \sqsubseteq s_2}{\lambda x : \sigma_1. s_1 \sqsubseteq \lambda x : \sigma_2. s_2}$	$\frac{s_1 \sqsubseteq s_2}{\Lambda a. s_1 \sqsubseteq \Lambda a. s_2}$
B-LE-APP		B-LE-CAST	
$\frac{s_1 \sqsubseteq s_3 \quad s_2 \sqsubseteq s_4}{s_1 s_2 \sqsubseteq s_3 s_4}$		$\frac{\sigma_1 \sqsubseteq \sigma_3 \quad \sigma_2 \sqsubseteq \sigma_4 \quad s_1 \sqsubseteq s_2}{\langle \sigma_1 \hookrightarrow \sigma_2 \rangle s_1 \sqsubseteq \langle \sigma_3 \hookrightarrow \sigma_4 \rangle s_2}$	

Figure 4.8: Less Precision

simulation relation as the precision, and a subset of it that is used in our system is given at the bottom of Figure 4.8.

**THE DYNAMIC GRADUAL GUARANTEE.** Besides the static criteria, there is also a criterion concerning the dynamic semantics, known as *the dynamic gradual guarantee* [Siek et al. 2015].

**Definition 9** (Dynamic Gradual Guarantee). Suppose  $e' \sqsubseteq e$ , and  $\bullet \vdash^G e : \sigma \rightsquigarrow s$  and  $\bullet \vdash^G e' : \sigma' \rightsquigarrow s'$ ,

- if  $s \Downarrow v$ , then  $s' \Downarrow v'$  and  $v' \sqsubseteq v$ . If  $s \Uparrow$  then  $s' \Uparrow$ .
- if  $s' \Downarrow v'$ , then  $s \Downarrow v$  where  $v' \sqsubseteq v$ , or  $s \Downarrow \text{blame}$ . If  $s' \Uparrow$  then  $s \Uparrow$  or  $s \Downarrow \text{blame}$ .

The first part of the dynamic gradual guarantee says that if a gradually typed program evaluates to a value, then making type annotations less precise always produces a program that evaluates to an less precise value. Unfortunately, coherence up to cast errors in the declarative system breaks the dynamic gradual guarantee. For instance:

$$(\lambda f : \forall a. a \rightarrow \text{Int}. \lambda x : \text{Int}. f x) (\lambda x. 1) 3 \quad (\lambda f : \forall a. a \rightarrow \text{Int}. \lambda x : ?. f x) (\lambda x. 1) 3$$

The left one evaluates to 1, whereas its less precise version (right) will give a cast error if  $a$  is instantiated to  $\text{Bool}$  for example. In Chapter 5, we will present an extension of the declarative system that will alleviate the issue.

#### 4.4 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 DK [Dunfield and Krishnaswami 2013]. However our algorithmic system differs from theirs in three aspects: (1) the addition of the unknown type  $?$ ; (2) the use of the matching judgment; and 3) the approach of *gradual inference only producing static types* [Garcia and Cimini 2015]. We then prove that our algorithm is both sound and complete with respect to the declarative type system. We also provide an implementation.

**ALGORITHMIC CONTEXTS.** The top of Figure 4.9 shows the syntax of the algorithmic system. A noticeable difference are the algorithmic contexts  $\Gamma$ , which are represented as an *ordered* list containing declarations of type variables  $a$  and term variables  $x : \sigma$ . Unlike declarative contexts, algorithmic contexts also contain declarations of existential type variables  $\hat{a}$ ,

Expressions	$e ::= x \mid n \mid \lambda x : \sigma. e \mid \lambda x. e \mid e_1 e_2 \mid e : \sigma \mid \text{let } x = e_1 \text{ in } e_2$
Types	$\sigma ::= \text{Int} \mid a \mid \hat{\alpha} \mid \sigma_1 \rightarrow \sigma_2 \mid \forall a. \sigma \mid ?$
Monotypes	$\tau ::= \text{Int} \mid a \mid \hat{\alpha} \mid \tau_1 \rightarrow \tau_2$
Algorithmic Contexts	$\Gamma, \Delta, \Theta ::= \bullet \mid \Gamma, x : \sigma \mid \Gamma, a \mid \Gamma, \hat{\alpha} \mid \Gamma, \hat{\alpha} = \tau \mid \Gamma, \blacktriangleright_{\hat{\alpha}}$
Complete Contexts	$\Omega ::= \bullet \mid \Omega, x : \sigma \mid \Omega, a \mid \Omega, \hat{\alpha} = \tau \mid \Omega, \blacktriangleright_{\hat{\alpha}}$

$\boxed{\Gamma \vdash^G \sigma}$	(Well-formedness of types)		
GPC-AD-INT	GPC-AD-UNKNOWN	GPC-AD-TVAR	GPC-AD-EVAR
$\overline{\Gamma \vdash^G \text{Int}}$	$\overline{\Gamma \vdash^G ?}$	$\overline{\Gamma[a] \vdash^G a}$	$\overline{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha}}$
GPC-AD-SOLVED	GPC-AD-ARROW	GPC-AD-FORALL	
$\overline{\Gamma[\hat{\alpha} = \tau] \vdash^G \hat{\alpha}}$	$\frac{\Gamma \vdash^G \sigma_1 \quad \Gamma \vdash^G \sigma_2}{\Gamma \vdash^G \sigma_1 \rightarrow \sigma_2}$	$\frac{\Gamma, a \vdash^G \sigma}{\Gamma \vdash^G \forall a. \sigma}$	

$\boxed{\vdash^G \Gamma}$	(Well-formedness of algorithmic contexts)		
GPC-WF-EMPTY	GPC-WF-VAR	GPC-WF-TVAR	
$\overline{\vdash^G \bullet}$	$\frac{\vdash^G \Gamma \quad x \notin \text{FV}(\Gamma) \quad \Gamma \vdash^G \sigma}{\vdash^G \Gamma, x : \sigma}$	$\frac{\vdash^G \Gamma \quad a \notin \text{FV}(\Gamma)}{\vdash^G \Gamma, a}$	
GPC-WF-EVAR	GPC-WF-SOLVED	GPC-WF-MARKER	
$\frac{\vdash^G \Gamma \quad \hat{\alpha} \notin \text{FV}(\Gamma)}{\vdash^G \Gamma, \hat{\alpha}}$	$\frac{\vdash^G \Gamma \quad \hat{\alpha} \notin \text{FV}(\Gamma) \quad \Gamma \vdash^G \tau}{\vdash^G \Gamma, \hat{\alpha} = \tau}$	$\frac{\vdash^G \Gamma \quad \blacktriangleright_{\hat{\alpha}} \notin \text{FV}(\Gamma)}{\vdash^G \Gamma, \blacktriangleright_{\hat{\alpha}}}$	

Figure 4.9: Syntax and well-formedness of the algorithmic GPC

which can be either unsolved (written  $\hat{\alpha}$ ) or solved to some monotype (written  $\hat{\alpha} = \tau$ ). Finally, algorithmic contexts include a *marker*  $\blacktriangleright_{\hat{\alpha}}$  (read “marker  $\hat{\alpha}$ ”), which is used to delineate existential variables created by the algorithm. We will have more to say about markers when we examine the rules. Complete contexts  $\Omega$  are the same as contexts, except that they contain no unsolved variables.

Apart from expressions in the declarative system, we add annotated expressions  $e : \sigma$ . The well-formedness judgments for types and contexts are shown in Figure 4.9.

**NOTATIONAL CONVENIENCE.** Following DK system, we use contexts as substitutions on types. We write  $[\Gamma]\sigma$  to mean  $\Gamma$  applied as a substitution to type  $\sigma$ . We also use a hole notation, which is useful when manipulating contexts by inserting and replacing declarations in the middle. The hole notation is used extensively in proving soundness and completeness. For example,  $\Gamma[\Theta]$  means  $\Gamma$  has the form  $\Gamma_L, \Theta, \Gamma_R$ ; if we have  $\Gamma[\hat{\alpha}] = (\Gamma_L, \hat{\alpha}, \Gamma_R)$ , then  $\Gamma[\hat{\alpha} = \tau] = (\Gamma_L, \hat{\alpha} = \tau, \Gamma_R)$ . Occasionally, we will see a context with two *ordered* holes, e.g.,  $\Gamma = \Gamma_0[\Theta_1][\Theta_2]$  means  $\Gamma$  has the form  $\Gamma_L, \Theta_1, \Gamma_M, \Theta_2, \Gamma_R$ .

**INPUT AND OUTPUT CONTEXTS.** The algorithmic system, compared with the declarative system, includes similar judgment forms, except that we replace the declarative context  $\Psi$  with an algorithmic context  $\Gamma$  (the *input context*), and add an *output context*  $\Delta$  after a backward turnstile, e.g.,  $\Gamma \vdash^G \sigma_1 \lesssim \sigma_2 \dashv \Delta$  is the judgment form for the algorithmic consistent subtyping. All algorithmic rules manipulate input and output contexts in a way that is consistent with the notion of *context extension*, which will be described in Section 4.4.5.

We start with the explanation of the algorithmic consistent subtyping as it involves manipulating existential type variables explicitly (and solving them if possible).

#### 4.4.1 ALGORITHMIC CONSISTENT SUBTYPING

Figure 4.10 presents the rules of algorithmic consistent subtyping  $\Gamma \vdash^G \sigma_1 \lesssim \sigma_2 \dashv \Delta$ , which says that under input context  $\Gamma$ ,  $\sigma_1$  is a consistent subtype of  $\sigma_2$ , with output context  $\Delta$ . The first five rules do not manipulate contexts, but illustrate how contexts are propagated.

Rule **GPC-AS-TVAR** and rule **GPC-AS-INT** do not involve existential variables, so the output contexts remain unchanged. Rule **GPC-AS-EVAR** says that any unsolved existential variable is a consistent subtype of itself. The output is still the same as the input context as the rule gives no clue as to what is the solution of that existential variable. Rules **GPC-AS-UNKNOWNL** and **AS-UNKNOWNR** are the counterparts of rule **GPC-CS-UNKNOWNL** and rule **GPC-CS-UNKNOWNR**.

Rule **GPC-AS-ARROW** 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

$\boxed{\Gamma \vdash^G \sigma_1 \lesssim \sigma_2 \dashv \Delta} \quad (\text{Under input context } \Gamma, \sigma_1 \text{ is a consistent subtype of } \sigma_2, \text{ with output context } \Delta)$		
$\frac{\text{GPC-AS-TVAR}}{\Gamma[a] \vdash^G a \lesssim a \dashv \Gamma[a]}$	$\frac{\text{GPC-AS-INT}}{\Gamma \vdash^G \text{Int} \lesssim \text{Int} \dashv \Gamma}$	$\frac{\text{GPC-AS-EVAR}}{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim \hat{\alpha} \dashv \Gamma[\hat{\alpha}]}$
$\frac{\text{GPC-AS-UNKNOWNL}}{\Gamma \vdash^G ? \lesssim \sigma \dashv \Gamma}$		$\frac{\text{GPC-AS-UNKNOWNR}}{\Gamma \vdash^G \sigma \lesssim ? \dashv \Gamma}$
$\frac{\text{GPC-AS-ARROW} \quad \Gamma \vdash^G \sigma_3 \lesssim \sigma_1 \dashv \Theta \quad \Theta \vdash^G [\Theta]\sigma_2 \lesssim [\Theta]\sigma_4 \dashv \Delta}{\Gamma \vdash^G \sigma_1 \rightarrow \sigma_2 \lesssim \sigma_3 \rightarrow \sigma_4 \dashv \Delta}$		$\frac{\text{GPC-AS-FORALLR} \quad \Gamma, a \vdash^G \sigma_1 \lesssim \sigma_2 \dashv \Delta, a, \Theta}{\Gamma \vdash^G \sigma_1 \lesssim \forall a. \sigma_2 \dashv \Delta}$
$\frac{\text{GPC-AS-FORALLL} \quad \Gamma, \blacktriangleright_{\hat{\alpha}}, \hat{\alpha} \vdash^G \sigma_1[a \mapsto \hat{\alpha}] \lesssim \sigma_2 \dashv \Delta, \blacktriangleright_{\hat{\alpha}}, \Theta}{\Gamma \vdash^G \forall a. \sigma_1 \lesssim \sigma_2 \dashv \Delta}$		$\frac{\text{GPC-AS-INSTL} \quad \hat{\alpha} \notin \text{FV}(\sigma) \quad \Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta}{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta}$
$\frac{\text{GPC-AS-INSTR} \quad \hat{\alpha} \notin \text{FV}(\sigma) \quad \Gamma[\hat{\alpha}] \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta}{\Gamma[\hat{\alpha}] \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta}$		

Figure 4.10: Algorithmic consistent subtyping

premise is the output context of the conclusion. Note that we do not simply check  $\sigma_2 \lesssim \sigma_4$ , but apply  $\Theta$  (the input context of the second premise) to both types (e.g.,  $[\Theta]\sigma_2$ ). This is to maintain an important invariant: whenever  $\Gamma \vdash^G \sigma_1 \lesssim \sigma_2 \dashv \Delta$  holds, the types  $\sigma_1$  and  $\sigma_2$  are fully applied under input context  $\Gamma$  (they contain no existential variables already solved in  $\Gamma$ ). The same invariant applies to every algorithmic judgment.

Rule **GPC-AS-FORALLR**, similar to the declarative rule **GPC-CS-FORALLR**, adds  $a$  to the input context. Note that the output context of the premise allows additional existential variables to appear after the type variable  $a$ , in a trailing context  $\Theta$ . These existential variables could depend on  $a$ ; since  $a$  goes out of scope in the conclusion, we need to drop them from the concluding output, resulting in  $\Delta$ . The next rule is essential to eliminating the guessing work. Instead of guessing a monotype  $\tau$  out of thin air, rule **GPC-AS-FORALLL** generates a fresh existential variable  $\hat{\alpha}$ , and replaces  $a$  with  $\hat{\alpha}$  in the body  $\sigma$ . The new existential variable  $\hat{\alpha}$  is then added to the input context, just before the marker  $\blacktriangleright_{\hat{\alpha}}$ . The output context  $(\Delta, \blacktriangleright_{\hat{\alpha}}, \Theta)$  allows additional existential variables to appear after  $\blacktriangleright_{\hat{\alpha}}$  in  $\Theta$ . For the same reasons as in rule **GPC-AS-FORALLR**, we drop them from the output context. A central idea behind these two rules is that we defer the decision of picking a monotype for a type variable, and hope that it could be solved later when we have more information at hand. As a side note, when both types are universal quantifiers, then either rule **GPC-AS-FORALLR** or rule **GPC-AS-FORALLL** could be tried. In practice, one can apply rule **GPC-AS-FORALLR** eagerly as it is invertible.

The last two rules (rule **GPC-AS-INSTL** and rule **GPC-AS-INSTR**) are specific to the algorithm, thus having no counterparts in the declarative version. They both check consistent subtyping with an unsolved existential variable on one side and an arbitrary type on the other side. Apart from checking that the existential variable does not occur in the type  $\sigma$ , both rules do not directly solve the existential variables, but leave the real work to the instantiation judgment.

#### 4.4.2 INSTANTIATION

Two symmetric judgments  $\Gamma \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta$  and  $\Gamma \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta$  defined in Figure 4.11 instantiate unsolved existential variables. They read “under input context  $\Gamma$ , instantiate  $\hat{\alpha}$  to a consistent subtype (or supertype) of  $\sigma$ , with output context  $\Delta$ ”. The judgments are extended naturally from DK system, whose original inspiration comes from Cardelli [1993]. Since these two judgments are mutually defined, we discuss them together.

Rule **GPC-INSTL-SOLVE** is the simplest one – when an existential variable meets a monotype – where we simply set the solution of  $\hat{\alpha}$  to the monotype  $\tau$  in the output context. We also need to check that the monotype  $\tau$  is well-formed under the prefix context  $\Gamma$ .



$\boxed{\Gamma \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta}$  (Under input context  $\Gamma$ , instantiate  $\hat{\alpha}$  such that  $\hat{\alpha} \lesssim \sigma$ , with output context  $\Delta$ )

$$\begin{array}{c}
\text{GPC-INSTL-SOLVE} \\
\frac{\Gamma \vdash^G \tau}{\Gamma, \hat{\alpha}, \Gamma' \vdash^G \hat{\alpha} \lesssim \tau \dashv \Gamma, \hat{\alpha} = \tau, \Gamma'} \\
\\
\text{GPC-INSTL-SOLVEU} \\
\frac{}{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim ? \dashv \Gamma[\hat{\alpha}]} \\
\\
\text{GPC-INSTL-REACH} \\
\frac{}{\Gamma[\hat{\alpha}][\hat{\beta}] \vdash^G \hat{\alpha} \lesssim \hat{\beta} \dashv \Gamma[\hat{\alpha}][\hat{\beta} = \hat{\alpha}]} \\
\\
\text{GPC-INSTL-FORALLR} \\
\frac{\Gamma[\hat{\alpha}], b \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta, b, \Theta}{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim \forall b. \sigma \dashv \Delta} \\
\\
\text{GPC-INSTL-ARR} \\
\frac{\Gamma[\hat{\alpha}_2, \hat{\alpha}_1, \hat{\alpha} = \hat{\alpha}_1 \rightarrow \hat{\alpha}_2] \vdash^G \sigma_1 \lesssim \hat{\alpha}_1 \dashv \Theta \quad \Theta \vdash^G \hat{\alpha}_2 \lesssim [\Theta]\sigma_2 \dashv \Delta}{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim \sigma_1 \rightarrow \sigma_2 \dashv \Delta}
\end{array}$$

$\boxed{\Gamma \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta}$  (Under input context  $\Gamma$ , instantiate  $\hat{\alpha}$  such that  $\sigma \lesssim \hat{\alpha}$ , with output context  $\Delta$ )

$$\begin{array}{c}
\text{GPC-INSTR-SOLVE} \\
\frac{\Gamma \vdash^G \tau}{\Gamma, \hat{\alpha}, \Gamma' \vdash^G \tau \lesssim \hat{\alpha} \dashv \Gamma, \hat{\alpha} = \tau, \Gamma'} \\
\\
\text{GPC-INSTR-SOLVEU} \\
\frac{}{\Gamma[\hat{\alpha}] \vdash^G ? \lesssim \hat{\alpha} \dashv \Gamma[\hat{\alpha}]} \\
\\
\text{GPC-INSTR-REACH} \\
\frac{}{\Gamma[\hat{\alpha}][\hat{\beta}] \vdash^G \hat{\beta} \lesssim \hat{\alpha} \dashv \Gamma[\hat{\alpha}][\hat{\beta} = \hat{\alpha}]} \\
\\
\text{GPC-INSTR-FORALLL} \\
\frac{\Gamma[\hat{\alpha}], \blacktriangleright_{\hat{\beta}}, \hat{\beta} \vdash^G \sigma[b \mapsto \hat{\beta}] \lesssim \hat{\alpha} \dashv \Delta, \blacktriangleright_{\hat{\beta}}, \Theta}{\Gamma[\hat{\alpha}] \vdash^G \forall b. \sigma \lesssim \hat{\alpha} \dashv \Delta} \\
\\
\text{GPC-INSTR-ARR} \\
\frac{\Gamma[\hat{\alpha}_2, \hat{\alpha}_1, \hat{\alpha} = \hat{\alpha}_1 \rightarrow \hat{\alpha}_2] \vdash^G \hat{\alpha}_1 \lesssim \sigma_1 \dashv \Theta \quad \Theta \vdash^G [\Theta]\sigma_2 \lesssim \hat{\alpha}_2 \dashv \Delta}{\Gamma[\hat{\alpha}] \vdash^G \sigma_1 \rightarrow \sigma_2 \lesssim \hat{\alpha} \dashv \Delta}
\end{array}$$

Figure 4.11: Algorithmic instantiation

Rule **GPC-INSTL-SOLVEU** is similar to rule **GPC-AS-UNKNOWNR** in that we put no constraint<sup>7</sup> on  $\hat{\alpha}$  when it meets the unknown type  $?$ . This design decision reflects the point that type inference only produces static types [Garcia and Cimini 2015].

Rule **GPC-INSTL-REACH** deals with the situation where two existential variables meet. Recall that  $\Gamma[\hat{\alpha}][\hat{\beta}]$  denotes a context where some unsolved existential variable  $\hat{\alpha}$  is declared before  $\hat{\beta}$ . In this situation, the only logical thing we can do is to set the solution of one existential variable to the other one, depending on which one is declared before. For example, in the output context of rule **GPC-INSTL-REACH**, we have  $\hat{\beta} = \hat{\alpha}$  because in the input context,  $\hat{\alpha}$  is declared before  $\hat{\beta}$ .

Rule **GPC-INSTL-FORALLR** is the instantiation version of rule **GPC-AS-FORALLR**. Since our system is predicative,  $\hat{\alpha}$  cannot be instantiated to  $\forall b. \sigma$ , but we can decompose  $\forall b. \sigma$  in the same way as in rule **GPC-AS-FORALLR**. Rule **GPC-INSTL-FORALLL** is the instantiation version of rule **GPC-AS-FORALLL**.

Rule **GPC-INSTL-ARR** applies when  $\hat{\alpha}$  meets an arrow type. It follows that the solution must also be an arrow type. This is why, in the first premise, we generate two fresh existential variables  $\hat{\alpha}_1$  and  $\hat{\alpha}_2$ , and insert them just before  $\hat{\alpha}$  in the input context, so that we can solve  $\hat{\alpha}$  to  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$ . Note that the first premise  $\sigma_1 \lesssim \hat{\alpha}_1$  switches to the other instantiation judgment.

