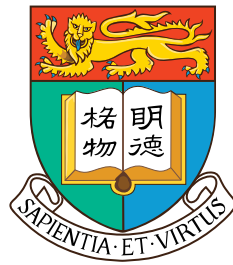


# Formalized Higher-Ranked Polymorphic Type Inference Algorithms

*by*

**Jinxu Zhao**  
(赵锦煦)



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

June 2021



Abstract of thesis entitled  
“Formalized Higher-Ranked Polymorphic Type Inference Algorithms”

Submitted by  
**Jinxu Zhao**

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

Modern functional programming languages, such as Haskell or OCaml, use sophisticated forms of type inference. In the meantime, more and more object-oriented programming languages implement advanced type inference algorithms, which greatly reduces the number of redundant and uninteresting types written by programmers, including C++11, Java 10, and C# 3.0. While being an important topic in the Programming Languages research, there is little work on the mechanization of the metatheory of type inference in theorem provers.

In particular, we are unaware of any complete formalization of the type inference algorithms that are the backbone of modern programming languages. This thesis presents the *first* full mechanical formalizations of the metatheory for three higher-ranked polymorphic type inference algorithms. Higher-ranked polymorphism is an advanced feature that enables more code reuse and has numerous applications, which is already implemented in languages like Haskell. All three systems are based on the bidirectional type system by Dunfield and Krishnaswami (DK). The DK type system has two variants, a declarative and an algorithmic one, that have been manually proven sound, complete, and decidable. While DK’s original formalization comes with very well-written manual proofs, there are several missing details and some incorrect proofs, which motivates us to fully formalize the metatheory in proof assistants.

Our first system focuses on the core problem in higher-ranked type inference algorithms—the subtyping relation. Our algorithm differs from those currently in the literature by using a novel approach based on *worklist judgments*. Worklist judgments simplify the propagation of information required by the unification process during subtyping. Furthermore, they enable a simple formulation of the meta-theoretical properties, which can be easily encoded in theorem provers. We formally prove soundness, completeness, and decidability of the subtyping algorithm w.r.t DK’s declarative specification.

The second system extends the first one with a type system. We present a mechanical formalization of DK's declarative type system with a novel algorithmic system. This system further unifies contexts with judgments, which precisely captures the scope of variables and simplifies the formalization of scoping in a theorem prover. Besides, the use of continuation-passing-style judgments allows simple formalization for inference judgments. We formally prove soundness, completeness, and decidability of the type inference algorithm. Despite the use of a different algorithm, we prove the same results as DK, although with significantly different proofs and proof techniques.

The third system is based on the second one and extended with object-oriented subtyping. In presence of object-oriented subtyping, meta-variables usually have multiple different solutions. Therefore we present a backtracking-based algorithm that non-deterministically checks against each possibility. We prove soundness w.r.t our specification, and also completeness under the rank-1 restriction.

Since such type inference algorithms are quite subtle and have a complex metatheory, mechanical formalizations are an important advance in type-inference research. With machine-checked proofs, there is little chance that any logical derivation goes wrong. In this thesis, all the properties we declare are fully formalized in the Abella theorem prover.

---

*An abstract of exactly 483 words*

# 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.

.....

**Jinxu Zhao**

June 2021

## ACKNOWLEDGMENTS

First, I would like to give my most sincere thanks to my supervisor Prof. Bruno C. d. S. Oliveira. Six years ago, I was lucky to get a chance exchanging to HKU and took two courses taught by Bruno, and both of them are quite challenging and interesting at the same time. Since then, I decide to explore more on programming languages and therefore applied for Bruno's Ph.D. student. It turns out to be a wise choice! Regular meetings with him is usually quite informative, encouraging and fun. Bruno is a very patient person; he always provide professional suggestions on how to proceed the research when I get stuck or move in the wrong direction. In short, I would never expect a better supervisor.

Collaborations in research have been unforgettable experiences. Prof. Tom Schrijvers offered help on my type inference project, giving me invaluable ideas and suggestions on both research, writing, and presentation. This piece of work contributes to the most recognized publications during my Ph.D., which won the distinguished paper award in ICFP. I also participated in projects of my colleagues Xuejing Huang and Yaoda Zhu, which are both challenging and interesting. I learnt a lot from you as well. I would like to thank everyone in the group, for your enlightening seminars and discussions: Tomas Tauber, Huang Li, Xuan Bi, Haoyuan Zhang, Yanlin Wang, Yanpeng Yang, Weixin Zhang, Ningning Xie, Xuejing Huang, Yaoda Zhou, Baber Rehman, Mingqi Xue, Yaozhu Sun, Wenjia Ye, Xu Xue, Chen Cui, Jinhao Tan.

HKU, CS Dr. François Pottier

Parents and family

Wife

# CONTENTS

DECLARATION	I
ACKNOWLEDGMENTS	II
LIST OF FIGURES	VII
LIST OF TABLES	IX
I PROLOGUE	1
1 INTRODUCTION	2
1.1 Type Systems and Type Inference Algorithms . . . . .	2
1.1.1 Functional Programming and System F . . . . .	4
1.1.2 Hindley-Milner Type System . . . . .	5
1.1.3 Higher-Ranked Polymorphism . . . . .	6
1.1.4 Bidirectional Typing . . . . .	7
1.1.5 Subtyping . . . . .	8
1.2 Mechanical Formalizations and Theorem Provers . . . . .	9
1.3 Contributions and Outline . . . . .	10
2 BACKGROUND	13
2.1 Hindley-Milner Type System . . . . .	13
2.1.1 Declarative System . . . . .	13
2.1.2 Algorithmic System and Principality . . . . .	15
2.2 Odersky-Läufer Type System . . . . .	16
2.2.1 Higher-Ranked Types . . . . .	17
2.2.2 Declarative System . . . . .	18
2.2.3 Relating to HM . . . . .	20
2.3 Dunfield-Krishnaswami Bidirectional Type System . . . . .	21
2.3.1 Declarative System . . . . .	21

2.4	MLsub . . . . .	23
2.4.1	Types and Polar Types . . . . .	24
2.4.2	Biunification . . . . .	25
II	HIGHER-RANKED TYPE INFERENCE ALGORITHMS . . . . .	27
3	HIGHER-RANKED POLYMORPHISM SUBTYPING ALGORITHM . . . . .	28
3.1	Overview: Polymorphic Subtyping . . . . .	28
3.1.1	Declarative Polymorphic Subtyping . . . . .	28
3.1.2	Finding Solutions for Variable Instantiation . . . . .	29
3.1.3	The Worklist Approach . . . . .	31
3.2	A Worklist Algorithm for Polymorphic Subtyping . . . . .	32
3.2.1	Syntax and Well-Formedness of the Algorithmic System . . . . .	32
3.2.2	Algorithmic Subtyping . . . . .	33
3.3	Metatheory . . . . .	36
3.3.1	Transfer to the Declarative System . . . . .	37
3.3.2	Soundness . . . . .	38
3.3.3	Completeness . . . . .	38
3.3.4	Decidability . . . . .	38
3.4	The Choice of Abella . . . . .	40
3.4.1	Statistics and Discussion . . . . .	42
4	A TYPE-INFERENCE ALGORITHM FOR HIGHER-RANKED POLYMORPHISM . . . . .	43
4.1	Overview . . . . .	44
4.1.1	DK's Declarative System . . . . .	45
4.1.2	DK's Algorithm . . . . .	47
4.1.3	Judgment Lists . . . . .	50
4.1.4	Single-Context Worklist Algorithm for Subtyping . . . . .	51
4.1.5	Algorithmic Type Inference for Higher-Ranked Types: Key Ideas . . . . .	51
4.2	Algorithmic System . . . . .	53
4.2.1	Syntax and Well-Formedness . . . . .	53
4.2.2	Algorithmic System . . . . .	55
4.3	Metatheory . . . . .	61
4.3.1	Declarative Worklist and Transfer . . . . .	61
4.3.2	Non-Overlapping Declarative System . . . . .	63
4.3.3	Soundness . . . . .	66
4.3.4	Completeness . . . . .	67



4.3.5	Decidability . . . . .	68
4.3.6	Abella and Proof Statistics . . . . .	71
4.4	Discussion . . . . .	71
4.4.1	Contrasting Our Scoping Mechanisms with DK's . . . . .	71
4.4.2	Elaboration . . . . .	75
4.4.3	Lexically-Scoped Type Variables . . . . .	76
5	HIGHER-RANKED POLYMORPHISM WITH OBJECT-ORIENTED SUBTYPING . . . . .	78
5.1	Overview . . . . .	78
5.1.1	Type Inference in Presence of Subtyping . . . . .	78
5.1.2	Judgment List and Eager Substitution . . . . .	80
5.1.3	Our Solution: Backtracking Algorithm . . . . .	81
5.2	Declarative System . . . . .	82
5.3	Backtracking Algorithm . . . . .	84
5.3.1	Syntax . . . . .	84
5.3.2	Algorithmic Subtyping . . . . .	84
5.3.3	Algorithmic Typing . . . . .	87
5.4	Metatheory . . . . .	89
5.4.1	Declarative Properties . . . . .	89
5.4.2	Transfer . . . . .	92
5.4.3	Soundness . . . . .	93
5.4.4	Partial Completeness of Subtyping: Rank-1 Restriction . . . . .	93
5.4.5	Algorithmic Rank-1 Restriction (Partial Completeness) . . . . .	93
5.4.6	Termination . . . . .	94
5.4.7	Formalization in the Abella Proof Assistant . . . . .	95
5.5	Discussion . . . . .	96
5.5.1	A Complete Algorithm Under Monotype Guessing Restrictions . . . . .	96
5.5.2	Lazy Substitution and Non-terminating Loops . . . . .	97
III	RELATED WORK . . . . .	100
6	RELATED WORK . . . . .	101
6.1	Higher-Ranked Polymorphic Type Inference Algorithms . . . . .	101
6.1.1	Predicative Algorithms . . . . .	101
6.1.2	Impredicative Algorithms . . . . .	102
6.2	Type Inference Algorithms with Subtyping . . . . .	104

## Contents

6.3	Techniques Used in Type Inference Algorithms . . . . .	105
6.3.1	Ordered Contexts in Type Inference . . . . .	105
6.3.2	The Essence of ML Type Inference . . . . .	105
6.3.3	Lists of Judgments in Unification . . . . .	105
6.4	Mechanical Formalization of Polymorphic Type Systems . . . . .	106
IV	EPILOGUE . . . . .	107
7	CONCLUSION AND FUTURE WORK . . . . .	108
7.1	Conclusion . . . . .	108
7.2	Future Work . . . . .	109
	BIBLIOGRAPHY . . . . .	113

## LIST OF FIGURES

2.1	HM Syntax . . . . .	14
2.2	HM Type System . . . . .	14
2.3	Syntax of Odersky-Läufer System . . . . .	18
2.4	Well-formedness of types in the Odersky-Läufer System . . . . .	18
2.5	Subtyping of the Odersky-Läufer System . . . . .	19
2.6	Typing of the Odersky-Läufer System . . . . .	20
2.7	Syntax of Declarative System . . . . .	21
2.8	Declarative Well-formedness and Subtyping . . . . .	22
2.9	Declarative Typing . . . . .	23
2.10	Types of MLsub . . . . .	24
3.1	Syntax of Declarative System . . . . .	29
3.2	Well-formedness of Declarative Types and Declarative Subtyping . . . . .	29
3.3	Syntax and Well-Formedness judgment for the Algorithmic System. . . . .	33
3.4	Algorithmic Subtyping . . . . .	34
3.5	A Success Derivation for the Algorithmic Subtyping Relation . . . . .	35
3.6	A Failing Derivation for the Algorithmic Subtyping Relation . . . . .	35
3.7	Transfer Rules . . . . .	37
3.8	Statistics for the proof scripts . . . . .	42
4.1	Syntax of Declarative System (Extends Figure 3.1) . . . . .	45
4.2	Declarative Well-formedness and Subtyping . . . . .	45
4.3	Declarative Typing . . . . .	46
4.4	Extended Syntax and Well-Formedness for the Algorithmic System . . . . .	54
4.5	Algorithmic Typing . . . . .	56
4.6	A Sample Derivation for Algorithmic Typing . . . . .	61
4.7	Declarative Worklists and Instantiation . . . . .	62
4.8	Declarative Transfer . . . . .	63
4.9	Context Subtyping . . . . .	65
4.10	Worklist measures . . . . .	69
4.11	Worklist Update . . . . .	70

## *List of Figures*

5.1	Declarative Syntax . . . . .	82
5.2	Declarative Subtyping . . . . .	82
5.3	Declarative Typing . . . . .	83
5.4	Algorithmic Syntax . . . . .	84
5.5	Algorithmic Garbage Collection and Subtyping . . . . .	85
5.6	Algorithmic Typing . . . . .	88
5.7	Declarative Worklists and Instantiation . . . . .	92
5.8	Declarative Transfer . . . . .	92

## LIST OF TABLES

4.1	Statistics for the proof scripts . . . . .	72
4.2	Translation Table for the Proof Scripts . . . . .	73
5.1	Statistics for the proof scripts . . . . .	95

## PART I

### PROLOGUE

# 1 INTRODUCTION

## 1.1 TYPE SYSTEMS AND TYPE INFERENCE ALGORITHMS

Statically typed programming languages are widely used nowadays. Programs are categorized by various types before they are compiled and executed. Type errors caught before execution usually indicate potential bugs, letting the programmers realize and correct such errors in advance.

In the early stages, programming languages like Java (before 1.5) were built on a simple type system, where only features like primitive types, explicitly-typed functions, and non-generic classes are supported. People soon realized the need to generalize similar programs that have different types when used. For example, a simple way to define a function that swaps the first two items in an integer array is

```
void swap2(int[] arr) {  
    int t = arr[0];  
    arr[0] = arr[1];  
    arr[1] = t;  
}
```

Similarly, the swap function for a float array can be defined as

```
void swap2(float[] arr) {  
    float t = arr[0];  
    arr[0] = arr[1];  
    arr[1] = t;  
}
```

which mostly shares the same body with the above definition for integer array, except that the type of element in the array changes from `int` to `float`. If such functionality is a commonly used one, such as the sorting function, we definitely want to define it once and use it on many types, such as `int`, `float`, `double`, `String`, etc. Luckily, later versions of Java (1.5 and above) provides a way to define a *generic* function that accepts input of different types:

```

<T> void swap2_generic(T[] arr) {
    T t = arr[0];
    arr[0] = arr[1];
    arr[1] = t;
}

```

where  $T$  denotes a generic type that programmers may arbitrarily pick. The `swap2_generic` function utilizes the feature of Java's type system, generics, to improve modularity. The following program invokes the generic function (suppose we defined the `swap2_generic` function in a `Utils` class as a static method)

```

Double[] arr = new Double[2];
arr[0] = 1.0; arr[1] = 2.0;
Utils.<Double>swap2_generic(arr);
System.out.println(arr[0] + " " + arr[1]);

```

However, the type parameter `<Double>` seems a bit redundant: given the type of `arr` is `Double[]`, the generic variable  $T$  can only be `Double`. In fact, with the help of *type inference*, we can simply write

```

Utils.swap2_generic(arr);

```

to call the function. When compiling the above code, type inference algorithms help programmers to fill in the missing type automatically, thus saving them from writing redundant code.

From the example above, we learn that the generics of Java together with type inference algorithms used in the compiler help programmers to write generic functions, given that the functionality is generic in nature. In other words, good features introduced to type systems accept more meaningful programs.

On the other hand, being able to accept all syntactically correct programs, as dynamically-typed programming languages do, is not desirable as well. Ill-typed programs are easy to write by mistake, like `"3" / 3`; or even well-typed programs with ambiguous/unexpected meaning like `"3" + 4 + 5` (for someone who is not familiar with the conventions, she might think this evaluates to "39"). Or even more commonly seen in practice, an accidentally misspelled variable name does not cause a runtime error until the line of code is actually executed. Type systems are designed to prevent such problems from happening, therefore statically-typed languages ensure type-safety, or "well-typed programs cannot go wrong". Type inference algorithms that come along with modern type systems help eliminate trivial or redundant parts of programs to improve the conciseness.



## 1.1.1 FUNCTIONAL PROGRAMMING AND SYSTEM F

Nowadays, more and more programming languages support the functional programming paradigm, where functions are first-class citizens, and programs are mostly constructed with function applications. Functional programming originates from the lambda calculus [Church 1932]. The simply-typed lambda calculus [Church 1941] extends the lambda calculus with a simple static type checking algorithm, preventing ill-typed programs before actual execution. However, the system does not allow polymorphic functions and thus is quite tedious to express higher-level logic.

In order to improve the expressiveness of functional programming languages, System F [Girard 1971, 1972; Reynolds 1974] introduces polymorphism via the universal quantifier  $\forall$  for types and the  $\Lambda$  binder (to introduce type-level functions) for expressions. For example, the identity function that can range over any type of the form  $A \rightarrow A$  can be encoded in System F:

$$\text{id} = \Lambda a. \lambda x : a. x : \forall a. a \rightarrow a$$

To enjoy the polymorphism of such a function, one needs to first supply a type argument like `id @Int` (we use the `@` sign to denote a type-level application), so that the polymorphic function is *instantiated* with a concrete type.

**IMPLICIT PARAMETRIC POLYMORPHISM** Although being much more expressive than simply-typed systems, plain System F is still tedious to use. It feels a bit silly to write `id @Int 3` compared with `id 3`, because the missing type is quite easy to figure out, in presence of the argument, which is of type `Int` and should represent  $a$  for the function application at the same time. With *implicit parametric polymorphism* [Milner 1978], type arguments like `@Int` are not written by the programmer explicitly, in contrast, it is the type inference algorithm's responsibility to guess them.

Practically speaking, in Java we can define the identity function as well,

```
<T> T id(T x) {
    return x;
}
```

And we typically use the function without specifying the type parameters, because the type system already supports implicit parametric polymorphism:

```
int n = id(3);
String s = id("3");
```

Theoretically speaking, it is unfortunate that there does not exist such a perfect algorithm that can automatically guess missing type applications for every possible System F

program [Tiuryn and Urzyczyn 1996]. For example, the following expression is ambiguous when the implicit type argument is not given:

$$f = (\text{choose} : \forall a. a \rightarrow a \rightarrow a) (\text{id} : \forall b. b \rightarrow b)$$

It is unclear how the type variable is instantiated during the polymorphic application: one possibility is that  $a$  is chosen to be the type of  $\text{id}$ , or  $\forall b. b \rightarrow b$ , resulting in the type  $(\forall b. b \rightarrow b) \rightarrow (\forall b. b \rightarrow b)$  (Note that we cannot express this type in programming languages like Java, simply because the  $\forall$  quantifiers do not appear at the top level, and its type system is already a restricted version of System F). However, another valid type is  $\forall b. (b \rightarrow b) \rightarrow (b \rightarrow b)$ , which is obtained by first instantiating  $\text{id}$  with a fresh type variable  $b$ , and generalizing the type after calculating the type of the application. Furthermore, between the two possible types, neither one is better: there exist programs that type check under either one of them and fail to type check under the other.

The fact that implicit parametric algorithm for full System F is impossible motivates people to discover restrictions on the type system under which type inference algorithms are capable of guessing the best types.

### 1.1.2 HINDLEY-MILNER TYPE SYSTEM

The Hindley-Milner (henceforth denoted as HM) type system [Damas and Milner 1982; Hindley 1969; Milner 1978] restricts System F types to *type schemes*, or *first-order polymorphism*, where polymorphic types can only have universal quantifiers in the top level. For example,  $\forall a b. (a \rightarrow b) \rightarrow a \rightarrow b$  is allowed, but not  $(\forall b. b \rightarrow b) \rightarrow (\forall b. b \rightarrow b)$ . The type system of many programming languages like Java adopts the idea from HM, where generics is a syntax to express polymorphic types in HM: all the generic variables must be declared at the top level before the function return type. An important property of the HM type inference algorithm is *principality*, where any unannotated program can be inferred to a most general type within its type system. This supports full type-inference without any type annotations.

For example, the following function

$$g = \lambda x. \lambda y. x$$

can be assigned to types  $\text{Int} \rightarrow \text{Bool} \rightarrow \text{Int}$ ,  $\text{Int} \rightarrow \text{Int} \rightarrow \text{Int}$  or infinitely many others. The HM inference algorithm will infer a more general type  $\forall a. \forall b. a \rightarrow b \rightarrow a$ . In order to

use the function as the types mentioned above, a more-general-than relation  $\leq$  is used to describe that the polymorphic type can be instantiated to more concrete types:

$$\begin{aligned} \forall a. \forall b. a \rightarrow b \rightarrow a &\leq \text{Int} \rightarrow \text{Bool} \rightarrow \text{Int} \\ \forall a. \forall b. a \rightarrow b \rightarrow a &\leq \text{Int} \rightarrow \text{Int} \rightarrow \text{Int} \\ \forall a. \forall b. a \rightarrow b \rightarrow a &\leq \forall a. a \rightarrow a \rightarrow a \end{aligned}$$

**PREDICATIVITY** In the HM system,  $\forall$  quantifiers can appear only on the top level, type instantiations will always be *monotypes*, i.e. types without the  $\forall$  quantifier. We refer to such a system as *predicative*. In contrast, System F does not restrict the types to instantiate, thus being an impredicative system. An important challenge is that full type inference for impredicative polymorphism is known to be undecidable [Wells 1999]. There are works that focus on practical inference of impredicative systems [Emrich et al. 2020; Le Botlan and Rémy 2003; Leijen 2008; Serrano et al. 2020, 2018; Vytiniotis et al. 2008]. However, throughout this work, we study predicative type systems only.

### 1.1.3 HIGHER-RANKED POLYMORPHISM

As functional languages evolved, the need for more expressive power has motivated language designers to look beyond HM, where there is still one obvious weakness that prevents some useful programs to type check: HM only have types of rank-1, since all the  $\forall$ 's appear on the top level. Thus one expected feature is to allow *higher-ranked polymorphism* where polymorphic types can occur anywhere in a type signature. This enables more code reuse and more expressions to type check, and has numerous applications [Gill et al. 1993; Jones 1995; Lämmel and Jones 2003; Launchbury and Peyton Jones 1995].

One of the interesting examples is the ST monad [Launchbury and Peyton Jones 1994] of Haskell, where the `runST` function is only possible to express in a rank-2 type system:

$$\text{runST} :: \forall a. (\forall s. \text{ST } s \ a) \rightarrow a$$

The type is rank-2 because of the inner  $\forall$  quantifier in the argument position of the type. Such a type encapsulates the state and makes sure that program states from different computation threads do not escape their scopes, otherwise the type checker should reject in advance.

In order to support higher-ranked types, we need to extend the type system of HM, but not taking a too big step since type inference for full System F would be impossible. A simple polymorphic subtyping relation proposed by Odersky and Läufer [1996] extends the HM system by allowing higher-ranked types, but instantiations are still limited to monotypes, thus the system remains predicative.

## 1.1.4 BIDIRECTIONAL TYPING

In order to improve the expressiveness for higher-ranked systems, some type annotations are necessary to guide type inference. In response to this challenge, several decidable type systems requiring some annotations have been proposed [Dunfield and Krishnaswami 2013; Le Botlan and Rémy 2003; Leijen 2008; Peyton Jones et al. 2007; Serrano et al. 2018; Vytiniotis et al. 2008].

As an example,

$$\text{hpoly} = \lambda(f : \forall a. a \rightarrow a). (f\ 1, f\ 'c')$$

the type of `hpoly` is  $(\forall a. a \rightarrow a) \rightarrow (\text{Int}, \text{Char})$ , which is a rank-2 type and is not typeable in HM. Notably (and unlike Hindley-Milner) the lambda argument  $f$  requires a *polymorphic* type annotation. This annotation is needed because the single universal quantifier does not appear at the top-level. Instead, it is used to quantify a type variable  $a$  used in the first argument of the function.

One of the key features that improve type inference algorithms with type annotations is bidirectional typing [Pierce and Turner 2000], a technique that combines two modes of typing: type checking, which checks an expression against a given type, and type synthesis (inference), which infers a type from an expression. Bidirectional typing is quite useful when the language supports type annotations, because those “hints” are handled with the checking mode. With bidirectional type-checking the typing judgment has two modes. In the checking mode  $\Gamma \vdash e \Leftarrow A$  both  $e$  and  $A$  are inputs: i.e. we check whether expression  $e$  has some type  $A$ . In the synthesis mode  $\Gamma \vdash e \Rightarrow A$  only  $e$  is an input: we need to calculate the output type  $A$  from the input expression  $e$ . It is clear that, with less information, the synthesis mode is more difficult to implement than the checking mode. Therefore, bidirectional type checking with type annotations provides a way for the programmer to guide the type inference algorithm when the expression is tricky to analyse, especially in the case where higher-ranked types are involved. In addition, bidirectional typing algorithms improve the quality of error messages in practice, due to the fact that they report errors in a relatively local range, compared with global unification algorithms.

Two closely related type systems that support *predicative* higher-ranked type inference were proposed by Peyton Jones et al. [Peyton Jones et al. 2007] and Dunfield and Krishnaswami [Dunfield and Krishnaswami 2013] (henceforth denoted as DK). These type systems are popular among language designers and their ideas have been adopted by several modern functional languages, including Haskell, Unison [Chiusano and Bjarnason 2015] and PureScript [Freeman 2017] among others.

DK developed a higher-ranked global bidirectional type system based on the declarative system by Odersky and Läufer [Odersky and Läufer 1996]. Beyond the existing works, they

introduce a third form of judgment, the application inference  $A \bullet e \Rightarrow C$ , where the function type  $A$  and argument expression  $e$  are input, and type  $C$  is the output representing the result type of the function application  $(f :: A) e$ . Note that  $f$  does not appear in this relation, and one can tell the procedure for type checking application expressions from the shape of the judgment — firstly, the type of the function is inferred to  $A$ ; then the application inference judgment is reduced to a checking judgment to verify if the argument  $e$  checks against the argument part of  $A$ ; finally, output the return part of  $A$ . Formally, the rules implement the procedure we described above:

$$\frac{\Psi \vdash e_1 \Rightarrow A \quad \Psi \vdash A \bullet e_2 \Rightarrow C}{\Psi \vdash e_1 e_2 \Rightarrow C} \text{Dec1} \rightarrow \text{E} \qquad \frac{\Psi \vdash e \Leftarrow A}{\Psi \vdash A \rightarrow C \bullet e \Rightarrow C} \text{Dec1} \rightarrow \text{App}$$

The use of application inference judgment avoids implicit instantiations of types like HM, instead, when the function type  $A$  is a polymorphic type, it is explicitly instantiated by the application inference until it becomes a function type:

$$\frac{\Psi \vdash \tau \quad \Psi \vdash [\tau/a]A \bullet e \Rightarrow C}{\Psi \vdash \forall a. A \bullet e \Rightarrow C} \text{Dec1} \forall \text{App}$$

As a result, DK is in a more syntax-directed system compared with HM-like systems.

DK also provided an elegant formalization of their sound and complete algorithm, which has also inspired implementations of type inference in some polymorphic programming languages (such as PureScript [Freeman 2017] or DDC [Disciple Development Team 2017]).

The focus of this thesis is also on predicative implicit higher-ranked bidirectional type inference algorithms.

### 1.1.5 SUBTYPING

The term “subtyping” is used as two slightly different concepts in this thesis. One of them refers to the polymorphic subtyping relation used in polymorphic type systems, which compares the degree of polymorphism between types, i.e. the more-general-than relation. Chapters 3 and 4 only focus on this type of subtyping relation. The other one is the subtyping usually seen in object-oriented programming, where there are some built-in or user-defined type conversion rules.

Type system in presence of object-oriented-style subtyping allows programmers to abstract functionalities inside a class and thus benefit from another form of polymorphism. Introduced by Cardelli [1988]; Mitchell [1984]; Reynolds [1985], subtyping is studied for explicitly typed programs. Compared with functional programming languages, mainstream

object-oriented languages lacks competitive type inference algorithms. We will further introduce object-oriented subtyping as a feature in Chapter 5, where we also introduce the top and bottom types as the minimal support for subtyping.

## 1.2 MECHANICAL FORMALIZATIONS AND THEOREM PROVERS

Although type inference is important in practice and receives a lot of attention in academic research, there is little work on mechanically formalizing such advanced forms of type inference in theorem provers. The remarkable exception is work done on the formalization of certain parts of Hindley-Milner type inference [Dubois 2000; Dubois and Menissier-Morain 1999; Garrigue 2015; Naraschewski and Nipkow 1999; Urban and Nipkow 2008]. However, there is still no formalization of the higher-ranked type systems that are employed by modern languages like Haskell. This is at odds with the current trend of mechanical formalizations in programming language research. In particular, both the POPLMark challenge [Aydemir et al. 2005] and CompCert [Leroy et al. 2012] have significantly promoted the use of theorem provers to model various aspects of programming languages. Today, papers in various programming language venues routinely use theorem provers to mechanically formalize: *dynamic and static semantics* and their correctness properties [Aydemir et al. 2008], *compiler correctness* [Leroy et al. 2012], *correctness of optimizations* [Bertot et al. 2006], *program analysis* [Chang et al. 2006] or proofs involving *logical relations* [Abel et al. 2018].

**MOTIVATIONS FOR MECHANICAL FORMALIZATIONS.** The main argument for mechanical formalizations is a simple one. Proofs for programming languages tend to be *long*, *tedious*, and *error-prone*. In such proofs, it is very easy to make mistakes that may invalidate the whole development. Furthermore, readers and reviewers often do not have time to look at the proofs carefully to check their correctness. Therefore errors can go unnoticed for a long time. In fact, manual proofs are commonly observed to have flaws that invalidate the claimed properties. For instance, Klein et al. [2012] reproduced the proofs of nine ICFP 2009 papers in Redex, and found problems in each one of them. We also found false lemmas and incorrect proofs in DK’s manual proof [Dunfield and Krishnaswami 2013]. Mechanical formalizations provide, in principle, a natural solution for these problems. Theorem provers can automatically check and validate the proofs, which removes the burden of checking from both the person doing the proofs as well as readers or reviewers.

Moreover, extending type-inference algorithms with new programming language features is often quite delicate. Studying the meta-theory for such extensions would be greatly aided by the existence of a mechanical formalization of the base language, which could then be reused and extended by the language designer. Compared with manual proofs which may

take a long time before one can fully understand every detail, theorem provers can quickly point out proofs that are invalidated after extensions.

**CHALLENGES IN VARIABLE BINDING AND ABELLA.** Handling variable binding is particularly challenging in type inference, because the algorithms typically do not rely simply on local environments, but instead propagate information across judgments. Yet, there is little work on how to deal with these complex forms of binding in theorem provers. We believe that this is the primary reason why theorem provers have still not been widely adopted for formalizing type-inference algorithms.

The Abella theorem prover [Gacek 2008] is one that specifically eases formalization on-binders. Two common treatments of binding are to use the De Bruijn indices [de Bruijn 1972] and the nominal logic framework of Pitts [Pitts 2003]. In contrast, Abella uses the abstraction operator in a typed  $\lambda$ -calculus to encode binding. Its  $\lambda$ -tree syntax, or HOAS, and features including the  $\nabla$  quantifier and higher-order unification, allow for better experiences than using Coq libraries utilizing other approaches. In practice, Abella uses the  $\nabla$  (nabla) quantifier and nominal constants to help quantify a “fresh” variable during formalization. For example, the common type checking rule

$$\frac{e \Leftarrow A \quad a \text{ fresh}}{e \Leftarrow \forall a. A}$$

is encoded as

`check E (all A) := nabla a, check E (A a)`

in Abella, where the  $\nabla$  quantifier introduces a fresh type variable  $a$  and later use it to “open” the body of  $\forall a. A$ .

Throughout the thesis, all the type systems and declared properties are mechanically formalized in Abella.

### 1.3 CONTRIBUTIONS AND OUTLINE

**CONTRIBUTIONS** In this thesis, we propose variants of type inference algorithms for higher-ranked polymorphic type systems and formalize each of them in the Abella theorem prover. It is the first work on higher-ranked type inference that comes with mechanical formalizations. In summary, the main contributions of this thesis are:

- **Chapter 3** presents a predicative polymorphic subtyping algorithm.

- We proved that our algorithm is *sound*, *complete*, and *decidable* with respect to OLs higher-ranked subtyping specification in the Abella theorem prover. And we are the first to formalize the meta-theoretical results of a polymorphic higher-ranked subtyping algorithm.
- Similar to DK’s algorithm, we employ an ordered context to collect type variables and existential variables (used as placeholders for guessing monotypes). However, our unification process is novel. DK’s algorithm solves variables on-the-fly and communicates the partial solutions through an output context. In contrast, our algorithm collects a list of judgments and propagate partial solutions across them via eager substitutions. Such technique eliminates the use of output contexts, and thus simplifies the metatheory and makes mechanical formalizations easier. Besides, using only a single context keeps the definition of well-formedness simple, resulting in an easy and elegant algorithm.
- Chapter 4 presents a new bidirectional higher-ranked typing inference algorithm based on DK’s declarative system.
  - We are the first to present a full mechanical formalization for a type inference algorithm of higher-ranked type system. The *soundness*, *completeness*, and *decidability* are shown in the Abella theorem prover, with respect to DK’s declarative system.
  - We propose *worklist judgments*, a new technique that unifies ordered contexts and judgments. This enables precise scope tracking of variables in judgments and avoids the duplication of context information across judgments in worklists. Similar to the previous algorithm, partial solutions are propagated across judgments in a single list consist of both variable bindings and judgments. Nevertheless, the unification of worklists and contexts exploits the fact that judgments are usually sharing a large part of common information. And one can easily tell when a variable is no longer referred to.
  - Furthermore, we support inference judgments so that bidirectional typing can be encoded as worklist judgments. The idea is to use a continuation passing style to enable the transfer of inferred information across judgments.
- Chapter 5 further extends the higher-ranked system with object-oriented subtyping.
  - We propose a bidirectional declarative system extended with the top and bottom types and relevant subtyping and typing rules. Several desirable properties are satisfied and mechanically proven.



- A new backtracking-based worklist algorithm is presented and proven to be *sound* with respect to our declarative specification in the Abella theorem prover. Extended with subtyping relations of the top and bottom types, simple solutions such as  $\top$  or  $\perp$  satisfies subtyping constraints in parallel with other solutions which does not involve object-oriented subtyping. Our backtracking technique is specifically well-suited for the non-deterministic trial without missing any of them.
- We also formalize the rank-1 restriction of subtyping relation, and proved that our algorithmic subtyping is *complete* under such restriction.

PRIOR PUBLICATIONS. This thesis is partially based on the publications by the author [Zhao et al. 2018, 2019], as indicated below.

Chapter 3: Jinxu Zhao, Bruno C. d. S. Oliveira, and Tom Schrijvers. 2018. “Formalization of a Polymorphic Subtyping Algorithm”. In *International Conference on Interactive Theorem Proving (ITP)*.

Chapter 4: Jinxu Zhao, Bruno C. d. S. Oliveira, and Tom Schrijvers. 2019. “A Mechanical Formalization of Higher-Ranked Polymorphic Type Inference”. In *International Conference on Functional Programming (ICFP)*.

MECHANIZED PROOFS. All proofs in this thesis is mechanically formalized in the Abella theorem prover and are available online:

Chapter 3: <https://github.com/JimmyZJX/Abella-subtyping-algorithm>

Chapter 4: <https://github.com/JimmyZJX/TypingFormalization>

Chapter 5: <https://github.com/JimmyZJX/Dissertation/tree/main/src>

## 2 BACKGROUND

In this chapter, we introduce some highly related type systems. They are basic concepts and should help the reader understand the rest of the thesis better. Section 2.1 introduces the Hindley-Milner type system [Damas and Milner 1982; Hindley 1969; Milner 1978], a widely used system that supports rank-1 (prenex) polymorphism. Section 2.2 presents the Odersky-Läufer type system, which is an extension of the Hindley-Milner by allowing higher-ranked types. Section 2.3 describes the Dunfield-Krishnaswami bidirectional type system, a bidirectional type system that further extends Odersky-Läufer’s. Finally, in Section 2.4 we introduce one of the recent advancements in the ML family, namely the MLsub type system, which integrates object-oriented subtyping with the Hindley-Milner type inference.

### 2.1 HINDLEY-MILNER TYPE SYSTEM

The Hindley-Milner type system, hereafter called “HM” for short, is a classical lambda calculus with parametric polymorphism. Thanks to its simple formulation and powerful inference algorithm, many modern functional programming languages are still using HM as their base, including the ML family and Haskell. The system is also known as Damas–Milner or Damas–Hindley–Milner. Hindley [Hindley 1969] and Milner [Milner 1978] independently discovered algorithms for the polymorphic typing problem and also proved the soundness of their algorithms, despite that there are slight differences in the expressiveness— Milner’s system supports let-generalization. Later on, Damas and Milner [1982] proved the completeness of their algorithm.

