

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 | SYSTEM AP | 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 | 30 |
| 3.2.2 | Type System | 31 |
| 3.2.3 | Subtyping | 34 |
| 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 | 40 |
| 3.3.4 | Type Safety | 40 |
| 3.3.5 | Coherence | 41 |
| 3.4 | Type Inference Algorithm | 42 |
| 3.5 | Discussion | 43 |
| 3.5.1 | Combining Application and Checking Modes | 43 |
| 3.5.2 | Additional Constructs | 44 |
| 3.5.3 | More Expressive Type Applications | 45 |
| 3.5.4 | Dependent Type Systems | 47 |
| III | HIGHER-RANKED POLYMORPHISM AND GRADUAL TYPING | 49 |
| 4 | HIGHER RANK GRADUAL TYPES | 51 |
| IV | UNIFICATION AND TYPE-INFERENC FOR DEPENDENT TYPES | 53 |
| 5 | UNIFICATION WITH PROMOTION | 55 |
| 6 | DEPENDENT TYPES | 57 |

| | | |
|-----|-------------------------|----|
| V | RELATED AND FUTURE WORK | 59 |
| 7 | RELATED WORK | 61 |
| 8 | FUTURE WORK | 63 |
| VI | EPILOGUE | 65 |
| 9 | CONCLUSION | 67 |
| | BIBLIOGRAPHY | 69 |
| VII | TECHNICAL APPENDIX | 73 |

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. | 30 |
| 3.2 | Typing rules of System AP. | 32 |
| 3.3 | Syntax and typing rules of System F. | 37 |
| 3.4 | Subtyping translation rules of System AP. | 38 |
| 3.5 | Typing translation rules of System AP. | 39 |
| 3.6 | Type erasure and eta-id equality of System F. | 41 |

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 5 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 [2007] 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 λB [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 6 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.¹

¹For 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 5: 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 6: 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 σ and monomorphic types (monotypes) τ . A polymorphic type is a sequence of universal quantifications (which can be empty) followed by a monotype τ , which can be the integer type Int , type variables a and function types $\tau_1 \rightarrow \tau_2$.

A context Ψ 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 σ of the expression e under the context Ψ . 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* τ_1 as the type of x , and derives the type τ_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 τ_1 , and the whole application returns type τ_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 σ , 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 Ψ . In rule [HM-INST](#), we can instantiate the type variables with arbitrary *monotypes*.

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

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

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 σ_1 is a subtype of σ_2 . The subtyping relation indicates that σ_1 is more *general* than σ_2 : for any instantiation of σ_2 , we can find an instantiation of σ_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 instantiation of a , we *skolemize* a , by making sure a does not appear in σ_1 (up to α -renaming). In this case, if σ_1 is still a subtype of σ_2 , we are sure then whatever a can be instantiated to, σ_1 can be instantiated to match σ_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 σ' such that $\Psi \vdash^{HM} e : \sigma'$, and for all σ 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 parser. 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 ODESKY-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 program¹:

```
(\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 (σ in Figure 2.1) has at most rank 1.

¹For 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 σ is a polymorphic type. Thus unlike HM, the argument type for a function is not limited to a monotype.

Second, the polymorphic types σ now include the integer type 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 Ψ 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. As before, we assume all items in a context are distinct.

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)

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

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

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

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

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

(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 σ is then generalized over a , returning $\forall a. \sigma$.

The subsumption rule **OL-SUB** is crucial for OL, which allows an expression of type σ_1 to have type σ_2 with σ_1 being a subtype of σ_2 ($\Psi \vdash \sigma_1 <: \sigma_2$). 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-FORALLL** 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-FORALLL**. 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 Ψ' is Ψ extended with type variables in Ψ and σ , 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 are 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.

“... are local in the sense 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.

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

Figure 2.6: Syntax of the Dunfield-Krishnaswami Type System

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:

```
(\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 → a) → (Int, 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-binding 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 inward 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 sys-

| | |
|---|---|
| $\Psi \vdash^{DK} e \Rightarrow \sigma$ | <i>(Type Inference)</i> |
| $\frac{\text{DK-INF-VAR} \quad (x : \sigma) \in \Psi}{\Psi \vdash^{DK} x \Rightarrow \sigma}$ | $\frac{\text{DK-INF-INT}}{\Psi \vdash^{DK} n \Rightarrow \text{Int}}$ |
| $\frac{\text{DK-INF-LAM} \quad \Psi \vdash^{DK} \tau_1 \rightarrow \tau_2 \quad \Psi, x : \tau_1 \vdash^{DK} e \Rightarrow \tau_2}{\Psi \vdash^{DK} \lambda x. e \Rightarrow \tau_1 \rightarrow \tau_2}$ | |
| $\frac{\text{DK-INF-APP} \quad \Psi \vdash^{DK} e_1 \Rightarrow \sigma \quad \Psi \vdash^{DK} e_2 \Leftarrow \sigma_1}{\Psi \vdash^{DK} e_1 e_2 \Rightarrow \sigma_2}$ | $\frac{\text{DK-INF-ANNO} \quad \Psi \vdash^{DK} e \Leftarrow \sigma}{\Psi \vdash^{DK} e : \sigma \Rightarrow \sigma}$ |
| $\Psi \vdash^{DK} e \Leftarrow \sigma$ | <i>(Type Checking)</i> |
| $\frac{\text{DK-CHK-INT}}{\Psi \vdash^{DK} n \Leftarrow \text{Int}}$ | $\frac{\text{DK-CHK-LAM} \quad \Psi, x : \sigma_1 \vdash^{DK} e \Leftarrow \sigma_2}{\Psi \vdash^{DK} \lambda x. e \Leftarrow \sigma_1 \rightarrow \sigma_2}$ |
| $\frac{\text{DK-CHK-GEN} \quad \Psi, a \vdash^{DK} e \Leftarrow \sigma}{\Psi \vdash^{DK} e \Leftarrow \forall a. \sigma}$ | |
| $\frac{\text{DK-CHK-SUB} \quad \Psi \vdash^{DK} e \Rightarrow \sigma_1 \quad \Psi \vdash^{DK} \sigma_1 <: \sigma_2}{\Psi \vdash^{DK} e \Leftarrow \sigma_2}$ | |
| $\Psi \vdash^{DK} \sigma_1 \triangleright \sigma_2$ | <i>(Matching)</i> |
| $\frac{\text{DK-M-FORALL} \quad \Psi \vdash^{DK} \tau \quad \Psi \vdash^{DK} \sigma[a \mapsto \tau] \triangleright \sigma_1 \rightarrow \sigma_2}{\Psi \vdash^{DK} \forall a. \sigma \triangleright \sigma_1 \rightarrow \sigma_2}$ | $\frac{\text{DK-M-ARR}}{\Psi \vdash^{DK} \sigma_1 \rightarrow \sigma_2 \triangleright \sigma_1 \rightarrow \sigma_2}$ |

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

2 Background

tems: the inference mode $\Psi \vdash^{DK} e \Rightarrow \sigma$, which takes a term e and produces a type σ , 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 σ as input, and ensures that the term e checks against σ . We first discuss rules in the inference mode.

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 σ of the expression e_1 . Then, because e_1 is applied to an argument, the type σ 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 σ_1 , the rule proceeds by checking e_2 against σ_1 . Similarly, for an annotated expression $e : \sigma$, rule **DK-INF-ANNO** simply checks e against σ . 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 against the integer type `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 σ_1 . Then the rule proceeds by checking e with σ_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 σ . Rule **DK-CHK-SUB** switches the mode from checking to inference: an expression e can be checked against σ_2 , if e infers the type σ_1 and σ_1 is a subtype of σ_2 .

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

Matching has two straightforward rules: rule **DK-M-FORALL** instantiates a polymorphic type, by substituting a with a well-formed monotype τ , 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 σ_1 to the expression e , we get a return type σ_2 .

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

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

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 σ_1 .

We prefer the presentation of rule **DK-INF-APP** with matching, as matching is a simple and standalone 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 SYSTEM AP

This section first presents a declarative, *syntax-directed* type system, System AP, for a lambda calculus with implicit higher-ranked polymorphism. The interesting aspects about the new type system are: 1) the typing rules, which employ a combination of inference and a so-called *application* modes; 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 σ_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**:

DK-CHK-INT

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

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](#).

DK-INF-INT

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

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 required¹. 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,

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

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 Σ . The context Σ 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 σ_1 , which then is pushed into the application context Σ 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.