#### 4.4.3 ALGORITHMIC TYPING

We now turn to the algorithmic typing rules in Figure 4.12. Because general type inference for System F is undecidable [Wells 1999], our algorithmic system uses bidirectional type checking to accommodate (first-class) polymorphism. Traditionally, two modes are employed in bidirectional systems: the checking mode  $\Gamma \vdash^G e \Leftarrow \sigma \dashv \Theta$ , which takes a term  $e$  and a type  $\sigma$  as input, and ensures that the term  $e$  checks against  $\sigma$ ; the inference mode  $\Gamma \vdash^G e \Rightarrow \sigma \dashv \Theta$ , which takes a term  $e$  and produces a type  $\sigma$ . We first discuss rules in the inference mode.

Rule **GPC-INF-VAR** and rule **GPC-INF-INT** do not generate any new information and simply propagate the input context. Rule **GPC-INF-ANNO** is standard, switching to the checking mode in the premise.

In rule **GPC-INF-LAMANN**, we generate a fresh existential variable  $\hat{\beta}$  for the function codomain, and check the function body against  $\hat{\beta}$ . Note that it is tempting to write  $\Gamma, x : \sigma \vdash^G e \Rightarrow \sigma_2 \dashv \Delta, x : \sigma, \Theta$  as the premise (in the hope of better matching its declarative counterpart rule **GPC-LAMANN**), which has a subtle consequence. Consider the expression  $\lambda x : \text{Int}. \lambda y. y$ .

<sup>7</sup>As we will see in Chapter 5 where we present a more refined system, the “no constraint” statement is not entirely true.

$\boxed{\Gamma \vdash^G e \Rightarrow \sigma \dashv \Delta}$  (Under input context  $\Gamma$ ,  $e$  infers output type  $\sigma$ , with output context  $\Delta$ )

$$\begin{array}{c} \text{GPC-INF-VAR} \\ \frac{(x : \sigma) \in \Gamma}{\Gamma \vdash^G x \Rightarrow \sigma \dashv \Gamma} \end{array} \quad \begin{array}{c} \text{GPC-INF-INT} \\ \frac{}{\Gamma \vdash^G n \Rightarrow \text{Int} \dashv \Gamma} \end{array} \quad \begin{array}{c} \text{GPC-INF-ANNO} \\ \frac{\Gamma \vdash^G \sigma \quad \Gamma \vdash^G e \Leftarrow \sigma \dashv \Delta}{\Gamma \vdash^G e : \sigma \Rightarrow \sigma \dashv \Delta} \end{array}$$

$$\begin{array}{c} \text{GPC-INF-LAMANN} \\ \frac{\Gamma \vdash^G \sigma \quad \Gamma, \hat{\beta}, x : \sigma \vdash^G e \Leftarrow \hat{\beta} \dashv \Delta, x : \sigma, \Theta}{\Gamma \vdash^G \lambda x : \sigma. e \Rightarrow \sigma \rightarrow \hat{\beta} \dashv \Delta} \end{array}$$

$$\begin{array}{c} \text{GPC-INF-LAM} \\ \frac{\Gamma, \hat{\alpha}, \hat{\beta}, x : \hat{\alpha} \vdash^G e \Leftarrow \hat{\beta} \dashv \Delta, x : \hat{\alpha}, \Theta}{\Gamma \vdash^G \lambda x. e \Rightarrow \hat{\alpha} \rightarrow \hat{\beta} \dashv \Delta} \end{array}$$

$$\begin{array}{c} \text{GPC-INF-LET} \\ \frac{\Gamma \vdash^G e_1 \Rightarrow \sigma \dashv \Theta_1 \quad \Theta_1, \hat{\alpha}, x : \sigma \vdash^G e_2 \Leftarrow \hat{\alpha} \dashv \Delta, x : \sigma, \Theta_2}{\Gamma \vdash^G \text{let } x = e_1 \text{ in } e_2 \Rightarrow \hat{\alpha} \dashv \Delta} \end{array}$$

$$\begin{array}{c} \text{GPC-INF-APP} \\ \frac{\Gamma \vdash^G e_1 \Rightarrow \sigma \dashv \Theta_1 \quad \Theta_1 \vdash^G [\Theta_1] \sigma \triangleright \sigma_1 \rightarrow \sigma_2 \dashv \Theta_2 \quad \Theta_2 \vdash^G e_2 \Leftarrow [\Theta_2] \sigma_1 \dashv \Delta}{\Gamma \vdash^G e_1 e_2 \Rightarrow \sigma_2 \dashv \Delta} \end{array}$$

$\boxed{\Gamma \vdash^G e \Leftarrow \sigma \dashv \Delta}$  (Under input context  $\Gamma$ ,  $e$  checks against input type  $\sigma$ , with output context  $\Delta$ )

$$\begin{array}{c} \text{GPC-CHK-LAM} \\ \frac{\Gamma, x : \sigma_1 \vdash^G e \Leftarrow \sigma_2 \dashv \Delta, x : \sigma_1, \Theta}{\Gamma \vdash^G \lambda x. e \Leftarrow \sigma_1 \rightarrow \sigma_2 \dashv \Delta} \end{array} \quad \begin{array}{c} \text{GPC-CHK-GEN} \\ \frac{\Gamma, a \vdash^G e \Leftarrow \sigma \dashv \Delta, a, \Theta}{\Gamma \vdash^G e \Leftarrow \forall a. \sigma \dashv \Delta} \end{array}$$

$$\begin{array}{c} \text{GPC-CHK-SUB} \\ \frac{\Gamma \vdash^G e \Rightarrow \sigma_1 \dashv \Theta \quad \Theta \vdash^G [\Theta] \sigma_1 \lesssim [\Theta] \sigma_2 \dashv \Delta}{\Gamma \vdash^G e \Leftarrow \sigma_2 \dashv \Delta} \end{array}$$

$\boxed{\Gamma \vdash^G \sigma \triangleright \sigma_1 \rightarrow \sigma_2 \dashv \Delta}$  (Under input context  $\Gamma$ ,  $\sigma$  matches output type  $\sigma_1 \rightarrow \sigma_2$ , with output context  $\Delta$ )

$$\begin{array}{c} \text{GPC-AM-FORALL} \\ \frac{\Gamma, \hat{\alpha} \vdash^G \sigma[a \mapsto \hat{\alpha}] \triangleright \sigma_1 \rightarrow \sigma_2 \dashv \Delta}{\Gamma \vdash^G \forall a. \sigma \triangleright \sigma_1 \rightarrow \sigma_2 \dashv \Delta} \end{array} \quad \begin{array}{c} \text{GPC-AM-ARR} \\ \frac{}{\Gamma \vdash^G \sigma_1 \rightarrow \sigma_2 \triangleright \sigma_1 \rightarrow \sigma_2 \dashv \Gamma} \end{array}$$

$$\begin{array}{c} \text{GPC-AM-UNKNOWN} \\ \frac{}{\Gamma \vdash^G ? \triangleright ? \rightarrow ? \dashv \Gamma} \end{array} \quad \begin{array}{c} \text{GPC-AM-VAR} \\ \frac{}{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \triangleright \hat{\alpha}_1 \rightarrow \hat{\alpha}_2 \dashv \Gamma[\hat{\alpha}_1, \hat{\alpha}_2, \hat{\alpha} = \hat{\alpha}_1 \rightarrow \hat{\alpha}_2]} \end{array}$$

Figure 4.12: Algorithmic typing

Under the new premise, this is untypable because of  $\bullet \vdash^G \lambda x : \text{Int}. \lambda y. y \Rightarrow \text{Int} \rightarrow \hat{\alpha} \rightarrow \hat{\alpha} \dashv$   $\bullet$  where  $\hat{\alpha}$  is not found in the output context. This explains why we put  $\hat{\beta}$  before  $x : \sigma$  so that it remains in the output context  $\Delta$ . Rule **GPC-INF-LAM**, which corresponds to rule **GPC-LAM**, one of the guessing rules, is similar to rule **GPC-INF-LAMANN**. As with the other algorithmic rules that eliminate guessing, we create new existential variables  $\hat{\alpha}$  (for function domain) and  $\hat{\beta}$  (for function codomain) and check the function body against  $\hat{\beta}$ . Rule **GPC-INF-LET** is similar to rule **GPC-INF-LAMANN**.

**ALGORITHMIC MATCHING.** Rule **GPC-INF-APP** deserves attention. It relies on the algorithmic matching judgment  $\Gamma \vdash^G \sigma \triangleright \sigma_1 \rightarrow \sigma_2 \dashv \Delta$ . The matching judgment algorithmically synthesizes an arrow type from an arbitrary type. Rule **GPC-AM-FORALL** replaces  $a$  with a fresh existential variable  $\hat{\alpha}$ , thus eliminating guessing. Rule **GPC-AM-ARR** and rule **GPC-AM-UNKNOWN** correspond directly to the declarative rules. Rule **GPC-AM-VAR**, which has no corresponding declarative version, is similar to rule **GPC-INSTL-ARR**/**GPC-INSTL-ARR**: we create  $\hat{\alpha}_1$  and  $\hat{\alpha}_2$  and solve  $\hat{\alpha} \rightarrow \hat{\alpha}_1 \rightarrow \hat{\alpha}_2$  in the output context.

Back to the rule **GPC-INF-APP**. This rule first infers the type of  $e_1$ , producing an output context  $\Theta_1$ . Then it applies  $\Theta_1$  to  $A$  and goes into the matching judgment, which delivers an arrow type  $\sigma_1 \rightarrow \sigma_2$  and another output context  $\Theta_2$ .  $\Theta_2$  is used as the input context when checking  $e_2$  against  $[\Theta_2]\sigma_1$ , where we go into the checking mode.

Rules in the checking mode are quite standard. Rule **GPC-CHK-LAM** checks against  $\sigma_1 \rightarrow \sigma_2$ . Rule **GPC-CHK-GEN**, like the declarative rule **GPC-GEN**, adds a type variable  $a$  to the input context. Rule **GPC-CHK-SUB** uses the algorithmic consistent subtyping judgment.

#### 4.4.4 DECIDABILITY

Our algorithmic system is decidable. It is not at all obvious to see why this is the case, as many rules are not strictly structural (e.g., many rules have  $[\Gamma]\sigma$  in the premises). This implies that we need a more sophisticated measure to support the argument. Since the typing rules (Figure 4.12) depend on the consistent subtyping rules (Figure 4.10), which in turn depends on the instantiation rules (Figure 4.11), to show the decidability of the typing judgment, we need to show that the instantiation and consistent subtyping judgments are decidable. The proof strategy mostly follows that of the DK system. Here only highlights of the proofs are given.

**DECIDABILITY OF INSTANTIATION.** The basic idea is that we need to show  $\sigma$  in the instantiation judgments  $\Gamma \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta$  and  $\Gamma \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta$  always gets smaller. Most of the rules are structural and thus easy to verify (e.g., rule **INSTL-FORALLR**); the non-trivial cases

are rule **INSTL-ARR** and rule **INSTR-ARR** where context applications appear in the premises. The key observation there is that the instantiation rules preserve the size of (substituted) types. The formal statement of decidability of instantiation needs a few pre-conditions: assuming  $\hat{\alpha}$  is unsolved in the input context  $\Gamma$ , that  $\sigma$  is well-formed under the context  $\Gamma$ , that  $\sigma$  is fully applied under the input context  $\Gamma$  ( $[\Gamma]\sigma = \sigma$ ), and that  $\hat{\alpha}$  does not occur in  $\sigma$ . Those conditions are actually met when instantiation is invoked: rule **CHK-SUB** applies the input context, and the subtyping rules apply input context when needed.

**Theorem 4.11** (Decidability of Instantiation). *If  $\Gamma = \Gamma_0[\hat{\alpha}]$  and  $\Gamma \vdash^G \sigma$  such that  $[\Gamma]\sigma = \sigma$  and  $\hat{\alpha} \notin \text{FV}(\sigma)$  then:*

1. *Either there exists  $\Delta$  such that  $\Gamma \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta$ , or not.*
2. *Either there exists  $\Delta$  such that  $\Gamma \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta$ , or not.*

**DECIDABILITY OF ALGORITHMIC CONSISTENT SUBTYPING.** Proving decidability of the algorithmic consistent subtyping is a bit more involved, as the induction measure consists of several parts. We measure the judgment  $\Gamma \vdash^G \sigma_1 \lesssim \sigma_2 \dashv \Delta$  lexicographically by

- (M1) the number of  $\forall$ -quantifiers in  $\sigma_1$  and  $\sigma_2$ ;
- (M2) the number of unknown types in  $\sigma_1$  and  $\sigma_2$ ;
- (M3)  $|\text{UNSOLVED}(\Gamma)|$ : the number of unsolved existential variables in  $\Gamma$ ;
- (M4)  $|\Gamma \vdash^G \sigma_1| + |\Gamma \vdash^G \sigma_2|$ .

Notice that because of our gradual setting, we also need to measure the number of unknown types (M2). This is a key difference from the DK system. For (M4), we use *contextual size*—the size of well-formed types under certain contexts, which penalizes solved variables (\*).

**Definition 10** (Contextual Size).

$$\begin{aligned}
|\Gamma \vdash^G \text{Int}| &= 1 \\
|\Gamma \vdash^G ?| &= 1 \\
|\Gamma \vdash^G a| &= 1 \\
|\Gamma \vdash^G \hat{\alpha}| &= 1 \\
|\Gamma[\hat{\alpha} = \tau] \vdash^G \hat{\alpha}| &= 1 + |\Gamma[\hat{\alpha} = \tau] \vdash^G \tau| \quad (*) \\
|\Gamma \vdash^G \forall a. \sigma| &= 1 + |\Gamma, a \vdash^G \sigma| \\
|\Gamma \vdash^G \sigma_1 \rightarrow \sigma_2| &= 1 + |\Gamma \vdash^G \sigma_1| + |\Gamma \vdash^G \sigma_2|
\end{aligned}$$

**Theorem 4.12** (Decidability of Algorithmic Consistent Subtyping). *Given a context  $\Gamma$  and types  $\sigma_1, \sigma_2$  such that  $\Gamma \vdash^G \sigma_1$  and  $\Gamma \vdash^G \sigma_2$  and  $[\Gamma]\sigma_1 = \sigma_1$  and  $[\Gamma]\sigma_2 = \sigma_2$ , it is decidable whether there exists  $\Delta$  such that  $\Gamma \vdash^G \sigma_1 \lesssim \sigma_2 \dashv \Delta$ .*

**DECIDABILITY OF ALGORITHMIC TYPING.** Similar to proving decidability of algorithmic consistent subtyping, the key is to come up with a correct measure. Since the typing rules depend on the matching judgment, we first show decidability of the algorithmic matching.

**Lemma 4.13** (Decidability of Algorithmic Matching). *Given a context  $\Gamma$  and a type  $\sigma$  it is decidable whether there exist types  $\sigma_1, \sigma_2$  and a context  $\Delta$  such that  $\Gamma \vdash^G \sigma \triangleright \sigma_1 \rightarrow \sigma_2 \dashv \Delta$ .*

Now we are ready to show decidability of typing. The proof is obtained by induction on the lexicographically ordered triple: size of  $e$ , typing judgment (where the inference mode  $\Rightarrow$  is considered smaller than the checking mode  $\Leftarrow$ ) and contextual size.

$$\left\langle e, \begin{array}{c} \Rightarrow \\ \Leftarrow \end{array}, |\Gamma \vdash^G \sigma| \right\rangle$$

The above measure is much simpler than the corresponding one in the DK system, where they also need to consider the application judgment together with the inference and checking judgments. This shows another benefit (besides the independence from typing) of adopting the matching judgment.

**Theorem 4.14** (Decidability of Algorithmic Typing).

1. *Inference: Given a context  $\Gamma$  and a term  $e$ , it is decidable whether there exist a type  $\sigma$  and a context  $\Delta$  such that  $\Gamma \vdash^G e \Rightarrow \sigma \dashv \Delta$ .*
2. *Checking: Given a context  $\Gamma$ , a term  $e$  and a type  $\sigma$  such that  $\Gamma \vdash^G \sigma$ , it is decidable whether there exists a context  $\Delta$  such that  $\Gamma \vdash^G e \Leftarrow \sigma \dashv \Delta$ .*

#### 4.4.5 CONTEXT EXTENSION

To be confident that our algorithmic type system and the declarative type system agree with each other, we need to prove that the algorithmic rules are sound and complete with respect to the declarative specification. Before we give the formal statements of the soundness and completeness theorems, we need a meta-theoretical device, called *context extension* [Dunfield and Krishnaswami 2013], to capture a notion of information increase from input contexts to output contexts.

A context extension judgment  $\Gamma \longrightarrow \Delta$  reads “ $\Gamma$  is extended by  $\Delta$ ”. Intuitively, this judgment says that  $\Delta$  has at least as much information as  $\Gamma$ : some unsolved existential variables in  $\Gamma$  may be solved in  $\Delta$ . The full inductive definition can be found Figure 4.13.

$\boxed{\Gamma \longrightarrow \Delta}$				(Context extension)
GPC-EXT-ID $\frac{}{\bullet \longrightarrow \bullet}$	GPC-EXT-VAR $\frac{\Gamma \longrightarrow \Delta \quad [\Delta]\sigma = [\Delta]\sigma'}{\Gamma, x : \sigma \longrightarrow \Delta, x : \sigma'}$	GPC-EXT-TVAR $\frac{\Gamma \longrightarrow \Delta}{\Gamma, a \longrightarrow \Delta, a}$	GPC-EXT-EVAR $\frac{\Gamma \longrightarrow \Delta}{\Gamma, \hat{\alpha} \longrightarrow \Delta, \hat{\alpha}}$	
GPC-EXT-SOLVED $\frac{\Gamma \longrightarrow \Delta \quad [\Delta]\tau = [\Delta]\tau'}{\Gamma, \hat{\alpha} = \tau \longrightarrow \Delta, \hat{\alpha} = \tau'}$	GPC-EXT-SOLVE $\frac{\Gamma \longrightarrow \Delta}{\Gamma, \hat{\alpha} \longrightarrow \Delta, \hat{\alpha} = \tau}$	GPC-EXT-ADD $\frac{\Gamma \longrightarrow \Delta}{\Gamma \longrightarrow \Delta, \hat{\alpha}}$		
	GPC-EXT-ADDSOLVE $\frac{\Gamma \longrightarrow \Delta}{\Gamma \longrightarrow \Delta, \hat{\alpha} = \tau}$	GPC-EXT-MARKER $\frac{\Gamma \longrightarrow \Delta}{\Gamma, \blacktriangleright_{\hat{\alpha}} \longrightarrow \Delta, \blacktriangleright_{\hat{\alpha}}}$		

Figure 4.13: Context extension

## 4.4.6 SOUNDNESS

Roughly speaking, soundness of the algorithmic system says that given a derivation of an algorithmic judgment with input context  $\Gamma$ , output context  $\Delta$ , and a complete context  $\Omega$  that extends  $\Delta$ , applying  $\Omega$  throughout the given algorithmic judgment should yield a derivable declarative judgment. For example, let us consider an algorithmic typing judgment  $\bullet \vdash^G \lambda x. x \Rightarrow \hat{\alpha} \rightarrow \hat{\alpha} \dashv \hat{\alpha}$ , and any complete context, say,  $\Omega = (\hat{\alpha} = \text{Int})$ , then applying  $\Omega$  to the above judgment yields  $\bullet \vdash^G \lambda x. x : \text{Int} \rightarrow \text{Int}$ , which is derivable in the declarative system.

However there is one complication: applying  $\Omega$  to the algorithmic expression does not necessarily yield a typable declarative expression. For example, by rule [GPC-CHK-LAM](#) we have  $\lambda x. x \Leftarrow (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$ , but  $\lambda x. x$  itself cannot have type  $(\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a)$  in the declarative system. To circumvent that, we add an annotation to the lambda abstraction, resulting in  $\lambda x : (\forall a. a \rightarrow a). x$ , which is typeable in the declarative system with the same type. To relate  $\lambda x. x$  and  $\lambda x : (\forall a. a \rightarrow a). x$ , we erase all annotations on both expressions.

**Definition 11** (Type annotation erasure). The erasure function is denoted as  $|\cdot|$ , and defined as follows:

$$\begin{array}{ll}
|x| = x & |n| = n \\
|\lambda x : \sigma. e| = \lambda x. |e| & |\lambda x. e| = \lambda x. |e| \\
|e_1 e_2| = |e_1| |e_2| & |e : \sigma| = |e|
\end{array}$$

**Theorem 4.15** (Instantiation Soundness). *Given  $\Delta \longrightarrow \Omega$  and  $[\Gamma]\sigma = \sigma$  and  $\hat{\alpha} \notin \text{FV}(\sigma)$ :*

1. *If  $\Gamma \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta$  then  $[\Omega]\Delta \vdash^G [\Omega]\hat{\alpha} \lesssim [\Omega]\sigma$ .*

2. If  $\Gamma \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta$  then  $[\Omega]\Delta \vdash^G [\Omega]\sigma \lesssim [\Omega]\hat{\alpha}$ .

Notice that the declarative judgment uses  $[\Omega]\Delta$ , an operation that applies a complete context  $\Omega$  to the algorithmic context  $\Delta$ , essentially plugging in all known solutions and removing all declarations of existential variables (both solved and unsolved), resulting in a declarative context.