#### 2.1.1 DECLARATIVE SYSTEM

**SYNTAX** The declarative syntax is shown in Figure 2.1. The HM types are consist of polymorphic types (or type schemes) and monomorphic types. A polymorphic type contains zero or more universal quantifiers only at the top level. When no universal quantifier occurs, the type belongs to a mono-type. Mono-types are constructed by a unit type 1, a type variable  $a$ , or a function type  $\tau_1 \rightarrow \tau_2$ .

## 2 Background

Type variables	$a, b$
Types	$\sigma ::= \tau \mid \forall a. \sigma$
Monotypes	$\tau ::= 1 \mid a \mid \tau_1 \rightarrow \tau_2$
Expressions	$e ::= x \mid () \mid \lambda x. e \mid e_1 e_2 \mid \mathbf{let} \ x = e_1 \ \mathbf{in} \ e_2$
Contexts	$\Psi ::= \cdot \mid \Psi, x : \sigma$

Figure 2.1: HM Syntax

$\sigma_1 \sqsubseteq \sigma_2$  HM Type Instantiation

$$\frac{\tau' = [\bar{\tau}/\bar{a}]\tau \quad \bar{b} \notin \text{FV}(\forall \bar{a}. \tau)}{\forall \bar{a}. \tau \sqsubseteq \forall \bar{b}. \tau'} \text{HM-TInst}$$

$\Psi \vdash_{HM} e : \sigma$  HM Typing

$$\begin{array}{c}
\frac{(x : \sigma) \in \Psi}{\Psi \vdash_{HM} x : \sigma} \text{HM-Var} \quad \frac{}{\Psi \vdash_{HM} () : 1} \text{HM-Unit} \quad \frac{\Psi, x : \tau_1 \vdash_{HM} e : \tau_2}{\Psi \vdash_{HM} \lambda x. e : \tau_1 \rightarrow \tau_2} \text{HM-Abs} \\
\\
\frac{\Psi \vdash_{HM} e_1 : \tau_1 \rightarrow \tau_2 \quad \Psi \vdash_{HM} e_2 : \tau_1}{\Psi \vdash_{HM} e_1 e_2 : \tau_2} \text{HM-App} \\
\\
\frac{\Psi \vdash_{HM} e_1 : \sigma \quad \Psi, x : \sigma \vdash_{HM} e_2 : \tau}{\Psi \vdash_{HM} \mathbf{let} \ x = e_1 \ \mathbf{in} \ e_2 : \tau} \text{HM-Let} \\
\\
\frac{\Psi \vdash_{HM} e : \sigma \quad \bar{a} \notin \text{FV}(\Psi)}{\Psi \vdash_{HM} e : \forall \bar{a}. \sigma} \text{HM-Gen} \quad \frac{\Psi \vdash_{HM} e : \sigma_1 \quad \sigma_1 \sqsubseteq \sigma_2}{\Psi \vdash_{HM} e : \sigma_2} \text{HM-Inst}
\end{array}$$

Figure 2.2: HM Type System

Expressions  $e$  includes variables  $x$ , literals  $()$ , lambda abstractions  $\lambda x. e$ , applications  $e_1 e_2$  and the let expression  $\mathbf{let} \ x = e_1 \ \mathbf{in} \ e_2$ . A context  $\Psi$  is a collection of type bindings for variables.

**TYPE INSTANTIATION** The relations between types are described via type instantiations. The rule shown to the top of Figure 2.2 checks if  $\forall \bar{a}. \tau$  is a *generic instance* of  $\forall \bar{b}. \tau'$ . This relation is valid when  $\tau' = [\bar{\tau}/\bar{a}]\tau$  for a series of mono-types  $\bar{\tau}$  and each variable in  $\bar{b}$  is not free in  $\forall \bar{a}. \tau$ .

For example,

$$\forall a. a \rightarrow a \sqsubseteq 1 \rightarrow 1$$

## 2 Background

is obtained by the substitution  $[1/a]$ , and

$$\forall a. a \rightarrow a \sqsubseteq \forall b. (b \rightarrow b) \rightarrow (b \rightarrow b)$$

substitutes  $a$  by  $b \rightarrow b$ , and generalizes  $b$  after the substitution.

**TYPING** The typing relation  $\Psi \vdash_{HM} e : \sigma$  synthesizes a type  $\sigma$  for an expression  $e$  under the context  $\Psi$ . Rule HM-Var looks up the binding of a variable  $x$  in the context. Rule HM-Unit always gives the unit type 1 to the unit expression  $()$ . For a lambda abstraction  $\lambda x. e$ , Rule HM-Abs guesses its input type ( $\tau_1$ ) and computes the type of its body ( $\tau_2$ ) as the return type. Rule HM-App eliminates a function type by an application  $e_1 e_2$ , where the argument type must be the same as the input type of the function, and the type of the whole application is  $\tau_2$ .

Rule HM-Let is also referred to as let-polymorphism. In (untyped) lambda calculus, **let**  $x = e_1$  **in**  $e_2$  behaves the same as  $(\lambda x. e_2) e_1$ . However, the HM let rule derives the type of  $e_1$  first, and binds the polymorphic type into the context before  $e_2$ . This enables polymorphic expressions to be reused multiple times in different instantiated types.

Rules HM-Gen and HM-Inst change the type of an expression at any time during the derivation. Rule HM-Gen generalizes over fresh type variables  $\bar{a}$ . Rule HM-Inst, as opposed to generalization, specializes a type according to the type instantiation relation.

The type system of HM supports *implicit instantiation* through Rule HM-Inst. This means that any expression (function) that has a polymorphic type can be automatically instantiated with a proper monotype for any reasonable application. The fact that only monotypes are guessed indicates that the system is *predicative*. In contrast, an *impredicative* system might guess polymorphic types. Unfortunately, type inference on impredicative systems is undecidable [Wells 1999]. In this thesis, we focus on predicative systems only.

### 2.1.2 ALGORITHMIC SYSTEM AND PRINCIPALITY

**SYNTAX-DIRECTED SYSTEM** The declarative system is not fully syntax-directed due to Rules HM-Gen and HM-Inst, which can be applied to any expression. A syntax-directed system can be obtained by replacing Rules HM-Var and HM-Let by the following rules:

$$\frac{(x : \sigma) \in \Psi \quad \sigma \sqsubseteq \tau}{\Psi \vdash_{HM}^S x : \tau} \text{ HM-Var-Inst} \qquad \frac{\begin{array}{l} \Psi \vdash_{HM} e_1 : \sigma \\ \bar{a} = \text{FV}(\sigma) - \text{FV}(\Psi) \\ \Psi, x : \forall \bar{a}. \sigma \vdash_{HM} e_2 : \tau \end{array}}{\Psi \vdash_{HM}^S \text{let } x = e_1 \text{ in } e_2 : \tau} \text{ HM-Let-Gen}$$

## 2 Background

A generalization on  $\sigma$ , the synthesized type of  $e_1$ , is added to Rule HM-Let, since it is the source place where a polymorphic type is generated. However, a too generalized type might reject applications due to its shape, therefore, an instantiation procedure is added to eliminate all the universal quantifiers on Rule HM-Var. We omit Rules HM-Unit, HM-Abs, and HM-App for the syntax-directed system  $\Psi \vdash_{HM}^S$ . The following property shows that the new system is (almost) equivalent to the original declarative system.

**Theorem 2.1** (Equivalence of Syntax-Directed System).

1. If  $\Psi \vdash_{HM}^S e : \sigma$  then  $\Psi \vdash_{HM} e : \sigma$
2. If  $\Psi \vdash_{HM} e : \sigma$  then  $\Psi \vdash_{HM} e : \tau$ , and  $\forall \bar{a}. \tau \sqsubseteq \sigma$ , where  $\bar{a} = FV(\tau) - FV(\Psi)$ .

**TYPE INFERENCE ALGORITHM** Although being syntax-directed solves some problems, the rules still require some guessings, including Rules HM-Abs and HM-Var-Inst. Proposed by Milner [1978], Algorithm W, an inference algorithm of the HM type system based on unification, is proven to be sound and complete w.r.t the declarative specifications.

**Theorem 2.2** (Algorithmic Completeness (Principality)). *If  $\Psi \vdash_{HM} e : \sigma$ , then W computes a principal type scheme  $\sigma_p$ , i.e.*

1.  $\Psi \vdash_{HM} e : \sigma_p$
2.  $\sigma_p \sqsubseteq \sigma$ .

Note that both the above theorems are not formalized by the author. The reader may refer to an existing formalization of algorithm W, for example, [Naraschewski and Nipkow 1999].

### 2.2 Odersky-Läufer Type System

The HM type system is simple, powerful, and easy-to-use. However, it only accepts types of rank-1, i.e. the  $\forall$  quantifier can only appear in the top-level. In practice, there are cases where *higher-ranked* types are needed. The rank-1 limitation prevents those programs from type check and thus loses expressiveness. The problem is that full type inference for System F is proven to be undecidable [Wells 1999]. Odersky and Läufer [1996] then proposed a system where programmers can make use of type annotations to guide the type system, especially on higher-ranked types. This extension of HM preserves nice properties of HM, while accepting higher-ranked types to be checked with the help of the programmers.

For example, consider the following function definition

$$\lambda f. (f \ 1, f \ 'c')$$

## 2 Background

This is not typeable in HM, because the argument of the lambda abstraction,  $f$ , is applied to both an integer and a character, which means that it should be of a polymorphic unknown type, thus the type of the lambda abstraction cannot be inferred by the HM type system. This seems reasonable, since there are several polymorphic types that fit the function, for example,

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

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

The solution is also natural: if the programmer may pick the type of argument she wants, the type system can figure out the rest. By adding type annotation on  $f$ , OL now accepts the definition

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

and infers type  $(\forall a. a \rightarrow a) \rightarrow (\text{Int}, \text{Char})$ .

In what follows, we will first formally define the rank of a type, and then introduce the declarative system of OL, finally discuss the relationship between OL and HM.

### 2.2.1 HIGHER-RANKED TYPES

The *rank* of a type represents how deep a universal quantifier appears at the contravariant position [Kfoury and Tiuryn 1992]. Formally speaking,

$$\begin{array}{ll} \text{Rank 0 / Monotypes} & \tau, \sigma^0 ::= 1 \mid a \mid \tau_1 \rightarrow \tau_2 \\ \text{Rank } k (k \geq 1), \text{ Polytypes} & \sigma^k ::= \sigma^{k-1} \mid \sigma^{k-1} \rightarrow \sigma^k \mid \forall a. \sigma^k \end{array}$$

The following example illustrates what rank a type belongs to:

$$\begin{array}{ll} 1 \rightarrow 1 & \text{Rank 0} \\ \forall a. a \rightarrow a & \text{Rank 1} \\ 1 \rightarrow \forall a. a \rightarrow a & \text{Rank 1} \\ (\forall a. a \rightarrow a) \rightarrow (\forall a. a \rightarrow a) & \text{Rank 2} \end{array}$$

According to the definition, monotypes are types that does not contain any universal quantifier. In the HM type system, all polymorphic types have rank 1.

Type variables	$a, b$
Types	$\sigma ::= 1 \mid a \mid \sigma_1 \rightarrow \sigma_2 \mid \forall a. \sigma$
Monotypes	$\tau ::= 1 \mid a \mid \tau_1 \rightarrow \tau_2$
Expressions	$e ::= x \mid () \mid \lambda x : \sigma. e \mid e : \sigma \mid \lambda x. e \mid e_1 e_2 \mid \text{let } x = e_1 \text{ in } e_2$
Contexts	$\Psi ::= \cdot \mid \Psi, x : \sigma \mid \Psi, a$

Figure 2.3: Syntax of Odersky-Läufer System

$$\begin{array}{c}
 \frac{}{\Psi \vdash_{OL} 1} \text{OL-WF-Unit} \qquad \frac{a \in \Psi}{\Psi \vdash_{OL} a} \text{OL-WF-TVar} \\
 \frac{\Psi \vdash_{OL} \sigma_1 \quad \Psi \vdash_{OL} \sigma_2}{\Psi \vdash_{OL} \sigma_1 \rightarrow \sigma_2} \text{OL-WF-Arr} \qquad \frac{\Psi, a \vdash_{OL} \sigma}{\Psi \vdash_{OL} \forall a. \sigma} \text{OL-WF-Forall}
 \end{array}$$

Figure 2.4: Well-formedness of types in the Odersky-Läufer System

### 2.2.2 DECLARATIVE SYSTEM

The syntax of Odersky-Läufer system is shown in Figure 2.3. There are several differences compared to the HM system.

First, polymorphic types can be of arbitrary rank, i.e. the forall quantifier may occur at any part of a type. Yet, mono-type remains the same definition as HM's.

Second, expressions now allows annotations  $e : \sigma$  and (argument) annotated lambda functions  $\lambda x : \sigma. e$ . Annotations on expressions help guide the type system properly, acting as a machine-checked document by the programmers. By annotating the argument of a lambda function with a polymorphic type  $\sigma$ , one may encode a function of higher rank in this system compared to HM's.

Finally, contexts consist of not only variable bindings, but also type variable declarations. Here we adopt a slightly different approach than the original work [Odersky and Läufer 1996], which does not track type variables explicitly in a context. Such explicit declarations reduce formalization difficulties, especially when dealing with freshness conditions or variable name encodings. This also enables us to formally define the well-formedness of types, shown in Figure 2.4.

**SUBTYPING** The subtyping relation, defined in Figure 2.5, is more powerful than that (type instantiation) of HM. It compares the degree of polymorphism between two types. Essentially, if  $A$  can always be instantiated to match any instantiation of  $B$ , then  $A$  is “at

## 2 Background

$$\begin{array}{c}
\frac{}{\Psi \vdash_{OL} 1 \leq 1} \text{OL-SUB-Unit} \qquad \frac{a \in \Psi}{\Psi \vdash_{OL} a \leq a} \text{OL-SUB-Var} \\
\frac{\Psi \vdash_{OL} \sigma'_1 \leq \sigma_1 \quad \Psi \vdash_{OL} \sigma_2 \leq \sigma'_2}{\Psi \vdash_{OL} \sigma_1 \rightarrow \sigma_2 \leq \sigma'_1 \rightarrow \sigma'_2} \text{OL-SUB-Arr} \\
\frac{\Psi \vdash_{OL} \tau \quad \Psi \vdash_{OL} [\tau/a]\sigma \leq \sigma'}{\Psi \vdash_{OL} \forall a. \sigma \leq \sigma'} \text{OL-SUB-}\forall\text{L} \qquad \frac{\Psi, a \vdash_{OL} \sigma \leq \sigma'}{\Psi \vdash_{OL} \sigma \leq \forall a. \sigma'} \text{OL-SUB-}\forall\text{R}
\end{array}$$

Figure 2.5: Subtyping of the Odersky-Läufer System

least as polymorphic as”  $B$ . We also say that  $A$  is “more polymorphic than”  $B$  and write  $A \leq B$ . In contrast to HM’s subtyping, higher-ranked types can be compared thanks to Rule OL-SUB-Arr; functions are contravariant on argument types and covariant on return types.

Subtyping rules OL-SUB-Var, OL-SUB-Unit handle simple cases that do not involve universal quantifiers. Rule OL-SUB- $\forall$ R states that if  $A$  is a subtype of  $B$  in the context  $\Psi, a$ , where  $a$  is fresh in  $A$ , then  $A \leq \forall a. B$ . Intuitively, if  $A$  is more general than  $\forall a. B$  (where the universal quantifier already indicates that  $\forall a. B$  is a general type), then  $A$  must instantiate to  $[\tau/a]B$  for every  $\tau$ .

The most interesting rule is OL-SUB- $\forall$ L. If some instantiation of  $\forall a. A$ ,  $[\tau/a]A$ , is a subtype of  $B$ , then  $\forall a. A \leq B$ . The monotype  $\tau$  we used to instantiate  $a$  is *guessed* in this declarative rule, but the algorithmic system does not guess and defers the instantiation until it can determine the monotype deterministically. The fact that  $\tau$  is a monotype rules out the possibility of polymorphic (or impredicative) instantiation. However this restriction ensures that the subtyping relation remains decidable. Allowing an arbitrary type (rather than a monotype) in rule OL-SUB- $\forall$ L is known to give rise to an undecidable subtyping relation [Chrzyszcz 1998]. Peyton Jones et al. [2007] also impose the restriction of predicative instantiation in their type system. Both systems are adopted by several practical programming languages.

Note that when we introduce a new binder in the premise, we implicitly pick a fresh one, which is made possible by renaming according to alpha-equivalence. This applies to rules OL-SUB- $\forall$ R here. We follow this implicit freshness condition and omit it throughout the whole thesis.

**TYPING** The type system of Odersky-Läufer, shown in Figure 2.6, extends HM’s type system in the following aspects.

Rule OL-Lam now accepts polymorphic return type, because such type is well-formed. The guess on parameter types is still limited to monotypes like HM’s. However, if a parameter



$$\begin{array}{c}
 \frac{(x : \sigma) \in \Psi}{\Psi \vdash_{OL} x : \sigma} \text{OL-Var} \quad \frac{}{\Psi \vdash_{OL} () : 1} \text{OL-Unit} \quad \frac{\Psi \vdash_{OL} e : \sigma}{\Psi \vdash_{OL} (e : \sigma) : \sigma} \text{OL-Anno} \\
 \frac{\Psi \vdash_{OL} \tau \quad \Psi, x : \tau \vdash_{OL} e : \sigma}{\Psi \vdash_{OL} \lambda x. e : \tau \rightarrow \sigma} \text{OL-Lam} \quad \frac{\Psi, x : \sigma_1 \vdash_{OL} e : \sigma_2}{\Psi \vdash_{OL} \lambda x : \sigma_1. e : \sigma_1 \rightarrow \sigma_2} \text{OL-LamAnno} \\
 \frac{\Psi \vdash_{OL} e_1 : \sigma_1 \rightarrow \sigma_2 \quad \Psi \vdash_{OL} e_2 : \sigma_1}{\Psi \vdash_{OL} e_1 e_2 : \sigma_2} \text{OL-App} \quad \frac{\Psi, a \vdash_{OL} e : \sigma}{\Psi \vdash_{OL} e : \forall a. \sigma} \text{OL-Gen} \\
 \frac{\Psi \vdash_{OL} e_1 : \sigma_1 \quad \Psi, x : \sigma_1 \vdash_{OL} e_2 : \sigma_2}{\Psi \vdash_{OL} \text{let } x = e_1 \text{ in } e_2 : \sigma_2} \text{OL-Let} \\
 \frac{\Psi \vdash_{OL} e : \sigma_1 \quad \Psi \vdash_{OL} \sigma_1 \leq \sigma_2}{\Psi \vdash_{OL} e : \sigma_2} \text{OL-Sub}
 \end{array}$$

Figure 2.6: Typing of the Odersky-Läufer System

type is specified in advance, the type system accepts polymorphic argument type with rule OL-LamAnno. Functions of arbitrary rank can be encoded through proper annotations. The application and let-generalization rules also accept polymorphic return types.

Rule OL-Gen encodes the generalization rule of HM in a different way under explicit type variable declarations. A fresh type variable is introduced into the context before the type of expression  $e$  is calculated. Then we conclude that  $e$  has a polymorphic type by generalizing the type variable. For example, the type of the identity function is derived as follows

$$\begin{array}{c}
 \frac{(x : a) \in (\cdot, a, x : a)}{\cdot, a \vdash_{OL} a} \text{OL-Var} \\
 \frac{\cdot, a \vdash_{OL} a \quad \cdot, a, x : a \vdash_{OL} x : a}{\cdot, a \vdash_{OL} \lambda x. x : a \rightarrow a} \text{OL-Lam} \\
 \frac{\cdot, a \vdash_{OL} \lambda x. x : a \rightarrow a}{\cdot \vdash_{OL} \lambda x. x : \forall a. a \rightarrow a} \text{OL-Gen}
 \end{array}$$

The subsumption rule OL-Sub converts the type of an expression with the help of the subtyping relation.

### 2.2.3 RELATING TO HM

The OL type system accepts higher-ranked types, but it only tries to instantiate monotypes like HM. Therefore, conservatively extends HM, such that every typed expression in HM is also typed in OL. In the meantime, all the “guessing” jobs OL needs to do remains in instantiating monotypes, thus the algorithm can be extended directly from any one for HM. Formally speaking, type inference for both systems can be reduced to a problem of solving

## 2 Background

Type variables	$a, b$		
Types	$A, B, C$	$::=$	$1 \mid a \mid \forall a. A \mid A \rightarrow B$
Monotypes	$\tau, \sigma$	$::=$	$1 \mid a \mid \tau \rightarrow \sigma$
Expressions	$e$	$::=$	$x \mid () \mid \lambda x. e \mid e_1 e_2 \mid (e : A)$
Contexts	$\Psi$	$::=$	$\cdot \mid \Psi, a \mid \Psi, x : A$

Figure 2.7: Syntax of Declarative System

certain forms of unification constraints, including equalities, conjunctions, universal and existential quantifiers, and let-generalizations.

### 2.3 DUNFIELD-KRISHNASWAMI BIDIRECTIONAL TYPE SYSTEM

Bidirectional typing is popular among new type systems. Compared with the ML-style systems, bidirectional typing additionally makes use of checking mode, which checks an expression against a known type. This is especially helpful in dealing with unannotated lambda functions, and when the type of the function can be inferred from the neighbor nodes in the syntax tree. For example,

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

is typeable in higher-ranked bidirectional systems, as the outer type annotation may act as if both the argument type and return type of the lambda is given. The Dunfield-Krishnaswami type system [Dunfield and Krishnaswami 2013], hereafter referred to as DK, extends the OL system by exploiting bidirectional typing. In this section, we only introduce the declarative system and leave the discussion of their algorithmic system to Chapters 3 and 4.

#### 2.3.1 DECLARATIVE SYSTEM

**SYNTAX.** The syntax of DK's declarative system [Dunfield and Krishnaswami 2013] is shown in Figure 2.7. A declarative type  $A$  is either the unit type  $1$ , a type variable  $a$ , a universal quantification  $\forall a. A$  or a function type  $A \rightarrow B$ . Nested universal quantifiers are allowed for types, but monotypes  $\tau$  do not have any universal quantifier. Terms include a unit term  $()$ , variables  $x$ , lambda-functions  $\lambda x. e$ , applications  $e_1 e_2$  and annotations  $(e : A)$ . Contexts  $\Psi$  are sequences of type variable declarations and term variables with their types declared  $x : A$ .

**WELL-FORMEDNESS** Well-formedness of types and terms is shown at the top of Figure 2.8. The rules are standard and simply ensure that variables in types and terms are well-scoped.

## 2 Background

$$\begin{array}{c}
\boxed{\Psi \vdash A} \text{ Well-formed declarative type} \\
\\
\frac{}{\Psi \vdash 1} \text{ wf}_{\text{dunit}} \quad \frac{a \in \Psi}{\Psi \vdash a} \text{ wf}_{\text{dvar}} \quad \frac{\Psi \vdash A \quad \Psi \vdash B}{\Psi \vdash A \rightarrow B} \text{ wf}_{\text{d}\rightarrow} \quad \frac{\Psi, a \vdash A}{\Psi \vdash \forall a. A} \text{ wf}_{\text{d}\forall} \\
\\
\boxed{\Psi \vdash e} \text{ Well-formed declarative expression} \\
\\
\frac{x : A \in \Psi}{\Psi \vdash x} \text{ wf}_{\text{d}\text{tmvar}} \quad \frac{}{\Psi \vdash ()} \text{ wf}_{\text{d}\text{tmunit}} \quad \frac{\Psi, x : A \vdash e}{\Psi \vdash \lambda x. e} \text{ wf}_{\text{d}\text{abs}} \\
\\
\frac{\Psi \vdash e_1 \quad \Psi \vdash e_2}{\Psi \vdash e_1 e_2} \text{ wf}_{\text{d}\text{app}} \quad \frac{\Psi \vdash A \quad \Psi \vdash e}{\Psi \vdash (e : A)} \text{ wf}_{\text{d}\text{anno}} \\
\\
\boxed{\Psi \vdash A \leq B} \text{ Declarative subtyping} \\
\\
\frac{a \in \Psi}{\Psi \vdash a \leq a} \leq_{\text{Var}} \quad \frac{}{\Psi \vdash 1 \leq 1} \leq_{\text{Unit}} \quad \frac{\Psi \vdash B_1 \leq A_1 \quad \Psi \vdash A_2 \leq B_2}{\Psi \vdash A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2} \leq_{\rightarrow} \\
\\
\frac{\Psi \vdash \tau \quad \Psi \vdash [\tau/a]A \leq B}{\Psi \vdash \forall a. A \leq B} \leq_{\forall L} \quad \frac{\Psi, b \vdash A \leq B}{\Psi \vdash A \leq \forall b. B} \leq_{\forall R}
\end{array}$$

Figure 2.8: Declarative Well-formedness and Subtyping

**DECLARATIVE SUBTYPING** The bottom of Figure 2.8 shows DK’s declarative subtyping judgment  $\Psi \vdash A \leq B$ , which was adopted from Odersky and Läufer [1996]. Since it is exactly the same with OL’s subtyping relation, we refer to Section 2.2 for detailed discussion.

**DECLARATIVE TYPING** The bidirectional type system, shown in Figure 2.9, has three judgments. The checking judgment  $\Psi \vdash e \Leftarrow A$  checks expression  $e$  against the type  $A$  in the context  $\Psi$ . The synthesis judgment  $\Psi \vdash e \Rightarrow A$  synthesizes the type  $A$  of expression  $e$  in the context  $\Psi$ . The application judgment  $\Psi \vdash A \bullet e \Rightarrow C$  synthesizes the type  $C$  of the application of a function of type  $A$  (which could be polymorphic) to the argument  $e$ .

Many rules are standard. Rule `Dec1Var` looks up term variables in the context. Rules `Dec11I` and `Dec11I $\Rightarrow$`  respectively check and synthesize the unit type. Rule `Dec1Anno` synthesizes the annotated type  $A$  of the annotated expression  $(e : A)$  and checks that  $e$  has type  $A$ . Checking an expression  $e$  against a polymorphic type  $\forall a. A$  in the context  $\Psi$  succeeds if  $e$  checks against  $A$  in the extended context  $(\Psi, a)$ . The subsumption rule `Dec1Sub` depends on the subtyping relation, and changes mode from checking to synthesis: if  $e$  synthesizes type  $A$  and  $A \leq B$  ( $A$  is more polymorphic than  $B$ ), then  $e$  checks against  $B$ . If a checking problem does not match any other rules, this rule can be applied to synthesize a type instead and

$\boxed{\Psi \vdash e \Leftarrow A}$	$e$ checks against input type $A$ .
$\boxed{\Psi \vdash e \Rightarrow A}$	$e$ synthesizes output type $A$ .
$\boxed{\Psi \vdash A \bullet e \Rightarrow C}$	Applying a function of type $A$ to $e$ synthesizes type $C$ .

$\frac{(x : A) \in \Psi}{\Psi \vdash x \Rightarrow A} \text{DeclVar}$	$\frac{\Psi \vdash e \Rightarrow A \quad \Psi \vdash A \leq B}{\Psi \vdash e \Leftarrow B} \text{DeclSub}$
$\frac{\Psi \vdash A \quad \Psi \vdash e \Leftarrow A}{\Psi \vdash (e : A) \Rightarrow A} \text{DeclAnno}$	$\frac{}{\Psi \vdash () \Leftarrow 1} \text{Decl1I} \quad \frac{}{\Psi \vdash () \Rightarrow 1} \text{Decl1I}\Rightarrow$
$\frac{\Psi, a \vdash e \Leftarrow A}{\Psi \vdash e \Leftarrow \forall a. A} \text{Decl}\forall\text{I}$	$\frac{\Psi \vdash \tau \quad \Psi \vdash [\tau/a]A \bullet e \Rightarrow C}{\Psi \vdash \forall a. A \bullet e \Rightarrow C} \text{Decl}\forall\text{App}$
$\frac{\Psi, x : A \vdash e \Leftarrow B}{\Psi \vdash \lambda x. e \Leftarrow A \rightarrow B} \text{Decl}\rightarrow\text{I}$	$\frac{\Psi \vdash \sigma \rightarrow \tau \quad \Psi, x : \sigma \vdash e \Leftarrow \tau}{\Psi \vdash \lambda x. e \Rightarrow \sigma \rightarrow \tau} \text{Decl}\rightarrow\text{I}\Rightarrow$
$\frac{\Psi \vdash e_1 \Rightarrow A \quad \Psi \vdash A \bullet e_2 \Rightarrow C}{\Psi \vdash e_1 e_2 \Rightarrow C} \text{Decl}\rightarrow\text{E}$	$\frac{\Psi \vdash e \Leftarrow A}{\Psi \vdash A \rightarrow C \bullet e \Rightarrow C} \text{Decl}\rightarrow\text{App}$

Figure 2.9: Declarative Typing

then check whether the synthesized type entails the checked type. Lambda abstractions are the hardest construct of the bidirectional type system to deal with. Checking  $\lambda x. e$  against function type  $A \rightarrow B$  is easy: we check the body  $e$  against  $B$  in the context extended with  $x : A$ . However, synthesizing a lambda-function is a lot harder, and this type system only synthesizes monotypes  $\sigma \rightarrow \tau$ .

Application  $e_1 e_2$  is handled by Rule  $\text{Decl}\rightarrow\text{E}$ , which first synthesizes the type  $A$  of the function  $e_1$ . If  $A$  is a function type  $B \rightarrow C$ , Rule  $\text{Decl}\rightarrow\text{App}$  is applied; it checks the argument  $e_2$  against  $B$  and returns type  $C$ . The synthesized type of function  $e_1$  can also be polymorphic, of the form  $\forall a. A$ . In that case, we instantiate  $A$  to  $[\tau/a]A$  with a monotype  $\tau$  using according to Rule  $\text{Decl}\rightarrow\text{I}\Rightarrow$ . If  $[\tau/a]A$  is a function type, Rule  $\text{Decl}\rightarrow\text{App}$  proceeds; if  $[\tau/a]A$  is another universal quantified type, Rule  $\text{Decl}\rightarrow\text{I}\Rightarrow$  is recursively applied.

To conclude, DK employs a bidirectional declarative type system. The type system is mostly syntax-directed, but there are still some guesses of monotypes that need to be resolved by an algorithm. We will continue to discuss DK's algorithm in Chapters 3 and 4.

## 2.4 MLSUB

Besides extending the HM type system with higher-ranked types and bidirectional type checking, another invaluable direction is to support object-oriented subtyping. MLsub [Dolan and Mycroft 2017] is one of such systems. In presence of subtyping, type inference does not sim-

## 2 Background

Types	$\tau$	$::=$	$1 \mid a \mid \top \mid \perp \mid \tau_1 \rightarrow \tau_2 \mid \tau_1 \sqcup \tau_2 \mid \tau_1 \sqcap \tau_2$
Positive Types	$\tau^+$	$::=$	$1 \mid a \mid \perp \mid \tau_1^- \rightarrow \tau_2^+ \mid \tau_1^+ \sqcup \tau_2^+$
Negative Types	$\tau^-$	$::=$	$1 \mid a \mid \top \mid \tau_1^+ \rightarrow \tau_2^- \mid \tau_1^- \sqcap \tau_2^-$

Figure 2.10: Types of MLsub

ply handle equality during unification. Therefore, types are extended with lattice operations to express bounds properly. Furthermore, polar types are introduced to help separate input and output types, which simplifies the type inference algorithm. Like HM’s type inference, MLsub always infers a principal type.

### 2.4.1 TYPES AND POLAR TYPES

In comparison to the type system of HM, types (Figure 2.10) now include  $\top$  and  $\perp$ , as minimal components to support subtyping. Besides, the least-upper-bound ( $\sqcup$ ) and greatest-lower-bound ( $\sqcap$ ) lattice operations are used to represent a bound expressed by two types. For finite types, a distributive lattice can be defined via a set of equivalence classes of  $\equiv$  [Dolan and Mycroft 2017]. The most interesting equations are the distributivity rule and rules for function types:

$$\begin{aligned}
\tau_1 \sqcup (\tau_2 \sqcap \tau_3) &\equiv (\tau_1 \sqcup \tau_2) \sqcap (\tau_1 \sqcup \tau_3) \\
\tau_1 \sqcap (\tau_2 \sqcup \tau_3) &\equiv (\tau_1 \sqcap \tau_2) \sqcup (\tau_1 \sqcap \tau_3) \\
(\tau_1 \rightarrow \tau_2) \sqcup (\tau_1' \rightarrow \tau_2') &\equiv (\tau_1 \sqcap \tau_1') \rightarrow (\tau_2 \sqcup \tau_2') \\
(\tau_1 \rightarrow \tau_2) \sqcap (\tau_1' \rightarrow \tau_2') &\equiv (\tau_1 \sqcup \tau_1') \rightarrow (\tau_2 \sqcap \tau_2')
\end{aligned}$$

The partial order  $\tau_1 \leq \tau_2$  is defined as  $\tau_1 \sqcup \tau_2 \equiv \tau_2$  or  $\tau_1 \sqcap \tau_2 \equiv \tau_1$ .  $\top$  and  $\perp$  are the least and greatest types. The above rules on function types imply the usual subtyping rule for function types, considering the definition of partial order:

$$\frac{\tau_1' \leq \tau_1 \quad \tau_2 \leq \tau_2'}{\tau_1 \rightarrow \tau_2 \leq \tau_1' \rightarrow \tau_2'}$$

Type schemes are not defined as the usual  $\sigma$ . Instead, a monotype  $\tau$  already represents a type scheme by omitting the  $\forall$  quantifiers—all the free type variables are implicitly generalized.

Recursive types play an important role regarding the principality of type inference, but we omit them for simplicity.

**POLAR TYPES** Polar types are restrictions on the lattice operations; they should not occur arbitrarily in any position. Specifically, function outputs consist of types  $(\tau_1, \tau_2)$  from different branches, resulting in  $\tau_1 \sqcup \tau_2$ ; a function input might be used in various ways (under different constraints), thus  $\tau_1 \sqcap \tau_2$  is more suitable. In summary,  $\sqcup$  only arises in return types, while  $\sqcap$  only arises in argument types. Figure 2.10 formally defines the restriction, where positive types  $\tau^+$  describe return types, and negative types  $\tau^-$  describe argument types.

An important consequence is that all the constraints are of the form  $\tau^+ \leq \tau^-$ , which represents the subtyping relation when using an output expression in a function application as an argument. The following subtyping rules involving the lattice operations reflects their basic properties:

$$\frac{\tau_1^+ \leq \tau^- \quad \tau_2^+ \leq \tau^-}{\tau_1^+ \sqcup \tau_2^+ \leq \tau^-} \quad \frac{\tau^+ \leq \tau_1^- \quad \tau^+ \leq \tau_2^-}{\tau^- \leq \tau_1^- \sqcap \tau_2^-}.$$

Interestingly, the polar subtyping judgments avoids difficult judgments like  $\tau_1 \sqcap \tau_2 \leq \tau$  or  $\tau \leq \tau_1 \sqcup \tau_2$  through its syntactic restriction.

#### 2.4.2 BIUNIFICATION

Type inference for MLsub, similar to that for ML, is mainly a unification algorithm. However, in presence of subtyping, equality-based unification loses information about subtyping constraints.

For the atomic constraint  $\hat{\alpha} = \tau$  where  $\hat{\alpha} \notin \text{FV}(\tau)$ , the ML unification algorithm produces the substitution  $[\tau/\hat{\alpha}]$ . In contrast, an MLsub atomic constraint might be  $\hat{\alpha} \leq \tau$ , and the substitution  $[\tau/\hat{\alpha}]$  treat the subtyping constraint as an equality constraint, which eliminates a whole set of possibilities.

Luckily, lattices in MLsub helps express subtyping constraints on types directly. The constraint  $\hat{\alpha} \leq \tau$  ( $\hat{\alpha} \notin \text{FV}(\tau)$ ) may produce the substitution  $[\hat{\alpha} \sqcap \tau/\hat{\alpha}]$ , since  $\hat{\alpha} \sqcap \tau \leq \tau$ . In the meantime,  $\hat{\alpha} \sqcap \tau$  does not lose any expressiveness: for any  $\tau_0$  s.t.  $\tau_0 \leq \tau$ , picking  $\hat{\alpha} = \tau_0$  gives  $\hat{\alpha} \sqcap \tau = \tau_0$ , and the substitution  $[\hat{\alpha} \sqcap \tau/\hat{\alpha}]$  is equivalent to  $[\tau_0/\hat{\alpha}]$ .