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 σ 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 \mathbf{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

¹It type-checks in DK, because in DK rules for lambdas are duplicated for having both inference (integrated with type inference techniques) and checking 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 TODO 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 case².

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 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 σ_1 (we know the type of first argument), or an application context σ_1, σ_2 (we know the types of both arguments).

²Although 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.

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 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 Ψ and σ_1 are interpreted as inputs and σ_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:

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

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 ??:

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

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:

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

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 bidirectional 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:

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

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

$\mathbf{let}\ s = \backslash x. x\ \mathbf{in}\ \mathbf{let}\ t = \backslash y. s\ \mathbf{in}\ e$

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 7. 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:

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

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 applica-

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

Figure 3.1: Syntax of System AP.

tion 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.

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 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 Ψ are standard: they map a term variable x to its type σ . 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 Ψ are distinct.

The main novelty lies in the *application contexts* Σ , 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, 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 σ .

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 Ψ , and application context Σ , e has type σ . 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.

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 τ 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 σ 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 σ_1 is a subtype of σ_2 under application context Σ , then x has type σ_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 σ_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 and the application context has σ_2 , then rule **AP-APP-LAMANN** first checks that σ_2 is a subtype of σ_1 before putting $x : \sigma_1$

3 System AP

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

Figure 3.2: Typing rules of System AP.

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 σ_1 . Then the type σ_1 is generalized in the same way as the HM type system. The resulting generalized type is σ_2 . Thus the type of e_1 is now inferred under an application context extended with type σ_2 . The generalization step is important to infer higher ranked types: since σ_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:

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

The type σ_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 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 $\sigma_n \rightarrow \dots \rightarrow \sigma_2 \rightarrow \sigma_1 \rightarrow \sigma$.

3 System AP

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 (Σ Coincides with Typing Results). *If $\Psi; \Sigma \vdash^{AP} e \Rightarrow \sigma$, then for some σ' , 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 σ' , 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 σ_1 is subtype of $\forall a. \sigma_2$ only if σ_1 is a subtype of σ_2 , with the assumption a is a fresh variable. Rule **AP-S-FORALLL** says $\forall a. \sigma_1$ is a subtype of σ_2 if we can instantiate it with some τ and show the result is a subtype of σ_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 σ_2 is a *computed* output, rather than an input. In other words σ_2 is determined from Σ and σ_1 . This is unlike the judgment $\vdash^{AP} \sigma_1 <: \sigma_2$, where both σ_1 and σ_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 Σ 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 `Int` or type variables a cannot have a non-empty application context. In rule **AP-AS-FORALL**, we instantiate the polymorphic type with some τ , 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 σ_3 . So we check $\vdash^{AP} \sigma_3 <: \sigma_1$. Then we con-

3 System AP

tinue with σ_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 (Σ Coincides with Subtyping Results). *If $\Sigma \vdash^{AP} \sigma_1 <: \sigma_2$, then for some σ_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 $\Sigma \vdash^{AP}$). *If $\vdash^{AP} \sigma_1 <: \Sigma \rightarrow \sigma_2$, then for some σ_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 $\vdash^{AP} \text{Int} \rightarrow \forall a. a <: \text{Int} \rightarrow \text{Int}$:

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

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

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

| | |
|-------------|---|
| 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.

The former one, holds because we have $\vdash^{AP} \forall a. a <: \text{Int}$ in the return type. But in the latter one, after 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 *η -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