With instantiation soundness, next we show that the algorithmic consistent subtyping is sound:

**Theorem 4.16** (Soundness of Algorithmic Consistent Subtyping). *If  $\Gamma \vdash^G \sigma_1 \lesssim \sigma_2 \dashv \Delta$  where  $[\Gamma]\sigma_1 = \sigma_1$  and  $[\Gamma]\sigma_2 = \sigma_2$  and  $\Delta \longrightarrow \Omega$  then  $[\Omega]\Delta \vdash^G [\Omega]\sigma_1 \lesssim [\Omega]\sigma_2$ .*

Finally the soundness theorem of algorithmic typing is:

**Theorem 4.17** (Soundness of Algorithmic Typing). *Given  $\Delta \longrightarrow \Omega$ :*

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

#### 4.4.7 COMPLETENESS

Completeness of the algorithmic system is the reverse of soundness: given a declarative judgment of the form  $[\Omega]\Gamma \vdash^G [\Omega] \dots$ , we want to get an algorithmic derivation of  $\Gamma \vdash^G \dots \dashv \Delta$ . 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 4.18** (Instantiation Completeness). *Given  $\Gamma \longrightarrow \Omega$  and  $\sigma = [\Gamma]\sigma$  and  $\hat{\alpha} \in \text{UNSOLVED}(\Gamma)$  and  $\hat{\alpha} \notin \text{FV}(\sigma)$ :*

1. If  $[\Omega]\Gamma \vdash^G [\Omega]\hat{\alpha} \lesssim [\Omega]\sigma$  then there are  $\Delta, \Omega'$  such that  $\Omega \longrightarrow \Omega'$  and  $\Delta \longrightarrow \Omega'$  and  $\Gamma \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta$ .
2. If  $[\Omega]\Gamma \vdash^G [\Omega]\sigma \lesssim [\Omega]\hat{\alpha}$  then there are  $\Delta, \Omega'$  such that  $\Omega \longrightarrow \Omega'$  and  $\Delta \longrightarrow \Omega'$  and  $\Gamma \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta$ .

Next is the completeness of consistent subtyping:



**Theorem 4.19** (Generalized Completeness of Consistent Subtyping). *If  $\Gamma \longrightarrow \Omega$  and  $\Gamma \vdash^G \sigma_1$  and  $\Gamma \vdash^G \sigma_2$  and  $[\Omega]\Gamma \vdash^G [\Omega]\sigma_1 \lesssim [\Omega]\sigma_2$  then there exist  $\Delta$  and  $\Omega'$  such that  $\Delta \longrightarrow \Omega'$  and  $\Omega \longrightarrow \Omega'$  and  $\Gamma \vdash^G [\Gamma]\sigma_1 \lesssim [\Gamma]\sigma_2 \dashv \Delta$ .*

We prove that the algorithmic matching is complete with respect to the declarative matching:

**Theorem 4.20** (Matching Completeness). *Given  $\Gamma \longrightarrow \Omega$  and  $\Gamma \vdash^G \sigma$ , if  $[\Omega]\Gamma \vdash^G [\Omega]\sigma \triangleright \sigma_1 \rightarrow \sigma_2$  then there exist  $\Delta, \Omega', \sigma'_1$  and  $\sigma'_2$  such that  $\Gamma \vdash^G [\Gamma]\sigma \triangleright \sigma'_1 \rightarrow \sigma'_2 \dashv \Delta$  and  $\Delta \longrightarrow \Omega'$  and  $\Omega \longrightarrow \Omega'$  and  $\sigma_1 = [\Omega']\sigma'_1$  and  $\sigma_2 = [\Omega']\sigma'_2$ .*

Finally here is the completeness theorem of the algorithmic typing:

**Theorem 4.21** (Completeness of Algorithmic Typing). *Given  $\Gamma \longrightarrow \Omega$  and  $\Gamma \vdash^G \sigma$ , if  $[\Omega]\Gamma \vdash^G e : \sigma$  then there exist  $\Delta, \Omega', \sigma'$  and  $e'$  such that  $\Delta \longrightarrow \Omega'$  and  $\Omega \longrightarrow \Omega'$  and  $\Gamma \vdash^G e' \Rightarrow \sigma' \dashv \Delta$  and  $\sigma = [\Omega']\sigma'$  and  $|e| = |e'|$ .*

## 4.5 SIMPLE EXTENSIONS AND VARIANTS

This section considers two simple variations on the presented system. The first variation extends the system with a top type, while the second variation considers a more declarative formulation using a subsumption rule.

### 4.5.1 TOP TYPES

We argued that our definition of consistent subtyping (Definition 5) generalizes the original definition by Siek and Taha [2007b]. We have shown its applicability to polymorphic types, for which Siek and Taha [2007b] approach cannot be extended naturally. To strengthen our argument, we show how to extend our approach to  $\top$  types, which is also not supported by Siek and Taha [2007b] approach.

**CONSISTENT SUBTYPING WITH  $\top$ .** 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 [GPC-S-TOP](#). This rule might seem strange at first glance, since even if we remove the requirement  $\sigma$  static, the rule still 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  $? <: \top$  are not allowed. For consistency, we add the rule that  $\top$  is consistent with  $\top$ , which is actually included in the original reflexive rule  $\sigma \sim \sigma$ . For consistent subtyping, every type is a consistent subtype of  $\top$ , for example,  $\text{Nat} \rightarrow ? \lesssim \top$ .

$$\frac{\sigma \text{ static}}{\Psi \vdash^G \sigma <: \top} \text{GPC-S-TOP} \qquad \frac{}{\top \sim \top} \qquad \frac{}{\Psi \vdash^G \sigma \lesssim \top} \text{GPC-CS-TOP}$$

It is easy to verify that Definition 5 is still equivalent to that in Figure 4.6 extended with rule GPC-CS-TOP. That is, Theorem 4.4 holds:

**Proposition 4.22** (Extension with  $\top$ ).  $\Psi \vdash^G \sigma_1 \lesssim \sigma_2 \Leftrightarrow \Psi \vdash^G \sigma_1 <: \sigma', \sigma' \sim \sigma'', \Psi \vdash^G \sigma'' <: \sigma_2$  for some  $\sigma', \sigma''$ .

We extend the definition of concretization (Definition 6) with  $\top$  by adding another equation  $\gamma(\top) = \{\top\}$ . Note that Castagna and Lanvin [2017] also have this equation in their calculus. It is easy to verify that Corollary 4.2 still holds:

**Proposition 4.23** (Equivalent to AGT on  $\top$ ).  $\sigma_1 \lesssim \sigma_2$  if and only if  $\sigma_1 \widetilde{<} \sigma_2$ .

SIEK AND TAHA’S DEFINITION OF CONSISTENT SUBTYPING DOES NOT WORK FOR  $\top$ . As with the analysis in Section 4.2.2,  $\text{Nat} \rightarrow ? \lesssim \top$  only holds when we first apply consistency, then subtyping. However we cannot find a type  $\sigma$  such that  $\text{Nat} \rightarrow ? <: \sigma$  and  $\sigma \sim \top$ . The following diagram depicts the situation:

$$\begin{array}{ccc} \emptyset & \xrightarrow{\sim} & \top \\ \uparrow & & \uparrow \\ \text{Nat} \rightarrow ? & \xrightarrow{\sim} & \text{Nat} \rightarrow \text{Nat} \end{array}$$

Additionally we have a similar problem in extending the restriction operator: *non-structural* masking between  $\text{Int} \rightarrow ?$  and  $\top$  cannot be easily achieved.

Note that both the top and universally quantified types can be seen as special cases of intersection types. Indeed, top is the intersection of the empty set, while a universally quantified type is the intersection of the infinite set of its instantiations [Davies and Pfenning 2000]. Recall from Section 4.2.3 that Castagna and Lanvin [2017] shows that consistent subtyping from AGT works well for intersection types, and our definition coincides with AGT (Corollary 4.2 and Proposition 4.23). We thus believe that our definition is compatible with conventional binary intersection types as well. Yet, a rigorous formalization would be needed to substantiate this belief.

## 4.5.2 A MORE DECLARATIVE TYPE SYSTEM

In Section 4.3 we present our declarative system in terms of the matching and consistent subtyping judgments. The rationale behind this design choice is that the resulting declarative system combines subtyping and type consistency in the application rule, thus making it easier to design an algorithmic system accordingly. Still, one may wonder if it is possible to design a more declarative specification. For example, even though we mentioned that the subsumption rule is incompatible with consistent subtyping, it might be possible to accommodate a subsumption rule for normal subtyping (instead of consistent subtyping). In this section, we discuss an alternative for the design of the declarative system.

**WRONG DESIGN.** A naive design that does not work is to replace rule **GPC-APP** in Figure 4.7 with the following two rules:

$$\begin{array}{c}
 \text{GPC-V-SUB} \\
 \frac{\Psi \vdash^G e : \sigma \quad \sigma <: \sigma_2}{\Psi \vdash^G e : \sigma_2}
 \end{array}
 \qquad
 \begin{array}{c}
 \text{GPC-V-APP1} \\
 \frac{\Psi \vdash^G e_1 : \sigma \quad \Psi \vdash^G e_2 : \sigma_1 \quad \sigma \sim \sigma_1 \rightarrow \sigma_2}{\Psi \vdash^G e_1 e_2 : \sigma_2}
 \end{array}$$

Rule **GPC-V-SUB** is the standard subsumption rule: if an expression  $e$  has type  $\sigma$ , then it can be assigned some type  $\sigma_2$  that is a supertype of  $\sigma$ . Rule **GPC-V-APP1** first infers that  $e_1$  has type  $\sigma$ , and  $e_2$  has type  $\sigma_1$ , then it finds some  $\sigma_2$  so that  $\sigma$  is consistent with  $\sigma_1 \rightarrow \sigma_2$ .

There would be two obvious benefits of this variant if it did work: firstly this approach closely resembles the traditional declarative type systems for calculi with subtyping; secondly it saves us from discussing various forms of  $\sigma$  in rule **GPC-V-APP1**, leaving the job to the consistency judgment.

The design is wrong because of the information loss caused by the choice of  $\sigma_2$  in rule **GPC-V-APP1**. Suppose we have  $\Psi \vdash^G \text{plus} : \text{Nat} \rightarrow \text{Nat} \rightarrow \text{Nat}$ , then we can apply it to 1 to get

$$\frac{\Psi \vdash^G \text{plus} : \text{Nat} \rightarrow \text{Nat} \rightarrow \text{Nat} \quad \Psi \vdash^G 1 : \text{Nat} \quad \text{Nat} \rightarrow \text{Nat} \rightarrow \text{Nat} \sim \text{Nat} \rightarrow ? \rightarrow \text{Nat}}{\Psi \vdash^G \text{plus } 1 : ? \rightarrow \text{Nat}} \text{GPC-V-APP1}$$

Further applying it to true we get

$$\frac{\Psi \vdash^G \text{plus } 1 : ? \rightarrow \text{Nat} \quad \Psi \vdash^G \text{true} : \text{Bool} \quad ? \rightarrow \text{Nat} \sim \text{Bool} \rightarrow \text{Nat}}{\Psi \vdash^G \text{plus } 1 \text{ true} : \text{Nat}} \text{GPC-V-APP1}$$

which is obviously wrong! The type consistency in rule [GPC-V-APP1](#) causes information loss for both the argument type  $\sigma_1$  and the return type  $\sigma_2$ . The problem is that information of  $\sigma_2$  can get lost again if it appears in further applications. The moral of the story is that we should be very careful when using type consistency. We hypothesize that it is inevitable to do case analysis for the type of the function in an application (i.e.,  $\sigma$  in rule [GPC-V-APP1](#)).

**PROPER DECLARATIVE DESIGN.** The proper design refines the first variant by using a matching judgment to carefully distinguish two cases for the typing result of  $e_1$  in rule [GPC-V-APP1](#): (1) when it is an arrow type, and (2) when it is an unknown type. This variant replaces rule [GPC-APP](#) in Figure 4.7 with the following rules:

$$\begin{array}{c}
 \text{GPC-V-SUB} \\
 \frac{\Psi \vdash^G e : \sigma \quad \sigma <: \sigma_2}{\Psi \vdash^G e : \sigma_2} \\
 \\
 \text{GPC-V-APP2} \\
 \frac{\Psi \vdash^G e : \sigma \quad \Psi \vdash^G \sigma \triangleright \sigma_1 \rightarrow \sigma_2 \quad \Psi \vdash^G e_2 : \sigma_3 \quad \sigma_1 \sim \sigma_3}{\Psi \vdash^G e_1 e_2 : \sigma_2} \\
 \\
 \frac{}{\Psi \vdash^G \sigma_1 \rightarrow \sigma_2 \triangleright \sigma_1 \rightarrow \sigma_2} \qquad \frac{}{\Psi \vdash^G ? \triangleright ? \rightarrow ?}
 \end{array}$$

Rule [GPC-V-SUB](#) is the same as in the first variant. In rule [GPC-V-APP2](#), we infer that  $e_1$  has type  $\sigma$ , and use the matching judgment to get an arrow type  $\sigma_1 \rightarrow \sigma_2$ . Then we need to ensure that the argument type  $\sigma_3$  is *consistent with* (rather than a consistent subtype of)  $\sigma_1$ , and use  $\sigma_2$  as the result type of the application. The matching judgment only deals with two cases, as polymorphic types are handled by rule [GPC-V-SUB](#). These rules are closely related to the ones in Siek and Taha [2006] and Siek and Taha [2007b].

The more declarative nature of this system also implies that it is highly non-syntax-directed, and it does not offer any insight into combining subtyping and consistency. We have proved in Coq the following lemmas to establish soundness and completeness of this system with respect to our original system (to avoid ambiguity, we use the notation  $\vdash_m^G$  to indicate the more declarative version):

**Lemma 4.24** (Completeness of  $\vdash_m^G$ ). *If  $\Psi \vdash^G e : \sigma$ , then  $\Psi \vdash_m^G e : \sigma$ .*

**Lemma 4.25** (Soundness of  $\vdash_m^G$ ). *If  $\Psi \vdash_m^G e : \sigma$ , then there exists some  $B$ , such that  $\Psi \vdash^G e : \sigma_2$  and  $\Psi \vdash^G \sigma_2 <: \sigma$ .*

# 5

## RESTORING THE DYNAMIC GRADUAL GUARANTEE WITH TYPE PARAMETERS

In Section 4.3.2 we have seen an example where a single source expression could produce two different target expressions with different runtime behaviors. As we explained, this is due to the guessing nature of the declarative system, and, from the (source) typing point of view, no guessed type is particularly better than any other. As a consequence, this breaks the dynamic gradual guarantee as discussed in Section 4.3.3.

To alleviate this situation, we introduce *static type parameters*, which are placeholders for monotypes, and *gradual type parameters*, which are placeholders for monotypes that are consistent with the unknown type. The concept of static type parameters and gradual type parameters in the context of gradual typing was first introduced by Garcia and Cimini [2015], and later played a central role in the work of Igarashi et al. [2017]. In our type system, type parameters mainly help capture the notion of *representative translations*, and should not appear in a source program. With them we are able to recast the dynamic gradual guarantee in terms of representative translations, and to prove that every well-typed source expression possesses at least one representative translation. With a coherence conjecture regarding representative translations, the dynamic gradual guarantee of our extended source language now can be reduced to that of  $\lambda B$ , which, at the time of writing, is still an open question. TODO: not open anymore?

### 5.1 DECLARATIVE TYPE SYSTEM

The new syntax of types is given at the top of Figure 5.1, with the differences highlighted. In addition to the types of Figure 4.2, we add *static type parameters*  $S$ , and *gradual type parameters*  $G$ . Both kinds of type parameters are monotypes. The addition of type parameters, however, leads to two new syntactic categories of types. *Castable types*  $\mathbb{C}$  represent types that can be cast from or to  $?$ . It includes all types, except those that contain static type parameters. *Castable monotypes*  $t$  are those castable types that are also monotypes.

Types	$\sigma, \sigma_2 ::= \text{Int} \mid a \mid \sigma \rightarrow \sigma_2 \mid \forall a. \sigma \mid ? \mid \mathcal{S} \mid \mathcal{G}$
Monotypes	$\tau, \tau ::= \text{Int} \mid a \mid \tau \rightarrow \tau \mid \mathcal{S} \mid \mathcal{G}$
Castable Types	$\mathbb{C} ::= \text{Int} \mid a \mid \mathbb{C}_1 \rightarrow \mathbb{C}_2 \mid \forall a. \mathbb{C} \mid ? \mid \mathcal{G}$
Castable Monotypes	$t ::= \text{Int} \mid a \mid t_1 \rightarrow t_2 \mid \mathcal{G}$

  

$\Psi \vdash^G \sigma \lesssim \sigma_2$

(Consistent Subtyping)

$\frac{\text{GPC-CS-TVAR} \quad a \in \Psi}{\Psi \vdash^G a \lesssim a}$	$\frac{\text{GPC-CS-INT}}{\Psi \vdash^G \text{Int} \lesssim \text{Int}}$	$\frac{\text{GPC-CS-ARROW} \quad \Psi \vdash^G \sigma_3 \lesssim \sigma_1 \quad \Psi \vdash^G \sigma_2 \lesssim \sigma_4}{\Psi \vdash^G \sigma_1 \rightarrow \sigma_2 \lesssim \sigma_3 \rightarrow \sigma_4}$
$\frac{\text{GPC-CS-FORALLR} \quad \Psi, a \vdash^G \sigma_1 \lesssim \sigma_2}{\Psi \vdash^G \sigma_1 \lesssim \forall a. \sigma_2}$	$\frac{\text{GPC-CS-FORALLL} \quad \Psi \vdash^G \tau \quad \Psi \vdash^G \sigma_1[a \mapsto \tau] \lesssim \sigma_2}{\Psi \vdash^G \forall a. \sigma_1 \lesssim \sigma_2}$	$\frac{\text{GPC-CS-UNKNOWNLL}}{\Psi \vdash^G ? \lesssim \mathbb{C}}$
$\frac{\text{GPC-CS-UNKNOWNRR}}{\Psi \vdash^G \mathbb{C} \lesssim ?}$	$\frac{\text{GPC-CS-SPAR}}{\Psi \vdash^G \mathcal{S} \lesssim \mathcal{S}}$	$\frac{\text{GPC-CS-GPAR}}{\Psi \vdash^G \mathcal{G} \lesssim \mathcal{G}}$

Figure 5.1: Syntax of types, and consistent subtyping in the extended declarative system.

**CONSISTENT SUBTYPING.** The new definition of consistent subtyping is given at the bottom of Figure 5.1, again with the differences highlighted. Now the unknown type is only a consistent subtype of all castable types, rather than of all types (rule [GPC-CS-UNKNOWNLL](#)), and vice versa (rule [GPC-CS-UNKNOWNRR](#)). Moreover, the static type parameter  $\mathcal{S}$  is a consistent subtype of itself (rule [GPC-CS-SPAR](#)), and similarly for the gradual type parameter (rule [GPC-CS-GPAR](#)). From this definition it follows immediately that  $?$  is incomparable with types that contain static type parameters  $\mathcal{S}$ , such as  $\mathcal{S} \rightarrow \text{Int}$ .

**TYPING AND TRANSLATION.** Given these extensions to types and consistent subtyping, the typing process remains the same as in Figure 4.7. To account for the changes in the translation, if we extend  $\lambda B$  with type parameters as in Garcia and Cimini [2015], then the translation remains the same as well.

## 5.2 SUBSTITUTIONS AND REPRESENTATIVE TRANSLATIONS

As we mentioned, type parameters serve as placeholders for monotypes. As a consequence, wherever a type parameter is used, any *suitable* monotype could appear just as well. To formalize this observation, we define substitutions for type parameters as follows:

**Definition 12** (Substitution). Substitutions for type parameters are defined as:

1. Let  $S^{\mathcal{S}} : \mathcal{S} \rightarrow \tau$  be a total function mapping static type parameters to monotypes.
2. Let  $S^{\mathcal{G}} : \mathcal{G} \rightarrow t$  be a total function mapping gradual type parameters to castable monotypes.
3. Let  $S^{\mathcal{P}} = S^{\mathcal{G}} \cup S^{\mathcal{S}}$  be a union of  $S^{\mathcal{S}}$  and  $S^{\mathcal{G}}$  mapping static and gradual type parameters accordingly.

Note that since  $\mathcal{G}$  might be compared with  $?$ , only castable monotypes are suitable substitutes, whereas  $\mathcal{S}$  can be replaced by any monotypes. Therefore, we can substitute  $\mathcal{G}$  for  $\mathcal{S}$ , but not the other way around.

Let us go back to our example and its two translations in Section 4.3.2. The problem with those translations is that neither  $\text{Int} \rightarrow \text{Int}$  nor  $\text{Bool} \rightarrow \text{Int}$  is general enough. With type parameters, however, we can state a more *general* translation that covers both through substitution:

$$\begin{aligned} f : \forall a. a \rightarrow \text{Int} \vdash^G (\lambda x : ?.fx) : ? \rightarrow \text{Int} \\ \rightsquigarrow (\lambda x : ?. (\langle \forall a. a \rightarrow \text{Int} \hookrightarrow \mathcal{G} \rightarrow \text{Int} \rangle f) (\langle ? \hookrightarrow \mathcal{G} \rangle x)) \end{aligned}$$

The advantage of type parameters is that they help reasoning about the dynamic semantics. Now we are not limited to a particular choice, such as  $\text{Int} \rightarrow \text{Int}$  or  $\text{Bool} \rightarrow \text{Int}$ , which might or might not emit a cast error at runtime. Instead we have a general choice  $\mathcal{G} \rightarrow \text{Int}$ .

What does the more general choice with type parameters tell us? First, we know that in this case, there is no concrete constraint on  $a$ , so we can instantiate it with a type parameter. Second, the fact that the general choice uses  $\mathcal{G}$  rather than  $\mathcal{S}$  indicates that any chosen instantiation needs to be a castable type. It follows that any concrete instantiation will have an impact on the runtime behavior; therefore it is best to instantiate  $a$  with  $?$ . However, type inference cannot instantiate  $a$  with  $?$ , and substitution cannot replace  $\mathcal{G}$  with  $?$  either. This means that we need a syntactic refinement process of the translated programs in order to replace type parameters with allowed gradual types.

**SYNTACTIC REFINEMENT.** We define syntactic refinement of the translated expressions as follows. As  $\mathcal{S}$  denotes no constraints at all, substituting it with any monotype would work. Here we arbitrarily use  $\text{Int}$ . We interpret  $\mathcal{G}$  as  $?$  since any monotype could possibly lead to a cast error.