In presence of polar types, the biunification algorithm of MLsub produces a *bisubstitution*  $[\hat{\alpha} \sqcap \tau^-/\hat{\alpha}^-]$  against the constraint  $\hat{\alpha} \leq \tau^-$ , where only negative occurrences are substituted, keeping polar types properly “polarized”. For example, a positive type  $\hat{\alpha} \rightarrow \hat{\alpha}$  becomes  $(\hat{\alpha} \sqcap \tau^-) \rightarrow \hat{\alpha}$  under such substitution and remains a positive type. A more important fact is that this type is equivalent to the constrained type scheme  $\hat{\alpha} \rightarrow \hat{\alpha}$  with the constraint  $\hat{\alpha} \leq \tau$ . Similarly, a constraint like  $\tau^+ \leq \hat{\alpha}$  is reduced to a substitution  $[\hat{\alpha} \sqcup \tau^+/\hat{\alpha}^+]$ .

## 2 Background

For example, the *choose* function is typed  $\forall a. a \rightarrow a \rightarrow a$  in ML. However, MLsub might also infer an equivalent type  $\forall a b. a \rightarrow b \rightarrow a \sqcup b$ . One can easily read the MLsub type in a form where constraints are explicitly stated

$$\forall a b c. a \rightarrow b \rightarrow c \text{ where } a \leq c, b \leq c.$$

Therefore, MLsub encodes the constraints directly onto types with the help of the lattice operations. Furthermore, a simplification step is taken after the type inference algorithm, reducing the size and improving readability of the type inferred.

As a result, biunification for MLsub extends unification for ML, accepting subtyping in addition to type schemes, while maintaining principality.

## PART II

# HIGHER-RANKED TYPE INFERENCE ALGORITHMS



# 3 HIGHER-RANKED POLYMORPHISM

## SUBTYPING ALGORITHM

In this chapter, we present a new algorithm for polymorphic subtyping with mechanical formalizations in the Abella theorem prover. There is little work on formalizing type inference algorithms before, especially for higher-ranked systems, due to the fact that environments and variable bindings are tricky to mechanize in theorem provers. In order to overcome the difficulty in formalization, we propose the novel algorithm by means of *worklist judgments*. Worklist judgments turn complicated global propagation of unification constraints into simple local substitutions. Moreover, we exploit several ideas in the recent inductive formulation of a type-inference algorithm by Dunfield and Krishnaswami [2013], which turn out to be useful for mechanization in a theorem prover.

Building on these ideas we develop a complete formalization of polymorphic subtyping in the Abella theorem prover. Moreover, we show that the algorithm is *sound*, *complete*, and *decidable* with respect to the well-known declarative formulation of polymorphic subtyping by Odersky and Läufer [1996]. While these meta-theoretical results are not new, as far as we know our work is the first to mechanically formalize them.

### 3.1 OVERVIEW: POLYMORPHIC SUBTYPING

This section discusses Odersky and Läufer declarative subtyping rules further in depth, and identifies the challenges in formalizing a corresponding algorithmic version. Then the key ideas of our approach that address those challenges are introduced.

#### 3.1.1 DECLARATIVE POLYMORPHIC SUBTYPING

In implicitly polymorphic type systems, the subtyping relation compares the degree of polymorphism of types. In short, if a polymorphic type  $A$  can always be instantiated to any instantiation of  $B$ , then  $A$  is “at least as polymorphic as”  $B$ , or we just say that  $A$  is “more polymorphic than”  $B$ , or  $A \leq B$ .

There is a very simple declarative formulation of polymorphic subtyping due to Odersky and Läufer [1996]. The syntax of this declarative system is shown in Figure 3.1. Types, rep-

Type variables	$a, b$
Types	$A, B, C ::= 1 \mid a \mid \forall a. A \mid A \rightarrow B$
Monotypes	$\tau ::= 1 \mid a \mid \tau_1 \rightarrow \tau_2$
Contexts	$\Psi ::= \cdot \mid \Psi, a$

Figure 3.1: Syntax of Declarative System

$$\begin{array}{c}
 \boxed{\Psi \vdash A} \\
 \\
 \frac{}{\Psi \vdash 1} \text{wf}_{\text{dunit}} \quad \frac{a \in \Psi}{\Psi \vdash a} \text{wf}_{\text{dvar}} \quad \frac{\Psi \vdash A \quad \Psi \vdash B}{\Psi \vdash A \rightarrow B} \text{wf}_{\text{d}\rightarrow} \quad \frac{\Psi, a \vdash A}{\Psi \vdash \forall a. A} \text{wf}_{\text{d}\forall} \\
 \\
 \boxed{\Psi \vdash A \leq B} \\
 \\
 \frac{a \in \Psi}{\Psi \vdash a \leq a} \leq_{\text{Var}} \quad \frac{}{\Psi \vdash 1 \leq 1} \leq_{\text{Unit}} \quad \frac{\Psi \vdash B_1 \leq A_1 \quad \Psi \vdash A_2 \leq B_2}{\Psi \vdash A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2} \leq_{\rightarrow} \\
 \\
 \frac{\Psi \vdash \tau \quad \Psi \vdash [\tau/a]A \leq B}{\Psi \vdash \forall a. A \leq B} \leq_{\forall L} \quad \frac{\Psi, a \vdash A \leq B}{\Psi \vdash A \leq \forall a. B} \leq_{\forall R}
 \end{array}$$

Figure 3.2: Well-formedness of Declarative Types and Declarative Subtyping

resented by  $A, B, C$ , are the unit type  $1$ , type variables  $a, b$ , universal quantification  $\forall a. A$  and function type  $A \rightarrow B$ . We allow nested universal quantifiers to appear in types, but not in monotypes. Contexts  $\Psi$  collect a list of type variables.

In Figure 3.2, we give the well-formedness and subtyping relation for the declarative system, which is identical to the subtyping relation introduced in Subsection 2.3.1.

### 3.1.2 FINDING SOLUTIONS FOR VARIABLE INSTANTIATION

The declarative system specifies the behavior of subtyping relations, but is not directly implementable: Rule  $\leq_{\forall L}$  requires guessing a monotype  $\tau$ . The core problem that an algorithm for polymorphic subtyping needs to solve is to find an algorithmic way to compute the monotypes, instead of guessing them. An additional challenge is that the declarative rule  $\leq_{\rightarrow}$  splits one judgment into two, and the (partial) solutions found for existential variables when processing the first judgment should be transferred to the second judgment.

It is worth mentioning that the problem of deciding subtyping of OL's type system can be reduced to first-order unification under a mixed prefix Miller [1992], where universal and existential quantifiers appears at the top level of constraints (equations). We further discuss another simpler algorithm by DK and present ours afterwards.

DUNFIELD AND KRISHNASWAMI’S APPROACH An elegant algorithmic solution to computing the monotypes is presented by Dunfield and Krishnaswami [2013]. Their algorithmic subtyping judgment has the form:

$$\Psi \vdash A \leq B \dashv \Phi$$

A notable difference to the declarative judgment is the presence of a so-called *output context*  $\Phi$ , which refines the *input context*  $\Psi$  with solutions for existential variables found while processing the two types being compared for subtyping. Both  $\Psi$  and  $\Phi$  are *ordered contexts* with the same structure. Ordered contexts are particularly useful to keep track of the correct scoping for variables, and are a notable difference to older type-inference algorithms [Damas and Milner 1982] that use global unification variables or constraints collected in a set.

Output contexts are useful to transfer information across judgments in Dunfield and Krishnaswami’s approach. For example, the algorithmic rule corresponding to  $\leq \rightarrow$  in their approach is:

$$\frac{\Psi \vdash B_1 <: A_1 \dashv \Phi \quad \Phi \vdash [\Phi]A_2 <: [\Phi]B_2 \dashv \Phi'}{\Psi \vdash A_1 \rightarrow A_2 <: B_1 \rightarrow B_2 \dashv \Phi'} <:\rightarrow$$

The information gathered by the output context when comparing the input types of the functions for subtyping is transferred to the second judgment by becoming the new input context, while any solution derived from the first judgment is applied to the types of the second judgment.

EXAMPLE If we want to show that  $\forall a. a \rightarrow a$  is a subtype of  $1 \rightarrow 1$ , the declarative system will guess the proper  $\tau = 1$  for Rule  $\leq \forall L$ :

$$\frac{\cdot \vdash 1 \quad \cdot \vdash 1 \rightarrow 1 \leq 1 \rightarrow 1}{\cdot \vdash \forall a. a \rightarrow a \leq 1 \rightarrow 1} \leq \forall L$$

Dunfield and Krishnaswami introduce an *existential variable*—denoted with  $\hat{\alpha}, \hat{\beta}$ —whenever a monotype  $\tau$  needs to be guessed. Below is a sample derivation of their algorithm:

$$\frac{\frac{\frac{}{\hat{\alpha} \vdash 1 \leq \hat{\alpha} \dashv \hat{\alpha} = 1} \text{InstRSolve} \quad \frac{}{\hat{\alpha} = 1 \vdash 1 \leq 1 \dashv \hat{\alpha} = 1} <:\text{Unit}}{\hat{\alpha} \vdash \hat{\alpha} \rightarrow \hat{\alpha} \leq 1 \rightarrow 1 \dashv \hat{\alpha} = 1} <:\rightarrow}{\cdot \vdash \forall a. a \rightarrow a \leq 1 \rightarrow 1 \dashv \cdot} <:\forall L$$

The first step applies Rule  $<:\forall L$ , which introduces a fresh existential variable,  $\hat{\alpha}$ , and opens the left-hand-side  $\forall$ -quantifier with it. Next, Rule  $<:\rightarrow$  splits the judgment in two. For the first branch, Rule `InstRSolve` satisfies  $1 \leq \hat{\alpha}$  by solving  $\hat{\alpha}$  to 1, and stores the solution in its output context. The output context of the first branch is used as the input context of the second branch, and the judgment is updated according to current solutions. Finally, the second branch becomes a base case, and Rule  $<:\text{Unit}$  finishes the derivation, makes no change to the input context and propagates the output context back.

Dunfield and Krishnaswami's algorithmic specification is elegant and contains several useful ideas for a mechanical formalization of polymorphic subtyping. For example, *ordered contexts* and *existential variables* enable a purely inductive formulation of polymorphic subtyping. However, the binding/scoping structure of their algorithmic judgment is still fairly complicated and poses challenges when porting their approach to a theorem prover.

#### 3.1.3 THE WORKLIST APPROACH

We inherit Dunfield and Krishnaswami's ideas of ordered contexts, existential variables and the idea of solving those variables, but drop output contexts. Instead, our algorithmic rule has the form:

$$\Gamma \vdash \Omega$$

where  $\Omega$  is a list of judgments  $A \leq B$  instead of a single one. This judgment form, which we call *worklist judgment*, simplifies two aspects of Dunfield and Krishnaswami's approach.

Firstly, as already stated, there are no output contexts. Secondly, the form of the ordered contexts becomes simpler. The transfer of information across judgments is simplified because all judgments share the input context. Moreover, the order of the judgments in the list allows information discovered when processing the earlier judgments to be easily transferred to the later judgments. In the worklist approach the rule for function types is:

$$\frac{\Gamma \vdash B_1 \leq A_1; A_2 \leq B_2; \Omega}{\Gamma \vdash A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2; \Omega} \leq_{a \rightarrow}$$

The derivation of the previous example with the worklist approach is:

$$\begin{array}{c}
 \frac{}{\cdot \vdash \cdot} \text{a\_nil} \\
 \frac{}{\cdot \vdash 1 \leq 1; \cdot} \leq_{\text{aunit}} \\
 \frac{}{\hat{\alpha} \vdash 1 \leq \hat{\alpha}; \hat{\alpha} \leq 1; \cdot} \leq_{\text{a solve\_ex}} \\
 \frac{}{\hat{\alpha} \vdash \hat{\alpha} \rightarrow \hat{\alpha} \leq 1 \rightarrow 1; \cdot} \leq_{\text{a} \rightarrow} \\
 \frac{}{\cdot \vdash \forall a. a \rightarrow a \leq 1 \rightarrow 1; \cdot} \leq_{\text{a} \forall \text{L}}
 \end{array}$$

To derive  $\cdot \vdash \forall a. a \rightarrow a \leq 1 \rightarrow 1$  with the worklist approach, we first introduce an existential variable and change the judgment to  $\hat{\alpha} \vdash \hat{\alpha} \rightarrow \hat{\alpha} \leq 1 \rightarrow 1; \cdot$ . Then, we split the judgment in two for the function types and the derivation comes to  $\hat{\alpha} \vdash 1 \leq \hat{\alpha}; \hat{\alpha} \leq 1; \cdot$ . When the first judgment is solved with  $\hat{\alpha} = 1$ , we immediately remove  $\hat{\alpha}$  from the context, while propagating the solution as a substitution to the rest of the judgment list, resulting in  $\cdot \vdash 1 \leq 1; \cdot$ , which finishes the derivation in two trivial steps.

With this form of eager propagation, solutions no longer need to be recorded in contexts, simplifying the encoding and reasoning in a proof assistant.

**KEY RESULTS** Both the declarative and algorithmic systems are formalized in Abella. We have proven 3 important properties for this algorithm: *decidability*, ensuring that the algorithm always terminates; and *soundness* and *completeness*, showing the equivalence of the declarative and algorithmic systems.

### 3.2 A WORKLIST ALGORITHM FOR POLYMORPHIC SUBTYPING

This section presents our algorithm for polymorphic subtyping. A novel aspect of our algorithm is the use of worklist judgments: a form of judgment that facilitates the propagation of information.

#### 3.2.1 SYNTAX AND WELL-FORMEDNESS OF THE ALGORITHMIC SYSTEM

Figure 3.3 shows the syntax and the well-formedness judgment.

**EXISTENTIAL VARIABLES** In order to solve the unknown types  $\tau$ , the algorithmic system extends the declarative syntax of types with *existential variables*  $\hat{\alpha}$ . They behave like unification variables, but are not globally defined. Instead, the ordered *algorithmic context*, inspired by Dunfield and Krishnaswami [2013], defines their scope. Thus the type  $\tau$  represented by the

Type variables	$a, b$
Existential variables	$\hat{\alpha}, \hat{\beta}$
Algorithmic types	$A, B, C ::= 1 \mid a \mid \hat{\alpha} \mid \forall a. A \mid A \rightarrow B$
Algorithmic context	$\Gamma ::= \cdot \mid \Gamma, a \mid \Gamma, \hat{\alpha}$
Algorithmic judgments	$\Omega ::= \cdot \mid A \leq B; \Omega$
$\boxed{\Gamma \vdash A}$	
$\frac{}{\Gamma \vdash 1} \text{wf}_{\text{aunit}} \quad \frac{a \in \Gamma}{\Gamma \vdash a} \text{wf}_{\text{avar}} \quad \frac{\hat{\alpha} \in \Gamma}{\Gamma \vdash \hat{\alpha}} \text{wf}_{\text{exvar}}$ $\frac{\Gamma \vdash A \quad \Gamma \vdash B}{\Gamma \vdash A \rightarrow B} \text{wf}_{\text{a} \rightarrow} \quad \frac{\Gamma, a \vdash A}{\Gamma \vdash \forall a. A} \text{wf}_{\text{a} \forall}$	

Figure 3.3: Syntax and Well-Formedness judgment for the Algorithmic System.

corresponding existential variable is always bound in the corresponding declarative context  $\Psi$ .

**WORKLIST JUDGMENTS** The form of our algorithmic judgments is non-standard. Our algorithm keeps track of an explicit list of outstanding work: the list  $\Omega$  of (reified) *algorithmic judgments* of the form  $A \leq B$ , to which a substitution can be applied once and for all to propagate the solution of an existential variable.

**HOLE NOTATION** To facilitate context manipulation, we use the syntax  $\Gamma[\Gamma_M]$  to denote a context of the form  $\Gamma_L, \Gamma_M, \Gamma_R$  where  $\Gamma$  is the context  $\Gamma_L, \bullet, \Gamma_R$  with a hole ( $\bullet$ ). Hole notations with the same name implicitly share the same  $\Gamma_L$  and  $\Gamma_R$ . A multi-hole notation like  $\Gamma[\hat{\alpha}][\hat{\beta}]$  means  $\Gamma_1, \hat{\alpha}, \Gamma_2, \hat{\beta}, \Gamma_3$ .

### 3.2.2 ALGORITHMIC SUBTYPING

The algorithmic subtyping judgment, defined in Figure 3.4, has the form  $\Gamma \vdash \Omega$ , where  $\Omega$  collects multiple subtyping judgments  $A \leq B$ . The algorithm treats  $\Omega$  as a worklist. In every step it takes one task from the worklist for processing, possibly pushes some new tasks on the worklist, and repeats this process until the list is empty. This last and single base case is handled by Rule `a_nil`. The remaining rules all deal with the first task in the worklist. Logically we can discern 3 groups of rules.

Firstly, we have five rules that are similar to those in the declarative system, mostly just adapted to the worklist style. For instance, Rule  $\leq_{\text{a} \rightarrow}$  consumes one judgment and pushes

$$\boxed{\Gamma \vdash \Omega}$$

$$\frac{}{\Gamma \vdash \cdot} \leq_{\text{a nil}}$$

$$\frac{\Gamma \vdash \Omega}{\Gamma \vdash 1 \leq 1; \Omega} \leq_{\text{a unit}} \quad \frac{a \in \Gamma \quad \Gamma \vdash \Omega}{\Gamma \vdash a \leq a; \Omega} \leq_{\text{a var}} \quad \frac{\hat{\alpha} \in \Gamma \quad \Gamma \vdash \Omega}{\Gamma \vdash \hat{\alpha} \leq \hat{\alpha}; \Omega} \leq_{\text{a exvar}}$$

$$\frac{\Gamma \vdash B_1 \leq A_1; A_2 \leq B_2; \Omega}{\Gamma \vdash A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2; \Omega} \leq_{\text{a } \rightarrow}$$

$$\frac{\hat{\alpha} \text{ fresh} \quad \Gamma, \hat{\alpha} \vdash [\hat{\alpha}/a]A \leq B; \Omega}{\Gamma \vdash \forall a. A \leq B; \Omega} \leq_{\text{a } \forall \text{L}} \quad \frac{b \text{ fresh} \quad \Gamma, b \vdash A \leq B; \Omega}{\Gamma \vdash A \leq \forall b. B; \Omega} \leq_{\text{a } \forall \text{R}}$$

$$\frac{\hat{\alpha} \notin FV(A) \cup FV(B) \quad \Gamma[\hat{\alpha}_1, \hat{\alpha}_2] \vdash \hat{\alpha}_1 \rightarrow \hat{\alpha}_2 \leq A \rightarrow B; [\hat{\alpha}_1 \rightarrow \hat{\alpha}_2/\hat{\alpha}]\Omega}{\Gamma[\hat{\alpha}] \vdash \hat{\alpha} \leq A \rightarrow B; \Omega} \leq_{\text{a instL}}$$

$$\frac{\hat{\alpha} \notin FV(A) \cup FV(B) \quad \Gamma[\hat{\alpha}_1, \hat{\alpha}_2] \vdash A \rightarrow B \leq \hat{\alpha}_1 \rightarrow \hat{\alpha}_2; [\hat{\alpha}_1 \rightarrow \hat{\alpha}_2/\hat{\alpha}]\Omega}{\Gamma[\hat{\alpha}] \vdash A \rightarrow B \leq \hat{\alpha}; \Omega} \leq_{\text{a instR}}$$

$$\frac{\Gamma[\hat{\alpha}][] \vdash [\hat{\alpha}/\hat{\beta}]\Omega}{\Gamma[\hat{\alpha}][\hat{\beta}] \vdash \hat{\alpha} \leq \hat{\beta}; \Omega} \leq_{\text{a solve\_ex}} \quad \frac{\Gamma[\hat{\alpha}][] \vdash [\hat{\alpha}/\hat{\beta}]\Omega}{\Gamma[\hat{\alpha}][\hat{\beta}] \vdash \hat{\beta} \leq \hat{\alpha}; \Omega} \leq_{\text{a solve\_ex}'}$$

$$\frac{\Gamma[a][] \vdash [a/\hat{\beta}]\Omega}{\Gamma[a][\hat{\beta}] \vdash a \leq \hat{\beta}; \Omega} \leq_{\text{a solve\_var}} \quad \frac{\Gamma[a][] \vdash [a/\hat{\beta}]\Omega}{\Gamma[a][\hat{\beta}] \vdash \hat{\beta} \leq a; \Omega} \leq_{\text{a solve\_var}'}$$

$$\frac{\Gamma[] \vdash [1/\hat{\alpha}]\Omega}{\Gamma[\hat{\alpha}] \vdash \hat{\alpha} \leq 1; \Omega} \leq_{\text{a solve\_unit}} \quad \frac{\Gamma[] \vdash [1/\hat{\alpha}]\Omega}{\Gamma[\hat{\alpha}] \vdash 1 \leq \hat{\alpha}; \Omega} \leq_{\text{a solve\_unit}'}$$

Figure 3.4: Algorithmic Subtyping

two to the worklist. A notable difference with the declarative Rule  $\leq_{\forall \text{L}}$  is that Rule  $\leq_{\text{a } \forall \text{L}}$  requires no guessing of a type  $\tau$  to instantiate the polymorphic type  $\forall a. A$ , but instead introduces an existential variable  $\hat{\alpha}$  to the context and to  $A$ . In accordance with the declarative system, where the monotype  $\tau$  should be bound in the context  $\Psi$ , here  $\hat{\alpha}$  should only be solved to a monotype bound in  $\Gamma$ . More generally, for any algorithmic context  $\Gamma[\hat{\alpha}]$ , the algorithmic variable  $\hat{\alpha}$  can only be solved to a monotype that is well-formed with respect to  $\Gamma_L$ .

Secondly, Rules  $\leq_{\text{a instL}}$  and  $\leq_{\text{a instR}}$  partially instantiate existential types  $\hat{\alpha}$ , to function types. The domain and range of the new function type are undetermined: they are set to two fresh existential variables  $\hat{\alpha}_1$  and  $\hat{\alpha}_2$ . To make sure that  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$  has the same scope

$$\begin{array}{c}
 \frac{}{\hat{\alpha}_1 \vdash \cdot} \text{a\_nil} \\
 \frac{}{\hat{\alpha}_1 \vdash 1 \leq 1; \cdot} \leq_{\text{aunit}} \\
 \frac{}{\hat{\alpha}_1, \hat{\alpha}_2 \vdash \hat{\alpha}_1 \leq \hat{\alpha}_2; 1 \leq 1; \cdot} \leq_{\text{asolve\_ex}} \\
 \frac{}{\hat{\alpha}_1, \hat{\alpha}_2, \hat{\beta} \vdash \hat{\alpha}_1 \leq \hat{\beta}; \hat{\beta} \leq \hat{\alpha}_2; 1 \leq 1; \cdot} \leq_{\text{asolve\_ex}} \\
 \frac{}{\hat{\alpha}_1, \hat{\alpha}_2, \hat{\beta} \vdash \hat{\beta} \rightarrow \hat{\beta} \leq \hat{\alpha}_1 \rightarrow \hat{\alpha}_2; 1 \leq 1; \cdot} \leq_{\text{a}\rightarrow} \\
 \frac{}{\hat{\alpha}, \hat{\beta} \vdash \hat{\beta} \rightarrow \hat{\beta} \leq \hat{\alpha}; 1 \leq 1; \cdot} \leq_{\text{ainstR}} \\
 \frac{}{\hat{\alpha} \vdash \forall a. a \rightarrow a \leq \hat{\alpha}; 1 \leq 1; \cdot} \leq_{\text{a}\forall\text{L}} \\
 \frac{}{\hat{\alpha} \vdash \hat{\alpha} \rightarrow 1 \leq (\forall a. a \rightarrow a) \rightarrow 1; \cdot} \leq_{\text{a}\rightarrow} \\
 \frac{}{\cdot \vdash \forall a. a \rightarrow 1 \leq (\forall a. a \rightarrow a) \rightarrow 1; \cdot} \leq_{\text{a}\forall\text{L}}
 \end{array}$$

Figure 3.5: A Success Derivation for the Algorithmic Subtyping Relation

$$\begin{array}{c}
 \frac{}{\hat{\alpha}, b \vdash \hat{\alpha} \leq b; \cdot} \text{stuck} \quad ? \\
 \frac{}{\hat{\alpha} \vdash \hat{\alpha} \leq \forall b. b; \cdot} \leq_{\text{a}\forall\text{R}} \\
 \frac{}{\hat{\alpha} \vdash 1 \leq 1; \hat{\alpha} \leq \forall b. b; \cdot} \leq_{\text{aunit}} \\
 \frac{}{\hat{\alpha} \vdash 1 \rightarrow \hat{\alpha} \leq 1 \rightarrow \forall b. b; \cdot} \leq_{\text{a}\rightarrow} \\
 \frac{}{\cdot \vdash \forall a. 1 \rightarrow a \leq 1 \rightarrow \forall b. b; \cdot} \leq_{\text{a}\forall\text{L}}
 \end{array}$$

Figure 3.6: A Failing Derivation for the Algorithmic Subtyping Relation

as  $\hat{\alpha}$ , the new variables  $\hat{\alpha}_1$  and  $\hat{\alpha}_2$  are inserted in the same position in the context where the old variable  $\hat{\alpha}$  was. To propagate the instantiation to the remainder of the worklist,  $\hat{\alpha}$  is substituted for  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$  in  $\Omega$ . The *occurs-check* side-condition is necessary to prevent a diverging infinite instantiation. For example  $1 \rightarrow \hat{\alpha} \leq \hat{\alpha}$  would diverge with no such check. Note that the algorithm does not choose to instantiate  $\hat{\alpha}$  directly with  $A \rightarrow B$ , since the type is not guaranteed to be a monotype, and such instantiation will be inconsistent with our predicative declarative system.

Thirdly, in the remaining six rules an existential variable can be immediately solved. Each of the six similar rules removes an existential variable from the context, performs a substitution on the remainder of the worklist and continues.

The algorithm on judgment list is designed to share the context across all judgments. However, the declarative system does not share a single context in its derivation. This gap is filled by strengthening and weakening lemmas of both systems, where most of them are straightforward to prove, except for the strengthening lemma of the declarative system, which is a little trickier.



**EXAMPLE** We illustrate the subtyping rules through a sample derivation in Figure 3.5, which shows that  $\forall a. a \rightarrow 1 \leq (\forall a. a \rightarrow a) \rightarrow 1$ . Thus the derivation starts with an empty context and a judgment list with only one element.

In step 1, we have only one judgment, and that one has a top-level  $\forall$  on the left hand side. So the only choice is Rule  $\leq_a \forall L$ , which opens the universally quantified type with an unknown existential variable  $\hat{\alpha}$ . Variable  $\hat{\alpha}$  will be solved later to some monotype that is well-formed within the context before  $\hat{\alpha}$ . That is, the empty context  $\cdot$  in this case. In step 2, Rule  $\leq_a \rightarrow$  is applied to the worklist, splitting the first judgment into two. Step 3 is similar to step 1, where the left-hand-side  $\forall$  of the first judgment is opened according to Rule  $\leq_a \forall L$  with a fresh existential variable. In step 4, the first judgment has an arrow on the left hand side, but the right-hand-side type is an existential variable. It is obvious that  $\hat{\alpha}$  should be solved to a monotype of the form  $\sigma \rightarrow \tau$ . Rule `instR` implements this, but avoids guessing  $\sigma$  and  $\tau$  by “splitting”  $\hat{\alpha}$  into two existential variables,  $\hat{\alpha}_1$  and  $\hat{\alpha}_2$ , which will be solved to some  $\sigma$  and  $\tau$  later. Step 5 applies Rule  $\leq_a \rightarrow$  again. Notice that after the split,  $\hat{\beta}$  appears in two judgments. When the first  $\hat{\beta}$  is solved during any step of the derivation, the next  $\hat{\beta}$  will be substituted by that solution. This propagation mechanism ensures the consistent solution of the variables, while keeping the context as simple as possible. Steps 6 and 7 solve existential variables. The existential variable that is right-most in the context is always solved in terms of the other. Therefore in step 6,  $\hat{\beta}$  is solved in terms of  $\hat{\alpha}_1$ , and in step 7,  $\hat{\alpha}_2$  is solved in terms of  $\hat{\alpha}_1$ . Additionally, in step 6, when  $\hat{\beta}$  is solved, the substitution  $[\hat{\alpha}_1/\hat{\beta}]$  is propagated to the rest of the judgment list, and thus the second judgment becomes  $\hat{\alpha}_1 \leq \hat{\alpha}_2$ . Steps 8 and 9 trivially finish the derivation. Notice that  $\hat{\alpha}_1$  is not instantiated at the end. This means that any well-scoped instantiation is fine.

**A FAILING DERIVATION** We illustrate the role of ordered contexts through another example:  $\forall a. 1 \rightarrow a \leq 1 \rightarrow \forall b. b$ . From the declarative perspective,  $a$  should be instantiated to some  $\tau$  first, then  $b$  is introduced to the context, so that  $b \notin FV(\tau)$ . As a result, we cannot find  $\tau$  such that  $\tau \leq b$ . Figure 3.6 shows the algorithmic derivation, which also fails due to the scoping— $\hat{\alpha}$  is introduced earlier than  $b$ , thus it cannot be solved to  $b$ .

### 3.3 METATHEORY

This section presents the three main meta-theoretical results that we have proved in Abella. The first two are soundness and completeness of our algorithm with respect to Odersky and Läufer’s declarative subtyping. The third result is our algorithm’s decidability.

$$\begin{array}{c}
 \boxed{\Gamma \rightarrow \Psi} \\
 \frac{}{\cdot \rightarrow \cdot} \rightarrow \cdot \quad \frac{\Gamma \rightarrow \Psi}{\Gamma, a \rightarrow \Psi, a} \rightarrow \text{var} \quad \frac{\Gamma \rightarrow \Psi}{\Gamma, \hat{\alpha} \rightarrow \Psi} \rightarrow \text{exvar} \\
 \boxed{\Gamma \mid \Omega \rightsquigarrow \Omega'} \\
 \frac{}{\cdot \mid \Omega \rightsquigarrow \Omega} \rightsquigarrow \cdot \quad \frac{\Gamma \mid \Omega \rightsquigarrow \Omega'}{\Gamma, a \mid \Omega \rightsquigarrow \Omega'} \rightsquigarrow \text{var} \\
 \frac{\Gamma \rightarrow \Psi \quad \Psi \vdash \tau \quad \Gamma \mid [\tau/\hat{\alpha}]\Omega \rightsquigarrow \Omega'}{\Gamma, \hat{\alpha} \mid \Omega \rightsquigarrow \Omega'} \rightsquigarrow \text{exvar}
 \end{array}$$

Figure 3.7: Transfer Rules

### 3.3.1 TRANSFER TO THE DECLARATIVE SYSTEM

To state the correctness of the algorithmic subtyping rules, Figure 3.7 introduces two *transfer* judgments to relate the declarative and the algorithmic system. The first judgment, transfer of contexts  $\Gamma \rightarrow \Psi$ , removes existential variables from the algorithmic context  $\Gamma$  to obtain a declarative context  $\Psi$ . The second judgment, transfer of the judgment list  $\Gamma \mid \Omega \rightsquigarrow \Omega'$ , replaces all occurrences of existential variables in  $\Omega$  by well-scoped mono-types. Notice that this judgment is not decidable, i.e. a pair of  $\Gamma$  and  $\Omega$  may be related with multiple  $\Omega'$ . However, if there exists some substitution that transforms  $\Omega$  to  $\Omega'$ , and each subtyping judgment in  $\Omega'$  holds, we know that  $\Omega$  is potentially satisfiable.

The following two lemmas generalize Rule  $\rightsquigarrow \text{exvar}$  from substituting the first existential variable to substituting any existential variable.

**Lemma 3.1** (Insert). *If  $\Gamma \rightarrow \Psi$  and  $\Psi \vdash \tau$  and  $\Gamma, \Gamma_1 \mid [\tau/\hat{\alpha}]\Omega \rightsquigarrow \Omega'$ , then  $\Gamma, \hat{\alpha}, \Gamma_1 \mid \Omega \rightsquigarrow \Omega'$ .*

**Lemma 3.2** (Extract). *If  $\Gamma, \hat{\alpha}, \Gamma_1 \mid \Omega \rightsquigarrow \Omega'$ , then  $\exists \tau$  s.t.  $\Gamma \rightarrow \Psi, \Psi \vdash \tau$  and  $\Gamma, \Gamma_1 \mid [\tau/\hat{\alpha}]\Omega \rightsquigarrow \Omega'$ .*

In order to match the shape of algorithmic subtyping relation for the following proofs, we define a relation  $\Psi \vdash \Omega$  for the declarative system, meaning that all the declarative judgments hold under context  $\Psi$ .

**Definition 1** (Declarative Subtyping Worklist).

$$\Psi \vdash \Omega := \forall (A \leq B) \in \Omega, \Psi \vdash A \leq B$$

### 3.3.2 SOUNDNESS

Our algorithm is sound with respect to the declarative specification. For any derivation of a list of algorithmic judgments  $\Gamma \vdash \Omega$ , we can find a valid transfer  $\Gamma \mid \Omega \rightsquigarrow \Omega'$  such that all judgments in  $\Omega'$  hold in  $\Psi$ , with  $\Gamma \rightarrow \Psi$ .

**Theorem 3.3** (Soundness). *If  $\Gamma \vdash \Omega$  and  $\Gamma \rightarrow \Psi$ , then there exists  $\Omega'$ , s.t.  $\Gamma \mid \Omega \rightsquigarrow \Omega'$  and  $\Psi \vdash \Omega'$ .*

The proof proceeds by induction on the derivation of  $\Gamma \vdash \Omega$ , finished off by appropriate applications of the insertion and extraction lemmas.

### 3.3.3 COMPLETENESS

Completeness of the algorithm means that any declarative derivation has an algorithmic counterpart.

**Theorem 3.4** (Completeness). *If  $\Psi \vdash \Omega'$  and  $\Gamma \rightarrow \Psi$  and  $\Gamma \mid \Omega \rightsquigarrow \Omega'$ , then  $\Gamma \vdash \Omega$ .*

The proof proceeds by induction on the derivation of  $\Psi \vdash \Omega'$ . As the declarative system does not involve information propagation across judgments, the induction can focus on the subtyping derivation of the first judgment without affecting other judgments. The difficult cases correspond to the  $\leq_{\text{aInstL}}$  and  $\leq_{\text{aInstR}}$  rules. When the proof by induction on  $\Psi \vdash \Omega'$  reaches the  $\leq \rightarrow$  case, the first declarative judgment has a shape like  $A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2$ . One of the possible cases for the first corresponding algorithmic judgment is  $\hat{\alpha} \leq A \rightarrow B$ . However, the case analysis does not indicate that  $\hat{\alpha}$  is fresh in  $A$  and  $B$ . Thus we cannot apply Rule  $\leq_{\text{aInstL}}$  and make use of the induction hypothesis. The following lemma helps us out in those cases: it rules out subtypings with infinite types as solutions (e.g.  $\hat{\alpha} \leq 1 \rightarrow \hat{\alpha}$ ) and guarantees that  $\hat{\alpha}$  is free in  $A$  and  $B$ .

**Lemma 3.5** (Prune Transfer for Instantiation). *If  $\Psi \vdash A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2; \Omega'$  and  $\Gamma \rightarrow \Psi$  and  $\Gamma \mid (\hat{\alpha} \leq A \rightarrow B; \Omega) \rightsquigarrow (A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2; \Omega')$ , then  $\hat{\alpha} \notin FV(A) \cup FV(B)$ .*

A similar lemma holds for the symmetric case  $(A \rightarrow B \leq \hat{\alpha}; \Omega)$ .

### 3.3.4 DECIDABILITY

The third key result for our algorithm is decidability.

**Theorem 3.6** (Decidability). *Given any well-formed judgment list  $\Omega$  under  $\Gamma$ , it is decidable whether  $\Gamma \vdash \Omega$  or not.*

We have proven this theorem by means of a lexicographic group of induction measurements  $\langle |\Omega|_{\forall}, |\Gamma|_{\hat{\alpha}}, |\Omega|_{\rightarrow} \rangle$  on the worklist  $\Omega$  and algorithmic context  $\Gamma$ . The worklist measures  $|\cdot|_{\forall}$  and  $|\cdot|_{\rightarrow}$  count the number of universal quantifiers and function types respectively.