3 System AP

| | | |
|---|--|-------------------------|
| $\boxed{\vdash^{AP} \sigma_1 <: \sigma_2 \rightsquigarrow f}$ | | (Subtyping Translation) |
| $\text{AP-S-TVAR} \quad \frac{}{\vdash^{AP} a <: a \rightsquigarrow \lambda x : a. x}$ | $\text{AP-S-INT} \quad \frac{}{\vdash^{AP} \text{Int} <: \text{Int} \rightsquigarrow \lambda x : \text{Int}. x}$ | |
| $\text{AP-S-ARROW} \quad \frac{\vdash^{AP} \sigma_3 <: \sigma_1 \rightsquigarrow f_1 \quad \vdash^{AP} \sigma_2 <: \sigma_4 \rightsquigarrow f_2}{\vdash^{AP} \sigma_1 \rightarrow \sigma_2 <: \sigma_3 \rightarrow \sigma_4 \rightsquigarrow \lambda x : \sigma_1 \rightarrow \sigma_2. \lambda y : \sigma_3. f_2(x(f_1 y))}$ | | |
| $\text{AP-S-FORALLL} \quad \frac{\vdash^{AP} \sigma[a \mapsto \tau] <: \sigma_2 \rightsquigarrow f}{\vdash^{AP} \forall a. \sigma_1 <: \sigma_2 \rightsquigarrow \lambda x : \forall a. \sigma_1. f(x\tau)}$ | $\text{AP-S-FORALLR} \quad \frac{a \notin \text{FV}(\sigma_1) \quad \vdash^{AP} \sigma_1 <: \sigma_2 \rightsquigarrow f}{\vdash^{AP} \sigma_1 <: \forall a. \sigma_2 \rightsquigarrow \lambda x : \sigma_1. \Lambda a. f x}$ | |
| $\boxed{\Sigma \vdash^{AP} \sigma_1 <: \sigma_2 \rightsquigarrow f}$ | | (Application Subtyping) |
| $\text{AP-AS-EMPTY} \quad \frac{}{\bullet \vdash^{AP} \sigma <: \sigma \rightsquigarrow \lambda x : \sigma. x}$ | $\text{AP-AS-FORALL} \quad \frac{\Sigma, \sigma_3 \vdash^{AP} \sigma_1[a \mapsto \tau] <: \sigma_2 \rightsquigarrow f}{\Sigma, \sigma_3 \vdash^{AP} \forall a. \sigma_1 <: \sigma_2 \rightsquigarrow \lambda x : \forall a. \sigma_1. f(x\tau)}$ | |
| $\text{AP-AS-ARROW} \quad \frac{\vdash^{AP} \sigma_3 <: \sigma_1 \rightsquigarrow f_1 \quad \Sigma \vdash^{AP} \sigma_2 <: \sigma_4 \rightsquigarrow f_2}{\Sigma, \sigma_3 \vdash^{AP} \sigma_1 \rightarrow \sigma_2 <: \sigma_3 \rightarrow \sigma_4 \rightsquigarrow \lambda x : \sigma_1 \rightarrow \sigma_2. \lambda y : \sigma_3. f_2(x(f_1 y))}$ | | |

Figure 3.4: Subtyping translation rules of System AP.

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 ???. 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 $\bullet \vdash^F f : \sigma_1 \rightarrow \sigma_2$.

$$\boxed{\Psi \vdash^{AP} e \Rightarrow \sigma \rightsquigarrow s} \quad (\text{Typing Inference})$$

$$\begin{array}{c} \text{AP-INF-INT} \\ \hline \Psi \vdash^{AP} n \Rightarrow \text{Int} \rightsquigarrow n \end{array} \quad \begin{array}{c} \text{AP-INF-LAM} \\ \hline \Psi, x : \tau \vdash^{AP} e \Rightarrow \sigma \rightsquigarrow s \\ \hline \Psi \vdash^{AP} \lambda x. e \Rightarrow \tau \rightarrow \sigma \rightsquigarrow \lambda x : \tau. s \end{array}$$

$$\begin{array}{c} \text{AP-INF-LAMANN} \\ \hline \Psi, x : \sigma_1 \vdash^{AP} e \Rightarrow \sigma_2 \rightsquigarrow s \\ \hline \Psi \vdash^{AP} \lambda x : \sigma_1. e \Rightarrow \sigma_1 \rightarrow \sigma_2 \rightsquigarrow \lambda x : \sigma_1. s \end{array}$$

$$\boxed{\Psi; \Sigma \vdash^{AP} e \Rightarrow \sigma \rightsquigarrow s} \quad (\text{Typing Application Mode})$$

$$\begin{array}{c} \text{AP-APP-VAR} \\ \hline (x : \sigma_1) \in \Psi \quad \Sigma \vdash^{AP} \sigma_1 <: \sigma_2 \rightsquigarrow f \\ \hline \Psi; \Sigma \vdash^{AP} x \Rightarrow \sigma_2 \rightsquigarrow f x \end{array} \quad \begin{array}{c} \text{AP-APP-LAM} \\ \hline \Psi, x : \sigma_1 \vdash^{AP} e \Rightarrow \sigma_2 \rightsquigarrow s \\ \hline \Psi; \Sigma, \sigma_1 \vdash^{AP} \lambda x. e \Rightarrow \sigma_1 \rightarrow \sigma_2 \rightsquigarrow \lambda x : \sigma_1. s \end{array}$$

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

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

Figure 3.5: Typing translation rules of System AP.