**Definition 13** (Syntactic Refinement). The syntactic refinement of a translated expression  $s$  is denoted by  $\lceil s \rceil$ , and defined as follows:

$\lceil \text{Int} \rceil$	=	$\text{Int}$
$\lceil a \rceil$	=	$a$
$\lceil \sigma_1 \rightarrow \sigma_2 \rceil$	=	$\lceil \sigma_1 \rceil \rightarrow \lceil \sigma_2 \rceil$
$\lceil \forall a. \sigma \rceil$	=	$\forall a. \lceil \sigma \rceil$
$\lceil ? \rceil$	=	$?$
$\lceil \mathcal{S} \rceil$	=	$\text{Int}$
$\lceil \mathcal{G} \rceil$	=	$?$

Applying the syntactic refinement to the translated expression, we get

$$(\lambda x : ?. (\langle \forall a. a \rightarrow \text{Int} \hookrightarrow ? \rightarrow \text{Int} \rangle f) (\langle ? \hookrightarrow ? \rangle x))$$

where two  $\mathcal{G}$  are refined by  $?$  as highlighted. It is easy to verify that both applying this expression to 3 and to *true* now results in a translation that evaluates to a value.

**REPRESENTATIVE TRANSLATIONS.** To decide whether one translation is more general than the other, we define a preorder between translations.

**Definition 14** (Translation Pre-order). Suppose  $\Psi \vdash^G e : \sigma \rightsquigarrow s_1$  and  $\Psi \vdash^G e : \sigma \rightsquigarrow s_2$ , we define  $s_1 \leq s_2$  to mean  $s_2 \equiv_\alpha S^{\mathcal{P}}(s_1)$  for some  $S^{\mathcal{P}}$ .

**Proposition 5.1.** *If  $s_1 \leq s_2$  and  $s_2 \leq s_1$ , then  $s_1$  and  $s_2$  are  $\alpha$ -equivalent (i.e., equivalent up to renaming of type parameters).*

The preorder between translations gives rise to a notion of what we call *representative translations*:

**Definition 15** (Representative Translation). A translation  $s$  is said to be a representative translation of a typing derivation  $\Psi \vdash^G e : \sigma \rightsquigarrow s$  if and only if for any other translation  $\Psi \vdash^G e : \sigma \rightsquigarrow s'$  such that  $s' \leq s$ , we have  $s \leq s'$ . From now on we use  $r$  to denote a representative translation.



An important property of representative translations, which we conjecture for the lack of rigorous proof, is that if there exists any translation of an expression that (after syntactic refinement) can reduce to a value, so can a representative translation of that expression. Conversely, if a representative translation runs into a blame, then no translation of that expression can reduce to a value.

**Conjecture 5.2** (Property of Representative Translations). *For any expression  $e$  such that  $\Psi \vdash^G e : \sigma \rightsquigarrow s$  and  $\Psi \vdash^G e : \sigma \rightsquigarrow r$  and  $\forall \mathcal{C}. \mathcal{C} : (\Psi \vdash^B \sigma) \rightsquigarrow (\bullet \vdash^B \text{Int})$ , we have*

- *If  $\mathcal{C}\{[s]\} \Downarrow n$ , then  $\mathcal{C}\{[r]\} \Downarrow n$ .*
- *If  $\mathcal{C}\{[r]\} \Downarrow \text{blame}$ , then  $\mathcal{C}\{[s]\} \Downarrow \text{blame}$ .*

Given this conjecture, we can state a stricter coherence property (without the “up to casts” part) between any two representative translations. We first strengthen Definition 8 following Ahmed et al. [2009]:

**Definition 16** (Contextual Approximation à la Ahmed et al. [2009]).

$$\begin{aligned} \Psi \vdash s_1 \preceq_{ctx} s_2 : \sigma &\triangleq \Psi \vdash^B s_1 : \sigma \wedge \Psi \vdash^B s_2 : \sigma \wedge \\ &\text{for all } \mathcal{C}. \mathcal{C} : (\Psi \vdash^B \sigma) \rightsquigarrow (\bullet \vdash^B \text{Int}) \implies \\ &(\mathcal{C}\{[s_1]\} \Downarrow n \implies \mathcal{C}\{[s_2]\} \Downarrow n) \wedge \\ &(\mathcal{C}\{[s_1]\} \Downarrow \text{blame} \implies \mathcal{C}\{[s_2]\} \Downarrow \text{blame}) \end{aligned}$$

The only difference is that now when a program containing  $s_1$  reduces to a value, so does one containing  $s_2$ .

From Conjecture 5.2, it follows that coherence holds between two representative translations of the same expression.

**Corollary 5.3** (Coherence for Representative Translations). *For any expression  $e$  such that  $\Psi \vdash^G e : \sigma \rightsquigarrow r_1$  and  $\Psi \vdash^G e : \sigma \rightsquigarrow r_2$ , we have  $\Psi \vdash r_1 \preceq_{ctx} r_2 : \sigma$ .*

We have proved that for every typing derivation, at least one representative translation exists.

**Lemma 5.4** (Representative Translation for Typing). *For any typing derivation  $\Psi \vdash^G e : \sigma$  there exists at least one representative translation  $r$  such that  $\Psi \vdash^G e : \sigma \rightsquigarrow r$ .*

For our example,  $(\lambda x : ?. (\langle \forall a. a \rightarrow \text{Int} \hookrightarrow \mathcal{G} \rightarrow \text{Int} \rangle f) (\langle ? \hookrightarrow \mathcal{G} \rangle x))$  is a representative translation, while the other two are not.

### 5.3 DYNAMIC GRADUAL GUARANTEE, RELOADED

Given the above propositions, we are ready to revisit the dynamic gradual guarantee. The nice thing about representative translations is that the dynamic gradual guarantee of our source language is essentially that of  $\lambda B$ , our target language. However, the dynamic gradual guarantee for  $\lambda B$  is still an open question. According to Igarashi et al. [2017], the difficulty lies in the definition of term precision that preserves the semantics. We leave it here as a conjecture as well. From a declarative point of view, we cannot prevent the system from picking undesirable instantiations, but we know that some choices are better than the others, so we can restrict the discussion of dynamic gradual guarantee to representative translations.

**Conjecture 5.5** (Dynamic Gradual Guarantee in terms of Representative Translations). *Suppose  $e' \sqsubseteq e$ ,*

1. *If  $\bullet \vdash^G e : \sigma \rightsquigarrow r$ ,  $\lceil r \rceil \Downarrow v$ , then for some  $\sigma_2$  and  $r'$ , we have  $\bullet \vdash^G e' : \sigma_2 \rightsquigarrow r'$ , and  $\sigma_2 \sqsubseteq \sigma$ , and  $\lceil r' \rceil \Downarrow v'$ , and  $v' \sqsubseteq v$ .*
2. *If  $\bullet \vdash^G e' : \sigma_2 \rightsquigarrow r'$ ,  $\lceil r' \rceil \Downarrow v'$ , then for some  $\sigma$  and  $r$ , we have  $\bullet \vdash^G e : \sigma \rightsquigarrow r$ , and  $\sigma_2 \sqsubseteq \sigma$ . Moreover,  $\lceil r \rceil \Downarrow v$  and  $v' \sqsubseteq v$ , or  $\lceil r \rceil \Downarrow \text{blame}$ .*

For the example in Section 4.3.3, now we know that the representative translation of the right one will evaluate to 1 as well.

$$(\lambda f : \forall a. a \rightarrow \text{Int}. \lambda x : \text{Int}. fx) (\lambda x : \text{Int}. 1) 3 \quad (\lambda f : \forall a. a \rightarrow \text{Int}. \lambda x : \text{Int}. fx) (\lambda x : ?. 1) 3$$

More importantly, in what follows, we show that our extended algorithm is able to find those representative translations.

### 5.4 EXTENDED ALGORITHMIC TYPE SYSTEM

To understand the design choices involved in the new algorithmic system, we consider the following algorithmic typing example:

$$f : \forall a. a \rightarrow \text{Int}, x : ? \vdash^G fx : \text{Int} \dashv f : \forall a. a \rightarrow \text{Int}, x : ?, \hat{\alpha}$$

Compared with the declarative typing, where we have many choices (e.g.,  $\text{Int} \rightarrow \text{Int}$ ,  $\text{Bool} \rightarrow \text{Int}$ , and so on) to instantiate  $\forall a. a \rightarrow \text{Int}$ , the algorithm computes the instantiation  $\hat{\alpha} \rightarrow \text{Int}$  with  $\hat{\alpha}$  unsolved in the output context. What can we know from the algorithmic typing? First we know that, here  $\hat{\alpha}$  is *not constrained* by the typing problem. Second, and more importantly,  $\hat{\alpha}$  has been compared with an unknown type (when typing  $(fx)$ ). Therefore, it is

Types	$\sigma ::= \text{Int} \mid a \mid \hat{\alpha} \mid \sigma_1 \rightarrow \sigma_2 \mid \forall a. \sigma \mid ? \mid \mathcal{S} \mid \mathcal{G}$
Monotypes	$\tau ::= \text{Int} \mid a \mid \hat{\alpha} \mid \tau_1 \rightarrow \tau_2 \mid \mathcal{S} \mid \mathcal{G}$
Existential variables	$\hat{\alpha} ::= \hat{\alpha}_S \mid \hat{\alpha}_G$
Castable Types	$\mathbb{C} ::= \text{Int} \mid a \mid \hat{\alpha} \mid \mathbb{C}_1 \rightarrow \mathbb{C}_2 \mid \forall a. \mathbb{C} \mid ? \mid \mathcal{G}$
Castable Monotypes	$t ::= \text{Int} \mid a \mid \hat{\alpha} \mid t_1 \rightarrow t_2 \mid \mathcal{G}$
Algorithmic Contexts	$\Gamma, \Delta, \Theta ::= \bullet \mid \Gamma, x : \sigma \mid \Gamma, a \mid \Gamma, \hat{\alpha} \mid \Gamma, \hat{\alpha}_S = \tau \mid \Gamma, \hat{\alpha}_G = t \mid \Gamma, \blacktriangleright \hat{\alpha}$
Complete Contexts	$\Omega ::= \bullet \mid \Omega, x : \sigma \mid \Omega, a \mid \Omega, \hat{\alpha}_S = \tau \mid \Omega, \hat{\alpha}_G = t \mid \Omega, \blacktriangleright \hat{\alpha}$

Figure 5.2: Syntax of types, contexts and consistent subtyping in the extended algorithmic system.

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; while the second kind (as in this example) are those where the only constraint is that *the variable was once compared with an unknown type* [Garcia and Cimini 2015].

The syntax of types is shown in Figure 5.2. A notable difference, apart from the addition of static and gradual parameters, is that we further split existential variables  $\hat{\alpha}$  into static existential variables  $\hat{\alpha}_S$  and gradual existential variables  $\hat{\alpha}_G$ . Depending on whether an existential variable has been compared with  $?$  or not, its solution space changes. More specifically, static existential variables can be solved to a monotype  $\tau$ , whereas gradual existential variables can only be solved to a castable monotype  $t$ , as can be seen in the changes of algorithmic contexts and complete contexts. As a result, the typing result for the above example now becomes

$$f : \forall a. a \rightarrow \text{Int}, x : ? \vdash^G f x : \text{Int} \dashv f : \forall a. a \rightarrow \text{Int}, x : ?, \hat{\alpha}_G$$

since we can solve any unconstrained  $\hat{\alpha}_G$  to  $\mathcal{G}$ , it is easy to verify that the resulting translation is indeed a representative translation.

Our extended algorithm is novel in the following aspects. We naturally extend the concept of existential variables [Dunfield and Krishnaswami 2013] to deal with comparisons between existential variables and unknown types. Unlike Garcia and Cimini [2015], where they use an extra set to store types that have been compared with unknown types, our two kinds of existential variables emphasize the type distinction better, and correspond more closely to the two kinds of type parameters, as we can solve  $\hat{\alpha}_S$  to  $\mathcal{S}$  and  $\hat{\alpha}_G$  to  $\mathcal{G}$ .

$\Gamma \vdash^G \sigma_1 \lesssim \sigma_2 \dashv \Delta$

*(Algorithmic Consistent Subtyping)*

GPC-AS-TVAR

$$\frac{}{\Gamma[a] \vdash^G a \lesssim a \dashv \Gamma[a]}$$

GPC-AS-INT

$$\frac{}{\Gamma \vdash^G \text{Int} \lesssim \text{Int} \dashv \Gamma}$$

GPC-AS-EVAR

$$\frac{}{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim \hat{\alpha} \dashv \Gamma[\hat{\alpha}]}$$

GPC-AS-SPAR

$$\frac{}{\Gamma \vdash^G \mathcal{S} \lesssim \mathcal{S} \dashv \Gamma}$$

GPC-AS-GPAR

$$\frac{}{\Gamma \vdash^G \mathcal{G} \lesssim \mathcal{G} \dashv \Gamma}$$

GPC-AS-UNKNOWNLL

$$\frac{}{\Gamma \vdash^G ? \lesssim \mathbb{C} \dashv \text{contaminate}(\Gamma, \mathbb{C})}$$

GPC-AS-UNKNOWNRR

$$\frac{}{\Gamma \vdash^G \mathbb{C} \lesssim ? \dashv \text{contaminate}(\Gamma, \mathbb{C})}$$

GPC-AS-ARROW

$$\frac{\Gamma \vdash^G \sigma_3 \lesssim \sigma_1 \dashv \Theta \quad \Theta \vdash^G [\Theta]\sigma_2 \lesssim [\Theta]\sigma_4 \dashv \Delta}{\Gamma \vdash^G \sigma_1 \rightarrow \sigma_2 \lesssim \sigma_3 \rightarrow \sigma_4 \dashv \Delta}$$

GPC-AS-FORALLR

$$\frac{\Gamma, a \vdash^G \sigma_1 \lesssim \sigma_2 \dashv \Delta, a, \Theta}{\Gamma \vdash^G \sigma_1 \lesssim \forall a. \sigma_2 \dashv \Delta}$$

GPC-AS-FORALLL

$$\frac{\Gamma, \blacktriangleright_{\hat{a}_S}, \hat{\alpha}_S \vdash^G \sigma_1[a \mapsto \hat{\alpha}_S] \lesssim \sigma_2 \dashv \Delta, \blacktriangleright_{\hat{a}_S}, \Theta}{\Gamma \vdash^G \forall a. \sigma_1 \lesssim \sigma_2 \dashv \Delta}$$

GPC-AS-INSTL

$$\frac{\hat{\alpha} \notin \text{FV}(\sigma) \quad \Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta}{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta}$$

GPC-AS-INSTR

$$\frac{\hat{\alpha} \notin \text{FV}(\sigma) \quad \Gamma[\hat{\alpha}] \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta}{\Gamma[\hat{\alpha}] \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta}$$

Figure 5.3: Extended algorithmic consistent subtyping

## 5.4.1 EXTENDED ALGORITHMIC CONSISTENT SUBTYPING

While the changes in the syntax seem negligible, the addition of static and gradual type parameters changes the algorithmic judgments in a significant way. We first discuss the algorithmic consistent subtyping, which is shown in Figure 5.3. For notational convenience, when static and gradual existential variables have the same rule form, we compress them into one rule. For example, rule **GPC-AS-EVAR** is really two rules  $\Gamma[\hat{\alpha}_S] \vdash^G \hat{\alpha}_S \lesssim \hat{\alpha}_S \dashv \Gamma[\hat{\alpha}_S]$  and  $\Gamma[\hat{\alpha}_G] \vdash^G \hat{\alpha}_G \lesssim \hat{\alpha}_G \dashv \Gamma[\hat{\alpha}_G]$ ; same for rules **GPC-AS-INSTL** and **GPC-AS-INSTR**.

Rules **GPC-AS-SPAR** and **GPC-AS-GPAR** are direct analogies of rules **GPC-CS-SPAR** and **GPC-CS-GPAR**. Though looking simple, rules **GPC-AS-UNKNOWNLL** and **GPC-AS-UNKNOWNRR** deserve much explanation. To understand what the output context  $\text{contaminate}(\Gamma, \mathbb{C})$  is for, let us first see why this seemingly intuitive rule  $\Gamma \vdash^G ? \lesssim \mathbb{C} \dashv \Gamma$  (like rule **GPC-AS-UNKNOWNL** in the original algorithmic system) is wrong. Consider the judgment  $\hat{\alpha}_S \vdash^G ? \lesssim \hat{\alpha}_S \rightarrow \hat{\alpha}_S \dashv \hat{\alpha}_S$ , which seems fine. If this holds, then – since  $\hat{\alpha}_S$  is unsolved in the output context – we can solve it to  $\mathcal{S}$  for example (recall that  $\hat{\alpha}_S$  can be solved to some monotype), resulting in  $? \lesssim \mathcal{S} \rightarrow \mathcal{S}$ . However, this is in direct conflict with rule **GPC-CS-UNKNOWNLL** in the declarative system precisely because  $\mathcal{S} \rightarrow \mathcal{S}$  is not a castable type! A possible solution would be to transform all static existential variables to gradual existential variables within  $\mathbb{C}$  whenever it is being compared to  $?$ : while  $\hat{\alpha}_S \vdash^G ? \lesssim \hat{\alpha}_S \rightarrow \hat{\alpha}_S \dashv \hat{\alpha}_S$  does not hold,  $\hat{\alpha}_G \vdash^G ? \lesssim \hat{\alpha}_G \rightarrow \hat{\alpha}_G \dashv \hat{\alpha}_G$  does. While substituting static existential variables with gradual existential variables seems to be intuitively correct, it is rather hard to formulate—not only do we need to perform substitution in  $\mathbb{C}$ , we also need to substitute accordingly in both the input and output contexts in order to ensure that no existential variables become unbound. However, making such changes is at odds with the interpretation of input contexts: they are “input”, which evolve into output contexts with more variables solved. Therefore, in line with the use of input contexts, a simple solution is to generate a new gradual existential variable and solve the static existential variable to it in the output context, without touching  $\mathbb{C}$  at all. So we have  $\hat{\alpha}_S \vdash^G ? \lesssim \hat{\alpha}_S \rightarrow \hat{\alpha}_S \dashv \hat{\alpha}_G, \hat{\alpha}_S = \hat{\alpha}_G$ .

Based on the above discussion, the following defines  $\text{contaminate}(\Gamma, \sigma)$ :

**Definition 17.**  $\text{contaminate}(\Gamma, \sigma)$  is defined inductively as follows

---

$\text{contaminate}(\bullet, \sigma)$	$=$	$\bullet$
$\text{contaminate}((\Gamma, x : \sigma), \sigma)$	$=$	$\text{contaminate}(\Gamma, \sigma), x : \sigma$
$\text{contaminate}((\Gamma, a), \sigma)$	$=$	$\text{contaminate}(\Gamma, \sigma), a$
$\text{contaminate}((\Gamma, \hat{\alpha}_S), \sigma)$	$=$	$\text{contaminate}(\Gamma, \sigma), \hat{\alpha}_G, \hat{\alpha}_S = \hat{\alpha}_G$ if $\hat{\alpha}_S$ occurs in $\sigma$
$\text{contaminate}((\Gamma, \hat{\alpha}_S), \sigma)$	$=$	$\text{contaminate}(\Gamma, \sigma), \hat{\alpha}_S$ if $\hat{\alpha}_S$ does not occur in $\sigma$
$\text{contaminate}((\Gamma, \hat{\alpha}_G), \sigma)$	$=$	$\text{contaminate}(\Gamma, \sigma), \hat{\alpha}_G$
$\text{contaminate}((\Gamma, \hat{\alpha} = \tau), \sigma)$	$=$	$\text{contaminate}(\Gamma, \sigma), \hat{\alpha} = \tau$
$\text{contaminate}((\Gamma, \blacktriangleright \hat{\alpha}), \sigma)$	$=$	$\text{contaminate}(\Gamma, \sigma), \blacktriangleright \hat{\alpha}$

---

$\text{contaminate}(\Gamma, \sigma)$  solves all static existential variables found within  $\sigma$  to fresh gradual existential variables in  $\Gamma$ . Notice the case for  $\text{contaminate}((\Gamma, \hat{\alpha}_S), \sigma)$  is exactly what we have just described.

Rule **GPC-AS-FORALLL** is slightly different from rule **GPC-AS-FORALL** in the original algorithmic system in that we replace  $a$  with a new static existential variable  $\hat{\alpha}_S$ . Note that  $\hat{\alpha}_S$  might be solved to a gradual existential variable later. The rest of the rules are the same as those in the original system.

#### 5.4.2 EXTENDED INSTANTIATION

The instantiation judgments shown in Figure 5.4 also change significantly. The complication comes from the fact that now we have two different kinds of existential variables, and the relative order they appear in the context affects their solutions.