**Definition 2** (Worklist Measures).

$$\begin{aligned}
 |1|_{\forall} &= 0 \\
 |a|_{\forall} &= 0 \\
 |\hat{\alpha}|_{\forall} &= 0 \\
 |A \rightarrow B|_{\forall} &= |A|_{\forall} + |B|_{\forall} \\
 |\forall x.A|_{\forall} &= |A|_{\forall} + 1 \\
 |\Omega|_{\forall} &= \sum_{A \leq B \in \Omega} |A|_{\forall} + |B|_{\forall} \\
 |1|_{\rightarrow} &= 0 \\
 |a|_{\rightarrow} &= 0 \\
 |\hat{\alpha}|_{\rightarrow} &= 0 \\
 |A \rightarrow B|_{\rightarrow} &= |A|_{\rightarrow} + |B|_{\rightarrow} + 1 \\
 |\forall x.A|_{\rightarrow} &= |A|_{\rightarrow} \\
 |\Omega|_{\rightarrow} &= \sum_{A \leq B \in \Omega} |A|_{\rightarrow} + |B|_{\rightarrow}
 \end{aligned}$$

The context measure  $|\cdot|_{\hat{\alpha}}$  counts the number of unsolved existential variables.

**Definition 3** (Context Measure).

$$|\cdot|_{\hat{\alpha}} = 0 \quad |\Gamma, a|_{\hat{\alpha}} = |\Gamma|_{\hat{\alpha}} \quad |\Gamma, \hat{\alpha}|_{\hat{\alpha}} = |\Gamma|_{\hat{\alpha}} + 1$$

It is not difficult to see that all but two of the algorithm's rules decrease one of the three measures. The two exceptions are the Rules  $\leq_{\text{a inst L}}$  and  $\leq_{\text{a inst R}}$ ; both increment the number of existential variables and the number of function types without affecting the number of universal quantifiers. To handle these rules, we handle a special class of judgments, which we call *instantiation judgments*  $\Omega_i$ , separately. They take the form:

**Definition 4** ( $\Omega_i$ ).

$$\Omega_i := \cdot \mid \hat{\alpha} \leq A; \Omega'_i \mid A \leq \hat{\alpha}; \Omega'_i \quad \text{where } \hat{\alpha} \notin FV(A) \cup FV(\Omega'_i)$$

These instantiation judgments are these ones consumed and produced by Rules  $\leq_{\text{a inst L}}$  and  $\leq_{\text{a inst R}}$ . The following lemma handles their decidability.

**Lemma 3.7** (Instantiation Decidability). *For any context  $\Gamma$  and judgment list  $\Omega_i, \Omega$ , it is decidable whether  $\Gamma \vdash \Omega_i, \Omega$  if both of the conditions hold*

- 1)  $\forall \Gamma', \Omega' \text{ s.t. } |\Omega'|_{\forall} < |\Omega_i, \Omega|_{\forall}$ , it is decidable whether  $\Gamma' \vdash \Omega'$ .
- 2)  $\forall \Gamma', \Omega' \text{ s.t. } |\Omega'|_{\forall} = |\Omega_i, \Omega|_{\forall}$  and  $|\Gamma'|_{\hat{\alpha}} = |\Gamma|_{\hat{\alpha}} - |\Omega_i|$ , it is decidable whether  $\Gamma' \vdash \Omega'$ .

In other words, for any instantiation judgment prefix  $\Omega_i$ , the algorithm either reduces the number of  $\forall$ 's or solves one existential variable per instantiation judgment. The proof of this lemma is by induction on the measure  $2 * |\Omega_i|_{\rightarrow} + |\Omega_i|$  of the instantiation judgment list.

In summary, the decidability theorem can be shown through a lexicographic group of induction measurements  $\langle |\Omega|_{\forall}, |\Omega|_{\hat{\alpha}}, |\Omega|_{\rightarrow} \rangle$ . The critical case is that, whenever we encounter an instantiation judgment at the front of the worklist, we refer to Lemma 3.7, which reduces the number of unsolved variables by consuming that instantiation judgment, or reduces the number of  $\forall$ -quantifiers. Other cases are relatively straightforward.

### 3.4 THE CHOICE OF ABELLA

We have chosen the Abella (v2.0.5) proof assistant [Gacek 2008] to develop our formalization. Although Abella takes a two-level logic approach, where the specification logic can be expressed separately from the reasoning logic, we only make use of its reasoning logic, due to the difficulty of expressing our algorithmic rules with only the specification. Abella is particularly helpful due to its built-in support for variable bindings, and its  $\lambda$ -tree syntax [Miller 2000] is a form of HOAS, which helps with the encoding and reasoning about substitutions. For instance, the type  $\forall x. x \rightarrow a$  is encoded as `all (x \ arrow x a)`, where `x \ arrow x a` is a lambda abstraction in Abella. An opening  $[b/x](x \rightarrow a)$  is encoded as an application `(x \ arrow x a) b`, which can be simplified (evaluated) to `arrow b a`. Name supply and freshness conditions are controlled by the  $\nabla$ -quantifier. The expression `new x, F` means that `x` is a unique variable in `F`, i.e. it is different from any other names occurring elsewhere. Such variables are called nominal constants. They can be of any type, in other words, every type may contain an unlimited number of such atomic nominal constants.

**ENCODING OF THE DECLARATIVE SYSTEM** As a concrete example, our declarative context (well-formedness relation) and well-formedness rules are encoded as follows.

```
Kind ty      type.
Type i       ty.                % the unit type
Type all     (ty → ty) → ty.    % forall-quantifier
Type arrow   ty → ty → ty.      % function type
Type bound   ty → o.            % variable collection in contexts

Define env : olist → prop by
```

```
env nil;
nabla x, env (bound x :: E) := env E.
```

```
Define wft: olist → ty → prop by
wft E i;
nabla x, wft (E x) x := nabla x, member (bound x) (E x);
wft E (arrow A B) := wft E A ∧ wft E B;
wft E (all A) := nabla x, wft (bound x :: E) (A x).
```

In the above code, we first define the syntax of types and contexts in our type system. A type of our system has type `ty`, and the context is of type `olist`. We use the type `olist` just as normal list of `o` with two constructors, namely `nil: olist` and `(::) : o → olist → olist`, where `o` purely means “the element type of `olist`”, and both are built-in types of Abella. The `member: o → olist → prop` relation is also pre-defined.

Note that in Abella, the symbol `(:=)` used in definitions is similar to the `(:-)` symbol in Prolog, where we write the condition to its right and the conclusion to its left. The second case of the relation `wft` states Rule `wfdvar`. The encoding `(E x)` basically means that the context *may* contain `x`. If we write `(E x)` as `E`, then the context should not contain `x`, and both `wft E x` and `member (bound x) E` make no sense. Instead, we treat `E : ty → olist` as an *abstract structure* of a context, such as `x \ bound x :: bound a :: nil`. For the fourth case of the relation `wft`, the type  $\forall x. A$  in our target language is expressed as `(all A)`, and its opening `A`, `(A x)`.

**ENCODING OF THE ALGORITHMIC SYSTEM** In terms of the algorithmic system, notably, Abella handles the  $\leq_a \text{instL}$  and  $\leq_a \text{instR}$  rules in a nice way:

```
% sub_alg_list : enva → [subty_judgment] → prop
Define subal: olist → olist → prop by
subal E nil;
subal E (subt i i :: Exp) := subal E Exp;
% some cases omitted ...
% <: instL
nabla x, subal (E x) (subt x (arrow A B) :: Exp x) :=
  exists E1 E2 F, nabla x y z, append E1 (exvar x :: E2) (E x) ∧
    append E1 (exvar y :: exvar z :: E2) (F y z) ∧
    subal (F y z) (subt (arrow y z) (arrow A B) :: Exp (arrow y z));
% <: instR is symmetric to <: instL, omitted here
% other cases omitted ...
```

In this piece of code, we use `Exp` to denote the worklist  $\Omega$ . An algorithmic existential variable  $\hat{\alpha}$  is constructed by applying the Abella term `exvar` to an actual `nabla` quantified variable.

File(s)	SLOC	# of Theorems	Description
olist.thm, nat.thm	303	55	Basic data structures
higher.thm, order.thm	164	15	Declarative system
higher_alg.thm	618	44	Algorithmic system
trans.thm	411	46	Transfer
sound.thm	166	2	Soundness theorem
depth.thm	143	12	Definition of depth
complete.thm	626	28	Lemmas and Completeness theorem
decidable.thm	1077	53	Lemmas and Decidability theorem
Total	3627	267	(33 definitions in total)

Figure 3.8: Statistics for the proof scripts

Thanks to the way Abella deals with nominal constants, the pattern `subt x (arrow A B)` implicitly states that  $x \notin FV(A) \wedge x \notin FV(B)$ . If the condition were not required, we would have encoded the pattern as `subt x (arrow (A x) (B x))` instead.

#### 3.4.1 STATISTICS AND DISCUSSION

Some basic statistics on our proof script are shown in Figure 3.8. The proof consists of 3627 lines of code with a total of 33 definitions and 267 theorems. We have to mention that Abella provides few built-in tactics and does not support user-defined ones, and we would reduce significant lines of code if Abella provided more handy tactics. Moreover, the definition of natural numbers, the plus operation and less-than relation are defined within our proof due to Abella's lack of packages. However, the way Abella deals with name bindings is very helpful for type system formalizations and substitution-intensive formalizations, such as this one.

# 4 A TYPE-INFERENCING ALGORITHM FOR HIGHER-RANKED POLYMORPHISM

This chapter presents the first fully mechanized formalization of the metatheory for higher-ranked polymorphic type inference. Following the worklist subtyping algorithm in the previous chapter, we address several drawbacks and extend the technique to DK’s bidirectional type system [Dunfield and Krishnaswami 2013]. We chose DK’s type system because it is quite elegant, well-documented and it comes with detailed manually written proofs. Furthermore, the system is adopted in practice by a few real implementations of functional languages, including PureScript and Unison. The DK type system has two variants: a declarative and an algorithmic one. The two variants have been *manually* proved to be *sound*, *complete* and *decidable*. We present a mechanical formalization in the Abella theorem prover [Gacek 2008] for DK’s declarative type system using a different algorithm.

**CHALLENGES** While our initial goal was to formalize both DK’s declarative and algorithmic versions, we faced technical challenges with the latter, prompting us to find an alternative formulation.

The first challenge that we faced were missing details as well as a few incorrect proofs and lemmas in DK’s formalization. While DK’s original formalization comes with very well-written manual proofs, there are still several details missing. These complicate the task of writing a mechanically verified proof. Moreover, some proofs and lemmas are wrong and, in some cases, it is not clear to us how to fix them.

Despite the problems in DK’s manual formalization, we believe that these problems do not invalidate their work and that their results are still true. In fact we have nothing but praise for their detailed and clearly written metatheory and proofs, which provided invaluable help to our own work. We expect that for most non-trivial manual proofs similar problems exist, so this should not be understood as a sign of sloppiness on their part. Instead, it should be an indicator that reinforces the arguments for mechanical formalizations: manual formalizations are error-prone due to the multiple tedious details involved in them. There are several other examples of manual formalizations that were found to have similar problems. For ex-



ample, Klein et al. [2012] mechanized formalizations in Redex for nine ICFP 2009 papers and all were found to have mistakes.

Another challenge was variable binding. Type inference algorithms typically do not rely simply on local environments but instead propagate information across judgments. While local environments are well-studied in mechanical formalizations, there is little work on how to deal with the complex forms of binding employed by type inference algorithms in theorem provers. To keep track of variable scoping, DK’s algorithmic version employs input and output contexts to track information that is discovered through type inference. However, modeling output contexts in a theorem prover is non-trivial.

Due to those two challenges, our work takes a different approach by refining and extending the idea of *worklist judgments* in Chapter 3, where we mechanically formalized an algorithm for *polymorphic subtyping* [Odersky and Läufer 1996]. The algorithm eliminates output contexts compared to DK’s algorithm, and therefore the problem of variable binding is solved. Unfortunately, the subtyping algorithm cannot be naively extended to support the bidirectional type system. A key innovation in the new algorithm to be introduced in this chapter is how to adapt the idea of worklist judgments to *inference judgments*, which are not needed for polymorphic subtyping, but are necessary for type inference. The idea is to use a *continuation passing style* to enable the transfer of inferred information across judgments. A further refinement to the idea of worklist judgments is the *unification between ordered contexts* [Dunfield and Krishnaswami 2013; Gundry et al. 2010] and *worklists*. This enables precise scope tracking of free variables in judgments. Furthermore it avoids the duplication of context information across judgments in worklists that occurs in other techniques [Abel and Pientka 2011; Reed 2009]. Despite the use of a different algorithm, we prove the same results as DK, although with significantly different proofs and proof techniques. The calculus and its metatheory have been fully formalized in the Abella theorem prover [Gacek 2008].

## 4.1 OVERVIEW

This section starts with a discussion on DK’s declarative type system. Then it introduces several techniques that have been used in algorithmic formulations, and which have influenced our own algorithmic design. Finally we introduce the novelties of our new algorithm. In particular the support for inference judgments in worklist judgments, and a new form of worklist judgment that unifies *ordered contexts* and the worklists themselves.

Type variables	$a, b$		
Types	$A, B, C$	$::=$	$1 \mid a \mid \forall a. A \mid A \rightarrow B$
Monotypes	$\tau, \sigma$	$::=$	$1 \mid a \mid \tau \rightarrow \sigma$
Expressions	$e$	$::=$	$x \mid () \mid \lambda x. e \mid e_1 e_2 \mid (e : A)$
Contexts	$\Psi$	$::=$	$\cdot \mid \Psi, a \mid \Psi, x : A$

Figure 4.1: Syntax of Declarative System (Extends Figure 3.1)

$$\begin{array}{c}
 \boxed{\Psi \vdash A} \text{ Well-formed declarative type} \\
 \\
 \frac{}{\Psi \vdash 1} \text{ wf}_{\text{dunit}} \quad \frac{a \in \Psi}{\Psi \vdash a} \text{ wf}_{\text{dvar}} \quad \frac{\Psi \vdash A \quad \Psi \vdash B}{\Psi \vdash A \rightarrow B} \text{ wf}_{\text{d}\rightarrow} \quad \frac{\Psi, a \vdash A}{\Psi \vdash \forall a. A} \text{ wf}_{\text{d}\forall} \\
 \\
 \boxed{\Psi \vdash e} \text{ Well-formed declarative expression} \\
 \\
 \frac{x : A \in \Psi}{\Psi \vdash x} \text{ wf}_{\text{d}\text{tmvar}} \quad \frac{}{\Psi \vdash ()} \text{ wf}_{\text{d}\text{tmunit}} \quad \frac{\Psi, x : A \vdash e}{\Psi \vdash \lambda x. e} \text{ wf}_{\text{d}\text{abs}} \\
 \\
 \frac{\Psi \vdash e_1 \quad \Psi \vdash e_2}{\Psi \vdash e_1 e_2} \text{ wf}_{\text{d}\text{app}} \quad \frac{\Psi \vdash A \quad \Psi \vdash e}{\Psi \vdash (e : A)} \text{ wf}_{\text{d}\text{anno}} \\
 \\
 \boxed{\Psi \vdash A \leq B} \text{ Declarative subtyping} \\
 \\
 \frac{a \in \Psi}{\Psi \vdash a \leq a} \leq_{\text{Var}} \quad \frac{}{\Psi \vdash 1 \leq 1} \leq_{\text{Unit}} \quad \frac{\Psi \vdash B_1 \leq A_1 \quad \Psi \vdash A_2 \leq B_2}{\Psi \vdash A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2} \leq_{\rightarrow} \\
 \\
 \frac{\Psi \vdash \tau \quad \Psi \vdash [\tau/a]A \leq B}{\Psi \vdash \forall a. A \leq B} \leq_{\forall L} \quad \frac{\Psi, b \vdash A \leq B}{\Psi \vdash A \leq \forall b. B} \leq_{\forall R}
 \end{array}$$

Figure 4.2: Declarative Well-formedness and Subtyping

#### 4.1.1 DK'S DECLARATIVE SYSTEM

Subsection 2.3.1 introduces DK's declarative subtyping and typing systems. We also duplicate the rules here for the convenience of the reader.

**OVERLAPPING RULES** A problem that we found in the declarative system is that some of the rules overlap with each other. Declarative subtyping rules  $\leq_{\forall L}$  and  $\leq_{\forall R}$  both match the conclusion  $\Psi \vdash \forall a. A \leq \forall a. B$ . In such a case, choosing  $\leq_{\forall R}$  first is always better, since we introduce the type variable  $a$  to the context earlier, which gives more flexibility on the choice of  $\tau$ . The declarative typing rule  $\text{Dec1Sub}$  overlaps with both  $\text{Dec1}\forall I$  and  $\text{Dec1}\rightarrow I$ .

$\boxed{\Psi \vdash e \Leftarrow A}$	$e$ checks against input type $A$ .
$\boxed{\Psi \vdash e \Rightarrow A}$	$e$ synthesizes output type $A$ .
$\boxed{\Psi \vdash A \bullet e \Rightarrow C}$	Applying a function of type $A$ to $e$ synthesizes type $C$ .

$\frac{(x : A) \in \Psi}{\Psi \vdash x \Rightarrow A} \text{DeclVar}$	$\frac{\Psi \vdash e \Rightarrow A \quad \Psi \vdash A \leq B}{\Psi \vdash e \Leftarrow B} \text{DeclSub}$
$\frac{\Psi \vdash A \quad \Psi \vdash e \Leftarrow A}{\Psi \vdash (e : A) \Rightarrow A} \text{DeclAnno}$	$\frac{}{\Psi \vdash () \Leftarrow 1} \text{Decl1I} \quad \frac{}{\Psi \vdash () \Rightarrow 1} \text{Decl1I} \Rightarrow$
$\frac{\Psi, a \vdash e \Leftarrow A}{\Psi \vdash e \Leftarrow \forall a. A} \text{Decl} \forall \text{I}$	$\frac{\Psi \vdash \tau \quad \Psi \vdash [\tau/a]A \bullet e \Rightarrow C}{\Psi \vdash \forall a. A \bullet e \Rightarrow C} \text{Decl} \forall \text{App}$
$\frac{\Psi, x : A \vdash e \Leftarrow B}{\Psi \vdash \lambda x. e \Leftarrow A \rightarrow B} \text{Decl} \rightarrow \text{I}$	$\frac{\Psi \vdash \sigma \rightarrow \tau \quad \Psi, x : \sigma \vdash e \Leftarrow \tau}{\Psi \vdash \lambda x. e \Rightarrow \sigma \rightarrow \tau} \text{Decl} \rightarrow \text{I} \Rightarrow$
$\frac{\Psi \vdash e_1 \Rightarrow A \quad \Psi \vdash A \bullet e_2 \Rightarrow C}{\Psi \vdash e_1 e_2 \Rightarrow C} \text{Decl} \rightarrow \text{E}$	$\frac{\Psi \vdash e \Leftarrow A}{\Psi \vdash A \rightarrow C \bullet e \Rightarrow C} \text{Decl} \rightarrow \text{App}$

Figure 4.3: Declarative Typing

However, we argue that more specific rules are always the best choices, i.e.  $\text{Decl} \forall \text{I}$  and  $\text{Decl} \rightarrow \text{I}$  should have higher priority than  $\text{DeclSub}$ .

For example,  $\Psi \vdash \lambda x. x \Leftarrow \forall a. a \rightarrow a$  succeeds if derived from Rule  $\text{Decl} \forall \text{I}$ :

$$\frac{\Psi, a, x : a \vdash x \Leftarrow a}{\Psi, a \vdash \lambda x. x \Leftarrow a \rightarrow a} \text{Decl} \rightarrow \text{I} \quad \frac{}{\Psi \vdash \lambda x. x \Leftarrow \forall a. a \rightarrow a} \text{Decl} \forall \text{I},$$

but fails when applied to  $\text{DeclSub}$ :

$$\frac{\Psi \vdash \sigma \rightarrow \tau \quad \Psi, x : \sigma \vdash e \Leftarrow \tau}{\Psi \vdash \lambda x. x \Rightarrow \sigma \rightarrow \tau} \text{Decl} \rightarrow \text{I} \Rightarrow \quad \frac{\text{Impossible!} \quad a \notin \text{FV}(\sigma) \quad ?}{\Psi, a \vdash a \leq \sigma \quad \Psi, a \vdash \tau \leq a} \leq \rightarrow \quad \frac{\Psi, a \vdash \sigma \rightarrow \tau \leq a \rightarrow a}{\Psi \vdash \sigma \rightarrow \tau \leq \forall a. a \rightarrow a} \leq \forall \text{R} \quad \frac{}{\Psi \vdash \lambda x. x \Leftarrow \forall a. a \rightarrow a} \text{DeclSub}.$$

Rule  $\text{Decl} \rightarrow \text{I}$  is also better at handling higher-order types. When the lambda-expression to be inferred has a polymorphic input type, such as  $\forall a. a \rightarrow a$ ,  $\text{DeclSub}$  may not derive

some judgments. For example,  $\Psi, id : \forall a. a \rightarrow a \vdash \lambda f. f \text{ id } (f ()) \Leftarrow (\forall a. a \rightarrow a) \rightarrow 1$  requires the argument of the lambda-expression to be a polymorphic type, otherwise it could not be applied to both  $id$  and  $()$ . If Rule  $\text{Dec1Sub}$  was chosen for derivation, the type of its argument is restricted by Rule  $\text{Dec1} \rightarrow \text{I} \Rightarrow$ , which is not a polymorphic type. By contrast, Rule  $\text{Dec1} \rightarrow \text{I}$  keeps the polymorphic argument type  $\forall a. a \rightarrow a$ , and will successfully derive the judgment.

We will come back to this topic in Section 4.3.2 and formally derive a system without overlapping rules.

#### 4.1.2 DK'S ALGORITHM

DK's algorithm version revolves around their notion of *algorithmic context*.

Algorithmic Contexts  $\Gamma, \Delta, \Theta ::= \cdot \mid \Gamma, a \mid \Gamma, x : A \mid \Gamma, \hat{a} \mid \Gamma, \hat{a} = \tau \mid \Gamma, \blacktriangleright_{\hat{a}}$

In addition to the regular (universally quantified) type variables  $a$ , the algorithmic context also contains *existential* type variables  $\hat{a}$ . These are placeholders for monotypes  $\tau$  that are still to be determined by the inference algorithm. When the existential variable is “solved”, its entry in the context is replaced by the assignment  $\hat{a} = \tau$ . A context application on a type, denoted by  $[\Gamma]A$ , substitutes all solved existential type variables in  $\Gamma$  with their solutions on type  $A$ .

All algorithmic judgments thread an algorithmic context. They have the form  $\Gamma \vdash \dots \dashv \Delta$ , where  $\Gamma$  is the input context and  $\Delta$  is the output context:  $\Gamma \vdash A \leq B \dashv \Delta$  for subtyping,  $\Gamma \vdash e \Leftarrow A \dashv \Delta$  for type checking, and so on. The output context is a functional update of the input context that records newly introduced existentials and solutions.

Solutions are incrementally propagated by applying the algorithmic output context of a previous task as substitutions to the next task. This can be seen in the subsumption rule:

$$\frac{\Gamma \vdash e \Rightarrow A \dashv \Theta \quad \Theta \vdash [\Theta]A \leq [\Theta]B \dashv \Delta}{\Gamma \vdash e \Leftarrow B \dashv \Delta} \text{DK\_Sub}$$

The inference task yields an output context  $\Theta$  which is applied as a substitution to the types  $A$  and  $B$  before performing the subtyping check to propagate any solutions of existential variables that appear in  $A$  and  $B$ .

**MARKERS FOR SCOPING.** The sequential order of entries in the algorithmic context, in combination with the threading of contexts, does not perfectly capture the scoping of all existen-

tial variables. For this reason the DK algorithm uses scope markers  $\blacktriangleright_{\hat{\alpha}}$  in a few places. An example is given in the following rule:

$$\frac{\Gamma, \blacktriangleright_{\hat{\alpha}}, \hat{\alpha} \vdash [\hat{\alpha}/a]A \leq B \dashv \Delta, \blacktriangleright_{\hat{\alpha}}, \Theta}{\Gamma \vdash \forall a. A \leq B \dashv \Delta} \text{DK}_{\leq \forall L}$$

To indicate that the scope of  $\hat{\alpha}$  is local to the subtyping check  $[\hat{\alpha}/a]A \leq B$ , the marker is pushed onto its input stack and popped from the output stack together with the subsequent part  $\Theta$ , which may refer to  $\hat{\alpha}$ . (Remember that later entries may refer to earlier ones, but not vice versa.) This way  $\hat{\alpha}$  does not escape its scope.

At first sight, the DK algorithm would seem to be a good basis for mechanization. After all, it comes with a careful description and extensive manual proofs. Unfortunately, we ran into several obstacles that have prompted us to formulate a different, more mechanization-friendly algorithm.

**BROKEN METATHEORY** While going through the manual proofs of DK’s algorithm, we found several problems. Indeed, two proofs of lemmas—Lemma 19 (Extension Equality Preservation) and Lemma 14 (Subsumption)—wrongly apply induction hypotheses in several cases. Fortunately, we have found simple workarounds that fix these proofs without affecting the appeals to these lemmas.

More seriously, we have also found a lemma that simply does not hold: Lemma 29 (Parallel Admissibility)<sup>1</sup>. This lemma is a cornerstone of the two metatheoretical results of the algorithm, soundness, and completeness with respect to the declarative system. In particular, both instantiation soundness (i.e. a part of subtyping soundness) and typing completeness directly require the broken lemma. Moreover, Lemma 54 (Typing Extension) also requires the broken lemma and is itself used 13 times in the proof of typing soundness and completeness. Unfortunately, we have not yet found a way to fix this problem.

In what follows, we briefly discuss the problem through counterexamples. False lemmas are found in the manual proofs of DK’s two papers [Dunfield and Krishnaswami 2013] and [Dunfield and Krishnaswami 2019].

- In the first paper, Lemma 29 on page 9 of its appendix says:

**Lemma 4.1** (Parallel Admissibility of [Dunfield and Krishnaswami 2013]).

*If  $\Gamma_L \longrightarrow \Delta_L$  and  $\Gamma_L, \Gamma_R \longrightarrow \Delta_L, \Delta_R$  then:*

<sup>1</sup>Ningning Xie found the issue with Lemma 29 in 2016 on an earlier attempt to mechanically formalize DK’s algorithm. The authors acknowledged the problem after we contacted them through email. Although they briefly mentioned that it should be possible to use a weaker lemma instead they did not go into details.

1.  $\Gamma_L, \hat{\alpha}, \Gamma_R \longrightarrow \Delta_L, \hat{\alpha}, \Delta_R$
2. If  $\Delta_L \vdash \tau'$  then  $\Gamma_L, \hat{\alpha}, \Gamma_R \longrightarrow \Delta_L, \hat{\alpha} = \tau', \Delta_R$ .
3. If  $\Gamma_L \vdash \tau$  and  $\Delta_L \vdash \tau'$  and  $[\Delta_L]\tau = [\Delta_L]\tau'$ , then  $\Gamma_L, \hat{\alpha} = \tau, \Gamma_R \longrightarrow \Delta_L, \hat{\alpha} = \tau', \Delta_R$ .

We give a counter-example to this lemma:

Pick  $\Gamma_L = \cdot, \Gamma_R = \hat{\beta}, \Delta_L = \hat{\beta}, \Delta_R = \cdot$ , then both conditions  $\cdot \longrightarrow \hat{\beta}$  and  $\hat{\beta} \longrightarrow \hat{\beta}$  hold, but the first conclusion  $\hat{\alpha}, \hat{\beta} \longrightarrow \hat{\beta}, \hat{\alpha}$  does not hold.

- In the second paper, as an extended work to the first paper, Lemma 26 on page 22 of its appendix says:

**Lemma 4.2** (Parallel Admissibility of [Dunfield and Krishnaswami 2019]).

If  $\Gamma_L \longrightarrow \Delta_L$  and  $\Gamma_L, \Gamma_R \longrightarrow \Delta_L, \Delta_R$  then:

1.  $\Gamma_L, \hat{\alpha} : \kappa, \Gamma_R \longrightarrow \Delta_L, \hat{\alpha} : \kappa, \Delta_R$
2. If  $\Delta_L \vdash \tau' : \kappa$  then  $\Gamma_L, \hat{\alpha} : \kappa, \Gamma_R \longrightarrow \Delta_L, \hat{\alpha} : \kappa = \tau', \Delta_R$ .
3. If  $\Gamma_L \vdash \tau : \kappa$  and  $\Delta_L \vdash \tau'$  types and  $[\Delta_L]\tau = [\Delta_L]\tau'$ , then  $\Gamma_L, \hat{\alpha} : \kappa = \tau, \Gamma_R \longrightarrow \Delta_L, \hat{\alpha} : \kappa = \tau', \Delta_R$ .

A similar counter-example is given:

Pick  $\Gamma_L = \cdot, \Gamma_R = \hat{\beta} : \star, \Delta_L = \hat{\beta} : \star, \Delta_R = \cdot$ , then both conditions  $\cdot \longrightarrow \hat{\beta} : \star$  and  $\hat{\beta} : \star \longrightarrow \hat{\beta} : \star$  hold, but the first conclusion  $\hat{\alpha} : \kappa, \hat{\beta} : \star \longrightarrow \hat{\beta} : \star, \hat{\alpha} : \kappa$  does not hold.

**COMPLEX SCOPING AND PROPAGATION** Besides the problematic lemmas in DK's metatheory, there are other reasons to look for an alternative algorithmic formulation of the type system that is more amenable to mechanization. Indeed, two aspects that are particularly challenging to mechanize are the scoping of universal and existential type variables, and the propagation of the instantiation of existential type variables across judgments. DK is already quite disciplined on these accounts, in particular compared to traditional constraint-based type-inference algorithms like Algorithm  $\mathcal{W}$  [Milner 1978] which features an implicit global scope for all type variables. Indeed, DK uses its explicit and ordered context  $\Gamma$  that tracks the relative scope of universal and existential variables and it is careful to only instantiate existential variables in a well-scoped manner.

Moreover, DK's algorithm carefully propagates instantiations by recording them into the context and threading this context through all judgments. While this works well on paper, this approach is still fairly involved and thus hard to reason about in a mechanized setting.

Indeed, the instantiations have to be recorded in the context and are applied incrementally to each remaining judgment in turn, rather than immediately to all remaining judgments at once. Also, as we have mentioned above, the stack discipline of the ordered contexts does not mesh well with the use of output contexts; explicit marker entries are needed in two places to demarcate the end of an existential variable's scope. Actually, we found a scoping issue related to the subsumption rule  $\text{DK\_Sub}$ , which might cause existential variables to leak across judgments. In Section 4.4.1 we give a detailed discussion.

The complications of scoping and propagation are compelling reasons to look for another algorithm that is easier to reason about in a mechanized setting.

#### 4.1.3 JUDGMENT LISTS

To avoid the problem of incrementally applying a substitution to remaining tasks, we can find inspiration in the formulation of constraint solving algorithms. For instance, the well-known unification algorithm by Martelli and Montanari [1982] decomposes the problem of unifying two terms  $s \doteq t$  into a number of related unification problems between pairs of terms  $s_i \doteq t_i$ . These smaller problems are not tackled independently, but kept together in a set  $S$ . The algorithm itself is typically formulated as a small-step-style state transition system  $S \rightarrow S'$  that rewrites the set of unification problems until it is in solved form or until a contradiction has been found. For instance, the variable elimination rule is written as:

$$x \doteq t, S \rightarrow x \doteq t, [t/x]S \quad \text{if } x \notin t \text{ and } x \in S$$

Because the whole set is explicitly available, the variable  $x$  can be simultaneously substituted.

In the above unification problem, all variables are implicitly bound in the same global scope. Some constraint solving algorithms for Hindley-Milner type inference use similar ideas, but are more careful tracking the scopes of variables [Pottier and Rémy 2005]. Recent unification algorithms for dependently-typed languages are also more explicit about scopes. For instance, Reed [2009] represents a unification problem as  $\Delta \vdash P$  where  $P$  is a set of equations to be solved and  $\Delta$  is a (modal) context. Abel and Pientka [2011] even use multiple contexts within a unification problem. Such a problem is denoted  $\Delta \Vdash \mathcal{K}$  where the meta-context  $\Delta$  contains all the typings of meta-variables in the constraint set  $\mathcal{K}$ . The latter consists of constraints like  $\Psi \vdash M = N : C$  that are equipped with their individual context  $\Psi$ . While accurately tracking the scoping of regular and meta-variables, this approach is not ideal because it repeatedly copies contexts when decomposing a unification problem, and later it has to substitute solutions into these copies.

#### 4.1.4 SINGLE-CONTEXT WORKLIST ALGORITHM FOR SUBTYPING

As we have seen in Chapter 3, an algorithm based on *worklist judgments* is mechanized and shown to be correct with respect to DK’s declarative subtyping judgment. This approach overcomes some problems in DK’s algorithmic formulation by using a worklist-based judgment of the form

$$\Gamma \vdash \Omega$$

where  $\Omega$  is a worklist (or sequence) of subtyping problems of the form  $A \leq B$ . The judgment is defined by case analysis on the first element of  $\Omega$  and recursively processes the worklist until it is empty. Using both a single common ordered context  $\Gamma$  and a worklist  $\Omega$  greatly simplifies propagating the instantiation of type variables in one subtyping problem to the remaining ones.

Unfortunately, this work does not solve all problems. In particular, it has two major limitations that prevent it from scaling to the whole DK system.

**SCOPING GARBAGE** Firstly, the single common ordered context  $\Gamma$  does not accurately reflect the type and unification variables currently in scope. Instead, it is an overapproximation that steadily accrues variables, and only drops those unification variables that are instantiated. In other words,  $\Gamma$  contains “garbage” that is no longer in scope. This complicates establishing the relation with the declarative system.

**NO INFERENCE JUDGMENTS** Secondly, and more importantly, the approach only deals with a judgment for *checking* whether one type is the subtype of another. While this may not be so difficult to adapt to the *checking* mode of term typing  $\Gamma \vdash e \Leftarrow A$ , it clearly lacks the functionality to support the *inference* mode of term typing  $\Gamma \vdash e \Rightarrow A$ . Indeed, the latter requires not only the communication of unification variable instantiations from one typing problem to another, but also an inferred type.

#### 4.1.5 ALGORITHMIC TYPE INFERENCE FOR HIGHER-RANKED TYPES: KEY IDEAS

Our new algorithmic type system builds on the work above, but addresses the open problems.

**SCOPE TRACKING** We avoid scoping garbage by blending the ordered context and the worklist into a single syntactic sort  $\Gamma$ , our algorithmic worklist. This algorithmic worklist interleaves (type and term) variables with *work* like checking or inferring types of expressions. The interleaving keeps track of the variable scopes in the usual, natural way: each variable is in scope of anything in front of it in the worklist. If there is nothing in front, the variable is



no longer needed and can be popped from the worklist. This way, no garbage (i.e. variables out-of-scope) builds up.

$$\begin{array}{ll} \text{Algorithmic judgment chain} & \omega ::= A \leq B \mid e \Leftarrow A \mid e \Rightarrow_a \omega \mid A \bullet e \Rightarrow_a \omega \\ \text{Algorithmic worklist} & \Gamma ::= \cdot \mid \Gamma, a \mid \Gamma, \hat{\alpha} \mid \Gamma, x : A \mid \Gamma \Vdash \omega \end{array}$$

For example, suppose that the top judgment of the following worklist checks the identity function against  $\forall a. a \rightarrow a$ :

$$\Gamma \Vdash \lambda x. x \Leftarrow \forall a. a \rightarrow a$$

To proceed, two auxiliary variables  $a$  and  $x$  are introduced to help the type checking:

$$\Gamma, a, x : a \Vdash x \Leftarrow a$$

which will be satisfied after several steps, and the worklist becomes

$$\Gamma, a, x : a$$

Since the variable declarations  $a, x : a$  are only used for a judgment already processed, they can be safely removed, leaving the remaining worklist  $\Gamma$  to be further reduced.

Our worklist can be seen as an all-in-one stack, containing variable declarations and subtyping/ typing judgments. The stack is an enriched form of ordered context, and it has a similar variable scoping scheme.

**INFERENCE JUDGMENTS** To express the DK's inference judgments, we use a novel form of work entries in the worklist: our algorithmic judgment chains  $\omega$ . In its simplest form, such a judgment chain is just a check, like a subtyping check  $A \leq B$  or a term typecheck  $e \Leftarrow A$ . However, the non-trivial forms of chains capture an inference task together with the work that depends on its outcome. For instance, a type inference task takes the form  $e \Rightarrow_a \omega$ . This form expresses that we need to infer the type, say  $A$ , for expression  $e$  and use it in the chain  $\omega$  by substituting it for the placeholder type variable  $a$ . Notice that such  $a$  binds a fresh type variable for the inner chain  $\omega$ , which behaves similarly to the variable declarations in the context.