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 σ_1 . Then the resulting value becomes an argument to x , and a value of type σ_2 is obtained. Finally we apply f_2 to the value of type σ_2 , so that a value of type σ_4 is computed. In rule **A-PS-FORALLL**, the input argument is a polymorphic type, so we feed the type τ to it and apply the coercion function f from the precondition. Rule **AP-S-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 $\bullet \vdash^F f : \sigma \rightarrow \sigma_2$.

| | | |
|-----------------------|--|--------------------------|
| $f_1 =_{\eta id} f_2$ | | <i>(Eta-id Equality)</i> |
|-----------------------|--|--------------------------|

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

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

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 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 η 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 η -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

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

Comparing it to rule **AP-INF-LAM**, τ is replaced by a new meta type variable $\widehat{\beta}$ from name supply $N_0 \widehat{\beta}$. 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 $([], N_0); \Psi \vdash^{AP} e \Rightarrow \sigma \hookrightarrow (S_1, N_1)$, then for any substitution V with $\text{dom}(V) = \text{fmv}(S_1 \Psi, S_1 \sigma)$, we have $V S_1 \Psi \vdash^{AP} e \Rightarrow V S_1 \sigma$.*

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

3.5 DISCUSSION

3.5.1 COMBINING APPLICATION AND CHECKING MODES

Ningning: maybe move to future work

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:

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

3 System AP

Here, e is checking using its annotation σ_1 , and then we instantiate σ_1 to σ_2 using subtyping with application context Σ .

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 σ_1 serves as the type of x , and typing checks the body with σ_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, $\mathbf{fst}(\lambda x. x, 1)$, then it is possible to have a judgment with a non-empty application context. Therefore, we can use the application subtyping to account for the application contexts:

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

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

However, in polymorphic type systems, we need to take the subsumption rule into consideration. For example, in the expression $(\lambda x : \forall a. \forall b. (a, b). \mathbf{fst} x)$, \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:

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

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 σ_1 into the application context. Since Ψ tracks type equalities, we apply Ψ as a type substitution to σ_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 σ_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 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: maybe 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 **let** $a = \text{Int}$ **in** $\lambda x : a. x + 1$ into $(\lambda a. \lambda x : a. x + 1) \text{Int}$, which is ill-typed because the type of x in the abstraction body is a and not 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-RANKED POLYMORPHISM AND GRADUAL TYPING

4 HIGHER RANK GRADUAL TYPES

PART IV

UNIFICATION AND TYPE-INFERENCE FOR DEPENDENT TYPES

5 UNIFICATION WITH PROMOTION

6

DEPENDENT TYPES

PART V

RELATED AND FUTURE WORK

7 RELATED WORK

8 FUTURE WORK

PART VI

EPILOGUE

9 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 page 4]
- 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]
- Gang Chen. 2003. Coercive Subtyping for the Calculus of Constructions (*POPL '03*). 10. [cited on page 42]
- Thierry Coquand. 1996. An algorithm for type-checking dependent types. *Science of Computer Programming* 26, 1-3 (1996), 167–177. [cited on page 15]
- 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 and 10]
- 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 page 15]
- 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, and 19]

- 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, 23, and 26]
- 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 page 4]
- J. Roger Hindley. 1969. The Principal Type-Scheme of an Object in Combinatory Logic. *Trans. Amer. Math. Soc.* 146 (1969), 29–60. [cited on page 7]
- Assaf J Kfoury and Jerzy Tiuryn. 1992. Type reconstruction in finite rank fragments of the second-order λ -calculus. *Information and computation* 98, 2 (1992), 228–257. [cited on page 11]
- Andres Löh, 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]
- 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]
- 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, and 12]
- 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, 16, 27, 29, 35, 37, and 42]
- 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 45]
- Paula Severi and Erik Poll. 1994. Pure Type Systems with Definitions. *Logical Foundations of Computer Science* (1994), 316–328. [cited on page 48]
- Jeremy Siek and Walid Taha. 2007. 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, 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 page 4]
- 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 and 14]
- 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 47]
- 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 page 5]
- 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 46]

PART VII

TECHNICAL APPENDIX