Rules **GPC-INSTL-SOLVE**S and **GPC-INSTL-SOLVE**G are the refinement to rule **GPC-INSTL-SOLVE** in the original system. The next two rules deal with situations where one side is an existential variable and the other side is an unknown type. Rule **GPC-INSTL-SOLVE**US is a special case of rule **GPC-AS-UNKNOWNRR** where we create a new gradual existential variable  $\hat{\alpha}_G$  and set the solution of  $\hat{\alpha}_S$  to be  $\hat{\alpha}_G$  in the output context. Rule **GPC-INSTL-SOLVE**UG is the same as rule **GPC-INSTL-SOLVE**U in the original system and simply propagates the input context. The next two rules **GPC-INSTL-REACHSG1** and **GPC-INSTL-REACHSG2** are a bit involved, but they both answer to the same question: how to solve a gradual existential variable when it is declared after some static existential variable. More concretely, in rule **GPC-INSTL-REACHSG1**, we feel that we need to solve  $\hat{\beta}_G$  to another existential variable. However, simply setting  $\hat{\beta}_G = \hat{\alpha}_S$  and leaving  $\hat{\alpha}_S$  untouched in the output context is wrong. The reason is that  $\hat{\beta}_G$  could be a gradual existential variable created by rule **GPC-AS-UNKNOWNLL**/**GPC-AS-UNKNOWNRR** and solving  $\hat{\beta}_G$  to a static existential variable would result in the same problem as we have discussed. Instead, we create another new gradual existential variable  $\hat{\alpha}_G$  and set

$$\boxed{\Gamma \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta}$$

(Instantiation I)

$$\frac{\text{GPC-INSTL-SOLVE}\mathbf{S} \quad \Gamma \vdash^G \tau}{\Gamma, \hat{\alpha}_S, \Gamma' \vdash^G \hat{\alpha}_S \lesssim \tau \dashv \Gamma, \hat{\alpha}_S = \tau, \Gamma'}$$

$$\frac{\text{GPC-INSTL-SOLVE}\mathbf{G} \quad \Gamma \vdash^G t}{\Gamma, \hat{\alpha}_G, \Gamma' \vdash^G \hat{\alpha}_G \lesssim t \dashv \Gamma, \hat{\alpha}_G = t, \Gamma'}$$

$$\frac{\text{GPC-INSTL-SOLVE}\mathbf{US}}{\Gamma[\hat{\alpha}_S] \vdash^G \hat{\alpha}_S \lesssim ? \dashv \Gamma[\hat{\alpha}_G, \hat{\alpha}_S = \hat{\alpha}_G]}$$

$$\frac{\text{GPC-INSTL-SOLVE}\mathbf{UG}}{\Gamma[\hat{\alpha}_G] \vdash^G \hat{\alpha}_G \lesssim ? \dashv \Gamma[\hat{\alpha}_G]}$$

$$\frac{\text{GPC-INSTL-REACH}\mathbf{SG1}}{\Gamma[\hat{\alpha}_S][\hat{\beta}_G] \vdash^G \hat{\alpha}_S \lesssim \hat{\beta}_G \dashv \Gamma[\hat{\alpha}_G, \hat{\alpha}_S = \hat{\alpha}_G][\hat{\beta}_G = \hat{\alpha}_G]}$$

$$\frac{\text{GPC-INSTL-REACH}\mathbf{SG2}}{\Gamma[\hat{\beta}_S][\hat{\alpha}_G] \vdash^G \hat{\alpha}_G \lesssim \hat{\beta}_S \dashv \Gamma[\hat{\beta}_G, \hat{\beta}_S = \hat{\beta}_G][\hat{\alpha}_G = \hat{\beta}_G]}$$

$$\frac{\text{GPC-INSTL-REACH}\mathbf{OTHER}}{\Gamma[\hat{\alpha}][\hat{\beta}] \vdash^G \hat{\alpha} \lesssim \hat{\beta} \dashv \Gamma[\hat{\alpha}][\hat{\beta} = \hat{\alpha}]}$$

$$\frac{\text{GPC-INSTL-FORALL}\mathbf{R} \quad \Gamma[\hat{\alpha}], b \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta, b, \Theta}{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim \forall b. \sigma \dashv \Delta}$$

$$\frac{\text{GPC-INSTL-ARR} \quad \Gamma[\hat{\alpha}_2, \hat{\alpha}_1, \hat{\alpha} = \hat{\alpha}_1 \rightarrow \hat{\alpha}_2] \vdash^G \sigma_1 \lesssim \hat{\alpha}_1 \dashv \Theta \quad \Theta \vdash^G \hat{\alpha}_2 \lesssim [\Theta]\sigma_2 \dashv \Delta}{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim \sigma_1 \rightarrow \sigma_2 \dashv \Delta}$$

$$\boxed{\Gamma \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta}$$

(Instantiation II, excerpt)

$$\frac{\text{GPC-INSTR-FORALL}\mathbf{LL} \quad \Gamma[\hat{\alpha}], \blacktriangleright_{\hat{b}_S}, \hat{\beta}_S \vdash^G \sigma[b \mapsto \hat{\beta}_S] \lesssim \hat{\alpha} \dashv \Delta, \blacktriangleright_{\hat{b}_S}, \Theta}{\Gamma[\hat{\alpha}] \vdash^G \forall b. \sigma \lesssim \hat{\alpha} \dashv \Delta}$$

Figure 5.4: Instantiation in the extended algorithmic system

the solutions of both  $\hat{\alpha}_S$  and  $\hat{\beta}_G$  to it; similarly in rule [GPC-INSTL-REACHSG2](#). Rule [GPC-INSTL-REACHOTHER](#) deals with the other cases (e.g.,  $\hat{\alpha}_S \lesssim \hat{\beta}_S$ ,  $\hat{\alpha}_G \lesssim \hat{\beta}_G$  and so on). In those cases, we employ the same strategy as in the original system.

As for the other instantiation judgment, most of the rules are symmetric and thus omitted. The only interesting rule is [GPC-ISTR-FORALLLL](#), which is similar to what we did for rule [GPC-AS-FORALLLL](#).

### 5.4.3 ALGORITHMIC TYPING AND METATHEORY

Fortunately, the changes in the algorithmic bidirectional system are minimal: we replace every existential variable with a static existential variable. Furthermore, we proved that the extended algorithmic system is sound and complete with respect to the extended declarative system.

DO WE REALLY NEED TYPE PARAMETERS IN THE ALGORITHMIC SYSTEM? As we mentioned earlier, type parameters in the declarative system are merely an analysis tool, and in practice, type parameters are inaccessible to programmers. For the sake of proving soundness and completeness, we have to endow the algorithmic system with type parameters. However, the algorithmic system already has static and gradual existential variables, which can serve the same purpose. In that regard, we could directly solve every *unsolved* static and gradual existential variable in the output context to `Int` and `?`, respectively.

### 5.4.4 RESTRICTED GENERALIZATION

In Section 4.3.2, we discussed the issue that the translation produces multiple target expressions due to the different choices for instantiations, and those translations have different dynamic semantics. Besides that, there is another cause for multiple translations: redundant generalization during translation by rule [GEN](#). Consider the simple expression  $(\lambda x : \text{Int}. x) 1$ , the following shows two possible translations:

- $\vdash (\lambda x : \text{Int}. x) 1 : \text{Int} \rightsquigarrow (\lambda x : \text{Int}. x) (\langle \text{Int} \hookrightarrow \text{Int} \rangle 1)$
- $\vdash (\lambda x : \text{Int}. x) 1 : \text{Int} \rightsquigarrow (\lambda x : \text{Int}. x) (\langle \forall a. \text{Int} \hookrightarrow \text{Int} \rangle (\Lambda a. 1))$

The difference comes from the fact that in the second translation, we apply rule [GEN](#) while typing `1` to get  $\bullet \vdash 1 : \forall a. \text{Int}$ . As a consequence, the translation of `1` is accompanied by a cast from  $\forall a. \text{Int}$  to `Int` since the former is a consistent subtype of the latter. This difference is harmless, because obviously these two expressions will reduce to the same value in  $\lambda B$ , thus



preserving coherence (up to cast error). While it is not going to break coherence, it does result in multiple representative translations for one expression (e.g., the above two translations are both the representative translations).

There are several ways to make the translation process more deterministic. For example, we can restrict generalization to happen only in let expressions and require let expressions to include annotations, as  $\text{let } x : \sigma = e_1 \text{ in } e_2$ . Another feasible option would be to give a declarative, bidirectional system as the specification (instead of the type assignment one), in the same spirit of DK [Dunfield and Krishnaswami 2013]. Then we can restrict generalization to be performed through annotations in checking mode.

With restricted generalization, we hypothesize that now each expression has exactly one representative translation (up to renaming of fresh type parameters). Instead of calling it a *representative* translation, we can say it is a *principal* translation. Of course the above is only a sketch; we have not defined the corresponding rules, nor studied metatheory.



## PART IV

### TYPE INFERENCE WITH PROMOTION



# 6 HIGHER-RANK TYPE INFERENCE WITH TYPE PROMOTION

Gundry et al. [2010] propose type inference in context as a general foundation for unification/type inference algorithms. The key idea is based on the notion of information increase. However, their semantic definition of information increase makes it hard to prove metatheory formally. Following their work, a more syntactic foundation for information increase is presented by DK (Dunfield and Krishnaswami [2013]) to deal with higher-rank polymorphism. However, the DK approach produces duplication and cannot be easily generalized to handle more complicated types.

In this section, we propose a strategy called *promotion* that helps resolve the dependency of existential variables in the framework of type inference in context. As an illustration, Section 6.2 applies promotion to the unification algorithm for simply typed lambda calculus. Section 6.3 further proposes *polymorphic promotion* to deal with subtyping for higher-rank polymorphism. Finally, we briefly discuss how to promote dependent types and gradual types in Section 6.4.

## 6.1 INTRODUCTION AND MOTIVATION

### 6.1.1 BACKGROUND: TYPE INFERENCE IN CONTEXTS

Gundry et al. [2010] model unification and type inference from a general perspective of information increase. The problem context is a ML-style polymorphic system, based on the invariant that types can only depend on bindings appearing earlier in the context. Besides contexts being ordered, a key insight of the approach lies in how to solve existential variables with other types. In particular, it requires to resolve the dependency between existential variables. Consider unifying  $\hat{\alpha}$  with  $\hat{\beta} \rightarrow \text{Int}$  under context  $\hat{\alpha}, \hat{\beta}$ . Here we cannot simply set  $\hat{\alpha}$  to  $\hat{\beta} \rightarrow \text{Int}$ , as  $\hat{\beta}$  is out of the scope of  $\hat{\alpha}$ . The way Gundry et al. [2010] solve this problem is to examine variables in the context from the tail to the head, *moving segments of context to the left if necessary*. The process finishes when the existential variable being unified is found. In this case, Gundry et al. [2010] would return a solution context  $\hat{\beta}, \hat{\alpha} = \hat{\beta} \rightarrow \text{Int}$ ,

where  $\hat{\beta}$  is moved to the front of  $\hat{\alpha}$ . However, while moving type variables around is a feasible way to resolve the dependency between existential variables, the unpredictable context movements make the information increase hard to formalize and reason about. In their system, the information increase of contexts is defined in a semantic way: a context  $\Gamma_1$  is more informative than another context  $\Gamma_2$ , if there exists a substitution such that every item in  $\Gamma_2$  is, after context substitution, well-formed under  $\Gamma_1$ . This semantic definition makes it hard to prove metatheory formally, especially when advanced features are involved.

DK [Dunfield and Krishnaswami 2013] also uses ordered contexts as the input and output of the type inference algorithm for a higher-rank polymorphic type system (Section 2.3). Unlike Gundry et al. [2010], DK does it in a more syntactic way. In its instantiation rules, it decomposes type constructs so that unification between existential variables can only happen between single variables, which can then be solved by setting the one that appears later to the one that appears earlier. This way, the information increase of contexts is modeled as the intuitive and syntactic definition of *context extension* ( $\Gamma \longrightarrow \Delta$ ), which allows for inductive reasoning and proofs. This approach is adopted in GPC (Chapter 4), so let us consider how DK works in terms of the instantiation rules in GPC. Specifically, consider the derivation of  $\hat{\alpha}, \hat{\beta} \vdash^G \hat{\alpha} \lesssim \hat{\beta} \rightarrow \text{Int}$  in GPC:

$$\Delta = \hat{\alpha}_1, \hat{\alpha}_2, \hat{\alpha} = \hat{\alpha}_1 \rightarrow \hat{\alpha}_2, \hat{\beta}$$

$$\frac{\text{GPC-INSTL-REACH} \quad \Delta \vdash^G \hat{\beta} \lesssim \hat{\alpha}_1 \dashv \Delta[\hat{\beta} = \hat{\alpha}_1] \quad \text{GPC-INSTR-SOLVE} \quad \Delta[\hat{\beta} = \hat{\alpha}_1] \vdash^G \hat{\alpha}_2 \lesssim \text{Int} \dashv \Delta[\hat{\alpha}_2 = \text{Int}][\hat{\beta} = \hat{\alpha}_1]}{\hat{\alpha}, \hat{\beta} \vdash^G \hat{\alpha} \lesssim \hat{\beta} \rightarrow \text{Int} \dashv \Delta[\hat{\alpha}_2 = \text{Int}][\hat{\beta} = \hat{\alpha}_1]} \text{GPC-INSTL-ARR}$$

By rule **GPC-INSTL-ARR**, the variable  $\hat{\alpha}$  is solved by an arrow type  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$  consisting of two fresh existential variables. The two variables  $\hat{\alpha}_1$  and  $\hat{\alpha}_2$  are then instantiated with  $\hat{\beta}$  and  $\text{Int}$ , respectively. By rule **GPC-INSTL-REACH**, the variable  $\hat{\beta}$  is solved by  $\hat{\alpha}_1$ , as  $\hat{\alpha}_1$  appears in the context before  $\hat{\beta}$ , or otherwise we need to apply rule **GPC-INSTR-REACH** to set  $\hat{\alpha}_1$  to  $\hat{\beta}$  instead. The final solution context is  $\Delta[\hat{\alpha}_2 = \text{Int}][\hat{\beta} = \hat{\alpha}_1] = \hat{\alpha}_1, \hat{\alpha}_2 = \text{Int}, \hat{\alpha} = \hat{\alpha}_1 \rightarrow \hat{\alpha}_2, \hat{\beta} = \hat{\alpha}_1$ .

**CHALLENGES.** However, while the approach of decomposing type constructs works perfectly for DK and GPC, it has two drawbacks. The first drawback is that it produces duplication, in order to deal with both the case when the existential variable appears on the left ( $\Gamma \vdash^G \hat{\alpha} \lesssim \sigma \dashv \Delta$  in Figure 4.11) and the case when it appears on the right ( $\Gamma \vdash^G \sigma \lesssim \hat{\alpha} \dashv \Delta$  in Figure 4.11). For example, rule **GPC-INSTL-ARR** has its symmetric counterpart rule **GPC-INSTR-ARR**:

$$\begin{array}{c}
\text{GPC-INSTL-ARR} \\
\frac{\Gamma[\hat{\alpha}_2, \hat{\alpha}_1, \hat{\alpha} = \hat{\alpha}_1 \rightarrow \hat{\alpha}_2] \vdash^G \sigma_1 \lesssim \hat{\alpha}_1 \dashv \Theta \quad \Theta \vdash^G \hat{\alpha}_2 \lesssim [\Theta]\sigma_2 \dashv \Delta}{\Gamma[\hat{\alpha}] \vdash^G \hat{\alpha} \lesssim \sigma_1 \rightarrow \sigma_2 \dashv \Delta} \\
\\
\text{GPC-INSTR-ARR} \\
\frac{\Gamma[\hat{\alpha}_2, \hat{\alpha}_1, \hat{\alpha} = \hat{\alpha}_1 \rightarrow \hat{\alpha}_2] \vdash^G \hat{\alpha}_1 \lesssim \sigma_1 \dashv \Theta \quad \Theta \vdash^G [\Theta]\sigma_2 \lesssim \hat{\alpha}_2 \dashv \Delta}{\Gamma[\hat{\alpha}] \vdash^G \sigma_1 \rightarrow \sigma_2 \lesssim \hat{\alpha} \dashv \Delta}
\end{array}$$

Worse, this kind of “duplication” would scale up with the number of type constructs in the system.

Second, while decomposition works for function types, it may not work easily for more complicated types, e.g., dependent types. For example, consider that under the context  $\hat{\alpha}, \hat{\beta}$ , we want to instantiate  $\hat{\alpha}$  with a dependent type  $\Pi a : \hat{\beta}.a$ . Here because  $\hat{\beta}$  appears after  $\hat{\alpha}$ , we cannot directly set  $\hat{\alpha} = \Pi a : \hat{\beta}.a$  which is ill typed. However, if we try to decompose the type  $\Pi a : \hat{\beta}.a$  like in rule **GPC-INSTL-ARR**, it is obvious that  $\hat{\alpha}_2$  should be solved by  $a$ . Then, in order to make the solution well typed, we need to put  $a$  in the front of  $\hat{\alpha}_2$  in the context. However, this means that  $a$  would remain in the context, and it would be available for any later existential variables that should not have access to  $a$ .

### 6.1.2 OUR APPROACH: TYPE PROMOTION

We propose the *promotion* process, which helps resolve the dependency between existential variables. Promotion combines the advantages of Gundry et al. [2010] and DK: it is a simple and predictable process, so that information increase can still be modeled as the syntactic context extension; moreover, it does not cause any duplication.

To understand how promotion works, let us consider again the example  $\hat{\alpha}, \hat{\beta} \vdash^G \hat{\alpha} \lesssim \hat{\beta} \rightarrow \text{Int}$ . The problem here is that  $\hat{\beta}$  is out of the scope of  $\hat{\alpha}$  so we cannot directly set  $\hat{\alpha} = \hat{\beta} \rightarrow \text{Int}$ . Therefore, we first *promote* the type  $\hat{\beta} \rightarrow \text{Int}$ . At a high level, the promotion process looks for free existential variables in the type, and solves those out-of-scope existential variables with fresh existential variables added to the front of  $\hat{\alpha}$ , such that existential variables in the promoted type are all in the scope of  $\hat{\alpha}$ . In this case, we will solve  $\hat{\beta}$  with a fresh variable  $\hat{\alpha}_1$ , producing the context  $\hat{\alpha}_1, \hat{\alpha}, \hat{\beta} = \hat{\alpha}_1$ . Notice that  $\hat{\alpha}_1$  is inserted right before  $\hat{\alpha}$ . Now the instantiation example becomes  $\hat{\alpha}_1, \hat{\alpha}, \hat{\beta} = \hat{\alpha}_1 \vdash^G \hat{\alpha} \lesssim \hat{\alpha}_1 \rightarrow \text{Int}$ , and  $\hat{\alpha}_1 \rightarrow \text{Int}$  is a valid solution for  $\hat{\alpha}$ . Therefore, we get a final solution context  $\hat{\alpha}_1, \hat{\alpha} = \hat{\alpha}_1 \rightarrow \text{Int}, \hat{\beta} = \hat{\alpha}_1$ . Comparing the result with the solution context we get from DK ( $\hat{\alpha}_1, \hat{\alpha}_2 = \text{Int}, \hat{\alpha} = \hat{\alpha}_1 \rightarrow \hat{\alpha}_2, \hat{\beta} = \hat{\alpha}_1$ ), it is obvious that these two solutions are equivalent up to substitution.

INTERPRETATION OF PROMOTION. The approach taken by Gundry et al. [2010] and the approach used by DK are based on the same observation: *the relative order between existential variables does not matter for solving a constraint*. The promotion process captures precisely this observation. Its task is to “move” existential variables to suitable positions *indirectly*, by solving those out-of-scope existential variables with fresh in-scope existential variables.

This seems to go against the design principle that the contexts are ordered. However, ordering is still important for variables whose order matters. For instance, for polymorphic types, the order between existential variables  $\hat{\alpha}$  and type variables  $a$  is important, so we cannot set  $\hat{\alpha}$  to  $a$  under the context  $(\hat{\alpha}, a)$  as  $a$  is not in the scope of  $\hat{\alpha}$ . Moreover, ordering prevents invalid cyclic contexts, e.g.,  $\hat{\alpha} = \hat{\beta} \rightarrow \text{Int}, \hat{\beta} = \hat{\alpha} \rightarrow \text{Int}$ .

UNIFICATION FOR SIMPLY TYPED LAMBDA CALCULUS. As a first illustration of the utility of the promotion process, Section 6.2 recasts the unification process for simply typed lambda calculus (STLC) using the promotion process. This system illustrates the key idea of promotion.

### 6.1.3 POLYMORPHIC PROMOTION

Instead of unification, the instantiation relation in DK actually deals with the polymorphic subtyping relation between existential variables and other types. The promotion process we described above only works for unification. In this section, we discuss promotion for polymorphic subtyping.

The difficulty of subtyping is that it needs to take unification into account at the same time. For example, given that  $\hat{\alpha}$  is a subtype of  $\text{Int}$ , the only possible solution is  $\hat{\alpha} = \text{Int}$ . Now consider  $\hat{\alpha} \vdash \forall a. a \rightarrow a <: \hat{\alpha}$ . How can we promote the polymorphic type  $\forall a. a \rightarrow a$  into a monotype which can serve as a valid solution for  $\hat{\alpha}$ ? One possible answer is to set  $\hat{\alpha} = \text{Int} \rightarrow \text{Int}$ , or  $\hat{\alpha} = \text{Bool} \rightarrow \text{Bool}$ . In fact, the most general solution for this subtyping problem is  $\hat{\alpha} = \hat{\beta} \rightarrow \hat{\beta}$  with fresh  $\hat{\beta}$ . Namely, we remove the universal quantifier in  $\forall a. a \rightarrow a$  and replace the variable  $a$  with a fresh existential variable  $\hat{\beta}$  added to the front of  $\hat{\alpha}$ , resulting in the solution context  $\hat{\beta}, \hat{\alpha} = \hat{\beta} \rightarrow \hat{\beta}$ .

On the other hand, how can we promote the type  $\forall a. a \rightarrow a$  in  $\hat{\alpha} \vdash \hat{\alpha} <: \forall a. a \rightarrow a$ ? It turns out that this subtyping is actually unsolvable. Therefore, in this case, promoting  $\forall a. a \rightarrow a$  will directly add the type variable  $a$  to the tail of the context to promote  $a \rightarrow a$ . Since  $a$  is added to the tail, it means that  $a$  is out of the scope of  $\hat{\alpha}$  and promoting  $a \rightarrow a$  would fail, which is exactly what we want. In fact, the promotion would succeed if only the variable is not used in the body of the polymorphic type. For example,  $\forall a. \text{Int} \rightarrow \text{Int}$  can be promoted to  $\text{Int} \rightarrow \text{Int}$ , which is a valid solution of  $\hat{\alpha}$  in  $\hat{\alpha} \vdash \hat{\alpha} <: \forall a. \text{Int} \rightarrow \text{Int}$ .



From these observations, we extend promotion to *polymorphic promotion*, which is able to resolve the polymorphic subtyping relation for existential variables. Depending on whether the existential variable appears on the right or left, polymorphic promotion has two modes, which we call the *contravariant mode* and the *covariant mode* respectively.