Take the following worklist as an example

$$\hat{\alpha} \Vdash \underline{(\lambda x. x) () \Rightarrow_a a \leq \hat{\alpha}}, x : \hat{\alpha}, \hat{\beta} \Vdash \underline{\hat{\alpha} \leq \hat{\beta}} \Vdash \underline{\hat{\beta} \leq 1}$$

There are three (underlined) judgment chains in the worklist, where the first and second judgment chains (from the right) are two subtyping judgments, and the third judgment chain,  $(\lambda x. x) () \Rightarrow_a a \leq \hat{\alpha}$ , is a sequence of an inference judgment followed by a subtyping judgment.

The algorithm first analyses the two subtyping judgments and will find the best solutions  $\hat{\alpha} = \hat{\beta} = 1$  (please refer to Figure 4.5 for detailed derivations). Then we substitute every instance of  $\hat{\alpha}$  and  $\hat{\beta}$  with 1, so the variable declarations can be safely removed from the worklist. Now we reduce the worklist to the following state

$$\cdot \Vdash \underline{(\lambda x. x) () \Rightarrow_a a \leq 1, x : 1}$$

which has a term variable declaration as the top element. After removing the garbage term variable declaration from the worklist, we process the last remaining inference judgment  $(\lambda x. x) () \Rightarrow ?$ , with the unit type 1 as its result. Finally, the last judgment becomes  $1 \leq 1$ , a trivial base case.

## 4.2 ALGORITHMIC SYSTEM

This section introduces a novel algorithmic system that implements DK's declarative specification. The new algorithm extends the idea of worklists in Chapter 3 in two ways. Firstly, unlike its worklists, the scope of variables is precisely tracked and variables do not escape their scope. This is achieved by unifying algorithmic contexts and the worklists themselves. Secondly, our algorithm also accounts for the type system (and not just subtyping). To deal with inference judgments that arise in the type system we employ a *continuation passing style* to enable the transfer of inferred information across judgments in a worklist.

### 4.2.1 SYNTAX AND WELL-FORMEDNESS

Figure 4.4 shows the syntax and well-formedness judgments used by the algorithm. Similarly to the declarative system the well-formedness rules are unsurprising and merely ensure well-scopedness.

**EXISTENTIAL VARIABLES** The algorithmic system inherits the syntax of terms and types from the declarative system. It only introduces one additional feature. In order to find unknown types  $\tau$  in the declarative system, the algorithmic system extends the declarative types  $A$  with *existential variables*  $\hat{\alpha}, \hat{\beta}$ . They behave like unification variables, but their scope is restricted by their position in the algorithmic worklist rather than being global. Any existential variable  $\hat{\alpha}$  should only be solved to a type that is well-formed with respect to the worklist to

Existential variables	$\hat{\alpha}, \hat{\beta}$
Algorithmic types	$A, B, C ::= 1 \mid a \mid \forall a. A \mid A \rightarrow B \mid \hat{\alpha}$
Algorithmic judgment chain	$\omega ::= A \leq B \mid e \Leftarrow A \mid e \Rightarrow_a \omega \mid A \bullet e \Rightarrow_a \omega$
Algorithmic worklist	$\Gamma ::= \cdot \mid \Gamma, a \mid \Gamma, \hat{\alpha} \mid \Gamma, x : A \mid \Gamma \vdash \omega$

$\boxed{\Gamma \vdash A}$  Well-formed algorithmic type

$$\begin{array}{c}
 \frac{}{\Gamma \vdash 1} \text{wf\_unit} \quad \frac{a \in \Gamma}{\Gamma \vdash a} \text{wf\_var} \quad \frac{\hat{\alpha} \in \Gamma}{\Gamma \vdash \hat{\alpha}} \text{wf\_exvar} \\
 \frac{\Gamma \vdash A \quad \Gamma \vdash B}{\Gamma \vdash A \rightarrow B} \text{wf\_}\rightarrow \quad \frac{\Gamma, a \vdash A}{\Gamma \vdash \forall a. A} \text{wf\_}\forall
 \end{array}$$

$\boxed{\Gamma \vdash e}$  Well-formed algorithmic expression

$$\begin{array}{c}
 \frac{x : A \in \Gamma}{\Gamma \vdash x} \text{wf\_tmvar} \quad \frac{}{\Gamma \vdash ()} \text{wf\_tmunit} \quad \frac{\Gamma, x : A \vdash e}{\Gamma \vdash \lambda x. e} \text{wf\_abs} \\
 \frac{\Gamma \vdash e_1 \quad \Gamma \vdash e_2}{\Gamma \vdash e_1 e_2} \text{wf\_app} \quad \frac{\Gamma \vdash A \quad \Gamma \vdash e}{\Gamma \vdash (e : A)} \text{wf\_anno}
 \end{array}$$

$\boxed{\Gamma \vdash \omega}$  Well-formed algorithmic judgment

$$\begin{array}{c}
 \frac{\Gamma \vdash A \quad \Gamma \vdash B}{\Gamma \vdash A \leq B} \text{wf\_}\leq \quad \frac{\Gamma \vdash e \quad \Gamma \vdash A}{\Gamma \vdash e \Leftarrow A} \text{wf\_}\Leftarrow \\
 \frac{\Gamma \vdash e \quad \Gamma, a \vdash \omega}{\Gamma \vdash e \Rightarrow_a \omega} \text{wf\_}\Rightarrow \quad \frac{\Gamma \vdash A \quad \Gamma, a \vdash \omega \quad \Gamma \vdash e}{\Gamma \vdash A \bullet e \Rightarrow_a \omega} \text{wf\_}\Rightarrow\Rightarrow
 \end{array}$$

$\boxed{\text{wf } \Gamma}$  Well-formed algorithmic worklist

$$\begin{array}{c}
 \frac{}{\text{wf } \cdot} \text{wf\_}\cdot \quad \frac{\text{wf } \Gamma}{\text{wf } \Gamma, a} \text{wf\_}\text{a} \quad \frac{\text{wf } \Gamma}{\text{wf } \Gamma, \hat{\alpha}} \text{wf\_}\hat{\alpha} \\
 \frac{\text{wf } \Gamma \quad \Gamma \vdash A}{\text{wf } \Gamma, x : A} \text{wf\_}\text{of} \quad \frac{\text{wf } \Gamma \quad \Gamma \vdash \omega}{\text{wf } \Gamma \vdash \omega} \text{wf\_}\omega
 \end{array}$$

Figure 4.4: Extended Syntax and Well-Formedness for the Algorithmic System

which  $\hat{\alpha}$  has been added. The point is that the monotype  $\tau$ , represented by the corresponding existential variable, is always well-formed under the corresponding declarative context. In other words, the position of  $\hat{\alpha}$ 's reflects the well-formedness restriction of  $\tau$ 's.

**JUDGMENT CHAINS** Judgment chains  $\omega$ , or judgments for short, are the core components of our algorithmic type-checking. There are four kinds of judgments in our system: subtyping ( $A \leq B$ ), checking ( $e \Leftarrow A$ ), inference ( $e \Rightarrow_a \omega$ ) and application inference ( $A \bullet e \Rightarrow_a \omega$ ). Subtyping and checking are relatively simple, since their result is only success or failure. However both inference and application inference return a type that is used in subsequent judgments. We use a continuation-passing-style encoding to accomplish this. For example, the judgment chain  $e \Rightarrow_a (a \leq B)$  contains two judgments: first we want to infer the type of the expression  $e$ , and then check if that type is a subtype of  $B$ . The *unknown* type of  $e$  is represented by a type variable  $a$ , which is used as a placeholder in the second judgment to denote the type of  $e$ .

**WORKLIST JUDGMENTS** Our algorithm has a non-standard form. We combine traditional contexts and judgment(s) into a single sort, the *worklist*  $\Gamma$ . The worklist is an *ordered* collection of both variable bindings and judgments. The order captures the scope: only the objects that come after a variable's binding in the worklist can refer to it. For example,  $[\cdot, a, x : a \Vdash x \Leftarrow a]$  is a valid worklist, but  $[\cdot \Vdash \underline{x} \Leftarrow \underline{a}, x : \underline{a}, a]$  is not (the underlined symbols refer to out-of-scope variables).

**HOLE NOTATION** We use the syntax  $\Gamma[\Gamma_M]$  to denote the worklist  $\Gamma_L, \Gamma_M, \Gamma_R$ , where  $\Gamma$  is the worklist  $\Gamma_L, \bullet, \Gamma_R$  with a hole ( $\bullet$ ). Hole notations with the same name implicitly share the same structure  $\Gamma_L$  and  $\Gamma_R$ . A multi-hole notation splits the worklist into more parts. For example,  $\Gamma[\hat{\alpha}][\hat{\beta}]$  means  $\Gamma_1, \hat{\alpha}, \Gamma_2, \hat{\beta}, \Gamma_3$ .

#### 4.2.2 ALGORITHMIC SYSTEM

The algorithmic typing reduction rules, defined in Figure 4.5, have the form  $\Gamma \longrightarrow \Gamma'$ . The reduction process treats the worklist as a stack. In every step, it pops the first judgment from the worklist for processing and possibly pushes new judgments onto the worklist. The syntax  $\Gamma \longrightarrow^* \Gamma'$  denotes multiple reduction steps.

$$\frac{}{\Gamma \longrightarrow^* \Gamma} \longrightarrow^* \text{id} \qquad \frac{\Gamma \longrightarrow \Gamma_1 \quad \Gamma_1 \longrightarrow^* \Gamma'}{\Gamma \longrightarrow^* \Gamma'} \longrightarrow^* \text{step}$$

In the case that  $\Gamma \longrightarrow^* \cdot$  this corresponds to successful type checking.

$$\begin{array}{c}
 \boxed{\Gamma \longrightarrow \Gamma'} \qquad \qquad \qquad \Gamma \text{ reduces to } \Gamma'. \\
 \\
 \Gamma, a \longrightarrow_1 \Gamma \quad \Gamma, \hat{\alpha} \longrightarrow_2 \Gamma \quad \Gamma, x : A \longrightarrow_3 \Gamma \\
 \\
 \Gamma \Vdash 1 \leq 1 \longrightarrow_4 \Gamma \\
 \Gamma \Vdash a \leq a \longrightarrow_5 \Gamma \\
 \Gamma \Vdash \hat{\alpha} \leq \hat{\alpha} \longrightarrow_6 \Gamma \\
 \Gamma \Vdash A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2 \longrightarrow_7 \Gamma \Vdash A_2 \leq B_2 \Vdash B_1 \leq A_1 \\
 \Gamma \Vdash \forall a. A \leq B \longrightarrow_8 \Gamma, \hat{\alpha} \Vdash [\hat{\alpha}/a]A \leq B \quad \text{when } B \neq \forall a. B' \\
 \Gamma \Vdash A \leq \forall b. B \longrightarrow_9 \Gamma, b \Vdash A \leq B \\
 \\
 \Gamma[\hat{\alpha}] \Vdash \hat{\alpha} \leq A \rightarrow B \longrightarrow_{10} [\hat{\alpha}_1 \rightarrow \hat{\alpha}_2/\hat{\alpha}](\Gamma[\hat{\alpha}_1, \hat{\alpha}_2] \Vdash \hat{\alpha}_1 \rightarrow \hat{\alpha}_2 \leq A \rightarrow B) \\
 \qquad \qquad \qquad \text{when } \hat{\alpha} \notin FV(A) \cup FV(B) \\
 \Gamma[\hat{\alpha}] \Vdash A \rightarrow B \leq \hat{\alpha} \longrightarrow_{11} [\hat{\alpha}_1 \rightarrow \hat{\alpha}_2/\hat{\alpha}](\Gamma[\hat{\alpha}_1, \hat{\alpha}_2] \Vdash A \rightarrow B \leq \hat{\alpha}_1 \rightarrow \hat{\alpha}_2) \\
 \qquad \qquad \qquad \text{when } \hat{\alpha} \notin FV(A) \cup FV(B) \\
 \\
 \Gamma[\hat{\alpha}][\hat{\beta}] \Vdash \hat{\alpha} \leq \hat{\beta} \longrightarrow_{12} [\hat{\alpha}/\hat{\beta}](\Gamma[\hat{\alpha}][]) \\
 \Gamma[\hat{\alpha}][\hat{\beta}] \Vdash \hat{\beta} \leq \hat{\alpha} \longrightarrow_{13} [\hat{\alpha}/\hat{\beta}](\Gamma[\hat{\alpha}][]) \\
 \Gamma[a][\hat{\beta}] \Vdash a \leq \hat{\beta} \longrightarrow_{14} [a/\hat{\beta}](\Gamma[a][]) \\
 \Gamma[a][\hat{\beta}] \Vdash \hat{\beta} \leq a \longrightarrow_{15} [a/\hat{\beta}](\Gamma[a][]) \\
 \Gamma[\hat{\beta}] \Vdash 1 \leq \hat{\beta} \longrightarrow_{16} [1/\hat{\beta}](\Gamma[]) \\
 \Gamma[\hat{\beta}] \Vdash \hat{\beta} \leq 1 \longrightarrow_{17} [1/\hat{\beta}](\Gamma[]) \\
 \\
 \Gamma \Vdash e \Leftarrow B \longrightarrow_{18} \Gamma \Vdash e \Rightarrow_a a \leq B \quad \text{when } e \neq \lambda x. e' \text{ and } B \neq \forall a. B' \\
 \Gamma \Vdash e \Leftarrow \forall a. A \longrightarrow_{19} \Gamma, a \Vdash e \Leftarrow A \\
 \Gamma \Vdash \lambda x. e \Leftarrow A \rightarrow B \longrightarrow_{20} \Gamma, x : A \Vdash e \Leftarrow B \\
 \Gamma[\hat{\alpha}] \Vdash \lambda x. e \Leftarrow \hat{\alpha} \longrightarrow_{21} [\hat{\alpha}_1 \rightarrow \hat{\alpha}_2/\hat{\alpha}](\Gamma[\hat{\alpha}_1, \hat{\alpha}_2], x : \hat{\alpha}_1 \Vdash e \Leftarrow \hat{\alpha}_2) \\
 \\
 \Gamma \Vdash x \Rightarrow_a \omega \longrightarrow_{22} \Gamma \Vdash [A/a]\omega \quad \text{when } (x : A) \in \Gamma \\
 \Gamma \Vdash (e : A) \Rightarrow_a \omega \longrightarrow_{23} \Gamma \Vdash [A/a]\omega \Vdash e \Leftarrow A \\
 \Gamma \Vdash () \Rightarrow_a \omega \longrightarrow_{24} \Gamma \Vdash [1/a]\omega \\
 \Gamma \Vdash \lambda x. e \Rightarrow_a \omega \longrightarrow_{25} \Gamma, \hat{\alpha}, \hat{\beta} \Vdash [\hat{\alpha} \rightarrow \hat{\beta}/a]\omega, x : \hat{\alpha} \Vdash e \Leftarrow \hat{\beta} \\
 \Gamma \Vdash e_1 e_2 \Rightarrow_a \omega \longrightarrow_{26} \Gamma \Vdash e_1 \Rightarrow_b (b \bullet e_2 \Rightarrow_a \omega) \\
 \\
 \Gamma \Vdash \forall a. A \bullet e \Rightarrow_a \omega \longrightarrow_{27} \Gamma, \hat{\alpha} \Vdash [\hat{\alpha}/a]A \bullet e \Rightarrow_a \omega \\
 \Gamma \Vdash A \rightarrow C \bullet e \Rightarrow_a \omega \longrightarrow_{28} \Gamma \Vdash [C/a]\omega \Vdash e \Leftarrow A \\
 \Gamma[\hat{\alpha}] \Vdash \hat{\alpha} \bullet e \Rightarrow_a \omega \longrightarrow_{29} [\hat{\alpha}_1 \rightarrow \hat{\alpha}_2/\hat{\alpha}](\Gamma[\hat{\alpha}_1, \hat{\alpha}_2] \Vdash \hat{\alpha}_1 \rightarrow \hat{\alpha}_2 \bullet e \Rightarrow_a \omega)
 \end{array}$$

Figure 4.5: Algorithmic Typing

Please note that when a new variable is introduced in the right-hand side worklist  $\Gamma'$ , we implicitly pick a fresh one, since the right-hand side can be seen as the premise of the reduction.

Rules 1-3 pop variable declarations that are essentially garbage. Those variables are out of scope for the remaining judgments in the worklist. All other rules concern a judgment at the front of the worklist. Logically we can discern 6 groups of rules.

**1. Algorithmic subtyping** We have six subtyping rules (Rules 4-9) that are similar to their declarative counterparts. For instance, Rule 7 consumes a subtyping judgment and pushes two back to the worklist. Rule 8 differs from declarative Rule  $\leq \forall L$  by introducing an existential variable  $\hat{\alpha}$  instead of guessing the monotype  $\tau$  instantiation. Each existential variable is later solved to a monotype  $\tau$  with the same context, so the final solution  $\tau$  of  $\hat{\alpha}$  should be well-formed under  $\Gamma$ .

**WORKLIST VARIABLE SCOPING** Rules 8 and 9 involve variable declarations and demonstrate how our worklist treats variable scopes. Rule 8 introduces an existential variable  $\hat{\alpha}$  that is only visible to the judgment  $[\hat{\alpha}/a]A \leq B$ . Reduction continues until all the subtyping judgments in front of  $\hat{\alpha}$  are satisfied. Finally we can safely remove  $\hat{\alpha}$  since no occurrence of  $\hat{\alpha}$  could have leaked into the left part of the worklist. Moreover, the algorithm garbage-collects the  $\hat{\alpha}$  variable at the right time: it leaves the environment immediately after being unreferenced completely for sure.

**EXAMPLE** Consider the derivation of the subtyping judgment  $(1 \rightarrow 1) \rightarrow 1 \leq (\forall a. 1 \rightarrow 1) \rightarrow 1$ :

$$\begin{aligned}
 & \cdot \vdash (1 \rightarrow 1) \rightarrow 1 \leq (\forall a. 1 \rightarrow 1) \rightarrow 1 \\
 \longrightarrow_7 & \cdot \Vdash 1 \leq 1 \Vdash \forall a. 1 \rightarrow 1 \leq 1 \rightarrow 1 \\
 \longrightarrow_8 & \cdot \Vdash 1 \leq 1, \hat{\alpha} \Vdash 1 \rightarrow 1 \leq 1 \rightarrow 1 \\
 \longrightarrow_7 & \cdot \Vdash 1 \leq 1, \hat{\alpha} \Vdash 1 \leq 1 \Vdash 1 \leq 1 \\
 \longrightarrow_4 & \cdot \Vdash 1 \leq 1, \hat{\alpha} \Vdash 1 \leq 1 \\
 \longrightarrow_4 & \cdot \Vdash 1 \leq 1, \hat{\alpha} \\
 \longrightarrow_2 & \cdot \Vdash 1 \leq 1 \\
 \longrightarrow_4 & \cdot
 \end{aligned}$$

First, the subtyping of two function types is split into two judgments by Rule 7: a covariant subtyping on the return type and a contravariant subtyping on the argument type. Then we

apply Rule 8 to reduce the  $\forall$  quantifier on the left side. The rule introduces an existential variable  $\hat{\alpha}$  to the context (even though the type  $\forall a. 1 \rightarrow 1$  does not actually refer to the quantified type variable  $a$ ). In the following 3 steps we satisfy the judgment  $1 \rightarrow 1 \leq 1 \rightarrow 1$  by Rules 7, 4, and 4 (again).

Now the existential variable  $\hat{\alpha}$ , introduced before but still unsolved, is at the top of the worklist and Rule 2 garbage-collects it. The process is carefully designed within the algorithmic rules: when  $\hat{\alpha}$  is introduced earlier by Rule 8, we foresee the recycling of  $\hat{\alpha}$  after all the judgments (potentially) requiring  $\hat{\alpha}$  have been processed. Finally Rule 4 reduces one of the base cases and finishes the subtyping derivation.

**2. Existential decomposition.** Rules 10 and 11 are algorithmic versions of Rule  $\leq \rightarrow$ ; they both partially instantiate  $\hat{\alpha}$  to function types. The domain  $\hat{\alpha}_1$  and range  $\hat{\alpha}_2$  of the new function type are not determined: they are fresh existential variables with the same scope as  $\hat{\alpha}$ . We replace  $\hat{\alpha}$  in the worklist with  $\hat{\alpha}_1, \hat{\alpha}_2$ . To propagate the instantiation to the rest of the worklist and maintain well-formedness, every reference to  $\hat{\alpha}$  is replaced by  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$ . The *occurs-check* condition prevents divergence as usual. For example, without it  $\hat{\alpha} \leq 1 \rightarrow \hat{\alpha}$  would diverge.

**3. Solving existentials** Rules 12-17 are base cases where an existential variable is solved. They all remove an existential variable and substitute the variable with its solution in the remaining worklist. Importantly the rules respect the scope of existential variables. For example, Rule 12 states that an existential variable  $\hat{\alpha}$  can be solved with another existential variable  $\hat{\beta}$  only if  $\hat{\beta}$  occurs after  $\hat{\alpha}$ .

One may notice that the subtyping relation for simple types is just equivalence, which is true according to the declarative system. The DK's system works in a similar way.

**4. Checking judgments.** Rules 18-21 deal with checking judgments. Rule 18 is similar to  $\text{Dec1Sub}$ , but rewritten in the continuation-passing-style. The side conditions  $e \neq \lambda x. e'$  and  $B \neq \forall a. B'$  prevent overlap with Rules 19, 20 and 21; this is further discussed at the end of this section. Rules 19 and 20 adapt their declarative counterparts to the worklist style. Rule 21 is a special case of  $\text{Dec1} \rightarrow \text{I}$ , dealing with the case when the input type is an existential variable, representing a monotype *function* as in the declarative system (it must be a function type, since the expression  $\lambda x. e$  is a function). The same instantiation technique as in Rules 10 and 11 applies. The declarative checking rule  $\text{Dec11I}$  does not have a direct counterpart in the algorithm, because Rules 18 and 24 can be combined to give the same result.

**RULE 21 DESIGN CHOICE** The addition of Rule 21 is a design choice we have made to simplify the side condition of Rule 18 (which avoids overlap). It also streamlines the algorithm and the metatheory as we now treat all cases where we can see that an existential variable should be instantiated to a function type (i.e., Rules 10, 11, 21 and 29) uniformly.

The alternative would have been to omit Rule 21 and drop the condition on  $e$  in Rule 18. The modified Rule 18 would then handle  $\Gamma \Vdash \lambda x. e \Leftarrow \hat{\alpha}$  and yield  $\Gamma \Vdash \lambda x. e \Rightarrow_a a \leq \hat{\alpha}$ , which would be further processed by Rule 25 to  $\Gamma, \hat{\beta}_1, \hat{\beta}_2 \Vdash \hat{\beta}_1 \rightarrow \hat{\beta}_2 \leq \hat{\alpha}, x : \hat{\beta}_1 \Vdash e \Leftarrow \hat{\beta}_2$ . As a subtyping constraint between monotypes is simply equality,  $\hat{\beta}_1 \rightarrow \hat{\beta}_2 \leq \hat{\alpha}$  must end up equating  $\hat{\beta}_1 \rightarrow \hat{\beta}_2$  with  $\hat{\alpha}$  and thus have the same effect as Rule 21, but in a more roundabout fashion.

In comparison, DK's algorithmic subsumption rule has no restriction on the expression  $e$ , and they do not have a rule that explicitly handles the case  $\lambda x. e \Leftarrow \hat{\alpha}$ . Therefore the only way to check a lambda function against an existential variable is by applying the subsumption rule, which further breaks into type inference of a lambda function and a subtyping judgment.

**5. Inference judgments.** Inference judgments behave differently compared with subtyping and checking judgments: they *return* a type instead of only accepting or rejecting. For the algorithmic system, where guesses are involved, it may happen that the output type of an inference judgment refers to new existential variables, such as Rule 25. In comparison to Rule 8 and 9, where new variables are only referred to by the sub-derivation, Rule 25 introduces variables  $\hat{\alpha}, \hat{\beta}$  that affect the remaining judgment chain. This rule is carefully designed so that the output variables are bound by earlier declarations, thus the well-formedness of the worklist is preserved, and the garbage will be collected at the correct time. By making use of the continuation-passing-style judgment chain, inner judgments always share the context with their parent judgment.

Rules 22-26 deal with type inference judgments, written in continuation-passing-style. When an inference judgment succeeds with type  $A$ , the algorithm continues to work on the inner-chain  $\omega$  by assigning  $A$  to its placeholder variable  $a$ . Rule 23 infers an annotated expression by changing into checking mode, therefore another judgment chain is created. Rule 24 is a base case, where the unit type 1 is inferred and thus passed to its child judgment chain. Rule 26 infers the type of an application by firstly inferring the type of the function  $e_1$ , and then leaving the rest work to an application inference judgment, which passes  $a$ , representing the return type of the application, to the remainder of the judgment chain  $\omega$ .

Rule 25 infers the type of a lambda expression by introducing  $\hat{\alpha}, \hat{\beta}$  as the input and output types of the function, respectively. After checking the body  $e$  under the assumption  $x : \hat{\alpha}$ , the return type might reflect more information than simply  $\hat{\alpha} \rightarrow \hat{\beta}$  through propagation when



existential variables are solved or partially solved. The variable scopes are maintained during the process: the assumption of argument type  $(x : \hat{\alpha})$  is recycled after checking against the function body; the existential variables used by the function type  $(\hat{\alpha}, \hat{\beta})$  are only visible in the remaining chain  $[\hat{\alpha} \rightarrow \hat{\beta}/a]\omega$ . The recycling process of Rule 25 differs from DK's corresponding rule significantly, and we further discuss the differences in Section 4.4.1.

**6. Application inference judgments** Finally, Rules 27-29 deal with application inference judgments. Rules 27 and 28 behave similarly to declarative rules  $\text{Decl}\forall\text{App}$  and  $\text{Decl} \rightarrow \text{App}$ . Rule 29 instantiates  $\hat{\alpha}$  to the function type  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$ , just like Rules 10, 11 and 21.

**EXAMPLE** Figure 4.6 shows a sample derivation of the algorithm. It checks the application  $(\lambda x. x) ()$  against the unit type. According to the algorithm, we apply Rule 18 (subsumption), changing to inference mode. Type inference of the application breaks into two steps by Rule 26: first we infer the type of the function, and then the application inference judgment helps to determine the return type. In the following 5 steps the type of the identity function,  $\lambda x. x$ , is inferred to be  $\hat{\alpha} \rightarrow \hat{\alpha}$ : checking the body of the lambda function (Rule 25), switching from check mode to inference mode (Rule 18), inferring the type of a term variable (Rule 22), solving a subtyping judgment between existential variables (Rule 12) and garbage collecting the term variable  $x$  (Rule 3).

After that, Rule 28 changes the application inference judgment to a check of the argument against the input type  $\hat{\alpha}$  and returns the output type  $\hat{\alpha}$ . Checking  $()$  against the existential variable  $\hat{\alpha}$  solves  $\hat{\alpha}$  to the unit type 1 through Rules 18, 24 and 16. Immediately after  $\hat{\alpha}$  is solved, the algorithm replaces every occurrence of  $\hat{\alpha}$  with 1. Therefore the worklist remains  $1 \leq 1$ , which is finished off by Rule 4. Finally, the empty worklist indicates the success of the whole derivation.

In summary, our type checking algorithm accepts  $(\lambda x. x) () \Leftarrow 1$ .

**NON-OVERLAPPING AND DETERMINISTIC REDUCTION** An important feature of our algorithmic rules is that they are directly implementable. Indeed, although written in the form of reduction rules, they do not overlap and are thus deterministic.

Consider in particular Rules 8 and 9, which correspond to the declarative rules  $\leq\forall\text{L}$  and  $\leq\forall\text{R}$ . While those declarative rules both match the goal  $\forall a. A \leq \forall b. B$ , we have eliminated this overlap in the algorithm by restricting Rule 8 ( $B \neq \forall a. B'$ ) and thus always applying Rule 9 to  $\forall a. A \leq \forall b. B$ .

Similarly, the declarative rule  $\text{DeclSub}$  overlaps highly with the other checking rules. Its algorithmic counterpart is Rule 18. Yet, we have avoided the overlap with other algorithmic checking rules by adding side-conditions to Rule 18, namely  $e \neq \lambda x. e'$  and  $B \neq \forall a. B'$ .

$$\begin{aligned}
 & \cdot \Vdash (\lambda x. x) () \Leftarrow 1 \\
 \longrightarrow_{18} & \cdot \Vdash (\lambda x. x) () \Rightarrow_a a \leq 1 \\
 \longrightarrow_{26} & \cdot \Vdash (\lambda x. x) \Rightarrow_b (b \bullet ()) \Rightarrow_a a \leq 1 \\
 \longrightarrow_{25} & \cdot, \hat{\alpha}, \hat{\beta} \Vdash \hat{\alpha} \rightarrow \hat{\beta} \bullet () \Rightarrow_a a \leq 1, x : \hat{\alpha} \Vdash x \Leftarrow \hat{\beta} \\
 \longrightarrow_{18} & \cdot, \hat{\alpha}, \hat{\beta} \Vdash \hat{\alpha} \rightarrow \hat{\beta} \bullet () \Rightarrow_a a \leq 1, x : \hat{\alpha} \Vdash x \Rightarrow_b b \leq \hat{\beta} \\
 \longrightarrow_{22} & \cdot, \hat{\alpha}, \hat{\beta} \Vdash \hat{\alpha} \rightarrow \hat{\beta} \bullet () \Rightarrow_a a \leq 1, x : \hat{\alpha} \Vdash \hat{\alpha} \leq \hat{\beta} \\
 \longrightarrow_{12} & \cdot, \hat{\alpha} \Vdash \hat{\alpha} \rightarrow \hat{\alpha} \bullet () \Rightarrow_a a \leq 1, x : \hat{\alpha} \\
 \longrightarrow_3 & \cdot, \hat{\alpha} \Vdash \hat{\alpha} \rightarrow \hat{\alpha} \bullet () \Rightarrow_a a \leq 1 \\
 \longrightarrow_{28} & \cdot, \hat{\alpha} \Vdash \hat{\alpha} \leq 1 \Vdash () \Leftarrow \hat{\alpha} \\
 \longrightarrow_{18} & \cdot, \hat{\alpha} \Vdash \hat{\alpha} \leq 1 \Vdash () \Rightarrow_a a \leq \hat{\alpha} \\
 \longrightarrow_{24} & \cdot, \hat{\alpha} \Vdash \hat{\alpha} \leq 1 \Vdash 1 \leq \hat{\alpha} \\
 \longrightarrow_{16} & \cdot \Vdash 1 \leq 1 \\
 \longrightarrow_4 & \cdot
 \end{aligned}$$

Figure 4.6: A Sample Derivation for Algorithmic Typing

These restrictions have not been imposed arbitrarily: we formally prove that the restricted algorithm is still complete. In Section 4.3.2 we discuss the relevant metatheory, with the help of a non-overlapping version of the declarative system.

### 4.3 METATHEORY

This section presents the metatheory of the algorithmic system presented in the previous section. We show that three main results hold: *soundness*, *completeness* and *decidability*. These three results have been mechanically formalized and proved in the Abella theorem prover [Gacek 2008].

#### 4.3.1 DECLARATIVE WORKLIST AND TRANSFER

To aid formalizing the correspondence between the declarative and algorithmic systems, we introduce the notion of a declarative worklist  $\Omega$ , defined in Figure 4.7. A declarative worklist  $\Omega$  has the same structure as an algorithmic worklist  $\Gamma$ , but does not contain any existential variables  $\hat{\alpha}$ .

**WORKLIST INSTANTIATION.** The relation  $\Gamma \rightsquigarrow \Omega$  expresses that the algorithmic worklist  $\Gamma$  can be instantiated to the declarative worklist  $\Omega$ , by appropriately instantiating all existential

Declarative worklist  $\Omega ::= \cdot \mid \Omega, a \mid \Omega, x : A \mid \Omega \Vdash \omega$

$$\boxed{\Gamma \rightsquigarrow \Omega} \quad \Gamma \text{ instantiates to } \Omega.$$

$$\frac{}{\Omega \rightsquigarrow \Omega} \rightsquigarrow \Omega \quad \frac{\Omega \vdash \tau \quad \Omega, [\tau/\hat{\alpha}]\Gamma \rightsquigarrow \Omega}{\Omega, \hat{\alpha}, \Gamma \rightsquigarrow \Omega} \rightsquigarrow \hat{\alpha}$$

Figure 4.7: Declarative Worklists and Instantiation

variables  $\hat{\alpha}$  in  $\Gamma$  with well-scoped monotypes  $\tau$ . The rules of this instantiation relation are shown in Figure 4.7 too. Rule  $\rightsquigarrow \hat{\alpha}$  replaces the first existential variable with a well-scoped monotype and repeats the process on the resulting worklist until no existential variable remains and thus the algorithmic worklist has become a declarative one. In order to maintain well-scopedness, the substitution is applied to all the judgments and term variable bindings in the scope of  $\hat{\alpha}$ .

Observe that the instantiation  $\Gamma \rightsquigarrow \Omega$  is not deterministic. From left to right, there are infinitely many possibilities to instantiate an existential variable and thus infinitely many declarative worklists that one can get from an algorithmic one. In the other direction, any valid monotype in  $\Omega$  can be abstracted to an existential variable in  $\Gamma$ . Thus different  $\Gamma$ 's can be instantiated to the same  $\Omega$ .

Lemmas 4.3 and 4.4 generalize Rule  $\rightsquigarrow \hat{\alpha}$  from substituting the first existential variable to substituting any existential variable.

**Lemma 4.3 (Insert).** *If  $\Gamma_L, [\tau/\hat{\alpha}]\Gamma_R \rightsquigarrow \Omega$  and  $\Gamma_L \vdash \tau$ , then  $\Gamma_L, \hat{\alpha}, \Gamma_R \rightsquigarrow \Omega$ .*

**Lemma 4.4 (Extract).** *If  $\Gamma_L, \hat{\alpha}, \Gamma_R \rightsquigarrow \Omega$ , then there exists  $\tau$  s.t.  $\Gamma_L \vdash \tau$  and  $\Gamma_L, [\tau/\hat{\alpha}]\Gamma_R \rightsquigarrow \Omega$ .*

**DECLARATIVE TRANSFER.** Figure 4.8 defines a relation  $\Omega \longrightarrow \Omega'$ , which transfers all judgments in the declarative worklists to the declarative type system. This relation checks that every judgment entry in the worklist holds using a corresponding conventional declarative judgment. The typing contexts of declarative judgments are recovered using an auxiliary erasure function  $\|\Omega\|$ <sup>2</sup>. The erasure function simply drops all judgment entries from the worklist, keeping only variable and type variable declarations.

<sup>2</sup>In the proof script we do not use the erasure function, for the declarative system and well-formedness judgments automatically fit the non-erased declarative worklist just as declarative contexts.

$\boxed{\ \Omega\ }$	Judgment erasure.
$\ \cdot\  = \cdot$	
$\ \Omega, a\  = \ \Omega\ , a$	
$\ \Omega, x : A\  = \ \Omega\ , x : A$	
$\ \Omega \vdash \omega\  = \ \Omega\ $	
$\boxed{\Omega \longrightarrow \Omega'}$	Declarative transfer.
$\Omega, a \longrightarrow \Omega$	
$\Omega, x : A \longrightarrow \Omega$	
$\Omega \vdash A \leq B \longrightarrow \Omega$	when $\ \Omega\  \vdash A \leq B$
$\Omega \vdash e \Leftarrow A \longrightarrow \Omega$	when $\ \Omega\  \vdash e \Leftarrow A$
$\Omega \vdash e \Rightarrow_a \omega \longrightarrow \Omega \vdash [A/a]\omega$	when $\ \Omega\  \vdash e \Rightarrow A$
$\Omega \vdash A \bullet e \Rightarrow_a \omega \longrightarrow \Omega \vdash [C/a]\omega$	when $\ \Omega\  \vdash A \bullet e \Rightarrow C$

Figure 4.8: Declarative Transfer

#### 4.3.2 NON-OVERLAPPING DECLARATIVE SYSTEM

DK's declarative system, shown in Figures 2.8 and 2.9, has a few overlapping rules. In contrast, our algorithm has removed all overlap; at most one rule applies in any given situation. This discrepancy makes it more difficult to relate the two systems.

To simplify matters, we introduce an intermediate system that is still declarative in nature, but has no overlap. This intermediate system differs only in a few rules from DK's declarative system:

1. Restrict the shape of  $B$  in the rule  $\forall L$  subtyping rule:

$$\frac{B \neq \forall b. B' \quad \Psi \vdash \tau \quad \Psi \vdash [\tau/a]A \leq B}{\Psi \vdash \forall a. A \leq B} \forall L'$$

2. Drop the redundant rule  $\text{Decl1I}$ , which can be easily derived by a combination of  $\text{DeclSub}$ ,  $\text{Decl1I} \Rightarrow$  and  $\leq \text{Unit}$ :

$$\frac{\frac{}{\Psi \vdash () \Rightarrow 1} \text{Decl1I} \Rightarrow \quad \frac{}{\Psi \vdash 1 \leq 1} \leq \text{Unit}}{\Psi \vdash () \Leftarrow 1} \text{DeclSub}$$

3. Restrict the shapes of  $e$  and  $A$  in the subsumption rule  $\text{DeclSub}$ :

$$\frac{e \neq \lambda x. e' \quad A \neq \forall a. A' \quad \Psi \vdash e \Rightarrow A \quad \Psi \vdash A \leq B}{\Psi \vdash e \Leftarrow B} \text{DeclSub}'$$

The resulting declarative system has no overlapping rules and more closely resembles the algorithmic system, which contains constraints of the same shape.