The contravariant mode promotes types as  $\forall a. a \rightarrow a$  in the case of  $\hat{\alpha} \vdash \forall a. a \rightarrow a <: \hat{\alpha}$ , where the universal quantifier is removed and the type variable  $a$  is replaced by a fresh existential variable added to front of the existential variable being solved. This corresponds to rule [GPC-ISTR-FORALLL](#), except that with promotion, the new existential variable  $\hat{\beta}$  (in rule [GPC-ISTR-FORALLL](#)) will be added directly before  $\hat{\alpha}$  and there is no need to create a marker or to discard the context after  $\hat{\beta}$  anymore.

The covariant mode promotes types as  $\forall a. a \rightarrow a$  in the case of  $\hat{\alpha} \vdash \hat{\alpha} <: \forall a. a \rightarrow a$ . In this case, promoting  $\forall a. a \rightarrow a$  will directly add the type variable  $a$  to the tail of the context, which corresponds to rule [GPC-ISTR-FORALLR](#). Since the type variable is out of the scope of the existential variable being solved, and promotion will succeed if only the variable is not used in the body of the polymorphic type.

While promoting polymorphic types behaves differently according to the mode, the mode does not matter for monotypes, as in both  $\hat{\alpha} <: \text{Int}$  and  $\text{Int} <: \hat{\alpha}$ ,  $\hat{\alpha} = \text{Int}$  would be the only solution. Since function types are contravariant in codomains and covariant in domains, promoting a function type under a certain mode proceeds to promote its codomain under the other mode and promote its domain under the original mode. For example,  $\hat{\alpha} = (\hat{\beta} \rightarrow \hat{\beta}) \rightarrow (\text{Int} \rightarrow \text{Int})$  is a solution for  $\hat{\alpha} \vdash \hat{\alpha} <: (\forall a. a \rightarrow a) \rightarrow (\forall a. \text{Int} \rightarrow \text{Int})$ , where  $(\forall a. a \rightarrow a) \rightarrow (\forall a. \text{Int} \rightarrow \text{Int})$  is promoted under the covariant mode, which means  $\forall a. a \rightarrow a$  is promoted under the *contravariant* mode and  $\forall a. \text{Int} \rightarrow \text{Int}$  is promoted under the original covariant mode.

**POLYMORPHIC PROMOTION FOR SUBTYPING.** We illustrate polymorphic promotion by showing that the original instantiation relationship in DK can be replaced by our polymorphic promotion process. Furthermore, we show that subtyping, which was built upon instantiation but now uses polymorphic promotion, remains sound and complete.

## 6.2 UNIFICATION FOR SIMPLY TYPED LAMBDA CALCULUS

This section first introduces simple typed lambda calculus, and then present a unification algorithm that uses the novel promotion mechanism.

Monotypes	$\tau ::= \text{Int} \mid \tau_1 \rightarrow \tau_2 \mid \hat{\alpha}$
Algorithmic Contexts	$\Gamma, \Delta, \Theta ::= \bullet \mid \Gamma, \hat{\alpha} \mid \Gamma, \hat{\alpha} = \tau$
Complete Contexts	$\Omega ::= \bullet \mid \Omega, \hat{\alpha} = \tau$

  

$\Gamma \Vdash \tau_1 \approx \tau_2 \dashv \Delta$

(Unification)

  

$$\frac{\text{U-REFL}}{\Gamma \Vdash \tau \approx \tau \dashv \Gamma}$$

$$\frac{\text{U-ARROW} \quad \Gamma \Vdash \tau_1 \approx \tau_3 \dashv \Theta \quad \Theta \Vdash [\Theta]\tau_2 \approx [\Theta]\tau_4 \dashv \Delta}{\Gamma \Vdash \tau_1 \rightarrow \tau_2 \approx \tau_3 \rightarrow \tau_4 \dashv \Delta}$$

  

$$\frac{\text{U-EVARL} \quad \Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \tau_1 \rightsquigarrow \tau_2 \dashv \Delta[\hat{\alpha}]}{\Gamma \Vdash \hat{\alpha} \approx \tau_1 \dashv \Delta[\hat{\alpha} = \tau_2]}$$

$$\frac{\text{U-EVARR} \quad \Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \tau_1 \rightsquigarrow \tau_2 \dashv \Delta[\hat{\alpha}]}{\Gamma \Vdash \tau_1 \approx \hat{\alpha} \dashv \Delta[\hat{\alpha} = \tau_2]}$$

  

$\Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \tau_1 \rightsquigarrow \tau_2 \dashv \Delta$

(Promotion)

  

$$\frac{\text{PR-INT}}{\Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \text{Int} \rightsquigarrow \text{Int} \dashv \Gamma}$$

$$\frac{\text{PR-ARROW} \quad \Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \tau_1 \rightsquigarrow \tau_3 \dashv \Theta \quad \Theta \vdash_{\hat{\alpha}}^{\text{pr}} [\Theta]\tau_2 \rightsquigarrow \tau_4 \dashv \Delta}{\Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \tau_1 \rightarrow \tau_2 \rightsquigarrow \tau_3 \rightarrow \tau_4 \dashv \Delta}$$

  

$$\frac{\text{PR-EVARL}}{\Gamma[\hat{\beta}][\hat{\alpha}] \vdash_{\hat{\alpha}}^{\text{pr}} \hat{\beta} \rightsquigarrow \hat{\beta} \dashv \Gamma[\hat{\beta}][\hat{\alpha}]}$$

$$\frac{\text{PR-EVARR}}{\Gamma[\hat{\alpha}][\hat{\beta}] \vdash_{\hat{\alpha}}^{\text{pr}} \hat{\beta} \rightsquigarrow \hat{\beta}_1 \dashv \Gamma[\hat{\beta}_1, \hat{\alpha}][\hat{\beta}]}$$

Figure 6.1: Types, contexts, unification and promotion of algorithmic STLC

### 6.2.1 DECLARATIVE SYSTEM

The definition of types in STLC is given below. We have only monotypes  $\tau$ , which includes the integer type  $\text{Int}$  and function types  $\tau_1 \rightarrow \tau_2$ . In this section, we focus on the unification process. Hence, we do not elaborate the details of expressions' syntax or typing rules.

Monotypes	$\tau ::= \text{Int} \mid \tau_1 \rightarrow \tau_2$
-----------	--

### 6.2.2 ALGORITHMIC SYSTEM

The syntax of the algorithmic system is given in Figure 6.1. Following DK [Dunfield and Krishnaswami 2013] and GPC, algorithmic monotypes include existential type variables  $\hat{\alpha}$ . Algorithmic contexts also contain declarations of existential type variables, either unsolved ( $\hat{\alpha}$ ) or solved ( $\hat{\alpha} = \tau$ ). Complete contexts  $\Omega$  contain only solved variables. We use the judgment  $\Gamma \Vdash^{\text{wf}} \tau$  to indicate that all existential variables in  $\tau$  are well-scoped. Its definition

is standard and thus omitted. We also use  $\Gamma \longrightarrow \Delta$  for context extension, whose definition is essentially a simplified version of the one in GPC (Section 4.4.5).

**UNIFICATION.** Figure 6.1 defines the unification process. The judgment  $\Gamma \vdash \tau_1 \approx \tau_2 \dashv \Delta$  reads that under the input context  $\Gamma$ , unifying  $\tau_1$  with  $\tau_2$  results in the output context  $\Delta$ . Rule **U-REFL** is our base case, and rule **U-ARROW** unifies the components of the arrow types. When unifying  $\hat{\alpha} \approx \tau_1$  (rule **U-EVARL**), we cannot simply set  $\hat{\alpha}$  to  $\tau_1$ , as  $\tau_1$  might include variables bound to the right of  $\hat{\alpha}$ . Instead, we need to *promote* ( $\vdash^{\text{pr}}$ )  $\tau_1$ . After promoting  $\tau_1$  to  $\tau_2$ , we can directly set  $\hat{\alpha} = \tau_2$ . Rule **U-VARR** deals with the case when  $\hat{\alpha}$  appears on the right.

**PROMOTION.** The promotion relation  $\Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \tau_1 \rightsquigarrow \tau_2 \dashv \Delta$  given at the bottom of Figure 6.1 reads that under the input context  $\Delta$ , promoting type  $\tau_1$  yields type  $\tau_2$ , so that  $\tau_2$  is well-formed in the prefix context of  $\hat{\alpha}$ , while retaining  $[\Delta]\tau_1 = [\Delta]\tau_2$ . At a high-level,  $\vdash^{\text{pr}}$  looks for free variables in  $\tau_1$ . Integers are always well-formed (rule **PR-INT**). Promoting a function recursively promotes its components (rule **PR-ARROW**). Variables bound to the left of  $\hat{\alpha}$  in  $\Gamma$  are unaffected (rule **PR-EVARL**), as they are already well-formed. In rule **PR-EVARR**, a unification variable  $\hat{\beta}$  bound to the right of  $\hat{\alpha}$  in  $\Gamma$  is replaced by a fresh variable introduced to  $\hat{\alpha}$ 's left. Promotion is a partial operation, as it requires  $\hat{\beta}$  either to be to the right or to the left of  $\hat{\alpha}$ . There is yet another possibility: if  $\hat{\beta} = \hat{\alpha}$ , then no rule applies. This is a desired property, as the  $\hat{\beta} = \hat{\alpha}$  case exactly corresponds to the “occurs-check” in a more typical presentation of unification. By preventing promoting  $\hat{\alpha}$  to the left of  $\hat{\alpha}$ , we prevent the possibility of an infinite substitution when applying an algorithmic context. Note that rule **U-REFL** solves the unification case  $\hat{\alpha} \approx \hat{\alpha}$ .

**EXAMPLE.** Below we give the derivation of  $\hat{\alpha}, \hat{\beta} \vdash \hat{\alpha} \approx \hat{\beta} \rightarrow \text{Int}$  discussed in Section 6.1.1.

$$\begin{array}{c}
 \text{PR-EVARL} \qquad \qquad \qquad \text{PR-INT} \\
 \hline
 \hat{\alpha}, \hat{\beta} \vdash_{\hat{\alpha}}^{\text{pr}} \hat{\beta} \rightsquigarrow \hat{\alpha} \dashv \hat{\alpha}_1, \hat{\alpha}, \hat{\beta} = \hat{\alpha}_1 \quad \hat{\alpha}_1, \hat{\alpha}, \hat{\beta} = \hat{\alpha}_1 \vdash_{\hat{\alpha}}^{\text{pr}} \text{Int} \rightsquigarrow \text{Int} \dashv \hat{\alpha}_1, \hat{\alpha}, \hat{\beta} = \hat{\alpha}_1 \\
 \hline
 \hat{\alpha}, \hat{\beta} \vdash_{\hat{\alpha}}^{\text{pr}} \hat{\beta} \rightarrow \text{Int} \rightsquigarrow \hat{\alpha}_1 \rightarrow \text{Int} \dashv \hat{\alpha}_1, \hat{\alpha}, \hat{\beta} = \hat{\alpha}_1 \quad \text{PR-ARROW} \\
 \hline
 \hat{\alpha}, \hat{\beta} \vdash \hat{\alpha} \approx \hat{\beta} \rightarrow \text{Int} \dashv \hat{\alpha}_1, \hat{\alpha} = \hat{\alpha}_1 \rightarrow \text{Int}, \hat{\beta} = \hat{\alpha}_1 \quad \text{U-EVAL-R}
 \end{array}$$

### 6.2.3 SOUNDNESS AND COMPLETENESS

We prove that our type promotion strategy and the unification algorithm are sound. First, we show that except for resolving the order problem, promotion will not change the type.

Namely, the input type and the output type are equivalent after substitution by the output context. Moreover, the promoted type is well-formed under the prefix context of  $\hat{\alpha}$ .

**Theorem 6.1** (Soundness of Promotion). *If  $\Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \tau_1 \rightsquigarrow \tau_2 \dashv \Delta$ , then  $[\Delta]\tau_1 = [\Delta]\tau_2$ . Moreover, given  $\Delta = \Delta_1, \hat{\alpha}, \Delta_2$ , we have  $\Delta_1 \vdash^{\text{wf}} \tau_2$ ,*

With soundness of promotion, we can prove that the unification algorithm is also sound:

**Theorem 6.2** (Soundness of Unification). *If  $\Gamma \vdash \tau_1 \approx \tau_2 \dashv \Delta$ , then  $[\Delta]\tau_1 = [\Delta]\tau_2$ .*

We can further prove that promotion is complete using the notion of context extension. Note that in the completeness statement we require  $\hat{\alpha} \notin \text{FV}(\tau_1)$ , or otherwise promotion would fail.

**Theorem 6.3** (Completeness of Promotion). *Given  $\Gamma \longrightarrow \Omega$ , and  $\Gamma \vdash^{\text{wf}} \hat{\alpha}$ , and  $\Gamma \vdash^{\text{wf}} \tau_1$ , and  $[\Gamma]\hat{\alpha} = \hat{\alpha}$ , and  $[\Gamma]\tau_1 = \tau_1$ , if  $\hat{\alpha} \notin \text{FV}(\tau_1)$ , there exist  $\tau_2, \Delta$  and  $\Omega'$  such that  $\Gamma \longrightarrow \Omega'$  and  $\Omega \longrightarrow \Omega'$  and  $\Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \tau_1 \rightsquigarrow \tau_2 \dashv \Delta$ .*

The completeness of unification is then built upon the completeness of promotion.

**Theorem 6.4** (Completeness of Unification). *Given  $\Gamma \longrightarrow \Omega$ , and  $\Gamma \vdash^{\text{wf}} \tau_1$ , and  $\Gamma \vdash^{\text{wf}} \tau_2$ , and  $[\Gamma]\tau_1 = \tau_1$ , and  $[\Gamma]\tau_2 = \tau_2$ , if  $[\Omega]\tau_1 = [\Omega]\tau_2$ , there exist  $\Delta$  and  $\Omega'$  such that  $\Gamma \longrightarrow \Omega'$  and  $\Omega \longrightarrow \Omega'$  and  $\Gamma \vdash \tau_1 \approx \tau_2 \dashv \Delta$ .*

### 6.3 SUBTYPING FOR HIGHER-RANK POLYMORPHISM

In this section, we adopt the type promotion strategy to a higher-rank polymorphic type system from DK [Dunfield and Krishnaswami 2013]. We show that promotion can be further extended to polymorphic promotion to deal with subtyping, which can be used to replace the instantiation relation in the original DK system while preserving completeness and soundness.

#### 6.3.1 DECLARATIVE SYSTEM

The definition of types in DK (Figure 2.6 in Section 2.3.2) is repeated below. Comparing to STLC, we have polymorphic types  $\forall a. \sigma$  and type variables  $a$ . Again, we omit the details about expressions since we focus on types in this section. Recall that DK shares the same subtyping relation as of OL (Figure 2.5), and we use the judgment  $\Psi \vdash^{DK} \sigma_1 <: \sigma_2$  to denote the subtyping relation in DK.

Types	$\sigma$	$::=$	$\text{Int} \mid a \mid \sigma_1 \rightarrow \sigma_2 \mid \forall a. \sigma$
Monotypes	$\tau$	$::=$	$\text{Int} \mid a \mid \tau_1 \rightarrow \tau_2$
Contexts	$\Psi$	$::=$	$\bullet \mid \Psi, x : \sigma \mid \Psi, a$

## 6.3.2 ALGORITHMIC SYSTEM

The syntax of the algorithmic system is given in Figure 6.2. The promotion mode  $\otimes$  is either covariant(+) or contravariant(-). We can use  $-\otimes$  to flip the promotion mode. Specifically,

$$\begin{aligned} -(+) &= - \\ -(-) &= + \end{aligned}$$

**SUBTYPING.** Figure 6.2 also includes the subtyping judgment  $\Gamma \vdash^{\text{sub}} \sigma_1 <: \sigma_2 \dashv \Delta$ , which reads that, under input context  $\Gamma$ , type  $\sigma_1$  is a subtype of  $\sigma_2$ , with output context  $\Delta$ . The rules except the last two are the same as the algorithmic subtyping rules in DK.

Rule **s-INSTL** and rule **INSTR** are specific to this system. Recall that in GPC (which follows DK), the consistent subtyping between  $\hat{\alpha}$  and  $\sigma$  relies on the instantiation rules, which are duplicated for the case when  $\hat{\alpha}$  is on the left and the case when  $\hat{\alpha}$  is on the right. Here, instead of instantiation, we directly use polymorphic promotion to promote the possibly polymorphic type  $\sigma$  into a monotype  $\tau$ . Specifically, rule **s-INSTL** uses polymorphic promotion under the covariant mode (+) and rule **s-INSTR** uses polymorphic promotion under the contravariant mode (-). If promotion succeeds, we can directly set  $\hat{\alpha}$  to  $\tau$ .

**POLYMORPHIC PROMOTION.** The judgment  $\Gamma \vdash_{\hat{\alpha}}^{\otimes} \sigma \rightsquigarrow \tau \dashv \Delta$  reads that under the input context  $\Gamma$ , promoting  $\sigma$  under promotion mode  $\otimes$  yields type  $\tau$ , so that  $\tau$  is well-formed in the prefix context of  $\hat{\alpha}$ .

The only difference between these two promotion modes is how to promote polymorphic types. Under the contravariant mode (rule **P-PR-FORALL**), a monotype would make the final type more polymorphic. Therefore, we replace the universal binder  $a$  with a fresh existential variable  $\hat{a}$  and put it before  $\hat{\alpha}$ . Otherwise, in rule **P-PR-FORALLR**, we put  $a$  in the context and promote  $\sigma$ . Notice that since  $a$  is added to the tail of the context, it is not in the scope of  $\sigma$  and can actually never be used in  $\sigma$  or otherwise promotion would fail. We then discard  $a$  in the return context. Note that we can simplify the rule by directly requiring  $a \notin \text{FV}(\sigma)$ , as in rule **P-PR-FORALLRR** given below. This way we would not need to add  $a$  into the context and the rule would remain sound.

Types	$\sigma ::= \text{Int} \mid a \mid \sigma_1 \rightarrow \sigma_2 \mid \forall a. \sigma \mid \hat{\alpha}$
Monotypes	$\tau ::= \text{Int} \mid a \mid \hat{\alpha} \mid \tau_1 \rightarrow \tau_2$
Algorithmic Contexts	$\Gamma, \Delta, \Theta ::= \bullet \mid \Gamma, a \mid \Gamma, \hat{\alpha} \mid \Gamma, \hat{\alpha} = \tau$
Complete Contexts	$\Omega ::= \bullet \mid \Omega, a \mid \Omega, \hat{\alpha} = \tau$
Promotion Modes	$\otimes ::= + \mid -$

  

$\Gamma \vdash^{\text{sub}} \sigma_1 <: \sigma_2 \dashv \Delta$

(Subtyping)

  

<b>S-TVAR</b> $\frac{}{\Gamma[a] \vdash^{\text{sub}} a <: a \dashv \Gamma[a]}$	<b>S-INT</b> $\frac{}{\Gamma \vdash^{\text{sub}} \text{Int} <: \text{Int} \dashv \Gamma}$	<b>S-EVAR</b> $\frac{}{\Gamma[\hat{\alpha}] \vdash^{\text{sub}} \hat{\alpha} <: \hat{\alpha} \dashv \Gamma[\hat{\alpha}]}$
<b>S-ARROW</b> $\frac{\Gamma \vdash^{\text{sub}} \sigma_3 <: \sigma_1 \dashv \Theta \quad \Theta \vdash^{\text{sub}} [\Theta]\sigma_2 <: [\Theta]\sigma_4 \dashv \Delta}{\Gamma \vdash^{\text{sub}} \sigma_1 \rightarrow \sigma_2 <: \sigma_3 \rightarrow \sigma_4 \dashv \Delta}$	<b>S-FORALLR</b> $\frac{\Gamma, a \vdash^{\text{sub}} \sigma_1 <: \sigma_2 \dashv \Delta, a, \Theta}{\Gamma \vdash^{\text{sub}} \sigma_1 <: \forall a. \sigma_2 \dashv \Delta}$	
<b>S-FORALLL</b> $\frac{\Gamma, \blacktriangleright_{\hat{\alpha}}, \hat{\alpha} \vdash^{\text{sub}} \sigma_1[a \mapsto \hat{\alpha}] <: \sigma_2 \dashv \Delta, \blacktriangleright_{\hat{\alpha}}, \Theta}{\Gamma \vdash^{\text{sub}} \forall a. \sigma_1 <: \sigma_2 \dashv \Delta}$	<b>S-INSTL</b> $\frac{\Gamma[\hat{\alpha}] \vdash_{\hat{\alpha}}^+ \sigma \rightsquigarrow \tau \dashv \Delta[\hat{\alpha}]}{\Gamma[\hat{\alpha}] \vdash^{\text{sub}} \hat{\alpha} <: \sigma \dashv \Delta[\hat{\alpha} = \tau]}$	
<b>S-INSTR</b> $\frac{\Gamma[\hat{\alpha}] \vdash_{\hat{\alpha}}^- \sigma \rightsquigarrow \tau \dashv \Delta[\hat{\alpha}]}{\Gamma[\hat{\alpha}] \vdash^{\text{sub}} \sigma <: \hat{\alpha} \dashv \Delta[\hat{\alpha} = \tau]}$		

  

$\Gamma \vdash_{\hat{\alpha}}^{\otimes} \sigma \rightsquigarrow \tau \dashv \Delta$

(Polymorphic Promotion)

  

<b>P-PR-FORALLL</b> $\frac{\Gamma[\hat{\beta}, \hat{\alpha}] \vdash_{\hat{\alpha}}^- \sigma[a \mapsto \hat{\beta}] \rightsquigarrow \tau \dashv \Delta}{\Gamma[\hat{\alpha}] \vdash_{\hat{\alpha}}^- \forall a. \sigma \rightsquigarrow \tau \dashv \Delta}$	<b>P-PR-FORALLR</b> $\frac{\Gamma, a \vdash_{\hat{\alpha}}^+ \sigma \rightsquigarrow \tau \dashv \Delta, a}{\Gamma \vdash_{\hat{\alpha}}^+ \forall a. \sigma \rightsquigarrow \tau \dashv \Delta}$
<b>P-PR-ARROW</b> $\frac{\Gamma \vdash_{\hat{\alpha}}^{\otimes} \sigma_1 \rightsquigarrow \tau_1 \dashv \Theta \quad \Theta \vdash_{\hat{\alpha}}^{\otimes} [\Theta]\sigma_2 \rightsquigarrow \tau_2 \dashv \Delta}{\Gamma \vdash_{\hat{\alpha}}^{\otimes} \sigma_1 \rightarrow \sigma_2 \rightsquigarrow \tau_1 \rightarrow \tau_2 \dashv \Delta}$	<b>P-PR-MONO</b> $\frac{\Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \tau_1 \rightsquigarrow \tau_2 \dashv \Delta}{\Gamma \vdash_{\hat{\alpha}}^{\otimes} \tau_1 \rightsquigarrow \tau_2 \dashv \Delta}$

  

$\Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \tau_1 \rightsquigarrow \tau_2 \dashv \Delta$

(Promotion)

  

**PR-TVAR**  

$$\frac{}{\Gamma[a][\hat{\alpha}] \vdash_{\hat{\alpha}}^{\text{pr}} a \rightsquigarrow a \dashv \Gamma[a][\hat{\alpha}]}$$

Figure 6.2: Types, contexts, subtyping and (polymorphic) promotion of the algorithmic system

$$\frac{\text{P-PR-FORALLRR} \quad a \notin \text{FV}(\sigma) \quad \Gamma \vdash_{\hat{\alpha}}^+ \sigma \rightsquigarrow \tau \dashv \Delta}{\Gamma \vdash_{\hat{\alpha}}^+ \forall a. \sigma \rightsquigarrow \tau \dashv \Delta}$$

Rule **P-PR-ARROW** flips the mode for codomains, and uses the same mode for domains. When the type to be promoted is a monotype, rule **P-PR-MONO** uses the promotion judgment ( $\vdash^{\text{pr}}$ ) directly. Note that for a monotype the mode does not matter, so rule **P-PR-MONO** applies in both modes.