We have proven soundness and completeness of the non-overlapping declarative system with respect to the overlapping one to establish their equivalence. Thus the restrictions do not change the expressive power of the system. Modification (2) is relatively easy to justify, with the derivation given above: the rule is redundant and can be replaced by a combination of three other rules. Modifications (1) and (3) require inversion lemmas for the rules that overlap. Firstly, Rule  $\forall\text{L}$  overlaps with Rule  $\forall\text{R}$  for the judgment  $\Psi \vdash \forall a. A \leq \forall b. B$ . The following inversion lemma for Rule  $\forall\text{R}$  resolves the overlap:

**Lemma 4.5** (Invertibility of  $\forall\text{R}$ ). *If  $\Psi \vdash A \leq \forall a. B$  then  $\Psi, a \vdash A \leq B$ .*

The lemma implies that preferring Rule  $\forall\text{R}$  does not affect the derivability of the judgment. Therefore the restriction  $B \neq \forall b. B'$  in  $\forall\text{L}'$  is valid.

Secondly, Rule  $\text{DeclSub}$  overlaps with both  $\text{Decl}\forall\text{I}$  and  $\text{Decl}\rightarrow\text{I}$ . We have proven two inversion lemmas for these overlaps:

**Lemma 4.6** (Invertibility of  $\text{Decl}\forall\text{I}$ ). *If  $\Psi \vdash e \Leftarrow \forall a. A$  then  $\Psi, a \vdash e \Leftarrow A$ .*

**Lemma 4.7** (Invertibility of  $\text{Decl}\rightarrow\text{I}$ ). *If  $\Psi \vdash \lambda x. e \Leftarrow A \rightarrow B$  then  $\Psi, x : A \vdash e \Leftarrow B$ .*

These lemmas express that applying the more specific rules, rather than the more general rule  $\text{DeclSub}$ , does not reduce the expressive power. This is required by the completeness lemma, as the algorithmic system explicitly prioritizes the more specific rules.

The proofs of the above two lemmas rely on an important property of the declarative system, the *subsumption lemma*. To be able to formulate this lemma, Figure 4.9 introduces the *context subtyping relation*  $\Psi \leq \Psi'$ . Context  $\Psi$  subsumes context  $\Psi'$  if they bind the same variables in the same order, but the types  $A$  of the term variables  $x$  in the former are subtypes of types  $A'$  assigned to those term variables in the latter. Now we can state the subsumption lemma:

**Lemma 4.8** (Subsumption). *Given  $\Psi' \leq \Psi$ :*

1. *If  $\Psi \vdash e \Leftarrow A$  and  $\Psi \vdash A \leq A'$  then  $\Psi' \vdash e \Leftarrow A'$ ;*
2. *If  $\Psi \vdash e \Rightarrow B$  then there exists  $B'$  s.t.  $\Psi \vdash B' \leq B$  and  $\Psi' \vdash e \Rightarrow B'$ ;*

$$\boxed{\Psi' \leq \Psi}$$

$$\begin{array}{c}
 \frac{}{\cdot \leq \cdot} \text{CtxSubEmpty} \qquad \frac{\Psi' \leq \Psi}{\Psi', a \leq \Psi, a} \text{CtxSubTyVar} \\
 \frac{\Psi' \leq \Psi \quad \Psi \vdash A' \leq A}{\Psi', x : A' \leq \Psi, x : A} \text{CtxSubTmVar}
 \end{array}$$

Figure 4.9: Context Subtyping

3. If  $\Psi \vdash A \bullet e \Rightarrow C$  and  $\Psi \vdash A' \leq A$ , then there exists  $C'$  s.t.  $\Psi \vdash C' \leq C$  and  $\Psi' \vdash A' \bullet e \Rightarrow C'$ .

This lemma expresses that any derivation in a context  $\Psi$  has a corresponding derivation in any context  $\Psi'$  that it subsumes.

While being written in a clear format and providing enough details, some proof is not fully accepted by the proof assistant. Specifically for the subsumption lemma, we have tried to follow DK's manual proof, but we discovered several problems in their reasoning that we have been unable to address. Fortunately we have found a different way to prove the lemma. A description of the problem and our fix are discussed in-depth as follows.

In the appendix of DK's paper [Dunfield and Krishnaswami 2013], the first two applications of induction hypotheses on page 22 are not accepted. Either of them seems to use a slightly different “induction hypothesis” than the true one. In fact, we cannot think of simple ways to fix the problem, since the induction hypothesis seems not strong enough for these two cases.

To fix the proof, we propose a new induction scheme by making use of our worklist measures. Recall that  $|e|_e$  measures term size; and the judgment measure counts checking as 2, inference as 1 and application inference as 3; and  $|A|_\forall$  counts the number of  $\forall$ 's in a type.

*Proof.* The proof is by a nested mutual induction on the lexicographical order of the measures

$$\langle |e|_e, | \cdot |_{\Leftrightarrow}, |A|_\forall + |A'|_\forall \rangle,$$

where the second measure simply give an natural number for each type of judgment: (1) checking judgments ( $\Psi \vdash e \Leftarrow A$ ) count as 2; (2) inference judgments ( $\Psi \vdash e \Rightarrow A$ ) count as 1; and (3) application inference judgments ( $\Psi \vdash e \bullet A \Rightarrow C$ ) count as 3; and the third measure does not apply to (case 2) since no  $A$  is used. Compared with DK's, our third measure that counts the degree of polymorphism fixes problems that occurred in DK's proof: in both places, our new induction hypothesis is more generalized.

All but two cases can be finished easily following the declarative typing derivation, and the proof shares a similar structure to DK's. One tricky case related to Rule  $\text{Decl}\forall\text{I}$  indicates that  $A$  has the shape  $\forall a. A_0$ , thus the subtyping relation derives from either  $\leq\forall\text{L}$  or  $\leq\forall\text{R}$ . For each of the case, the third measure  $|A|_{\forall} + |A'|_{\forall}$  decreases (the  $\leq\forall\text{L}$  case requires a type substitution lemma obtaining  $\Psi \vdash e \Leftarrow [\tau/a]A_0$  from the typing derivation).

Another tricky case is  $\text{Decl}\rightarrow\text{I}$ . When the subtyping relation is derived from  $\leq\rightarrow$ , a simple application of induction hypothesis finishes the proof. When the subtyping relation is derived from  $\leq\forall\text{R}$ ,  $|A'|_{\forall}$  decreases, and thus the induction hypothesis finishes this case.  $\square$

In short, we found a specific problem when trying to prove the subsumption lemma following DK's manual proof, yet we addressed the problem by using a slightly different induction scheme.

**THREE-WAY SOUNDNESS AND COMPLETENESS THEOREMS** We now have three systems that can be related: DK's overlapping declarative system, our non-overlapping declarative system, and our algorithmic system. We have already established the first relation, that the two declarative systems are equivalent. In what follows, we will establish the soundness of our algorithm directly against the original overlapping declarative system. However, we have found that showing completeness of the algorithm is easier against the non-overlapping declarative system. Of course, as a corollary, it follows that our algorithm is also complete with respect to DK's declarative system.

### 4.3.3 SOUNDNESS

Our algorithm is sound with respect to DK's declarative system. For any worklist  $\Gamma$  that reduces successfully, there is a valid instantiation  $\Omega$  that transfers all judgments to the declarative system.

**Theorem 4.9** (Soundness). *If  $\text{wf } \Gamma$  and  $\Gamma \longrightarrow^*$ , then there exists  $\Omega$  s.t.  $\Gamma \rightsquigarrow \Omega$  and  $\Omega \longrightarrow^*$ .*

The proof proceeds by induction on the derivation of  $\Gamma \longrightarrow^*$ . Interesting cases are those involving existential variable instantiations, including Rules 10, 11, 21 and 29. Applications of Lemmas 4.3 and 4.4 help analyse the full instantiation of those existential variables. For example, when  $\hat{\alpha}$  is solved to  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$  in the algorithm, applying the Extract lemma gives two instantiations  $\hat{\alpha}_1 = \sigma$  and  $\hat{\alpha}_2 = \tau$ . It follows that  $\hat{\alpha} = \sigma \rightarrow \tau$ , which enables the induction hypothesis and finishes the corresponding case. Some immediate corollaries which show the soundness for specific judgment forms are:

**Corollary 4.10** (Soundness, single judgment form). *Given  $\text{wf } \Gamma$ :*

1. If  $\Gamma \Vdash A \leq B \longrightarrow^*$  .  
 then there exist  $A', B', \Omega$  s.t.  $\Gamma \Vdash A \leq B \rightsquigarrow \Omega \Vdash A' \leq B'$  and  $\|\Omega\| \vdash A' \leq B'$ ;
2. If  $\Gamma \Vdash e \Leftarrow A \longrightarrow^*$  .  
 then there exist  $A', \Omega$  s.t.  $\Gamma \Vdash e \Leftarrow A \rightsquigarrow \Omega \Vdash e \Leftarrow A'$  and  $\|\Omega\| \vdash e \Leftarrow A'$ ;
3. If  $\Gamma \Vdash e \Rightarrow_a \omega \longrightarrow^*$  for any  $\omega$   
 then there exists  $\Omega, \omega', A$  s.t.  $\Gamma \rightsquigarrow \Omega$  and  $\|\Omega\| \vdash e \Rightarrow A$ ;
4. If  $\Gamma \Vdash A \bullet e \Rightarrow_a \omega \longrightarrow^*$  for any  $\omega$   
 then there exists  $\Omega, \omega', A', C$  s.t.  $\Gamma \Vdash A \bullet e \Rightarrow_a \omega \rightsquigarrow \Omega \Vdash A' \bullet e \Rightarrow_a \omega'$  and  $\|\Omega\| \vdash A' \bullet e \Rightarrow C$ .

#### 4.3.4 COMPLETENESS

The completeness of our algorithm means that any derivation in the declarative system has an algorithmic counterpart. We explicitly relate between an algorithmic context  $\Gamma$  and a declarative context  $\Omega$  to avoid potential confusion.

**Theorem 4.11** (Completeness). *If wf  $\Gamma$  and  $\Gamma \rightsquigarrow \Omega$  and  $\Omega \longrightarrow^*$  , then  $\Gamma \longrightarrow^*$  .*

We prove completeness by induction on the derivation of  $\Omega \longrightarrow^*$  . and use the non-overlapping declarative system. Since the declarative worklist is reduced judgment by judgment (shown in Figure 4.8), the induction always analyses the first judgment by a small step. As the algorithmic system introduces existential variables, a declarative rule may correspond to multiple algorithmic rules, and thus we analyse each of the possible cases.

Most cases are relatively easy to prove. The Insert and Extract lemmas are applied when the algorithm uses existential variables, but transferred to a monotype for the declarative system, such as Rules 6, 8, 10, 11, 12-17, 21, 25, 27 and 29.

Algorithmic Rules 10 and 11 require special treatment. When the induction reaches the  $\leq \rightarrow$  case, the first judgment is of shape  $A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2$ . One of the corresponding algorithmic judgments is  $\hat{\alpha} \leq A \rightarrow B$ . However, the case analysis does not imply that  $\hat{\alpha}$  is fresh in  $A$  and  $B$ , therefore Rule 10 cannot be applied and the proof gets stuck. The following lemma helps us out in those cases: the success in declarative subtyping indicates the freshness of  $\hat{\alpha}$  in  $A$  and  $B$  in its corresponding algorithmic judgment. In other words, the declarative system does not accept infinite types. A symmetric lemma holds for  $A \rightarrow B \leq \hat{\alpha}$ .

**Lemma 4.12** (Prune Transfer for Instantiation). *If  $(\Gamma \Vdash \hat{\alpha} \leq A \rightarrow B) \rightsquigarrow (\Omega \Vdash C \leq A_1 \rightarrow B_1)$  and  $\|\Omega\| \vdash C \leq A_1 \rightarrow B_1$ , then  $\hat{\alpha} \notin FV(A) \cup FV(B)$ .*

The following corollary is derived immediately from Theorem 4.11.



**Corollary 4.13** (Completeness, single judgment form). *Given wf  $\Gamma$  containing no judgments:*

1. *If  $\Omega \vdash A' \leq B'$  and  $\Gamma \Vdash A \leq B \rightsquigarrow \Omega \vdash A' \leq B'$   
then  $\Gamma \Vdash A \leq B \longrightarrow^*$  .;*
2. *If  $\Omega \vdash e \Leftarrow A'$  and  $\Gamma \Vdash e \Leftarrow A \rightsquigarrow \Omega \vdash e \Leftarrow A'$   
then  $\Gamma \Vdash e \Leftarrow A \longrightarrow^*$  .;*
3. *If  $\Omega \vdash e \Rightarrow A$  and  $\Gamma \Vdash e \Rightarrow_a 1 \leq 1 \rightsquigarrow \Omega \vdash e \Rightarrow_a 1 \leq 1$   
then  $\Gamma \Vdash e \Rightarrow_a 1 \leq 1 \longrightarrow^*$  .;*
4. *If  $\Omega \vdash A' \bullet e \Rightarrow C$  and  $\Gamma \Vdash A \bullet e \Rightarrow_a 1 \leq 1 \rightsquigarrow \Omega \vdash A' \bullet e \Rightarrow_a 1 \leq 1$   
then  $\Gamma \Vdash A \bullet e \Rightarrow_a 1 \leq 1 \longrightarrow^*$  .*

#### 4.3.5 DECIDABILITY

Finally, we show that our algorithm is decidable:

**Theorem 4.14** (Decidability). *Given wf  $\Gamma$ , it is decidable whether  $\Gamma \longrightarrow^* \cdot$  or not.*

Our decidability proof is based on a lexicographic group of induction measures  $\langle |\Gamma|_e, |\Gamma|_{\Leftarrow}, |\Gamma|_{\forall}, |\Gamma|_{\hat{\alpha}}, |\Gamma|_{\rightarrow} + |\Gamma| \rangle$  on the worklist  $\Gamma$ . Formal definitions of these measures can be found in Figure 4.10. The first two,  $|\cdot|_e$  and  $|\cdot|_{\Leftarrow}$ , measure the total size of terms and the total difficulty of judgments, respectively. In the latter, check judgments count for 2, inference judgments for 1 and function inference judgments for 3. Another two measures,  $|\cdot|_{\forall}$  and  $|\cdot|_{\rightarrow}$ , count the total number of universal quantifiers and function types, respectively. Finally,  $|\cdot|_{\hat{\alpha}}$  counts the number of existential variables in the worklist, and  $|\cdot|$  is simply the length of the worklist.

It is not difficult to see that all but two algorithmic reduction rules decrease the group of measures. (The result of Rule 29 could be directly reduced by Rule 28, which decreases the measures.) The two exceptions are Rules 10 and 11. Both rules increase the number of existential variables without decreasing the number of universal quantifiers. However, they are both immediately followed by Rule 7, which breaks the subtyping problem into two smaller problems of the form  $\hat{\alpha} \leq A$  and  $A \leq \hat{\alpha}$  which we call *instantiation judgments*.

We now show that entirely reducing these smaller problems leaves the worklist in a state with an overall smaller measure. Our starting point is a worklist  $\Gamma, \Gamma_i$  where  $\Gamma_i$  are instantiation judgments.

$$\Gamma_i := \cdot \mid \Gamma_i, \hat{\alpha} \leq A \mid \Gamma_i, A \leq \hat{\alpha} \quad \text{where } \hat{\alpha} \notin FV(A) \cup FV(\Gamma_i)$$

$\boxed{ \Gamma _e,  \omega _e,  e _e}$ Size measure.	$\boxed{ \Gamma _\forall,  \omega _\forall,  A _\forall}$ Polymorphism measure.
$ \Gamma _e = \sum_{\omega \in \Gamma}  \omega _e$	$ \Gamma _{\Leftrightarrow} = \sum_{\omega \in \Gamma}  \omega _\forall$
$ A \leq B _e = 0$	$ A \leq B _\forall =  A _\forall +  B _\forall$
$ e \Leftarrow A _e =  e _e$	$ e \Leftarrow A _\forall =  A _\forall$
$ e \Rightarrow \omega _e =  e _e +  \omega _e$	$ e \Rightarrow \omega _\forall =  \omega _\forall$
$ A \bullet e \Rightarrow_a \omega _e =  e _e +  \omega _e$	$ A \bullet e \Rightarrow_a \omega _\forall =  A _\forall +  \omega _\forall$
$ x _e =  () _e = 1$	$ 1 _\forall =  a _\forall =  \hat{\alpha} _\forall = 1$
$ \lambda x. e _e =  e _e + 1$	$ A \rightarrow B _\forall =  A _\forall +  B _\forall$
$ e_1 e_2 _e =  e_1 _e +  e_2 _e + 1$	$ \forall a. A _\forall =  A _\forall + 1$
$ e : A _e =  e _e + 1$	
$\boxed{ \Gamma _{\Leftrightarrow},  \omega _{\Leftrightarrow}}$ Judgment measure.	$\boxed{ \Gamma _{\rightarrow},  \omega _{\rightarrow},  A _{\rightarrow}}$ Function measure.
$ \Gamma _{\Leftrightarrow} = \sum_{\omega \in \Gamma}  \omega _e$	$ \Gamma _{\rightarrow} = \sum_{\omega \in \Gamma}  \omega _\forall$
$ A \leq B _{\Leftrightarrow} = 0$	$ A \leq B _{\rightarrow} =  A _{\rightarrow} +  B _{\rightarrow}$
$ e \Leftarrow A _{\Leftrightarrow} = 2$	$ e \Leftarrow A _{\rightarrow} =  A _{\rightarrow}$
$ e \Rightarrow \omega _{\Leftrightarrow} =  \omega _{\Leftrightarrow} + 1$	$ e \Rightarrow \omega _{\rightarrow} =  \omega _{\rightarrow}$
$ A \bullet e \Rightarrow_a \omega _{\Leftrightarrow} =  \omega _{\Leftrightarrow} + 3$	$ A \bullet e \Rightarrow_a \omega _{\rightarrow} =  A _{\rightarrow} +  \omega _{\rightarrow}$
$\boxed{ \Gamma _{\hat{\alpha}}}$ Existential measure.	$ 1 _{\rightarrow} =  a _{\rightarrow} =  \hat{\alpha} _{\rightarrow} = 1$
$ \Gamma _{\hat{\alpha}} = \#\hat{\alpha} \in \Gamma$	$ A \rightarrow B _{\rightarrow} =  A _{\rightarrow} +  B _{\rightarrow} + 1$
	$ \forall a. A _{\rightarrow} =  A _{\rightarrow}$

Figure 4.10: Worklist measures

Fully reducing these instantiation judgments at the top of the worklist has a twofold impact. Firstly, new entries may be pushed onto the worklist which are not instantiation judgments. This only happens when  $\Gamma_i$  contains a universal quantifier that is reduced by Rule 8 or 9. The new entries then are of the form  $\Gamma_{\leq}$ :

$$\Gamma_{\leq} := \cdot \mid \Gamma_{\leq}, a \mid \Gamma_{\leq}, \hat{\alpha} \mid \Gamma_{\leq}, A \leq B$$

Secondly, reducing the instantiation judgments may also affect the remainder of the worklist  $\Gamma$ , by solving existing existentials and introducing new ones. This worklist update is captured in the update judgment  $\Gamma \rightarrow \Gamma'$  defined in Figure 4.11. For instance, an existential variable

$$\begin{array}{c}
 \boxed{\Gamma \rightarrow \Gamma'} \quad \Gamma \text{ updates to } \Gamma'. \\
 \frac{}{\Gamma \rightarrow \Gamma} \rightarrow \text{id} \quad \frac{|A|_{\forall} = 0 \quad \Gamma_L, [A/\hat{\alpha}] \Gamma_R \rightarrow \Gamma'}{\Gamma_L, \hat{\alpha}, \Gamma_R \rightarrow \Gamma'} \rightarrow \text{solve} \quad \frac{\Gamma_L, \hat{\alpha}, \Gamma_R \rightarrow \Gamma'}{\Gamma_L, \Gamma_R \rightarrow \Gamma'} \rightarrow \hat{\alpha}
 \end{array}$$

Figure 4.11: Worklist Update

instantiation,  $\Gamma_L, \hat{\alpha}, \Gamma_R \rightarrow \Gamma_L, \hat{\alpha}_1, \hat{\alpha}_2, [\hat{\alpha}_1 \rightarrow \hat{\alpha}_2/\hat{\alpha}] \Gamma_R$ , can be derived as a combination of the three rules that define the update relation.

The good news is that worklist updates do not affect the three main worklist measures:

**Lemma 4.15** (Measure Invariants of Worklist Extension). *If  $\Gamma \rightarrow \Gamma'$  then  $|\Gamma|_e = |\Gamma'|_e$  and  $|\Gamma|_{\Leftrightarrow} = |\Gamma'|_{\Leftrightarrow}$  and  $|\Gamma|_{\forall} = |\Gamma'|_{\forall}$ .*

Moreover, we can characterize the reduction of the instantiation judgments as follows.

**Lemma 4.16** (Instantiation Decidability). *For any well-formed algorithmic worklist  $(\Gamma, \Gamma_i)$ :*

- 1) *If  $|\Gamma_i|_{\forall} = 0$ , then there exists  $\Gamma'$   
 s.t.  $(\Gamma, \Gamma_i) \rightarrow^* \Gamma'$  and  $|\Gamma'|_{\hat{\alpha}} = |\Gamma|_{\hat{\alpha}} - |\Gamma_i|$  and  $\Gamma \rightarrow \Gamma'$ .*
- 2) *If  $|\Gamma_i|_{\forall} > 0$ , then there exist  $\Gamma', \Gamma_{\leq}$   
 s.t.  $(\Gamma, \Gamma_i) \rightarrow^* (\Gamma', \Gamma_{\leq})$  and  $|\Gamma_{\leq}|_{\forall} = |\Gamma_i|_{\forall} - 1$  and  $\Gamma \rightarrow \Gamma'$ .*

Hence, reducing the instantiation judgment prefix  $\Gamma_i$  either decreases the number of universal quantifiers or eliminates one existential variable per instantiation judgment. The proof of this lemma proceeds by induction on the measure  $2 * |\Gamma_i|_{\rightarrow} + |\Gamma_i|$  of the instantiation judgment list  $\Gamma_i$ .

Let us go back to the whole algorithm and summarize our findings. The decidability theorem is shown through a lexicographic group of induction measures

$$\langle |\Gamma|_e, |\Gamma|_{\Leftrightarrow}, |\Gamma|_{\forall}, |\Gamma|_{\hat{\alpha}}, |\Gamma|_{\rightarrow} + |\Gamma| \rangle$$

which is decreased by nearly all rules. In the exceptional case that the measure does not decrease immediately, we encounter an instantiation judgment at the top of the worklist. We can then make use of Lemma 4.16 to show that  $|\Gamma|_{\hat{\alpha}}$  or  $|\Gamma|_{\forall}$  decreases when that instantiation judgment is consumed or partially reduced. Moreover, Lemma 4.15 establishes that no higher-priority measure component increases. Hence, in the exceptional case we have an overall measure decrease too.

Combining all three main results (soundness, completeness and decidability), we conclude that the declarative system is decidable by means of our algorithm.

**Corollary 4.17** (Decidability of Declarative Typing). *Given wf  $\Omega$ , it is decidable whether  $\Omega \longrightarrow^* \cdot$  or not.*

#### 4.3.6 ABELLA AND PROOF STATISTICS

We have chosen the Abella (v2.0.7-dev<sup>3</sup>) proof assistant [Gacek 2008] to develop our formalization. Abella is designed to help with formalizations of programming languages, due to its built-in support for variable binding and the  $\lambda$ -tree syntax [Miller 2000], which is a form of HOAS. Nominal variables, or  $\nabla$ -quantified variables, are used as an unlimited name supply, which supports explicit freshness control and substitutions. Although Abella lacks packages, tactics and support for modules, its higher-order unification and the ease of formalizing substitution-intensive relations are very helpful.

While the algorithmic rules are in a small-step style, the proof script rewrites them into a big-step style for easier inductions. In addition, we do prove the equivalence of the two representations.

**STATISTICS OF THE PROOF** The proof script consists of 7,977 lines of Abella code with a total of 60 definitions and 596 theorems. Figure 4.1 briefly summarizes the contents of each file. The files are linearly dependent due to the limitations of Abella.

**TRANSLATION TABLE FOR THE PROOF SCRIPTS** In the proof scripts, we use textual relational definitions rather than the symbolic ones used in the paper. The mapping, shown in Table 4.2, should be helpful for anyone who reads the script.

## 4.4 DISCUSSION

This section discusses some insights that we gained from our work and contrasts the scoping mechanisms we have employed with those in DK’s algorithm. We also discuss a way to improve the precision of their scope tracking. Furthermore we discuss and sketch an extension of our algorithm with an elaboration to a target calculus, and discuss an extension of our algorithm with scoped type variables [Peyton Jones and Shields 2004].

### 4.4.1 CONTRASTING OUR SCOPING MECHANISMS WITH DK’S

A nice feature of our worklists is that, simply by interleaving variable declarations and judgment chains, they make the scope of variables precise. DK’s algorithm deals with garbage

<sup>3</sup>We use a development version because the developers just fixed a serious bug that accepts a simple proof of false, which also affects our proof. Specifically, our scripts compile against commit 92829a of Abella’s GitHub repository.

Table 4.1: Statistics for the proof scripts

File(s)	LOC	#Thm	Description
olist.thm, nat.thm	311	57	Basic data structures
typing.thm	245	7	Declarative & algorithmic system, debug examples
decl.thm	226	33	Basic declarative properties
order.thm	235	27	The $ \cdot _{\forall}$ measure; decl. subtyping strengthening
alg.thm	679	80	Basic algorithmic properties
trans.thm	616	53	Worklist instantiation and declarative transfer; Lemmas <a href="#">4.3</a> , <a href="#">4.4</a>
declTyping.thm	909	70	Non-overlapping declarative system; Lemmas <a href="#">4.5</a> , <a href="#">4.6</a> , <a href="#">4.7</a> , <a href="#">4.8</a>
soundness.thm	1,107	78	Soundness theorem; aux. lemmas on transfer
depth.thm	206	14	The $ \cdot _{\rightarrow}$ measure; Lemma <a href="#">4.12</a>
dcl.thm	380	12	Non-overlapping declarative worklist
completeness.thm	1,124	61	Completeness theorem; aux. lemmas and relations
inst_decidable.thm	837	45	Other worklist measures; Lemma <a href="#">4.16</a>
decidability.thm	983	57	Decidability theorem and corollary
smallStep.thm	119	2	The equivalence between big-step and small-step
<i>Total</i>	7,977	596	(60 definitions in total)

Table 4.2: Translation Table for the Proof Scripts

Symbol	Abella textual relation	File (.thm)	Description
$A$	<code>ty, wft, wfta</code>	<code>typing</code>	The "type" type, decl./alg. well-formedness
$e$	<code>tm, wftm</code>	<code>typing</code>	The "term" type, (alg.) well-formedness
$\Psi, \Gamma$	<code>olist, wfj</code>	<code>typing</code>	olist is Abella's built-in list type, i.e. [o]
$\omega$	<code>judgment, wfjg</code>	<code>typing</code>	Judgment chain and its well-formedness
$\Gamma \rightarrow \Gamma'$	<code>reduction</code>	<code>smallStep</code>	Algorithmic reduction: paper version
$\Gamma \rightarrow^* \cdot$	<code>judge</code>	<code>typing</code>	(judge $\Gamma$ ): a success reduction on $\Gamma$
$\Gamma \rightsquigarrow \Omega$	<code>tex</code>	<code>trans</code>	"transfer existential variables"
$\Omega \rightarrow^* \cdot$	<code>dc, dcl</code>	<code>trans, dcl</code>	Declarative chain representation
$\Psi' \leq \Psi$	<code>esub</code>	<code>declTyping</code>	Declarative context subtyping
$\Gamma \rightarrow \Gamma'$	<code>jExt</code>	<code>inst_decidable</code>	"Judgment extension"
$\Gamma_i$	<code>iexp</code>	<code>inst_decidable</code>	instantiation judgments ( $\Gamma \vdash_{wf} \Gamma_i$ )
$\Gamma \leq$	<code>subExp</code>	<code>inst_decidable</code>	subtyping judgments
$ \cdot _e$	<code>tmSize, tmSizeJ, tmSize1</code>	<code>declTyping, decidability</code>	Size measure for term, judgment chain or context
$ \cdot _{\Leftrightarrow}$	<code>m_judgeJ, m_judge</code>	<code>decidability</code>	Judgment measure for judgment chain or context
$ \cdot _{\hat{\alpha}}$	<code>nVar</code>	<code>inst_decidable</code>	Existential measure for context
$ \cdot _{\forall}$	<code>order, orderJ, order1</code>	<code>order, inst_decidable</code>	Polymorphism measure for type, judgment chain or context
$ \cdot _{\rightarrow}$	<code>depth, depthJ, depth1</code>	<code>order, decidability</code>	Function measure for type, judgment chain or context

collecting variables in a different way: through type variable or existential variable markers (as discussed in Section 4.1.2). Despite the sophistication employed in DK's algorithm to keep scoping precise, there is still a chance that unused existential variables leak their scope to an output context and accumulate indefinitely. For example, the derivation of the judgment  $(\lambda x. x) () \Leftarrow 1$  is as follows

$$\begin{array}{c}
 \frac{\dots x \Rightarrow \hat{\alpha} \dots \quad \dots \hat{\alpha} \leq \hat{\beta} \dots}{\Gamma, \hat{\alpha}, \hat{\beta}, x : \hat{\alpha} \vdash x \Leftarrow \hat{\beta} \dashv \Gamma_1, x : \hat{\alpha}} \quad \frac{\dots () \Leftarrow \hat{\alpha} \dots}{\Gamma_1 \vdash \hat{\alpha} \rightarrow \hat{\alpha} \bullet () \Rightarrow \hat{\alpha} \dashv \Gamma_2} \\
 \hline
 \Gamma \vdash \lambda x. x \Rightarrow \hat{\alpha} \rightarrow \hat{\beta} \dashv \Gamma_1 \quad \Gamma_1 \vdash \hat{\alpha} \rightarrow \hat{\alpha} \bullet () \Rightarrow \hat{\alpha} \dashv \Gamma_2 \\
 \hline
 \Gamma \vdash (\lambda x. x) () \Rightarrow \hat{\alpha} \dashv \Gamma_2 \\
 \hline
 \Gamma_2 \vdash 1 \leq 1 \dashv \Gamma_2 \\
 \hline
 \Gamma \vdash (\lambda x. x) () \Leftarrow 1 \dashv \Gamma, \hat{\alpha} = 1, \hat{\beta} = \hat{\alpha}
 \end{array}$$

where  $\Gamma_1 := (\Gamma, \hat{\alpha}, \hat{\beta} = \hat{\alpha})$  solves  $\hat{\beta}$ , and  $\Gamma_2 := (\Gamma, \hat{\alpha} = 1, \hat{\beta} = \hat{\alpha})$  solves both  $\hat{\alpha}$  and  $\hat{\beta}$ .

If the reader is not familiar with DK's algorithm, he/she might be confused about the inconsistent types across judgment. As an example,  $(\lambda x. x) ()$  synthesizes  $\hat{\alpha}$ , but the second premise of the subsumption rule uses 1 for the result. This is because a context application  $[\Gamma, \hat{\alpha} = 1, \hat{\beta} = \hat{\alpha}] \hat{\alpha} = 1$  happens between the premises.

The existential variables  $\hat{\alpha}$  and  $\hat{\beta}$  are clearly not used after the subsumption rule, but according to the algorithm, they are kept in the context until some parent judgment recycles a block of variables, or to the very end of a type inference task. In that sense, DK's algorithm does not control the scoping of variables precisely.

Two rules we may blame for not garbage collecting correctly are the inference rule for lambda functions and an application inference rule:

$$\frac{\Gamma, \hat{\alpha}, \hat{\beta}, x : \hat{\alpha} \vdash e \Leftarrow \hat{\beta} \dashv \Delta, x : \hat{\alpha}, \Theta}{\Gamma \vdash \lambda x. e \Rightarrow \hat{\alpha} \rightarrow \hat{\beta} \dashv \Delta} \text{DK}_{\rightarrow} \text{I} \Rightarrow \quad \frac{\Gamma, \hat{\alpha} \vdash [\hat{\alpha}/a] A \bullet e \Rightarrow C \dashv \Delta}{\Gamma \vdash \forall a. A \bullet e \Rightarrow C \dashv \Delta} \text{DK}_{\forall} \text{App}$$

In contrast, Rule 25 of our algorithm collects the existential variables right after the second judgment chain, and Rule 27 collects one existential variable similarly:

$$\Gamma \Vdash \lambda x. e \Rightarrow_a \omega \longrightarrow_{25} \Gamma, \hat{\alpha}, \hat{\beta} \Vdash [\hat{\alpha} \rightarrow \hat{\beta}/a] \omega, x : \hat{\alpha} \Vdash e \Leftarrow \hat{\beta}$$

$$\Gamma \Vdash \forall a. A \bullet e \Rightarrow_a \omega \longrightarrow_{27} \Gamma, \hat{\alpha} \Vdash [\hat{\alpha}/a] A \bullet e \Rightarrow_a \omega$$

It seems impossible to achieve a similar outcome in DK's system by only modifying these two rules. Taking  $\text{DK}_{\rightarrow} \text{I} \Rightarrow$  as an example, the declaration or solution for  $\hat{\alpha}$  and  $\hat{\beta}$  may be

referred to by subsequent judgments. Therefore leaving  $\hat{\alpha}$  and  $\hat{\beta}$  in the output context is the only choice, when the subsequent judgments cannot be consulted.

To fix the problem, one possible modification is on the algorithmic subsumption rule, as garbage collection for a checking judgment is always safe:

$$\frac{\Gamma, \blacktriangleright_{\hat{\alpha}} \vdash e \Rightarrow A \dashv \Theta \quad \Theta \vdash [\Theta]A \leq [\Theta]B \dashv \Delta, \blacktriangleright_{\hat{\alpha}}, \Delta'}{\Gamma \vdash e \Leftarrow B \dashv \Delta} \text{DK\_Sub}$$

Here we employ the markers in a way they were originally not intended for. We create a dummy fresh existential  $\hat{\alpha}$  and add a marker to the input context of the inference judgment. After the subtyping judgment is processed we look for the marker and drop everything afterwards. We pick this rule because it is the only one where a checking judgment calls an inference judgment. That is the point where our algorithm recycles variables—right after a judgment chain is satisfied. After applying this patch, to the best of our knowledge, DK's algorithm behaves equivalently to our algorithm in terms of variable scoping. However, this exploits markers in a way they were not intended to be used and seems ad-hoc.

#### 4.4.2 ELABORATION

Type-inference algorithms are often extended with an associated elaboration. For example, for languages with implicit polymorphism, it is common to have an elaboration to a variant of System F [Reynolds 1983], which recovers type information and explicit type applications. Therefore a natural question is whether our algorithm can also accommodate such elaboration. While our algorithmic reduction does not elaborate to System F, we believe that it is not difficult to extend the algorithm with a (type-directed) elaboration. We explain the rough idea as follows.

Since the judgment form of our algorithmic worklist contains a collection of judgments, elaboration expressions are also generated as a list of equal length to the number of judgments (*not judgment chains*) in the worklist. As usual, subtyping judgments translate to coercions (denoted by  $f$  and represented by System F functions), all three other types of judgments translate to terms in System F (denoted by  $t$ ).

Let  $\Phi$  be the elaboration list, containing translated type coercions and terms:

$$\Phi ::= \cdot \mid \Phi, f \mid \Phi, t$$

Then the form of our algorithmic judgment becomes:

$$\Gamma \hookrightarrow \Phi$$



We take Rule 18 as an example, rewriting small-step reduction in a relational style,

$$\frac{\Gamma \Vdash e \Rightarrow_a a \leq B \hookrightarrow \Phi, f, t}{\Gamma \Vdash e \Leftarrow B \hookrightarrow \Phi, f, t} \text{ Translation\_18}$$

As is shown in the conclusion of the rule, a checking judgment at the top of the worklist corresponds to a top element for elaboration. The premise shows that one judgment chain may relate to more than one elaboration elements, instead, one single judgment relates to one elaborated term. The outer judgment, being processed before inner ones, elaborates to the top element in the elaboration list.

Moreover, existential variables need special treatment, since they may be solved at any point, or be recycled if not solved within their scopes. If an existential variable is solved, we not only propagate the solution to the other judgments, but also replace occurrences in the elaboration list. If an existential variable is recycled, meaning that it is not constrained, we can pick any well-formed monotype as its solution. The unit type 1, as the simplest type in the system, is a good choice.

#### 4.4.3 LEXICALLY-SCOPED TYPE VARIABLES

We have further extended the type system with the support for lexically-scoped type variables [Peyton Jones and Shields 2004]. Our Abella formalization for this extension proves all the metatheory we discuss in Section 4.3.

From a practical point of view, this extension allows the implementation of a function to refer to type variables from its type signature. For example,

$$(\lambda x. \lambda y. (x : a)) : \forall a b. a \rightarrow b \rightarrow a$$

has an annotation  $(x : a)$  that refers to the type variable  $a$  in the type signature. This is not a surprising feature, since the declarative system already accepts similar programs

$$\frac{\Psi, a \vdash e \Leftarrow A}{\Psi \vdash e \Leftarrow \forall a. A} \text{ Decl}\forall\text{I} \quad \frac{\Psi \vdash \forall a. A \quad \Psi \vdash e \Leftarrow \forall a. A}{\Psi \vdash (e : \forall a. A) \Rightarrow \forall a. A} \text{ DeclAnno}$$

The main issue is the well-formedness condition. Normally  $\Psi \vdash (e : A)$  follows from  $\Psi \vdash e$  and  $\Psi \vdash A$ . However, when  $A = \forall a. A'$ , the type variable  $a$  is not in scope at  $e$ ,

therefore  $\Psi \vdash e$  is not derivable. To address the problem, we add a new syntactic form that explicitly binds a type variable simultaneously in a function and its annotation.

$$\text{Expressions} \quad e ::= \dots \mid \Lambda a. e : A$$

This new type-lambda syntax  $\Lambda a. e : A$  actually annotates its body  $e$  with  $\forall a. A$ , while making  $a$  visible inside the body of the function. The well-formedness judgments are extended accordingly:

$$\frac{\Psi, a \vdash e \quad \Psi, a \vdash A}{\Psi \vdash \Lambda a. e : A} \text{wf}_d\Lambda \qquad \frac{\Gamma, a \vdash e \quad \Gamma, a \vdash A}{\Gamma \vdash \Lambda a. e : A} \text{wf}_\Lambda$$

Corresponding rules are introduced for both the declarative system and the algorithmic system:

$$\frac{\Psi, a \vdash A \quad \Psi, a \vdash e \Leftarrow A}{\Psi \vdash \Lambda a. e : A \Rightarrow \forall a. A} \text{Decl}\Lambda$$

$$\Gamma \Vdash \Lambda a. e : A \Rightarrow_b \omega \longrightarrow_{30} \Gamma \Vdash [(\forall a. A)/b]\omega, a \Vdash e \Leftarrow A$$

In practice, programmers would not write the syntax  $\Lambda a. e : A$  directly. The `Scoped-TypeVariables` extension of Haskell is effective only when the type signature is explicitly universally quantified (which the compiler translates into an expression similar to  $\Lambda a. e : A$ ); otherwise the program means the normal syntax  $e : \forall a. A$  and may not later refer to the type variable  $a$ .

We have proven all three desired properties for the extended system, namely soundness, completeness and decidability.

# 5 HIGHER-RANKED POLYMORPHISM WITH OBJECT-ORIENTED SUBTYPING

## 5.1 OVERVIEW

In this chapter, we present a worklist algorithm that further supports subtyping, by introducing the top and bottom types. Such type inference algorithms are known to be hard to design, due to the complication of unifying variables under subtyping inequality rather than equality. Therefore, algorithmic existential variables that appear in subtyping judgments might have more solutions than before, and the eager substitution rules miss some possibilities. Our new *backtracking* algorithm extends the existing one by some overlapping rules that non-deterministically try different sorts of solutions instead of a single sort, improving the rate of success guessing.

Formalization of the algorithm in the Abella theorem prover shows soundness of the algorithm with respect to our declarative specification. Although the backtracking algorithm is still incomplete in some corner cases, the proof script indicates that the only source of incompleteness comes from higher-ranked subtyping relations. Following that discovery, we formally proved partial completeness of our subtyping algorithm under the rank-1 restriction, and is at least comparable to local type inference algorithms.

### 5.1.1 TYPE INFERENCE IN PRESENCE OF SUBTYPING

In all previous chapters, “subtyping” refers to a relationship that compares the degree of polymorphism between two types. In addition to that, we also include the top and bottom types. Both of the types have practical uses, especially in object-oriented programming languages. The top type,  $\top$ , is the super type of any type, i.e. any type is more general than  $\top$  and thus can be considered as an instance of type  $\top$ . In typical object-oriented programming languages, the `Object` class, as the base class of any class, is the  $\top$  type. In contrast, the bottom type,  $\perp$ , is dual to  $\top$ . An instance of  $\perp$  can be cast to a value of any type, which is also viewed as the uninhabited type therefore one cannot find any value that has the  $\perp$  type. However, there are few exceptions. One of them is a `null` pointer value (`(void *)` in C++, for example). Another practical use for bottom types is for exceptions. The type `Exception`  $\rightarrow \perp$ .

given to the `raise` function may pass type checkers naturally. The type  $\text{Exception} \rightarrow \forall a. a$  is also a reasonable choice, which in fact reveals that  $\perp$  behaves almost identically as  $\forall a. a$ . From the theoretical point of view, both of them represent “falsity”.

A number of systems and algorithms are proposed for type inference with subtyping. Since subtyping constraints cannot be easily reduced in presence of subtyping, it is natural to extend the syntax of types so that unsolved constraints are carried with them, and the idea is studied by previous work [Eifrig et al. 1995b; Trifonov and Smith 1996]. For example, the following function

```
select p v d = if (p v) then v else d
```

is inferred to have the type

$$(\alpha \rightarrow \text{bool}) \rightarrow \alpha \rightarrow \beta \rightarrow \gamma \mid \alpha \leq \gamma, \beta \leq \gamma$$

However, this representation usually contains a large set of constraints, which might confuse programmers. Besides, subtyping comparison between constraint types may introduce additional difficulty for programmers to understand, which weakens its practical use.

**MLsub** MLsub [Dolan and Mycroft 2017], as introduced in Section 2.4, is based on a type system with lattices and polar types. Inspired by [Pottier 1998], they separated input and output types into polar types, which greatly simplifies the subtyping judgment. With the further help of their biunification algorithm, constraints are no longer carried with types, instead they are expressed directly on types. For example, the `select` function defined above has type

$$(\alpha \rightarrow \text{bool}) \rightarrow \alpha \rightarrow \beta \rightarrow \alpha \sqcup \beta$$

in MLsub. The  $\sqcup$  symbol in the output type indicates that any of the types might be returned when executed. On the other hand, the  $\sqcap$  symbol might occur on an input type, when that argument is used as different types, meaning that the type must be able to convert to all these types.

Nevertheless, there are still drawbacks in practice. Firstly, the type inference algorithm produces types that have comparable sizes to previous systems with constraint types. Although a simplification algorithm is provided, it is not guaranteed to produce the most simple form.

Secondly, MLsub always infers an expression and returns its principal type, but some of the types are not easy to understand. For example, the `twice` function defined as

```
twice f x = f (f x)
```

has type  $(\alpha \sqcup \beta \rightarrow \alpha) \rightarrow \beta \rightarrow \alpha$ , where normally people might expect  $(\alpha \rightarrow \alpha) \rightarrow \alpha \rightarrow \alpha$ .

Thirdly, MLsub tries to infer on almost every possible expression, even those that might just be a mistake made by the programmer. For instance, the function

```
positive f x = if x > 0 then f x else f 0
```

always checks the positivity of  $x$  before applying to  $f$ . However, a slight mistake is made in the following definition

```
positive' f x = if x > 0 then f x else f
```

then the algorithm might still accept the program by giving the type:

$$\alpha \sqcap (\beta \rightarrow \alpha) \rightarrow \beta \sqcap \text{int} \rightarrow \alpha$$

Type inference algorithms, as logical tools that help people program, should accept good programs and reject bad ones. Typically, such a function is considered to be badly-written, but unfortunately, MLsub accepts it without giving any warnings.

### 5.1.2 JUDGMENT LIST AND EAGER SUBSTITUTION

The judgment list algorithm we discussed in the last chapter treats subtyping judgments mainly as unification. However, the algorithm may lose some solutions when top and bottom types are introduced. For example, the judgment  $\hat{\alpha} \leq 1$  has the best solution  $\hat{\alpha} := 1$  in the previous system, but now  $\hat{\alpha} := \perp$  should also be allowed. As a result, several subtyping judgments cannot be reduced easily with a single eager substitution.

Similar incompleteness also affects the instantiation and existential variable solving rules

$$\begin{aligned} \Gamma[\hat{\alpha}] \Vdash \hat{\alpha} \leq A \rightarrow B &\longrightarrow_{10} [\hat{\alpha}_1 \rightarrow \hat{\alpha}_2 / \hat{\alpha}] (\Gamma[\hat{\alpha}_1, \hat{\alpha}_2] \Vdash \hat{\alpha}_1 \rightarrow \hat{\alpha}_2 \leq A \rightarrow B) \\ \Gamma[\hat{\alpha}][\hat{\beta}] \Vdash \hat{\alpha} \leq \hat{\beta} &\longrightarrow_{12} [\hat{\alpha} / \hat{\beta}] (\Gamma[\hat{\alpha}][\hat{\beta}]) \end{aligned}$$

The instantiation judgment  $\hat{\alpha} \leq A \rightarrow B$  immediately split  $\hat{\alpha}$  into  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$ , which is no longer the case —  $\hat{\alpha}$  can also be  $\perp$ .

The rule that solves one existential variable against another,  $\hat{\alpha} \leq \hat{\beta}$ , is even more problematic, because the shapes of the instantiations of both existential variables are unknown, there are infinitely many possibilities if no additional constraints present. For example, if  $\hat{\alpha}$  is finally instantiated to  $\text{Int} \rightarrow \text{Int}$ , then  $\hat{\beta}$  may be solved to the same type as  $\hat{\alpha}$ , or other types that satisfy the subtyping relationship, including  $\text{Int} \rightarrow \top$ ,  $\perp \rightarrow \text{Int}$  or  $\perp \rightarrow \top$ . What's worse, given only this single subtyping judgment  $\hat{\alpha} \leq \hat{\beta}$ , we have no assumption on both existential variables.

Theoretically speaking, if we extend the type system with  $\top$  and  $\perp$  types, the definition of monotypes becomes

$$\tau ::= 1 \mid \perp \mid \top \mid \tau \rightarrow \tau \mid a$$

, which causes the subtyping problem to be undecidable Su et al. [2002]. Therefore, it is impossible to derive a complete algorithm for the declarative system, and we try to keep the algorithm complete for common usages of polymorphism, but pose restrictions in some higher-ranked cases.

### 5.1.3 OUR SOLUTION: BACKTRACKING ALGORITHM

After examining the examples where the worklist algorithm behaves incompletely, we propose some improvements to the worklist algorithm to accept some simple cases. We have observed that subtyping judgments like  $\hat{\alpha} \leq A$  and  $A \leq \hat{\alpha}$  have trivial solutions  $\hat{\alpha} := \perp$  and  $\hat{\alpha} := \top$  respectively, thus we can try these simple solutions before the more complicated analysis.

For example, for the judgment list  $\Gamma \Vdash \hat{\alpha} \leq A \rightarrow B$ , we first assign  $\hat{\alpha} := \perp$  and continue to check other judgments. If the judgment reduction succeeds, then we know that  $\hat{\alpha} := \perp$  is a possible solution. If the reduction failed, we continue to try the other possibility according to the previous instantiation rule by splitting  $\hat{\alpha}$  into two existential variables  $\hat{\alpha}_1$  and  $\hat{\alpha}_2$ . The procedure must keep track of all backtracking points until type-checking the whole program. Note that for the reduced judgment  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2 \leq A \rightarrow B$ ,  $\hat{\alpha}_1$  and  $\hat{\alpha}_2$  might be solved to  $\top$  or  $\perp$ .

For the type system, there are also interesting changes to both the declarative system and algorithmic system. Two algorithmic typing rules also adopt the backtracking approach:

- The checking judgment  $e \Leftarrow \hat{\alpha}$  now has a trivial solution  $\hat{\alpha} := \top$ , because  $e \Leftarrow \top$  should be accepted for any (well-formed) expression.
- The introduction of top and bottom types comes with another new declarative rule  $\perp \bullet e \Rightarrow \perp$ . Therefore the application inference judgment  $\hat{\alpha} \bullet e \Rightarrow_a \omega$  can be satisfied by  $\hat{\alpha} := \perp$ , with output  $\perp$  applied to the rest of the judgment chain  $\omega$ .

By applying the above modifications, the algorithm remains sound w.r.t our new declarative system, and at the same time, we address most issues regarding the top and bottom types, obtaining a more complete algorithm. Furthermore, we formally verified that the subtyping algorithm is complete under the rank-1 restriction. However, for rank-2 judgments like  $\hat{\alpha} \leq \hat{\beta}$ , a naive complete algorithm would traverse an infinitely large tree of possibilities and diverge for some set of constraints.

Type variables	$a, b$
Types	$A, B, C ::= 1 \mid \top \mid \perp \mid a \mid \forall a. A \mid A \rightarrow B$
Monotypes	$\tau ::= 1 \mid \top \mid \perp \mid a \mid \tau_1 \rightarrow \tau_2$
Expressions	$e ::= x \mid () \mid \lambda x. e \mid e_1 e_2 \mid (e : A)$
Context	$\Psi ::= \cdot \mid \Psi, a \mid \Psi, x : A$

Figure 5.1: Declarative Syntax

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

$$\begin{array}{c}
 \frac{a \in \Psi}{\Psi \vdash a \leq a} \leq \text{Var} \quad \frac{}{\Psi \vdash 1 \leq 1} \leq \text{Unit} \quad \frac{\Psi \vdash B_1 \leq A_1 \quad \Psi \vdash A_2 \leq B_2}{\Psi \vdash A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2} \leq \rightarrow \\
 \frac{\Psi \vdash \tau \quad \Psi \vdash [\tau/a]A \leq B}{\Psi \vdash \forall a. A \leq B} \leq \forall L \quad \frac{\Psi, b \vdash A \leq B}{\Psi \vdash A \leq \forall b. B} \leq \forall R \\
 \frac{}{A \leq \top} \leq \text{Top} \quad \frac{}{\perp \leq A} \leq \text{Bot}
 \end{array}$$

Figure 5.2: Declarative Subtyping

Following the formal statements, we conclude that the backtracking algorithm is complete under a rank-1 restriction, but not always on the general higher-ranked settings. In those settings, our algorithm only consider the equality case and ignore any other possibilities. In another perspective, our algorithm does not try to solve completely when parametric polymorphism and subtyping polymorphism occur at the same time. This is not theoretically satisfactory, but we argue that this handles most common programs, and the programmer may use annotation to guide the inference algorithm with the real type he intends to instantiate.

## 5.2 DECLARATIVE SYSTEM

**SYNTAX** The syntax of the declarative system, shown in Figure 5.1, is similar to the previous systems by having a primitive type 1, type variables  $a$ , polymorphic types  $\forall a. A$  and function types  $A \rightarrow B$ . Additionally, top and bottom types are introduced to the type system.

The well-formedness formalization of the system is standard and almost identical to the previous systems, therefore we omit the formal definitions.

$\boxed{\Psi \vdash e \Leftarrow A}$	$e$ checks against input type $A$ .
$\boxed{\Psi \vdash e \Rightarrow A}$	$e$ synthesizes output type $A$ .
$\boxed{\Psi \vdash A \bullet e \Rightarrow C}$	Applying a function of type $A$ to $e$ synthesizes type $C$ .

$$\begin{array}{c}
 \frac{(x : A) \in \Psi}{\Psi \vdash x \Rightarrow A} \text{DeclVar} \quad \frac{\Psi \vdash e \Rightarrow A \quad \Psi \vdash A \leq B}{\Psi \vdash e \Leftarrow B} \text{DeclSub} \\
 \frac{\Psi \vdash A \quad \Psi \vdash e \Leftarrow A}{\Psi \vdash (e : A) \Rightarrow A} \text{DeclAnno} \quad \frac{}{\Psi \vdash () \Rightarrow 1} \text{Decl1I} \Rightarrow \\
 \frac{}{\Psi \vdash () \Leftarrow 1} \text{Decl1I} \quad \frac{\Psi \vdash e}{\Psi \vdash e \Leftarrow \top} \text{Decl}\top \quad \frac{\Psi \vdash e}{\Psi \vdash \perp \bullet e \Rightarrow \perp} \text{Decl}\perp\text{App} \\
 \frac{\Psi, a \vdash e \Leftarrow A}{\Psi \vdash e \Leftarrow \forall a. A} \text{Decl}\forall\text{I} \quad \frac{\Psi \vdash \tau \quad \Psi \vdash [\tau/a]A \bullet e \Rightarrow C}{\Psi \vdash \forall a. A \bullet e \Rightarrow C} \text{Decl}\forall\text{App} \\
 \frac{\Psi, x : A \vdash e \Leftarrow B}{\Psi \vdash \lambda x. e \Leftarrow A \rightarrow B} \text{Decl}\rightarrow\text{I} \quad \frac{\Psi \vdash \sigma \rightarrow \tau \quad \Psi, x : \sigma \vdash e \Leftarrow \tau}{\Psi \vdash \lambda x. e \Rightarrow \sigma \rightarrow \tau} \text{Decl}\rightarrow\text{I} \Rightarrow \\
 \frac{\Psi \vdash e_1 \Rightarrow A \quad \Psi \vdash A \bullet e_2 \Rightarrow C}{\Psi \vdash e_1 e_2 \Rightarrow C} \text{Decl}\rightarrow\text{E} \quad \frac{\Psi \vdash e \Leftarrow A}{\Psi \vdash A \rightarrow C \bullet e \Rightarrow C} \text{Decl}\rightarrow\text{App}
 \end{array}$$

Figure 5.3: Declarative Typing

**DECLARATIVE SUBTYPING** Shown in Figure 5.2, the declarative subtyping extends the polymorphic subtyping relation originally proposed by Odersky and Läufer [1996] by adding Rules  $\leq\text{Top}$  and  $\leq\text{Bot}$ , defining the properties of the  $\top$  and  $\perp$  types, respectively. Although the new rules seem quite simple, they may increase the uncertainty of polymorphic instantiations. For example, the subtyping judgment

$$\forall a. a \rightarrow a \leq \perp \rightarrow \top$$

accepts any well-formed instantiation on the polymorphic type  $\forall a. a \rightarrow a$ .

**DECLARATIVE TYPING** The declarative typing rules, shown in Figure 5.3, extend DK’s higher-ranked type system in order to support the top and bottom types. Rule  $\text{Decl}\top$  allows any well-formed expression to check against  $\top$ . Rule  $\text{Decl}\perp\text{App}$  returns the  $\perp$  type when a function of  $\perp$  type is applied to any argument. All other rules remain exactly the same as the system of Chapters 3 and 4.

It’s worth mentioning that the design of the two new rules is driven by the subsumption property described in Section 5.4.1. They maintain the property in presence of a more powerful declarative subtyping, and we will discuss further later in that part.



### 5.3 BACKTRACKING ALGORITHM

#### 5.3.1 SYNTAX

Existential variables	$\hat{\alpha}, \hat{\beta}$
Types	$A, B, C ::= 1 \mid \top \mid \perp \mid a \mid \forall a. A \mid A \rightarrow B \mid \hat{\alpha}$
Algorithmic judgment chain	$\omega ::= A \leq B \mid e \Leftarrow A \mid e \Rightarrow_a \omega \mid A \bullet e \Rightarrow_a \omega$
Algorithmic worklist	$\Gamma ::= \cdot \mid \Gamma, a \mid \Gamma, \hat{\alpha} \mid \Gamma \vdash \omega$

Figure 5.4: Algorithmic Syntax

The algorithmic syntax is shown in Figure 5.4. Compared with the declarative system, existential variables  $\hat{\alpha}, \hat{\beta}$  are used as placeholders for unsolved mono-types. The judgment chain  $\omega$  and worklist context  $\Gamma$  are defined in the same way as that of Chapter 4.

The well-formedness definition for worklist context remains the same as that of Chapter 4. The hole notation is also inherited.

**THE SOLVE NOTATION** We define a set of auxiliary substitution functions in the form  $\{\hat{\alpha} := \tau_A\}$  to improve readability of our algorithmic definitions. Basically it refers to a global substitution when solving an existential variable.

A detailed definition of the notation is shown as follows:

$$\begin{aligned}
 \{\hat{\alpha} := \tau\} (\Gamma_L, \hat{\alpha}, \Gamma_R) &= \Gamma_L, [\tau/\hat{\alpha}] \Gamma_R & \tau &= 1, \top, \perp \text{ or } a \\
 \{\hat{\alpha} := \hat{\beta}\} (\Gamma_L, \hat{\alpha}, \Gamma_R) &= \Gamma_L, [\hat{\beta}/\hat{\alpha}] \Gamma_R & \hat{\beta} &\in \Gamma_L \\
 \{\hat{\alpha} := \hat{\alpha}_1 \rightarrow \hat{\alpha}_2\} (\Gamma_L, \hat{\alpha}, \Gamma_R) &= \Gamma_L, \hat{\alpha}_1, \hat{\alpha}_2, [\hat{\alpha}_1 \rightarrow \hat{\alpha}_2/\hat{\alpha}] \Gamma_R
 \end{aligned}$$

#### 5.3.2 ALGORITHMIC SUBTYPING

Figure 5.5 describes the algorithmic rules for subtyping. The relation is stated in a small-step “reduction” form, i.e. in each step, the worklist is analysed from the right-hand-side and reduced according to the top judgment. The overall procedure succeeds iff the worklist eventually reduces to  $\cdot$  (the empty worklist).

We categorize them into 7 groups according to their behavior:

1. Rules 1-3 are basic garbage collection rules. Given that the worklist  $\Gamma$  is well-formed, no reference of a variable should occur before its declaration. Therefore removing the declaration in the top position does not break well-formedness.

$\boxed{\Gamma \longrightarrow \Gamma'}$   $\Gamma$  reduces to  $\Gamma'$ .

$$\begin{aligned}
 & \Gamma, a \longrightarrow_1 \Gamma \\
 & \Gamma, \hat{\alpha} \longrightarrow_2 \Gamma \\
 & \Gamma, x : A \longrightarrow_3 \Gamma \\
 & \Gamma \Vdash 1 \leq 1 \longrightarrow_4 \Gamma \\
 & \Gamma \Vdash a \leq a \longrightarrow_5 \Gamma \\
 & \Gamma \Vdash A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2 \longrightarrow_6 \Gamma \Vdash A_2 \leq B_2 \Vdash B_1 \leq A_1 \\
 & \Gamma \Vdash \forall a. A \leq B \longrightarrow_7 \Gamma, \hat{\alpha} \Vdash [\hat{\alpha}/a]A \leq B \\
 & \Gamma \Vdash A \leq \forall b. B \longrightarrow_8 \Gamma, b \Vdash A \leq B \\
 & \Gamma \Vdash A \leq \top \longrightarrow_9 \Gamma \\
 & \Gamma \Vdash \perp \leq B \longrightarrow_{10} \Gamma \\
 & \Gamma[\hat{\alpha}] \Vdash \hat{\alpha} \leq \hat{\alpha} \longrightarrow_{11} \Gamma \\
 & \Gamma[\hat{\alpha}] \Vdash A \leq \hat{\alpha} \longrightarrow_{12} \{\hat{\alpha} := \top\} \Gamma[\hat{\alpha}] \\
 & \Gamma[\hat{\alpha}] \Vdash \hat{\alpha} \leq A \longrightarrow_{13} \{\hat{\alpha} := \perp\} \Gamma[\hat{\alpha}] \\
 & \Gamma[\hat{\alpha}] \Vdash \hat{\alpha} \leq A \rightarrow B \longrightarrow_{14} \{\hat{\alpha} := \hat{\alpha}_1 \rightarrow \hat{\alpha}_2\} (\Gamma[\hat{\alpha}] \Vdash \hat{\alpha} \leq A \rightarrow B) \\
 & \hspace{10em} \text{when } \hat{\alpha} \notin FV(A \rightarrow B) \\
 & \Gamma[\hat{\alpha}] \Vdash A \rightarrow B \leq \hat{\alpha} \longrightarrow_{15} \{\hat{\alpha} := \hat{\alpha}_1 \rightarrow \hat{\alpha}_2\} (\Gamma[\hat{\alpha}] \Vdash A \rightarrow B \leq \hat{\alpha}) \\
 & \hspace{10em} \text{when } \hat{\alpha} \notin FV(A \rightarrow B) \\
 & \Gamma[a][\hat{\beta}] \Vdash a \leq \hat{\beta} \longrightarrow_{16} \{\hat{\beta} := a\} \Gamma[a][\hat{\beta}] \\
 & \Gamma[a][\hat{\beta}] \Vdash \hat{\beta} \leq a \longrightarrow_{17} \{\hat{\beta} := a\} \Gamma[a][\hat{\beta}] \\
 & \Gamma[\hat{\beta}] \Vdash 1 \leq \hat{\beta} \longrightarrow_{18} \{\hat{\beta} := 1\} \Gamma[\hat{\beta}] \\
 & \Gamma[\hat{\beta}] \Vdash \hat{\beta} \leq 1 \longrightarrow_{19} \{\hat{\beta} := 1\} \Gamma[\hat{\beta}] \\
 & \Gamma[\hat{\alpha}][\hat{\beta}] \Vdash \hat{\alpha} \leq \hat{\beta} \longrightarrow_{20} \{\hat{\beta} := \hat{\alpha}\} \Gamma[\hat{\alpha}][\hat{\beta}] \\
 & \Gamma[\hat{\alpha}][\hat{\beta}] \Vdash \hat{\beta} \leq \hat{\alpha} \longrightarrow_{21} \{\hat{\beta} := \hat{\alpha}\} \Gamma[\hat{\alpha}][\hat{\beta}]
 \end{aligned}$$

Figure 5.5: Algorithmic Garbage Collection and Subtyping

An existential variable that is unsolved in the top position indicates that it is not constrained, thus picking any well-formed mono-type as its solution is acceptable. In our algorithmic formalization, we simply drop the existential variable.

2. Rules 4-10 directly correspond to the declarative subtyping rules. With no top-level existential variables, there is nothing to guess immediately, and thus the algorithm behaves just like the declarative system.
3. Rule 11 is a base case in the algorithmic system. The declarative reflexivity property suggests that any solution is acceptable, thus the judgment holds without any constraint.
4. Rules 12-13 are important rules that require backtracking techniques for implementation. These rules *overlap* with all the remaining rules when solving an existential variable. In other words, they simply try if  $\top$  or  $\perp$  satisfies the constraints in parallel with other possibilities.
5. Rules 14-15 compares an existential variable  $\hat{\alpha}$  with a function type, resulting in solving the  $\hat{\alpha}$  by  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$ . The freshness condition rules out the possibility when there is a cyclic dependency. For example, the judgment

$$\hat{\alpha} \leq 1 \rightarrow \hat{\alpha}$$

is satisfied with either of these solutions to  $\hat{\alpha}$ :

$$\perp, 1 \rightarrow \perp, \top \rightarrow \perp, 1 \rightarrow 1 \rightarrow \perp, \dots$$

However, we argue that comparing  $\hat{\alpha}$  with a function type that contains  $\hat{\alpha}$  itself is hardly useful in practice, and most of the solutions are meaningless. Therefore, in our algorithm, only the  $\perp$  solution is considered with Rule 13. The condition of Rule 14 rejects the judgment for further analysis and thus does not produce more solutions. This is one source of incompleteness of our algorithm with respect to the declarative specification.

6. Rules 16-19 solve existential variables against a type variable or the unit type. For instance, the judgment

$$\hat{\alpha} \leq 1$$

only has two solutions:  $\hat{\alpha} = 1$  or  $\hat{\alpha} = \perp$ . In similar cases, one of the solutions is produced by Rule 12 or 13, and the other one is given by one of the rules in this group. Additional well-formedness checks are performed when type variables are encountered; a solution to an existential variable must be well-formed in the context before the existential variable is defined.

7. Rules 20-21 deal with subtyping judgments that compare two different existential variables. The only difference between them is the variable order. Similar to the type variable case in the previous group, existential variables must solve to another one defined earlier. With Rules 12, 13, 20 and 21, a judgment like

$$\hat{\alpha} \leq \hat{\beta}$$

could possibly give any of the following solutions:

$$\hat{\alpha} = \hat{\beta} \text{ (or } \hat{\beta} = \hat{\alpha}) \text{ or } \hat{\alpha} = \perp \text{ or } \hat{\beta} = \top$$

Those are good attempts, but unfortunately, they do not cover the complete set of possibilities. The following example worklist

$$\hat{\alpha}, \hat{\beta} \Vdash \hat{\beta} \leq 1 \rightarrow 1 \Vdash \hat{\alpha} \leq \hat{\beta}$$

has a solution  $\hat{\alpha} = 1 \rightarrow \top, \hat{\beta} = 1 \rightarrow 1$  missed by our algorithm. Similar situations happen when the judgments are specifically ordered; if  $\hat{\beta} \leq 1 \rightarrow 1$  is the top-most judgment, the algorithm will not miss this solution.

Although such treatment for existential variable solving is incomplete in theory, the algorithm completely handles practical programs where unification is simply equality. For other programs that exploit complex guessing involving subtyping, the programmer may put type annotations when the type inference algorithm does not find the optimal solution.

### 5.3.3 ALGORITHMIC TYPING

The algorithmic typing rules are split into three groups, according to the category of the top-most judgment.

1. Checking mode. Rules 22-27 reduce top-level checking judgments. Rules 22, 23, 24 and 26 directly reflect how the declarative system behaves. Rule 25 splits  $\hat{\alpha}$  into  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$  to mimic the same logic as Rule 24, since a lambda expression must be of a function type. Rule 27 is another backtracking rule, which overlaps with all other checking rules. Solving  $\hat{\alpha}$  to  $\top$  in such case reflects the declarative rule  $\text{Dec1}\top$  when the unknown type is  $\top$ . Furthermore, we eliminate the algorithmic counterpart of declarative rule  $\text{Dec11I}$ , because a combination of Rules 22 and 30 already accepts the judgment  $() \Leftarrow 1$ .

$\boxed{\Gamma \longrightarrow \Gamma'}$   $\Gamma$  reduces to  $\Gamma'$  (continued).

$$\begin{aligned}
 & \Gamma \Vdash e \Leftarrow B \longrightarrow_{22} \Gamma \Vdash e \Rightarrow_a a \leq B \quad \text{when } e \neq \lambda x. e' \text{ and } B \neq \forall a. B' \\
 & \Gamma \Vdash e \Leftarrow \forall a. A \longrightarrow_{23} \Gamma, a \Vdash e \Leftarrow A \\
 & \Gamma \Vdash \lambda x. e \Leftarrow A \rightarrow B \longrightarrow_{24} \Gamma, x : A \Vdash e \Leftarrow B \\
 & \Gamma[\hat{\alpha}] \Vdash \lambda x. e \Leftarrow \hat{\alpha} \longrightarrow_{25} \{\hat{\alpha} := \hat{\alpha}_1 \rightarrow \hat{\alpha}_2\} \Gamma, x : \hat{\alpha}_1 \Vdash e \Leftarrow \hat{\alpha}_2 \\
 & \Gamma \Vdash e \Leftarrow \top \longrightarrow_{26} \Gamma \\
 & \Gamma[\hat{\alpha}] \Vdash e \Leftarrow \hat{\alpha} \longrightarrow_{27} \{\hat{\alpha} := \top\} \Gamma[\hat{\alpha}] \\
 & \Gamma \Vdash x \Rightarrow_a \omega \longrightarrow_{28} \Gamma \Vdash [A/a]\omega \quad \text{when } (x : A) \in \Gamma \\
 & \Gamma \Vdash (e : A) \Rightarrow_a \omega \longrightarrow_{29} \Gamma \Vdash [A/a]\omega \Vdash e \Leftarrow A \\
 & \Gamma \Vdash () \Rightarrow_a \omega \longrightarrow_{30} \Gamma \Vdash [1/a]\omega \\
 & \Gamma \Vdash \lambda x. e \Rightarrow_a \omega \longrightarrow_{31} \Gamma, \hat{\alpha}, \hat{\beta} \Vdash [\hat{\alpha} \rightarrow \hat{\beta}/a]\omega, x : \hat{\alpha} \Vdash e \Leftarrow \hat{\beta} \\
 & \Gamma \Vdash e_1 e_2 \Rightarrow_a \omega \longrightarrow_{32} \Gamma \Vdash e_1 \Rightarrow_b (b \bullet e_2 \Rightarrow_a \omega) \\
 & \Gamma \Vdash \forall a. A \bullet e \Rightarrow_a \omega \longrightarrow_{33} \Gamma, \hat{\alpha} \Vdash [\hat{\alpha}/a]A \bullet e \Rightarrow_a \omega \\
 & \Gamma \Vdash A \rightarrow C \bullet e \Rightarrow_a \omega \longrightarrow_{34} \Gamma \Vdash [C/a]\omega \Vdash e \Leftarrow A \\
 & \Gamma \Vdash \perp \bullet e \Rightarrow_a \omega \longrightarrow_{35} \Gamma \Vdash [\perp/a]\omega \\
 & \Gamma[\hat{\alpha}] \Vdash \hat{\alpha} \bullet e \Rightarrow_a \omega \longrightarrow_{36} \{\hat{\alpha} := \hat{\alpha}_1 \rightarrow \hat{\alpha}_2\} \Gamma \Vdash \hat{\alpha}_1 \rightarrow \hat{\alpha}_2 \bullet e \Rightarrow_a \omega \\
 & \Gamma[\hat{\alpha}] \Vdash \hat{\alpha} \bullet e \Rightarrow_a \omega \longrightarrow_{37} \{\hat{\alpha} := \perp\} \Gamma[\hat{\alpha}] \Vdash [\perp/a]\omega
 \end{aligned}$$

Figure 5.6: Algorithmic Typing

2. Inference mode. Rules 28-32 reduce inference judgments. Like the work presented in Chapter 4, the encoding of the return type is by an explicit substitution on the binder of a judgment chain. Rules 28, 29, 30 and 32 correspond to the declarative system straightforwardly. Rule 31 guesses an unannotated lambda with a mono function type  $\hat{\alpha} \rightarrow \hat{\beta}$ . Writing in a slightly different way may help improve readability:

$$\Gamma \Vdash \lambda x. e \Rightarrow_a \omega \longrightarrow_{31'} \Gamma, \hat{\alpha}, \hat{\beta} \Vdash [\hat{\alpha} \rightarrow \hat{\beta}/a]\omega \Vdash \lambda x. e \Leftarrow \hat{\alpha} \rightarrow \hat{\beta}$$

Rule 32 illustrates how the judgment chain works, with a continuation-passing-style encoding of the type inference task.

3. Application inference mode. Rules 33-37 reduce application inference judgments. Each of these judgments accepts an input function type and an argument expression and produces the expected return type. Rules 33, 34 and 35 are direct translations from the declarative rules. Rule 33, specifically, enables implicit parametric polymorphism via existential variable solving. Rules 36 and 37 deal with cases where a single existen-

tial variable  $\hat{\alpha}$  behaves as a function type. Rule 36 splits  $\hat{\alpha}$  into a function unknown type  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$ . Rule 37 tries the solution  $\hat{\alpha} = \perp$  and returns the  $\perp$  type.

## 5.4 METATHEORY

In this section we present several properties that are formally verified. For the declarative system, the typing subsumption and subtyping transitivity lemmas are discussed in detail. The algorithmic system is proven to be sound with respect to the declarative system via a transfer relation. A partial completeness theorem is shown under the rank-1 restriction. We then briefly describe the challenges we face when proving termination. Lastly, proof statistics of Abella are discussed.

### 5.4.1 DECLARATIVE PROPERTIES

**THE TYPING SUBSUMPTION LEMMA.** An important desired property for a type system is *checking subsumption*, which says that any expression  $e$  can check against  $B$  if  $e$  checks against  $A$  and  $A \leq B$ . Since our bidirectional type system defines the checking mode, inference mode and application inference mode mutually, we formalize the generalized *typing subsumption*.