**PROMOTION.** The promotion judgment is the same as before, and still only works for monotypes, except that now we have rule **PR-TVAR** for type variables  $a$ . Note again that promotion is a partial operation, as it requires  $a$  to be the left of  $\hat{\alpha}$ , since variable orders matter.

### 6.3.3 SOUNDNESS AND COMPLETENESS

The statement of soundness of promotion remains the same as before.

**Theorem 6.5** (Soundness of Promotion). *If  $\Gamma \vdash_{\hat{\alpha}}^{\text{pr}} \tau_1 \rightsquigarrow \tau_2 \dashv \Delta$ , and  $\Delta = \Delta_1, \hat{\alpha}, \Delta_2$ , then  $\Delta_1 \vdash^{\text{wf}} \tau_2$ , and  $[\Delta]\tau_1 = [\Delta]\tau_2$ .*

Based on soundness of promotion, we prove that after polymorphic promotion, the promoted type is also well-formed under the prefix context of  $\hat{\alpha}$ . Moreover, polymorphic promotion builds a subtyping relation according to the promotion mode: under the contravariant mode ( $-$ ), the original type is a subtype of the promoted type; under the covariant mode ( $+$ ), the promoted type is a subtype of the original type.

**Theorem 6.6** (Soundness of Polymorphic Promotion). *If  $\Gamma \vdash_{\hat{\alpha}}^{\otimes} \sigma \rightsquigarrow \tau \dashv \Delta$ , and  $\Delta = \Delta_1, \hat{\alpha}, \Delta_2$ , then  $\Delta_1 \vdash^{\text{wf}} \tau_2$ . Moreover, given  $\Delta \longrightarrow \Omega$ ,*

- *if  $\otimes = \neg$ , then  $[\Omega]\Gamma \vdash^{DK} [\Omega]\sigma <: [\Omega]\tau$ ; and*
- *if  $\otimes = +$ , then  $[\Omega]\Gamma \vdash^{DK} [\Omega]\tau <: [\Omega]\sigma$ .*

With soundness of polymorphic promotion, next we show that the new subtyping judgment using polymorphic promotion instead of instantiation remains sound.

**Theorem 6.7** (Soundness of Subtyping). *If  $\Gamma \vdash^{\text{sub}} \sigma_1 <: \sigma_2 \dashv \Delta$ , and  $\Delta \longrightarrow \Omega$ , then  $[\Omega]\Gamma \vdash^{DK} [\Omega]\sigma_1 <: [\Omega]\sigma_2$ .*

Now we turn to completeness. The completeness of promotion is the same as before.

**Theorem 6.8** (Completeness of Promotion). *Given  $\Gamma \longrightarrow \Omega$ , and  $\Gamma \vdash^{\text{wf}} \hat{\alpha}$ , and  $\Gamma \vdash^{\text{wf}} \tau$ , and  $[\Gamma]\hat{\alpha} = \hat{\alpha}$ , and  $[\Gamma]\tau = \tau$ , if  $\hat{\alpha} \notin \text{FV}(\tau)$ , there exist  $\tau_2$ ,  $\Delta$  and  $\Omega'$  such that  $\Gamma \longrightarrow \Omega'$  and  $\Omega \longrightarrow \Omega'$  and  $\Gamma \vdash_{\hat{\alpha}}^{\text{Pr}} \tau \rightsquigarrow \tau_2 \dashv \Delta$ .*

Completeness of polymorphic promotion has two parts. If the existential variable appears on the left, then we promote the type under the covariant mode (+), or otherwise the contravariant mode (-). Moreover, it also requires  $\hat{\alpha} \notin \text{FV}(\sigma)$ .

**Theorem 6.9** (Completeness of Polymorphic Promotion). *Given  $\Gamma \longrightarrow \Omega$ , and  $\Gamma \vdash^{\text{wf}} \hat{\alpha}$ , and  $\Gamma \vdash^{\text{wf}} \sigma$ , and  $[\Gamma]\hat{\alpha} = \hat{\alpha}$ , and  $[\Gamma]\tau = \sigma$ , and  $\hat{\alpha} \notin \text{FV}(\sigma)$ ,*

- *if  $[\Omega]\Gamma \vdash^{\text{DK}} [\Omega]\hat{\alpha} <: [\Omega]\sigma$ , then there exist  $\tau$ ,  $\Delta$  and  $\Omega'$  such that  $\Gamma \vdash_{\hat{\alpha}}^+ \sigma \rightsquigarrow \tau \dashv \Delta$ ; and*
- *if  $[\Omega]\Gamma \vdash^{\text{DK}} [\Omega]\sigma <: [\Omega]\hat{\alpha}$ , then there exist  $\tau$ ,  $\Delta$  and  $\Omega'$  such that  $\Gamma \vdash_{\hat{\alpha}}^- \sigma \rightsquigarrow \tau \dashv \Delta$ .*

Finally, we prove that our subtyping is complete. With this, we have proved our claim that the original instantiation relation in DK can be replaced by the polymorphic promotion process, as the subtyping algorithm using polymorphic promotion remains sound and complete.

**Theorem 6.10** (Completeness of Subtyping). *Given  $\Gamma \longrightarrow \Omega$ , and  $\Gamma \vdash^{\text{wf}} \sigma_1$ , and  $\Gamma \vdash^{\text{wf}} \sigma_2$ , if  $[\Omega]\Gamma \vdash^{\text{DK}} [\Omega]\tau_1 <: [\Omega]\tau_2$ , there exist  $\Delta$  and  $\Omega'$  such that  $\Delta \longrightarrow \Omega'$  and  $\Omega \longrightarrow \Omega'$  and  $\Gamma \vdash^{\text{sub}} [\Gamma]\sigma_1 <: [\Gamma]\sigma_2 \dashv \Delta$ .*

## 6.4 DISCUSSION

This section discusses two extensions of promotion. The first extension explores dependent types, while the second extension considers gradual types.

### 6.4.1 PROMOTING DEPENDENT TYPES

In Section 6.1.1 we mentioned the drawback of decomposing type constructs that it cannot be easily applied to more advanced types like dependent types. In this section, we discuss how we can apply promotion to dependent types.

Consider rule **PR-PI** given below that promotes a dependent type  $\Pi a : \tau_1. \tau_2$ .

$$\frac{\text{PR-PI} \quad \Gamma \vdash_{\hat{\alpha}}^{\text{Pr}} \tau_1 \rightsquigarrow \tau_3 \dashv \Theta \quad \Theta, a \vdash_{\hat{\alpha}}^{\text{Pr}} [\Theta]\tau_2 \rightsquigarrow \tau_4 \dashv \Delta, a}{\Gamma \vdash_{\hat{\alpha}}^{\text{Pr}} \Pi a : \tau_1. \tau_2 \rightarrow \text{Int} \rightsquigarrow \Pi a : \tau_3. \tau_4 \dashv \Delta}$$



Here we first promote  $\tau_1$ , returning  $\tau_3$ . Then we add  $a$  into the context to promote  $\tau_2$ . Finally, we return  $\Pi a : \tau_3. \tau_4$  and discard  $a$  in the output context.

Unfortunately, this design does not work. In particular, consider promoting  $\Pi a : \hat{\beta}. a$ .

$$\frac{\text{PR-EVARL} \quad \text{PR-EVARL}}{\frac{\widehat{\beta}, \widehat{\alpha} \vdash_{\widehat{\alpha}}^{\text{pr}} \widehat{\beta} \rightsquigarrow \widehat{\beta} \dashv \widehat{\beta}, \widehat{\alpha} \quad \widehat{\beta}, \widehat{\alpha}, a \vdash_{\widehat{\alpha}}^{\text{pr}} a \rightsquigarrow ???}{\widehat{\beta}, \widehat{\alpha} \vdash_{\widehat{\alpha}}^{\text{pr}} \Pi a : \widehat{\beta}. a \rightarrow \text{Int} \rightsquigarrow} \text{PR-PI}$$

We expect that the promotion would return  $\Pi a : \widehat{\beta}. a$ . However, after we add  $a$  into the context to promote  $a$ , rule **PR-TVARR** does not apply, as  $a$  is out of the scope of  $\widehat{\alpha}$ !

The issue can be fixed by changing rule **PR-TVARR** to rule **PR-TVARR** to not consider the order of type variables.

$$\text{PR-TVARR} \quad \frac{}{\Gamma \vdash_{\widehat{\alpha}}^{\text{pr}} a \rightsquigarrow a \dashv \Gamma}$$

Then, while promotion resolves the ordering of existential variables, since there is no constraint for type variables, it is not guaranteed anymore that the promoted type is well-formed in the prefix context of  $\widehat{\alpha}$ . Therefore, we need to adjust the rule of subtyping to check explicitly that the result is well-formed, i.e.,

$$\begin{array}{c} \text{s-INSTLL} \\ \frac{\Gamma[\widehat{\alpha}] \vdash_{\widehat{\alpha}}^+ \sigma \rightsquigarrow \tau \dashv \Delta_1, \widehat{\alpha}, \Delta_2 \quad \Delta_1 \vdash^{\text{wf}} \tau}{\Gamma[\widehat{\alpha}] \vdash^{\text{sub}} \widehat{\alpha} <: \sigma \dashv \Delta_1, \widehat{\alpha} = \tau, \Delta_2} \\ \\ \text{s-INSTRR} \\ \frac{\Gamma[\widehat{\alpha}] \vdash_{\widehat{\alpha}}^- \sigma \rightsquigarrow \tau \dashv \Delta_1, \widehat{\alpha}, \Delta_2 \quad \Delta_1 \vdash^{\text{wf}} \tau}{\Gamma[\widehat{\alpha}] \vdash^{\text{sub}} \sigma <: \widehat{\alpha} \dashv \Delta_1, \widehat{\alpha} = \tau, \Delta_2} \end{array}$$

Xie and Oliveira [2017] include a more detailed discussion and formalization of applying promotion to a dependently typed lambda calculus.

#### 6.4.2 PROMOTING GRADUAL TYPES

We have shown that polymorphic promotion works for DK. A natural extension is to also apply polymorphic promotion to GPC (Chapter 4). Then the key is to show how to promote the

unknown type. Since comparing with the unknown type does not impose any constraints, we can simply replace it with a fresh existential variable:

$$\frac{\text{P-PR-UNKNOWN}}{\Gamma[\hat{\alpha}] \vdash_{\hat{\alpha}}^{\otimes} ? \rightsquigarrow \hat{\beta} \dashv \Gamma[\hat{\beta}, \hat{\alpha}]}$$

For example, we have  $\hat{\alpha} \vdash_{\hat{\alpha}}^{\text{pr}} \text{Int} \rightarrow ? \rightsquigarrow \text{Int} \rightarrow \hat{\beta} \dashv \hat{\beta}, \hat{\alpha}$ .

For the extended GPC which restores the dynamic guarantee (Chapter 5), we can replace the unknown type with a fresh gradual existential variables instead.

$$\frac{\text{P-PR-UNKNOWN}_G}{\Gamma[\hat{\alpha}] \vdash_{\hat{\alpha}}^{\otimes} ? \rightsquigarrow \hat{\beta}_G \dashv \Gamma[\hat{\beta}_G, \hat{\alpha}]}$$

With these rules it would be possible to apply polymorphic promotion to GPC. Note this discussion is a sketch and we have not fully worked out the full algorithm yet.

# 7

## KIND INFERENCE FOR DATATYPES

In recent years, languages like Haskell have seen a dramatic surge of new features that significantly extends the expressive power of their type systems. With these features, the challenge of *kind inference* for datatype declarations has presented itself and become a worthy research problem on its own.

In this section, we apply promotion to kind inference for datatypes. Inspired by previous research on type-inference, we offer declarative specifications for what datatype declarations should be accepted, both for Haskell98 and for a more advanced system we call PolyKinds, based on the extensions in modern Haskell, including a limited form of dependent types. We believe these formulations to be novel and without precedent, even for Haskell98. These specifications are complemented with implementable algorithmic versions. We study *soundness*, *completeness* and the existence of *principal kinds* in these systems, proving the properties where they hold. This work can serve as a guide both to language designers who wish to formalize their datatype declarations and also to implementors keen to have principled inference of principal types.

### 7.1 INTRODUCTION AND MOTIVATION

#### 7.1.1 KIND INFERENCE FOR DATATYPES

Modern functional languages such as Haskell, ML, and OCaml come with powerful forms of type inference. The global type-inference algorithms employed in those languages are derived from the Hindley-Milner type system (HM) [Damas and Milner 1982; Hindley 1969], with multiple extensions. As the languages evolve, researchers also formalize the key aspects of type inference for the new extensions. Common extensions of HM include *higher-ranked polymorphism* [Odersky and Läufer 1996; Peyton Jones et al. 2007] and *type-inference for GADTs* [Peyton Jones et al. 2006], which have both been formally studied thoroughly.

Most research work for extensions of HM so far (including OL, DK, AP and GPC) has focused on forms of polymorphism, where type variables all have the same kind. In these systems, the type variables introduced by universal quantifiers and/or type declarations all stand for proper types (i.e., they have kind  $\star$ ). In such a simplified setting, datatype declara-

tions such as `data Maybe a = Nothing | Just a` pose no problem at all for type inference: with only one possible kind for `a`, there is nothing to infer.

However, real-world implementations for languages like Haskell support a non-trivial kind language, including kinds other than  $\star$ . Haskell98 accepts *higher-kinded polymorphism* [Jones 1995], enabling datatype declarations such as `data AppInt f = Mk (f Int)`. The type of constructor `Mk` applies the type variable `f` to an argument `Int`. Accordingly, `AppInt Bool` would not work, as the type `Bool Int` (in the instantiated type of `Mk`) is invalid. Instead, we must write something like `AppInt Maybe`: the argument to `AppInt` must be suitable for applying to `Int`. In Haskell98, `AppInt` has kind  $(\star \rightarrow \star) \rightarrow \star$ . For Haskell98-style higher-kinded polymorphism, Jones [1995] presents one of the few extensions of HM that deals with a non-trivial language of kinds. His work addresses the related problem of inference for *constructor type classes*, although he does not show directly how to do inference for datatype declarations.

Modern Haskell<sup>1</sup> has a much richer type and kind language compared to Haskell98. In recent years, Haskell has seen a dramatic surge of new features that extend the expressive power of algebraic datatypes. Such features include *GADTs*, *kind polymorphism* [Yorgey et al. 2012] with *implicit kind arguments*, and *dependent kinds* [Weirich et al. 2013], among others. With great power comes great responsibility: now we must be able to infer these kinds, too. For instance, consider these datatype declarations:

```
data App f a = MkApp (f a)
data Fix f    = In (f (Fix f))
data T        = MkT1 (App Maybe Int)
               | MkT2 (App Fix Maybe)  -- accept or reject?
```

Should the declaration for `T` be accepted or rejected? In a Haskell98 setting, the kind of `App` is  $(\star \rightarrow \star) \rightarrow \star \rightarrow \star$ . Therefore `T` should be rejected, because in `MkT2` the datatype `App` is applied to `Fix :: (\star \rightarrow \star) \rightarrow \star` and `Maybe :: \star \rightarrow \star`, which do not match the expected kinds of `App`. However, with kind polymorphism, `T` is accepted, because `App` has the more general kind  $\forall k. (k \rightarrow \star) \rightarrow k \rightarrow \star$ . With this kind, both uses of `App` in `T` are valid.

The questions we ask in this section are these: *Which datatype declarations should be accepted? What kinds do accepted datatypes have?* Surprisingly, the literature is essentially silent on these questions—we are unaware of any formal treatment of kind inference for datatype declarations.

Inspired by previous research on type inference, we offer declarative specifications for two languages: Haskell98, as standardized [Peyton Jones 2003] (?); and PolyKinds, a significant

<sup>1</sup>We consider the Glasgow Haskell Compiler’s implementation of Haskell, in version 8.8.

fragment of modern Haskell (??). These specifications are complemented with algorithmic versions that can guide implementations (Sections ?? and ??). To relate the declarative and algorithmic formulations we study various properties, including *soundness*, *completeness*, and the existence of *principal kinds* (Sections ??, ??, and ??).

### 7.1.2 KIND INFERENCE IN HASKELL98

Haskell98’s kind language contains a constant (the kind  $\star$ ) and kinds built from arrows ( $k_1 \rightarrow k_2$ ). Kind inference for Haskell98 datatypes is thus closely related to type inference for the simply typed  $\lambda$ -calculus (STLC). For example, consider a term  $+$  with type  $Int \rightarrow Int \rightarrow Int$  and a type constructor  $+$  with kind  $\star \rightarrow \star \rightarrow \star$ . At the term level, we infer that  $add\ a\ b = a + b$  yields  $add :: Int \rightarrow Int \rightarrow Int$ . Similarly, we can create a datatype  $Add\ a\ b = Add\ (a : + : b)$  and infer  $Add :: \star \rightarrow \star \rightarrow \star$ .

**NO PRINCIPAL TYPES.** Consider now the function definition  $k\ a = 1$ . In the STLC, there are infinitely many (incomparable) types that can be assigned to  $k$ , including  $k :: Int \rightarrow Int$  and  $k :: (Int \rightarrow Int) \rightarrow Int$ . Assuming there are no type variables, the STLC accordingly has no *principal types*. An analogous datatype declaration is  $data\ K\ a = K\ Int$ . As with  $k$ , there are infinitely many (incomparable) kinds that can be assigned to  $K$ , including  $K :: \star \rightarrow \star$  and  $K :: (\star \rightarrow \star) \rightarrow \star$ .

**DEFAULTING.** Definitions like  $k$  (in STLC) or  $K$  (in Haskell98) do not have a principal type/kind, which raises the immediate question of what type/kind to infer. Haskell98 solves this problem by using a *defaulting* strategy: *if the kind of a type variable cannot be inferred, then it is defaulted to  $\star$* . Therefore the kind of  $K$  in Haskell98 is  $\star \rightarrow \star$ . From the perspective of type inference, such defaulting strategy may seem somewhat ad-hoc, but due to the role that  $\star$  plays at the type level it seems a defensible design for kind inference. Defaulting brings complications in writing a declarative specification. We discuss this point further in Section ??.

### 7.1.3 KIND INFERENCE IN MODERN GHC HASKELL

The type and kind languages for modern GHC are *unified* (i.e., types and kinds are indistinguishable), *dependently typed*, and the kind system includes the  $\star :: \star$  axiom Cardelli [1986]; Weirich et al. [2013]. We informally use the word *type* or *kind* where we find it appropriate. Unlike Haskell98’s datatypes, whose inference problem is quite closely related to the well-studied inference problem for STLC, type inference for various features in modern Haskell

is not well-studied. While we are motivated concretely by Haskell, many of the challenges we face would be present in any dependently typed language seeking principled type inference. We use the term PolyKinds to refer to the fragment of modern Haskell we model.<sup>2</sup> We enumerate the key features of this fragment below.

**KIND POLYMORPHISM AND DEPENDENT TYPES** Global type inference, in the style of Damas and Milner [1982], allows polymorphic kinds to be assigned to datatype definitions. For instance, reconsider `data K a = K Int`. In PolyKinds, `K` can be given the kind  $K :: \forall\{k\}. k \rightarrow \star$ . This example shows one of the interesting new features of PolyKinds over Haskell98: *kind polymorphism* [Yorgey et al. 2012]. The polymorphic kind is obtained via *generalization*, which is a standard feature in Damas-Milner algorithms. Polymorphic types are helpful for recovering principal types, since they generalize many otherwise incomparable monomorphic types.

System-F-based languages do not have dependent types. In contrast, PolyKinds supports dependent kinds such as `data D ::  $\forall(k :: \star) (a :: k). K a \rightarrow \star$` . There are two noteworthy aspects about the kind of `D`. Firstly, kind and type variables are *typed*: different type variables may have different kinds. Secondly, the kinds of later variables can *depend* on earlier ones. In `D`, the kind of `a` depends on `k`. Both typed variables and dependent kinds bring technical complications that do not exist in many previous studies of type inference (e.g., Peyton Jones et al. [2007]; Vytiniotis et al. [2011]).

**FIRST-ORDER UNIFICATION WITH DEPENDENT KINDS AND TYPED VARIABLES.** Although PolyKinds is dependently typed, its unification problem is remarkably *first-order*. This is in contrast to many other dependently typed languages, where unification is usually *higher-order* [Andrews 1971; Huet 1973]. Since unification plays a central role in inference algorithms this is a crucial difference. Higher-order unification is well-known to be undecidable in the general case [Goldfarb 1981]. As a consequence, type-inference algorithms for most dependently typed languages make various trade-offs.

A key reason why unification can be kept as a first-order problem in PolyKinds is because the type language *does not include lambdas*. Type-level lambdas have been avoided since the start in Haskell, since they bring major challenges for (term-level) type inference [Jones 1995].

<sup>2</sup>Some of the features we model are slightly different in our presentation than they exist in GHC. TODO outlines the differences. These minor differences do not affect the applicability of our work to improving the GHC implementations, but they may affect the ability to test our examples in GHC.

The unification problem for PolyKinds is still challenging, compared to unification for System-F-like languages: unification must be *kind-directed*, as first observed at the term level by Jones [1995]. Consider the following (contrived) example:

```
data X :: ∀a (b :: ★ → ★). a b → ★      -- accepted
data Y :: ∀(c :: Maybe Bool). X c → ★    -- rejected
```

In  $X$ 's kind, we discover  $a :: (\star \rightarrow \star) \rightarrow \star$ . When checking  $Y$ 's kind, we must infer how to instantiate  $X$ : that is, we must choose  $a$  and  $b$  so that  $a\ b$  unifies with  $\text{Maybe Bool}$ , which is  $c$ 's kind. It is tempting to solve this with  $a \mapsto \text{Maybe}$  and  $b \mapsto \text{Bool}$ , but doing so would be ill-kinded, as  $a$  and  $\text{Maybe}$  have different kinds. Our unification thus features *heterogeneous constraints*?. When solving a unification variable, we need to first unify the kinds on both sides.

Because unification recurs into kinds, and because types are undifferentiated from kinds, it might seem that unification might not terminate. In Section ?? we show that the first-order unification with heterogeneous constraints employed in PolyKinds is guaranteed to terminate.

**MUTUAL AND POLYMORPHIC RECURSION** Recursion and mutual recursion are omnipresent in datatype declarations. In PolyKinds, mutually recursive definitions will be kinded together and then get generalized together. For example, both  $P$  and  $Q$  get kind  $\forall(k :: \star). k \rightarrow \star$ .

```
data P a = MkP (Q a)
data Q a = MkQ (P a)
```

The recursion is simple here: all recursive occurrences are at the same type. In existing type-inference algorithms, such recursive definitions are well understood and do not bring considerable complexity to type inference. However, we must also consider *polymorphic recursion* as in  $\text{Poly}$ :

```
data Poly :: ∀k. k → ★
data Poly k = C1 (Poly Int) | C2 (Poly Maybe)
```

This example includes a *kind signature*, meaning that we must *check* the kind of the datatype, not *infer* it. In the definition of  $\text{Poly}$ , the type  $\text{Poly Int}$  requires an instantiation  $k \mapsto \star$ , while the type  $\text{Poly Maybe}$  requires an instantiation of  $k \mapsto (\star \rightarrow \star)$ . These differing instantiations cause the declaration to be polymorphic recursive.

PolyKinds deals with such cases of polymorphic recursion, which also appear at the term level—for example, when writing recursive functions over GADTs or nested datatypes [Bird

and Meertens 1998]. Polymorphic recursion is known to render type-inference undecidable [Henglein 1993]. Furthermore, most existing formalizations of type inference avoid the question entirely, either by not modeling recursion at all or not allowing polymorphic recursion. Our PolyKinds system has full support for polymorphic recursion, implemented directly without the use of a *fix* operator. Polymorphic recursion is allowed only on datatypes with a kind signature; other datatypes are treated as monomorphic during inference.

**VISIBLE KIND APPLICATION** PolyKinds lifts *visible type application* (VTA) [Eisenberg et al. 2016], whereby we can explicitly instantiate a function call, as in *id @Bool True*, to kinds, giving us *visible kind application* (VKA). Following the design of VTA, we distinguish *specified variables* from *inferred variables*. As described by Eisenberg et al. [2016, Section 3.1], only specified variables can be instantiated via VKA. Instantiation of variables is inferred when no explicit kind application is given. To illustrate, consider *data T :: ∀a b. a b → \**. Here, *a* and *b* are specified variables. Because their order is given, explicit instantiation of *a* must happen before *b*. For example, *T @Maybe* instantiates *a* to *Maybe*. On the other hand, the kind of *a* and *b* can be generalized to *a :: k → \** and *b :: k*. Elaborating the kind of *T*, we write *T :: ∀{k :: \*} (a :: k → \*) (b :: k). a b → \**. The variable *k* is *inferred* and is not available for instantiation with VKA. This split between specified and inferred variables supports predictable type inference: if the variables invented by the compiler (e.g., *k*) were available for instantiation, then we have no way of knowing what order to instantiate them.

**OPEN KIND SIGNATURES AND GENERALIZATION ORDER** Echoing the design of Haskell, PolyKinds supports *open kind signatures*. We say a signature is *closed* if it contains no free variables (e.g., *cusk (T :: ∀a. a → \*)*). Otherwise, it is *open* (e.g., *cusk (Q :: ∀(a :: (f b)) (c :: k). f c → \*)*). Free variables (in this case, *f*, *b*, *k*) will be generalized over. We have a decision to make: in which order do we generalize the free variables? This question is non-trivial, as there can be dependency between the variables. We infer *k :: \**, *f :: k → \**, *b :: k*. Even though *f* and *b* appear before *k*, their kinds end up depending on *k* and we must quantify *k* before *f* and *b*. Inferring this order is a challenge: we cannot know the correct order before completing inference. We thus introduce *local scopes*, which are sets of variables that may be reordered. Since the ordering is not fixed by the programmer, these variables are considered *inferred*, not *specified*, with respect to VKA.

**EXISTENTIAL QUANTIFICATION.** PolyKinds supports existentially quantified variables on datatype constructors. This is useful, for example, to model GADTs. Given *data T1 = ∀a. MkT1 a*, we get *MkT1 :: ∀(a :: \*). a → T1*. The type of the data constructor declaration



can also be generalized. Given  $\text{data } P1 :: \forall(a :: \star). \star$ , from  $\text{data } T2 = \text{MkT2 } P1$ , we infer  $\text{MkT2} :: \forall\{a :: \star\}. P1 @a \rightarrow T2$ , where  $P1$  is elaborated to  $P1 @a$  with  $a$  generalized as an inferred variable.

#### 7.1.4 DESIRABLE PROPERTIES FOR KIND INFERENCE

The goal of this work is to provide concrete, principled guidance to implementors of dependently typed languages, such as GHC/Haskell. It is thus important to be able to describe our inference algorithm as sound and complete against a *declarative specification*. This declarative specification is what we might imagine a programmer to have in her head as she programs. This system should be designed with a minimum of low-level detail and a minimum of surprises. It is then up to an algorithm to live up to the expectations set by the specification. The algorithm is sound when all programs it accepts are also accepted by the specification; it is complete when all programs accepted by the specification are accepted by the algorithm.

Why choose the particular set of features described here? Because they lead to interesting kind inference challenges. We have found that the features above are sufficient in exploring kind inference in modern Haskell. We consider unformalized extensions in ??.



## PART V

### RELATED AND FUTURE WORK



## 8 RELATED WORK



## 9 FUTURE WORK





## PART VI

## EPILOGUE



# 10 CONCLUSION



# BIBLIOGRAPHY

[Citing pages are listed after each reference.]

- Amal Ahmed, Robert Bruce Findler, Jacob Matthews, and Philip Wadler. 2009. Blame for All. In *Proceedings for the 1st Workshop on Script to Program Evolution (STOP '09)*. Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/1570506.1570507> [cited on pages 4, 59, 61, 69, 71, 73, 74, and 97]
- P. B. Andrews. 1971. Resolution in type Theory. *Journal of Symbolic Logic* 36 (1971), 414–432. [cited on page 126]
- Andrea Asperti, Wilmer Ricciotti, Claudio Sacerdoti Coen, and Enrico Tassi. 2012. A Bi-Directional Refinement Algorithm for the Calculus of (Co) Inductive Constructions. *Logical Methods in Computer Science* 8 (2012), 1–49. [cited on page 15]
- Richard S. Bird and Lambert Meertens. 1998. Nested datatypes. In *LNCS 1422: Proceedings of Mathematics of Program Construction*, Johan Jeuring (Ed.). Springer-Verlag, Marstrand, Sweden, 52–67. <http://www.cs.ox.ac.uk/people/richard.bird/online/BirdMeertens98Nested.pdf> [cited on page 127]
- L. Cardelli. 1986. *A polymorphic lambda-calculus with Type:Type*. Technical Report 10. SRC. [cited on page 125]
- Luca Cardelli. 1993. *An implementation of FSub*. Technical Report. Research Report 97, Digital Equipment Corporation Systems Research Center. [cited on page 80]
- Giuseppe Castagna and Victor Lanvin. 2017. Gradual Typing with Union and Intersection Types. *Proc. ACM Program. Lang.* 1, ICFP, Article 41 (Aug. 2017), 28 pages. [cited on pages 65 and 90]
- Gang Chen. 2003. Coercive Subtyping for the Calculus of Constructions (*POPL '03*). 10. [cited on page 42]

- Alonzo Church. 1941. *The calculi of lambda-conversion*. Number 6. Princeton University Press. [cited on page 56]
- Matteo Cimini and Jeremy G. Siek. 2016. The Gradualizer: A Methodology and Algorithm for Generating Gradual Type Systems. In *Proceedings of the 43rd Symposium on Principles of Programming Languages*. [cited on pages 66 and 73]
- Thierry Coquand. 1996. An algorithm for type-checking dependent types. *Science of Computer Programming* 26, 1-3 (1996), 167–177. [cited on page 15]
- Haskell Brooks Curry, Robert Feys, William Craig, J Roger Hindley, and Jonathan P Seldin. 1958. *Combinatory logic*. Vol. 1. North-Holland Amsterdam. [cited on page 56]
- Luis Damas and Robin Milner. 1982. Principal Type-Schemes for Functional Programs. In *Proceedings of the 9th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '82)*. Association for Computing Machinery, New York, NY, USA, 207–212. <https://doi.org/10.1145/582153.582176> [cited on pages 7, 10, 123, and 126]
- Rowan Davies and Frank Pfenning. 2000. Intersection Types and Computational Effects. In *Proceedings of the Fifth ACM SIGPLAN International Conference on Functional Programming (ICFP '00)*. Association for Computing Machinery, New York, NY, USA, 198–208. <https://doi.org/10.1145/351240.351259> [cited on pages 15 and 90]
- Joshua Dunfield and Neelakantan R. Krishnaswami. 2013. Complete and Easy Bidirectional Typechecking for Higher-Rank Polymorphism. In *Proceedings of the 18th ACM SIGPLAN International Conference on Functional Programming (ICFP '13)*. Association for Computing Machinery, New York, NY, USA, 429–442. <https://doi.org/10.1145/2500365.2500582> [cited on pages 15, 16, 18, 19, 53, 72, 76, 86, 99, 105, 109, 110, 114, and 116]
- Joshua Dunfield and Frank Pfenning. 2004. Tridirectional Typechecking. *SIGPLAN Not.* 39, 1 (Jan. 2004), 281–292. <https://doi.org/10.1145/982962.964025> [cited on pages 15, 24, and 26]
- Richard A Eisenberg, Stephanie Weirich, and Hamidhasan G Ahmed. 2016. Visible type application. In *European Symposium on Programming*. Springer, 229–254. [cited on page 128]
- Ronald Garcia and Matteo Cimini. 2015. Principal Type Schemes for Gradual Programs. In *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '15)*. Association for Computing Machinery, New York,

- NY, USA, 303–315. <https://doi.org/10.1145/2676726.2676992> [cited on pages 4, 76, 82, 93, 94, and 99]
- Ronald Garcia, Alison M Clark, and Éric Tanter. 2016. Abstracting Gradual Typing. In *Proceedings of the 43rd Symposium on Principles of Programming Languages*. [cited on pages 55, 59, 65, and 66]
- Warren D Goldfarb. 1981. The undecidability of the second-order unification problem. *Theoretical Computer Science* 13, 2 (1981), 225–230. [cited on page 126]
- Adam Gundry, Conor McBride, and James McKinna. 2010. Type inference in context. In *Proceedings of the third ACM SIGPLAN workshop on Mathematically structured functional programming*. ACM, 43–54. [cited on pages 109, 110, 111, and 112]
- Fritz Henglein. 1993. Type inference with polymorphic recursion. *ACM Trans. Program. Lang. Syst.* 15, 2 (April 1993), 253–289. <https://doi.org/10.1145/169701.169692> [cited on page 128]
- J. Roger Hindley. 1969. The Principal Type-Scheme of an Object in Combinatory Logic. *Trans. Amer. Math. Soc.* 146 (1969), 29–60. [cited on pages 7 and 123]
- G. Huet. 1973. A unification algorithm for typed lambda calculus. *Theoretical Computer Science* 1, 1 (1973), 27–57. [cited on page 126]
- Yuu Igarashi, Taro Sekiyama, and Atsushi Igarashi. 2017. On Polymorphic Gradual Typing. In *Proceedings of the 22nd International Conference on Functional Programming*. [cited on pages 59, 61, 66, 93, and 98]
- Khurram A. Jafery and Joshua Dunfield. 2017. Sums of Uncertainty: Refinements Go Gradual. In *Proceedings of the 44th Symposium on Principles of Programming Languages*. 14. [cited on page 66]
- Mark P Jones. 1995. A system of constructor classes: overloading and implicit higher-order polymorphism. *Journal of functional programming* 5, 1 (1995), 1–35. [cited on pages 124, 126, and 127]
- Mark P Jones. 2000. Type classes with functional dependencies. In *European Symposium on Programming*. Springer, 230–244. [cited on page 58]
- Assaf J Kfoury and Jerzy Tiuryn. 1992. Type reconstruction in finite rank fragments of the second-order  $\lambda$ -calculus. *Information and computation* 98, 2 (1992), 228–257. [cited on page 11]

- Oleg Kiselyov, Ralf Lämmel, and Kean Schupke. 2004. Strongly typed heterogeneous collections. In *Proceedings of the 2004 ACM SIGPLAN workshop on Haskell*. ACM, 96–107. [cited on page 58]
- Andres Löb, Conor McBride, and Wouter Swierstra. 2010. A tutorial implementation of a dependently typed lambda calculus. *Fundamenta informaticae* 102, 2 (2010), 177–207. [cited on page 15]
- Jacob Matthews and Amal Ahmed. 2008. Parametric polymorphism through run-time sealing or, theorems for low, low prices!. In *European Symposium on Programming*. Springer, 16–31. [cited on page 71]
- Conor McBride. 2002. Faking it Simulating dependent types in Haskell. *Journal of functional programming* 12, 4-5 (2002), 375–392. [cited on page 58]
- Robin Milner. 1978. A theory of type polymorphism in programming. *Journal of computer and system sciences* 17, 3 (1978), 348–375. [cited on page 7]
- James Hiram Morris Jr. 1969. *Lambda-calculus models of programming languages*. Ph.D. Dissertation. Massachusetts Institute of Technology. [cited on page 73]
- Georg Neis, Derek Dreyer, and Andreas Rossberg. 2009. Non-parametric Parametricity. In *Proceedings of the 14th ACM SIGPLAN International Conference on Functional Programming (ICFP '09)*. ACM, New York, NY, USA, 135–148. <https://doi.org/10.1145/1596550.1596572> [cited on page 71]
- Martin Odersky and Konstantin Läufer. 1996. Putting Type Annotations to Work. In *Proceedings of the 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '96)*. Association for Computing Machinery, New York, NY, USA, 54–67. <https://doi.org/10.1145/237721.237729> [cited on pages 7, 10, 12, 59, and 123]
- Michel Parigot. 1992. Recursive programming with proofs. *Theoretical Computer Science* 94, 2 (1992), 335–356. [cited on page 56]
- Simon Peyton Jones. 2003. *Haskell 98 language and libraries: the revised report*. Cambridge University Press. [cited on page 124]
- Simon Peyton Jones, Mark Jones, and Erik Meijer. 1997. Type classes: exploring the design space. In *Haskell workshop*, Vol. 1997. [cited on page 58]



- Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Mark Shields. 2007. Practical type inference for arbitrary-rank types. *Journal of functional programming* 17, 1 (2007), 1–82. [cited on pages 10, 15, 27, 29, 35, 36, 37, 43, 72, 123, and 126]
- Simon Peyton Jones, Dimitrios Vytiniotis, Stephanie Weirich, and Geoffrey Washburn. 2006. Simple unification-based type inference for GADTs. In *Proceedings of the Eleventh ACM SIGPLAN International Conference on Functional Programming (ICFP '06)*. ACM, New York, NY, USA, 50–61. <https://doi.org/10.1145/1159803.1159811> [cited on page 123]
- Benjamin C Pierce. 2002. *Types and programming languages*. [cited on page 72]
- Benjamin C. Pierce and David N. Turner. 2000. Local Type Inference. *ACM Trans. Program. Lang. Syst.* 22, 1 (Jan. 2000), 1–44. <https://doi.org/10.1145/345099.345100> [cited on pages 15, 23, and 46]
- John C. Reynolds. 1983. Types, Abstraction and Parametric Polymorphism. In *Proceedings of the IFIP 9th World Computer Congress*. [cited on page 71]
- John C. Reynolds. 1991. The Coherence of Languages with Intersection Types. In *Proceedings of the International Conference on Theoretical Aspects of Computer Software*. [cited on page 72]
- Paula Severi and Erik Poll. 1994. Pure Type Systems with Definitions. *Logical Foundations of Computer Science* (1994), 316–328. [cited on page 49]
- Jeremy Siek and Walid Taha. 2007a. Gradual Typing for Objects. In *Proceedings of the 21st European Conference on Object-Oriented Programming (ECOOP'07)*. Springer-Verlag, Berlin, Heidelberg, 2–27. [cited on page 4]
- Jeremy G. Siek and Walid Taha. 2006. Gradual Typing for Functional Languages. In *Proceedings of the 2006 Scheme and Functional Programming Workshop*. [cited on page 92]
- Jeremy G. Siek and Walid Taha. 2007b. Gradual Typing for Objects. In *European Conference on Object-Oriented Programming*. [cited on pages 53, 54, 59, 62, 64, 65, 89, 90, and 92]
- Jeremy G Siek, Michael M Vitousek, Matteo Cimini, and John Tang Boyland. 2015. Refined criteria for gradual typing. In *1st Summit on Advances in Programming Languages (SNAPL 2015)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik. [cited on pages 4, 66, 69, 73, and 76]

- Dimitrios Vytiniotis, Simon Peyton Jones, Tom Schrijvers, and Martin Sulzmann. 2011. OutsideIn (X) Modular type inference with local assumptions. *Journal of functional programming* 21, 4-5 (2011), 333–412. [cited on page 126]
- Stephanie Weirich, Justin Hsu, and Richard A. Eisenberg. 2013. System FC with Explicit Kind Equality. In *Proceedings of the 18th ACM SIGPLAN International Conference on Functional Programming (ICFP '13)*. ACM, New York, NY, USA, 275–286. <https://doi.org/10.1145/2500365.2500599> [cited on pages 124 and 125]
- Joe B Wells. 1999. Typability and Type Checking in System F are Equivalent and Undecidable. *Annals of Pure and Applied Logic* 98, 1-3 (1999), 111–156. [cited on pages 10, 14, and 82]
- Hongwei Xi and Frank Pfenning. 1999. Dependent Types in Practical Programming. In *Proceedings of the 26th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '99)*. Association for Computing Machinery, New York, NY, USA, 214–227. <https://doi.org/10.1145/292540.292560> [cited on pages 15 and 48]
- Ningning Xie, Xuan Bi, and Bruno C d S Oliveira. 2018. Consistent Subtyping for All. In *European Symposium on Programming*. Springer, 3–30. [cited on page 5]
- Ningning Xie, Xuan Bi, Bruno C. D. S. Oliveira, and Tom Schrijvers. 2019a. Consistent Subtyping for All. *ACM Transactions on Programming Languages and Systems* 42, 1, Article 2 (Nov. 2019), 79 pages. <https://doi.org/10.1145/3310339> [cited on page 5]
- Ningning Xie, Richard A. Eisenberg, and Bruno C. d. S. Oliveira. 2019b. Kind Inference for Datatypes. *Proc. ACM Program. Lang.* 4, POPL, Article 53 (Dec. 2019), 28 pages. <https://doi.org/10.1145/3371121> [cited on page 5]
- Ningning Xie and Bruno C d S Oliveira. 2017. Towards Unification for Dependent Types. In *Draft Proceedings of the 18th Symposium on Trends in Functional Programming (TFP '18)*. Extended abstract. [cited on pages 5 and 121]
- Ningning Xie and Bruno C d S Oliveira. 2018. Let Arguments Go First. In *European Symposium on Programming*. Springer, 272–299. [cited on pages 5, 28, and 47]
- Brent A. Yorgey, Stephanie Weirich, Julien Cretin, Simon Peyton Jones, Dimitrios Vytiniotis, and José Pedro Magalhães. 2012. Giving Haskell a Promotion. In *Proceedings of the 8th ACM SIGPLAN Workshop on Types in Language Design and Implementation (TLDI '12)*. ACM, New York, NY, USA, 53–66. <https://doi.org/10.1145/2103786.2103795> [cited on pages 124 and 126]

## PART VII

## TECHNICAL APPENDIX