First of all, we give the definition of worklist subtyping, which is used to further generalize the typing subsumption lemma. This is necessary because rules like  $\text{Dec1} \rightarrow \text{I}$  will push the argument type  $A$  into the context, thus when checking against a super type of  $A \rightarrow B$ , say  $C \rightarrow D$ , will cause the bind of  $x$  in the context to a subtype of  $A$  (since  $C \leq A$ ).

**Definition 5** (Worklist Subtyping). Worklist subtyping compares the type of variables bound in the worklist.  $\Psi <: \Psi'$  iff each binding in  $\Psi$  is converted to one with a super type.

$$\begin{array}{c}
 \frac{}{\cdot <: \cdot} <: \text{nil} \qquad \frac{\Psi <: \Psi'}{\Psi, a <: \Psi', a} <: \text{ty} \\
 \\
 \frac{\Psi' \vdash A \leq B \quad \Psi <: \Psi'}{\Psi, x : A <: \Psi', x : B} <: \text{of} \qquad \frac{\Psi <: \Psi'}{\Psi \Vdash \omega <: \Psi' \Vdash \omega} <: \omega
 \end{array}$$

A basic property of worklist subtyping is that they acts similarly when dealing with subtyping between well-formed types.

**Lemma 5.1** (Worklist Subtyping Equivalence).

Given  $\Psi <: \Psi'$ ,  $\Psi \vdash A \leq B \iff \Psi' \vdash A \leq B$ .

Finally, we give the statement of typing subsumption lemma, which is generalized by the worklist subtyping relation.

**Lemma 5.2** (Typing Subsumption). *Given  $\Psi <: \Psi'$ ,*

- 1) *If  $\Psi' \vdash e \Leftarrow A$  and  $\Psi' \vdash A \leq B$ , then  $\Psi \vdash e \Leftarrow B$ ;*
- 2) *If  $\Psi' \vdash e \Rightarrow A$ , then  $\exists B$  s.t.  $\Psi' \vdash B \leq A$  and  $\Psi \vdash e \Rightarrow B$ .*
- 3) *If  $\Psi' \vdash C \bullet e \Rightarrow A$  and  $\Psi' \vdash D \leq C$ , then  $\exists B$  s.t.  $\Psi' \vdash B \leq A$  and  $\Psi \vdash D \bullet e \Rightarrow B$ .*

*Proof.* By induction on the following size measure (lexicographical order on a 3-tuple):

- Checking ( $e \Leftarrow A$ ):  $\langle |e|, 1, |A|_{\forall} + |B|_{\forall} \rangle$
- Inference ( $e \Rightarrow A$ ):  $\langle |e|, 0, 0 \rangle$
- Application inference ( $A \bullet e \Rightarrow C$ ):  $\langle |e|, 2, |C|_{\forall} + |D|_{\forall} \rangle$

Most of the cases are straightforward. When rule  $\leq \forall L$  is applied for the subtyping predicate like  $\Psi' \vdash A \leq B$ , a mono-type substitution is performed on  $\forall a. A$ , resulting in  $[\tau/a]A$ . Since  $\tau$  is a mono-type, the result type reduces the number of  $\forall$ 's, and thus reduces the size measure.  $\square$

Interestingly, the two new declarative rules  $\text{Dec1}\top$  and  $\text{Dec1}\perp\text{App}$  are discovered when we were trying to prove the property instead of before exploring the meta-theory. Given the typing and subtyping judgments  $\Psi \vdash e \Leftarrow A$  and  $\Psi \vdash A \leq \top$ , we should derive  $\Psi \vdash e \Leftarrow \top$  from the lemma, therefore Rule  $\text{Dec1}\top$  is required, saying that any expression can be checked against the top type. Note that here we do not include the condition  $e \Leftarrow A$ , since we may elaborate the term to the top term which no longer refers to the original term.

Similarly, the most general type  $\perp$ , being able to convert to any type due to Rule  $\leq \text{Bot}$ , can be converted to any function type, or simply the most general one  $\top \rightarrow \perp$ , which accepts any input and returns the  $\perp$  type, resulting in the derivation  $\Psi \vdash \perp \bullet e \Rightarrow \perp$ . From the lemma we can also derive that by  $\Psi \vdash C \bullet e \Rightarrow A$ ,  $\Psi \vdash \perp \leq C$  and  $\Psi \vdash \perp \leq A$ .

With the addition of Rules  $\text{Dec1}\top$  and  $\text{Dec1}\perp\text{App}$ , we can prove the typing subsumption lemma. To the best of the author's knowledge, they are the minimal set of rules that make the lemma hold.

**TRANSITIVITY OF SUBTYPING** The transitivity lemma for declarative subtyping is a commonly expected property. The proof depends on the following subtyping derivation size relation and an auxiliary lemma.

**Definition 6** (Subtyping Derivation Size).

$$\begin{aligned}
 |1 \leq 1| &= 0 \\
 |a \leq a| &= 0 \\
 |A \leq \top| &= 0 \\
 |\perp \leq B| &= 0 \\
 |A_1 \rightarrow A_2 \leq B_1 \rightarrow B_2| &= |B_1 \leq A_1| + |A_2 \leq B_2| + 1 \\
 |\forall a. A \leq B| &= |[\tau/a]A \leq B| + 1 \\
 |A \leq \forall a. B| &= |A \leq B| + 1
 \end{aligned}$$

**Lemma 5.3** (Monotype Subtyping Substitution). *If  $\Psi \vdash \tau$  and  $\Psi, a, \Psi_R \vdash A \leq B$ , then  $\Psi, [\tau/a]\Psi_R \vdash [\tau/a]A \leq [\tau/a]B$ .*

*Proof.* A routine induction on the subtyping relation  $\Psi, a, \Psi_R \vdash A \leq B$  finishes the proof.  $\square$

**Corollary 5.4** (Monotype Subtyping Substitution for Type Variables). *If  $\Psi \vdash \tau$  and  $\Psi, a \vdash A \leq B$ , then  $\Psi \vdash [\tau/a]A \leq [\tau/a]B$ .*

The above lemma and corollary reveal the fact that a type variable occurred in the subtyping relation represents an *arbitrary* well-formed monotype. And it also explains the difference in treatment of polymorphic types between Rules  $\leq \forall L$  and  $\leq \forall R$ : Rule  $\leq \forall R$  is in fact equivalent to:

$$\frac{\forall \tau \text{ s.t. } \Psi \vdash \tau \implies \Psi \vdash A \leq [\tau/b]B}{\Psi \vdash A \leq \forall b. B} \leq \forall R'$$

Finally, with the size measure defined and required lemma proven, we can obtain the transitivity lemma for declarative subtyping.

**Lemma 5.5** (Subtyping Transitivity). *If  $\Psi \vdash A \leq B$  and  $\Psi \vdash B \leq C$  then  $\Psi \vdash A \leq C$ .*

*Proof.* Induction on the lexicographical order defined by  $\langle |B|_{\forall}, |A \leq B| + |B \leq C| \rangle$ . Most cases preserve the first element of the size measures  $|B|_{\forall}$ , and are relatively easy to prove. The difficult case is when  $B$  is a polymorphic type, when the conditions are  $\Psi \vdash A \leq \forall a. B$  and  $\Psi \vdash \forall a. B \leq C$ . They are derived through rules  $\leq \forall L$  and  $\leq \forall R$ , respectively. Therefore, we have  $\Psi, a \vdash A \leq B$  and  $\Psi \vdash [\tau/a]B \leq C$ . To exploit the induction hypothesis, the contexts should be unified. By Corollary 5.4,  $\Psi \vdash A \leq [\tau/a]B$ . Notice that the freshness condition is implicit for rule  $\leq \forall L$ . Clearly,  $|[\tau/a]B|_{\forall} < |\forall a. B|_{\forall}$ , i.e. the first size measure decreases. By induction hypothesis we get  $\Psi \vdash A \leq C$  and finish this case.  $\square$



$$\begin{array}{l}
 \text{Declarative worklist} \quad \Omega ::= \cdot \mid \Omega, a \mid \Omega, x : A \mid \Omega \Vdash \omega \\
 \\
 \boxed{\Gamma \rightsquigarrow \Omega} \quad \Gamma \text{ instantiates to } \Omega. \\
 \\
 \frac{}{\Omega \rightsquigarrow \Omega} \rightsquigarrow \Omega \quad \frac{\Omega \vdash \tau \quad \Omega, [\tau/\hat{\alpha}]\Gamma \rightsquigarrow \Omega}{\Omega, \hat{\alpha}, \Gamma \rightsquigarrow \Omega} \rightsquigarrow \hat{\alpha}
 \end{array}$$

Figure 5.7: Declarative Worklists and Instantiation

#### 5.4.2 TRANSFER

Follow the approach of Section 4.3, the transfer relation and the declarative instantiation relation are defined in Figure 5.7.

Similarly, Lemmas 5.6 and 5.7 generalizing Rule  $\rightsquigarrow \hat{\alpha}$  hold as well.

**Lemma 5.6 (Insert).** *If  $\Gamma_L, [\tau/\hat{\alpha}]\Gamma_R \rightsquigarrow \Omega$  and  $\Gamma_L \vdash \tau$ , then  $\Gamma_L, \hat{\alpha}, \Gamma_R \rightsquigarrow \Omega$ .*

**Lemma 5.7 (Extract).** *If  $\Gamma_L, \hat{\alpha}, \Gamma_R \rightsquigarrow \Omega$ , then there exists  $\tau$  s.t.  $\Gamma_L \vdash \tau$  and  $\Gamma_L, [\tau/\hat{\alpha}]\Gamma_R \rightsquigarrow \Omega$ .*

$$\begin{array}{l}
 \boxed{\|\Omega\|} \quad \text{Judgment erasure.} \\
 \\
 \begin{array}{l}
 \|\cdot\| = \cdot \\
 \|\Omega, a\| = \|\Omega\|, a \\
 \|\Omega, x : A\| = \|\Omega\|, x : A \\
 \|\Omega \Vdash \omega\| = \|\Omega\|
 \end{array} \\
 \\
 \boxed{\Omega \longrightarrow \Omega'} \quad \text{Declarative transfer.} \\
 \\
 \begin{array}{ll}
 \Omega, a \longrightarrow \Omega & \\
 \Omega, x : A \longrightarrow \Omega & \\
 \Omega \Vdash A \leq B \longrightarrow \Omega & \text{when } \|\Omega\| \vdash A \leq B \\
 \Omega \Vdash e \Leftarrow A \longrightarrow \Omega & \text{when } \|\Omega\| \vdash e \Leftarrow A \\
 \Omega \Vdash e \Rightarrow_a \omega \longrightarrow \Omega \Vdash [A/a]\omega & \text{when } \|\Omega\| \vdash e \Rightarrow A \\
 \Omega \Vdash A \bullet e \Rightarrow_a \omega \longrightarrow \Omega \Vdash [C/a]\omega & \text{when } \|\Omega\| \vdash A \bullet e \Rightarrow C
 \end{array}
 \end{array}$$

Figure 5.8: Declarative Transfer

Figure 5.8 defines a relation  $\Omega \longrightarrow \Omega'$ , checking that every judgment entry in the worklist holds using a corresponding declarative judgment.

## 5.4.3 SOUNDNESS

Our algorithm is sound with respect to the declarative system. For any worklist  $\Gamma$  that reduces successfully, there is a valid instantiation  $\Omega$  that transfers all judgments to the declarative system.

**Theorem 5.8** (Soundness). *If  $\text{wf}\Gamma$  and  $\Gamma \longrightarrow^* \cdot$ , then there exists  $\Omega$  s.t.  $\Gamma \rightsquigarrow \Omega$  and  $\Omega \longrightarrow^* \cdot$ .*

Soundness is a basic desired property of a type inference algorithm, which ensures that the algorithm is always producing valid declarative derivations when the judgments are accepted.

## 5.4.4 PARTIAL COMPLETENESS OF SUBTYPING: RANK-1 RESTRICTION

The algorithm is incomplete due to the subtyping rules 14, 15, 20 and 21. However, subtyping is complete with respect to the declarative system in a rank-1 setting.

**DECLARATIVE RANK-1 RESTRICTION** Rank-1 types are also named type schemes in the Hindley-Milner type system.

$$\text{Declarative Type Schemes} \quad \sigma ::= \forall a. \sigma \mid \tau$$

In other words, the universal quantifiers only appear in the top level of all polymorphic types.

For declarative subtyping, a judgment must be of form  $\sigma_1 \leq \sigma_2$ .

## 5.4.5 ALGORITHMIC RANK-1 RESTRICTION (PARTIAL COMPLETENESS)

The algorithmic mono-types and type schemes are defined as following:

$$\begin{aligned} \text{Algorithmic Mono-types} \quad \tau_A &::= 1 \mid \top \mid \perp \mid a \mid A \rightarrow B \mid \hat{\alpha} \\ \text{Algorithmic Type Schemes} \quad \sigma_A &::= \forall a. \sigma_A \mid \tau_A \end{aligned}$$

Starting from the declarative judgment  $\sigma_1 \leq \sigma_2$ , the algorithmic derivation might involve different other kinds of judgments. The following derivation, as an example, shows how a rank-1 judgment derives.

$$\begin{aligned} &\cdot \Vdash \forall a. a \rightarrow a \leq \forall b. (b \rightarrow b) \rightarrow (b \rightarrow b) \\ \longrightarrow_8 &b \Vdash \forall a. a \rightarrow a \leq (b \rightarrow b) \rightarrow (b \rightarrow b) \\ \longrightarrow_7 &b, \hat{\alpha} \Vdash \hat{\alpha} \rightarrow \hat{\alpha} \leq (b \rightarrow b) \rightarrow (b \rightarrow b) \\ \longrightarrow_6 &b, \hat{\alpha} \Vdash \hat{\alpha} \leq b \rightarrow b \Vdash b \rightarrow b \leq \hat{\alpha} \\ \longrightarrow &\dots \end{aligned}$$

In this derivation, we begin from a judgment of the form  $\sigma \leq \sigma$ . After Rule 8 is applied, the judgment becomes  $\sigma \leq \tau$ , since the right-hand-side polymorphic type is reduced to a declarative mono-type. Then, Rule 7 introduces existential variables to the left-hand-side, resulting in a judgment like  $\tau_A \leq \tau$ , or  $\sigma_A \leq \tau$  in a more general case. Finally, Rule 6 breaks a judgment between functions into two sub-judgments, which swaps the positions of the argument types and creates a judgment like  $\tau \leq \tau_A$ . Notice that  $\sigma_A$  is not possible to occur to the right because the function type may not contain any polymorphic types as its argument type.

After a detailed analysis on the judgments derivations, we found that the only possible judgments that a rank-1 declarative subtyping judgment might step to belong to the following two categories:

$$\sigma_A \leq \sigma \quad \text{or} \quad \tau \leq \sigma_A$$

All the possible judgment types shown above fall into these categories. For example,  $\tau_A \leq \tau$  is a special form of  $\sigma_A \leq \sigma$ , and  $\tau \leq \tau_A$  belongs to  $\tau \leq \sigma_A$ .

An interesting observation is that  $\hat{\alpha} \leq \hat{\beta}$  does not belong to either category, neither does  $\hat{\alpha} \leq A \rightarrow B$  when  $\hat{\alpha} \in \text{FV}(A \rightarrow B)$ . Therefore, in the rank-1 setting, both cases of incompleteness never occur, and our algorithm is complete.

**Theorem 5.9** (Completeness of Rank-1 Subtyping). *Given  $\Psi \vdash \sigma_1 \leq \sigma_2$ ,*

- *If  $\Gamma \Vdash \sigma_A \leq \sigma \rightsquigarrow \Psi \vdash \sigma_1 \leq \sigma_2$   
then  $\Gamma \Vdash \sigma_A \leq \sigma \longrightarrow^* \cdot$ ;*
- *If  $\Gamma \Vdash \tau \leq \sigma_A \rightsquigarrow \Psi \vdash \sigma_1 \leq \sigma_2$   
then  $\Gamma \Vdash \tau \leq \sigma_A \longrightarrow^* \cdot$ .*

#### 5.4.6 TERMINATION

The measure used in Chapter 4 no longer works because subtyping judgments like

$$\hat{\alpha} \leq \perp \rightarrow \top$$

cause  $\hat{\alpha}$  to split into  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$ , without solving any part of it, resulting in an increased number of existential variables and possibly increased complexity of the worklist through the size-increasing substitution  $\{\hat{\alpha} := \hat{\alpha}_1 \rightarrow \hat{\alpha}_2\}$ .

We have performed a large set of tests on generated subtyping judgments that are consist of algorithmic monotypes, and all judgments terminated within a reasonable number of derivation depth. Unfortunately, we have not yet find any formal proof for the termination statement.

Table 5.1: Statistics for the proof scripts

File(s)	LOC	#Thm	Description
olist.thm, nat.thm	311	57	Basic data structures
typing.thm	273	7	Declarative & algorithmic system, debug examples
decl.thm	241	33	Basic declarative properties
order.thm	274	27	The $ \cdot _{\forall}$ measure; decl. subtyping strengthening
alg.thm	699	82	Basic algorithmic properties
trans.thm	635	53	Worklist instantiation and declarative transfer; Lemmas 4.3, 4.4
declTyping.thm	1,087	76	Non-overlapping declarative system; Lemmas 4.5, 4.6, 4.7, 4.8
soundness.thm	1,206	81	Soundness theorem; aux. lemmas on transfer
dcl.thm	417	12	Non-overlapping declarative worklist
scheme.thm	1,113	98	Type scheme (rank-1 restriction)
completeness.thm	1,045	63	Completeness theorem; aux. lemmas and relations
<i>Total</i>	7,301	592	(48 definitions in total)

#### 5.4.7 FORMALIZATION IN THE ABELLA PROOF ASSISTANT

We have chosen the Abella (v2.0.7-dev<sup>1</sup>) proof assistant [Gacek 2008] to develop our formalization. Equipped with HOAS, Abella eases the formalization and proof tasks a lot compared with various libraries in Coq. Additionally, our algorithm heavily uses eager substitutions, and Abella greatly simplifies relevant proofs thanks to its built-in substitution representation and higher-order unification algorithms.

The reader may find the source code of the proof at <https://github.com/JimmyZJX/Dissertation/tree/main/src/Subtyping>.

**STATISTICS OF THE PROOF** The proof script consists of 7,301 lines of Abella code with a total of 48 definitions and 592 theorems. Figure 5.1 briefly summarizes the contents of each file. The files are linearly dependent due to the limitations of Abella.

<sup>1</sup>We use a forked version <https://github.com/JimmyZJX/abella> by only enhancing the Abella prover with a handy “aplys” tactic.

## 5.5 DISCUSSION

### 5.5.1 A COMPLETE ALGORITHM UNDER MONOTYPE GUESSING RESTRICTIONS

The incompleteness of the algorithm is mainly due to incomplete guesses for existential variable instantiations when dealing with subtyping. For example, when reducing the algorithmic subtyping judgment

$$\hat{\alpha} \leq \hat{\alpha} \rightarrow 1$$

we only consider the possibility  $\hat{\alpha} := \perp$  and stop reducing because the judgment involves a self-reference. In fact, there are infinitely many valid solutions just in our type system, including  $\hat{\alpha} := \top \rightarrow \perp$  and  $\hat{\alpha} := \top \rightarrow 1$ . The problem can be solved by introducing recursive types, which may be used to represent a whole class of solutions, and recursive constraints can be solved with a single recursive type. However, the other type of incompleteness that comes from

$$\hat{\alpha} \leq \hat{\beta}$$

is still not addressed.

We propose an alternative approach that requires little change to the type system by restricting the declarative system. If the declarative system did not treat the bottom and top types as monotypes in the first place, the algorithm would not need to guess those types. Formally, if monotypes are defined as follows,

$$\text{Monotypes}' \quad \tau' ::= 1 \mid a \mid \tau_1 \rightarrow \tau_2$$

then subtyping relation between monotypes is simply equality.

The incomplete examples mentioned above can be addressed now, judgments like  $\hat{\alpha} \leq \hat{\alpha} \rightarrow 1$  can no longer instantiate  $\hat{\alpha}$  to a function type, otherwise it will infinitely loop without emitting any valid solution; judgments like  $\hat{\alpha} \leq \hat{\beta}$  can now be safely treated as  $\hat{\alpha} = \hat{\beta}$ , because the instantiations of  $\hat{\alpha}$  and  $\hat{\beta}$  must be equal to each other. Formally speaking, we proved the following lemma

**Lemma 5.10** (Subtyping between Monotypes is Equality).

*If  $\Psi \vdash \tau'_1 \leq \tau'_2$ , then  $\tau'_1 = \tau'_2$ .*

which also holds for the system in Chapter 4, yet it does not hold for the system with subtyping when monotypes are not restricted.

On the meantime, algorithmic subtyping rules 12 and 13, and typing rules 27 and 37 should be removed from the algorithm under the monotype restriction, since these rules solve existential variables to  $\top$  or  $\perp$ . Other than that, the algorithm remains unchanged.

Note that after removing these rules, the algorithm no longer has overlapping rules, therefore backtracking is not needed any more.

Following the proof techniques in Chapter 4, both *soundness* and *completeness* theorems are successfully proven and also mechanized. Given that the algorithm reduces judgments in a similar way to the one in the previous chapter, the result is not so surprising. However, it still reveals two facts: firstly, if we do not try to guess any type involving subtyping relation (the top and bottom types in our case), the algorithm should remain simple and complete; secondly, although the current algorithm for the system with subtyping is incomplete, all source of incompleteness comes from subtyping relations.

The source code of the proof can be found at [https://github.com/JimmyZJX/Dissertation/tree/main/src/Subtyping\\_mono\\_restriction](https://github.com/JimmyZJX/Dissertation/tree/main/src/Subtyping_mono_restriction).

### 5.5.2 LAZY SUBSTITUTION AND NON-TERMINATING LOOPS

In this subsection, we propose a possible way to fix the incompleteness issue for the subtyping part of the algorithm and discuss its impact to the existing formalization.

As we have analysed before, one major source of incompleteness comes from Rules 20 and 21, where two different existential variables are compared. A natural idea is to delay the solving procedure until all the corresponding constraints are collected. Constraint collection can be implemented by the following new worklist definition, where the existential variables also come with a set of upper bounds and a set of lower bounds. Both of the bounds consist of algorithmic mono-types that does not contain  $\hat{\alpha}$  (similar to the occurs-check condition).

$$\text{Algorithmic worklist} \quad \Gamma ::= \cdot \mid \Gamma, a \mid \Gamma, \overline{\tau_A} \leq \hat{\alpha} \leq \overline{\tau_A} \mid \Gamma \Vdash \omega$$

**DIFFERENCES OF ALGORITHMIC RULES.** In presence of the bound collections, the algorithm behaves slightly differently on variable solving rules. Firstly, rules that eagerly solve existential variables to the top or bottom type are discarded (Rules 12 and 13). Next, Rules 20 and 21 add the corresponding existential variable as part of a bound rather than performing a global substitution.

$$\begin{aligned} \Gamma[\hat{\alpha}][L \leq \hat{\beta} \leq U] \Vdash \hat{\alpha} \leq \hat{\beta} &\longrightarrow_{20'} \Gamma[\hat{\alpha}][L \cup \{\hat{\alpha}\} \leq \hat{\beta} \leq U] \\ \Gamma[\hat{\alpha}][L \leq \hat{\beta} \leq U] \Vdash \hat{\beta} \leq \hat{\alpha} &\longrightarrow_{21'} \Gamma[\hat{\alpha}][L \leq \hat{\beta} \leq U \cup \{\hat{\alpha}\}] \end{aligned}$$

Instantiation rules (Rules 14 and 15) also behave lazily in order not to lose completeness.

$$\begin{aligned}
 \Gamma[L \leq \hat{\alpha} \leq U] \Vdash \hat{\alpha} \leq A \rightarrow B &\longrightarrow_{14'} \\
 \Gamma[\hat{\alpha}_1, \hat{\alpha}_2, L \leq \hat{\alpha} \leq U \cup \{\hat{\alpha}_1 \rightarrow \hat{\alpha}_2\}] \Vdash \hat{\alpha}_1 \rightarrow \hat{\alpha}_2 \leq A \rightarrow B \\
 \Gamma[L \leq \hat{\alpha} \leq U] \Vdash A \rightarrow B \leq \hat{\alpha} &\longrightarrow_{15'} \\
 \Gamma[\hat{\alpha}_1, \hat{\alpha}_2, L \cup \{\hat{\alpha}_1 \rightarrow \hat{\alpha}_2\} \leq \hat{\alpha} \leq U] \Vdash A \rightarrow B \leq \hat{\alpha}_1 \rightarrow \hat{\alpha}_2
 \end{aligned}$$

These new rules differ from the original ones: the two fresh existential variables representing a function type  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2$  now acts as a bridge between  $\hat{\alpha}$  and  $A \rightarrow B$ . The algorithm breaks the judgment  $\hat{\alpha} \leq \hat{\alpha}_1 \rightarrow \hat{\alpha}_2$  into two sub-problems:  $\hat{\alpha} \leq \hat{\alpha}_1 \rightarrow \hat{\alpha}_2$  and  $\hat{\alpha}_1 \rightarrow \hat{\alpha}_2 \leq A \rightarrow B$ . And it is not difficult to prove that the problems before and after such transformation are equivalent to each other.

The change of existential variable declaration and solving mechanisms also result in the change of the algorithmic rule for garbage collecting them. When  $\hat{\alpha}$  is at the top of the worklist and thus going to be recycled, we need to further check if it can actually be solved, by ensuring that lower bounds are indeed subtypes of upper bounds: any type that is a subtype of  $\hat{\alpha}$  should be a subtype of any super type of  $\hat{\alpha}$ . Note that the bounds should not contain  $\hat{\alpha}$ .

$$\Gamma, \{l_i\}^n \leq \hat{\alpha} \leq \{u_j\}^m \longrightarrow_{2'} \Gamma \Vdash_{1 \leq i \leq n, 1 \leq j \leq m} l_i \leq u_j$$

The notion  $\{l_i\}^n$  indicates that the bound collection is of length  $n$ ; the notion

$$\Vdash_{1 \leq i \leq n, 1 \leq j \leq m} l_i \leq u_j$$

is similar to the syntax of list-comprehension, which creates a list of judgments

$$\Vdash l_1 \leq u_1 \cdots \Vdash l_1 \leq u_m \cdots \Vdash l_2 \leq u_1 \cdots \Vdash l_n \leq u_m$$

According to the transitivity of subtyping, this rule is sound: any valid instantiation of  $\hat{\alpha}$ , denoted  $\tau_{\hat{\alpha}}$ , indicates that the instantiation of  $l_i \leq \tau_{\hat{\alpha}}$  and  $\tau_{\hat{\alpha}} \leq u_j$  holds, therefore the instantiation of  $l_i \leq u_j$  holds for any  $1 \leq i \leq n, 1 \leq j \leq m$ .

On the other hand, such check is sufficient. In other words, we can always find a proper solution for  $\hat{\alpha}$  if any lower bound is a subtype of any upper bound. This is supported by the following property on declarative monotypes:

**Lemma 5.11** (Sufficient Bound Check for Monotypes).

*Given collections of monotypes  $\{l_i\}^n$  and  $\{u_j\}^m$ ,*

*if  $\Psi \vdash l_i \leq u_j$  holds for any  $1 \leq i \leq n, 1 \leq j \leq m$ ,*

*then there exists  $\tau$  s.t.  $\Psi \vdash l_i \leq \tau$  and  $\Psi \vdash \tau \leq u_i$  for any  $1 \leq i \leq n, 1 \leq j \leq m$ .*

IMPACT OF THE NEW ALGORITHM. The new algorithm addresses the incompleteness issue thoroughly by re-designing in a lazy-substitution style. However, the algorithm might loop forever due to the new version of instantiation rules. For example, the following derivation results in infinite instantiations:

$$\begin{aligned}
 & \hat{\alpha}, \hat{\beta} \Vdash \hat{\alpha} \leq 1 \rightarrow \hat{\beta} \Vdash \hat{\beta} \leq 1 \rightarrow \hat{\alpha} \\
 \longrightarrow & \hat{\alpha}, \hat{\beta}_1, \hat{\beta}_2, \{\} \leq \hat{\beta} \leq \{\hat{\beta}_1 \rightarrow \hat{\beta}_2\} \Vdash \hat{\alpha} \leq 1 \rightarrow \hat{\beta} \Vdash \hat{\beta}_1 \rightarrow \hat{\beta}_2 \leq 1 \rightarrow \hat{\alpha} \\
 \longrightarrow & \dots \\
 \longrightarrow & \hat{\alpha}, \{1\} \leq \hat{\beta}_1 \leq \{\}, \{\} \leq \hat{\beta}_2 \leq \{\hat{\alpha}\}, \{\} \leq \hat{\beta} \leq \{\hat{\beta}_1 \rightarrow \hat{\beta}_2\} \Vdash \hat{\alpha} \leq 1 \rightarrow \hat{\beta} \\
 \longrightarrow & \dots
 \end{aligned}$$

A further analyse on the problem suggests that the *loop dependencies* are not detected and reduced smartly by the lazy algorithm: it insist too much on not missing any possible solutions. By loop dependency we refer to the following pattern  $\hat{\alpha} \leq A \rightarrow B$  where  $\hat{\alpha} \in A \rightarrow B$ . In the above example, one can derive (by transitivity of subtyping)  $\hat{\alpha} \leq 1 \rightarrow \hat{\beta} \leq 1 \rightarrow 1 \rightarrow \hat{\alpha}$  which forms a (indirect) loop dependency on  $\hat{\alpha}$ .

By performing experimental implementations and randomly generated tests, we observe that the algorithm accepts more valid relations when there are no loop dependencies, but may be trapped in infinite loops otherwise. We have also tried to formalize how a loop detection procedure might be used to help prevent the non-termination problem, but most of them turn out to be rather complicated and hard to mechanically formalize.



## PART III

### RELATED WORK

# 6 RELATED WORK

Throughout the thesis, we have already discussed much of the closest related work. In this section we summarize the key differences and novelties, and discuss some other related work.

## 6.1 HIGHER-RANKED POLYMORPHIC TYPE INFERENCE ALGORITHMS

### 6.1.1 PREDICATIVE ALGORITHMS

Higher-ranked polymorphism is a convenient and practical feature of programming languages. Since full type-inference for System F is undecidable [Wells 1999], various decidable partial type-inference algorithms were developed. The declarative formulation of subtyping in Chapters 3 and 4 (and later extended in Chapter 5), originally proposed by Odersky and Läufer [1996], is *predicative*:  $\forall$ 's only instantiate to monotypes. The monotype restriction on instantiation is considered reasonable and practical for most programs, except for those that require sophisticated forms of higher-order polymorphism. In predicative higher-ranked settings, type annotations may guide the type system to accept lambda functions of higher-ranked types.

In addition to OL's type system, the bidirectional system proposed by DK [Dunfield and Krishnaswami 2013] accepts even better type annotations via bidirectional type checking, which also allow higher-ranked function types to check against an unannotated lambda. Such annotations also improve readability of the program, and are not much of a burden in practice. DK's algorithm is shown to be sound, complete and decidable in 70 pages of manual proofs. Though carefully written, some of the proofs are incorrect (see discussion in Section 4.1.2 and 4.3.2), which creates difficulties when formalizing them in a proof assistant. In their follow-up work Dunfield and Krishnaswami [2019] enrich the bidirectional higher-ranked system with existentials and indexed types. With a more complex declarative system, they developed a proof of over 150 pages. It is even more difficult to argue its correctness for every single detail within such a big development. Unfortunately, we find that their Lemma 26 (Parallel Admissibility) appears to have the same issue as lemma 29 in [Dunfield and Krishnaswami 2013]: the conclusion is false. We also discuss the issue in more detail in Section 4.1.2.

Peyton Jones et al. [2007] developed another higher-ranked predicative bidirectional type system. Their subtyping relation is enriched with *deep skolemisation*, which is more general than ours and allows more valid relations. The concept is also known as the distributivity rule [Mitchell 1984]. For example,

$$\forall a. \forall b. a \rightarrow b \rightarrow b \leq \forall a. a \rightarrow (\forall b. b \rightarrow b)$$

this subtyping relation does not hold in OL’s subtyping relation, because the instantiation of  $b$  in the left-hand-side type happens strictly before the introduction of the variable  $b$  in the right-hand-side type. Deep skolemisation can be achieved through a pre-processing that extracts all the  $\forall$ ’s in the return position of a function type to the top level. After such pre-processing, the right-hand-side type becomes  $\forall a. \forall b. a \rightarrow b \rightarrow b$  and the subtyping relation holds.

In terms of handling applications where type parameters are implicit, they do not use the application inference judgment as DK’s type system. Instead, they employ a complicated mechanism for implicit instantiation taken care by the unification process for the algorithm. A manual proof is given, showing that the algorithm is sound and complete with respect to their declarative specification.

In a more recent work, Xie and Oliveira [2018] proposed a variant of a bidirectional type inference system for a predicative system with higher-ranked types. Type information flows from arguments to functions with an additional *application* mode. This variant allows more higher-order typed programs to be inferred without additional annotations. Following the new mode, the let-generalization of the Hindley-Milner system is well supported as syntactic sugar. The formalization includes some mechanized proofs for the declarative type system, but all proofs regarding the algorithmic type system are manual.

### 6.1.2 IMPREDICATIVE ALGORITHMS

Impredicative System F allows instantiation with polymorphic types, but unfortunately its subtyping system is already undecidable [Chrzęszcz 1998]. Works on partial impredicative type-inference algorithms navigate a variety of design tradeoffs for a decidable algorithm. Generally speaking, such algorithms tend to be more complicated, and thus less adopted in real-world programming languages.

$\text{ML}^F$  [Le Botlan and Rémy 2003, 2009; Rémy and Yakobowski 2008] extends types of System F with a form of bounded quantification and proposes an impredicative system. The type inference algorithm always infers principal types given proper type annotations. Through the technique called “monomorphic abstraction of polymorphic types”, polymorphic instantiations are expressed by constraints on type variables of type schemes. An annotation is

only needed when an argument of a lambda function is used polymorphically. Moreover, their type system is robust against a wide range of program transformations, including let-expansion, let-reduction and  $\eta$ -expansion. However, the extended type structure of  $\text{ML}^F$  introduces non-compatible types with System F, and it complicates the metatheory and implementation of the type system.

HMF [Leijen 2008] takes a slightly different approach by not extending types with bounds, therefore programmers still work with plain System F types. The algorithm works similarly compared to  $\text{ML}^F$ , except that it does not output richer types. Instead, when there are ambiguity and no principal type can be inferred, being a conservative algorithm, HMF prefers predicative instantiations to maintain backward the compatibility of HM. Soon after HMF, HML [Leijen 2009] is proposed as an extension of HMF with flexible types. The system eliminates rigid quantifications ( $\alpha = \sigma$ ) of  $\text{ML}^F$  and only keeps the flexible quantifications ( $\alpha \geq \sigma$ ). Furthermore, a variant “Rigid HML” is presented by restricting let-bindings to System F types only, so that type annotations can still stay inside System F. Similar to HMF, they both require type annotations at higher-ranked types or ambiguous implicit instantiations.

FPH [Vytiniotis et al. 2008] employ a box notion  $\boxed{\sigma}$  to encapsulate polymorphic instantiations internally in the algorithm. A type inside a box indicates that the instantiation is impredicative. Therefore an inferred type with a box type in it means that incomparable System F types may be produced, and that is rejected before entering the environment. Although more annotations might be required compared with other approaches, the algorithm of FPH is much simpler, thanks to the simple syntax and subtyping relation of box types.

Guarded Impredicative Polymorphism [Serrano et al. 2018] was recently proposed as an improvement on GHC’s type inference algorithm with impredicative instantiation. They make use of local information in  $n$ -ary applications to infer polymorphic instantiations with a relatively simple specification and unification algorithm. Although not all impredicative instantiations can be handled well, their algorithm is already quite useful in practice.

A recent follow-up work, the Quick Look [Serrano et al. 2020], takes a more conservative approach. The algorithm will try its best to infer impredicative instantiations, and only use the instantiation when it is the best one. In order to do so, the algorithm also needs to analyse all the argument types during a function call, and make use of the *guardedness* property to decide the principality of the inferred instantiation. When Quick Look cannot give the best instantiation, it will instead try predicative inference as HM. This conservative approach makes sure that the algorithm does not infer bad types and leads the subsequent type checking in wrong directions. In the meantime, they treat the function arrow ( $\rightarrow$ ) as invariant like normal type constructors. This change significantly restricts the subtyping relation and thus simplifies guesses for implicit parameters. Manual  $\eta$ -expansions are required to regain the co- and contravariance of function types.

FreezeML [Emrich et al. 2020] is another recent work that extends the ML type system with impredicative instantiations. A special syntax of expression, the explicit freezing  $\lceil x \rceil$ , is added to guide the type system. Freezed variables are prevented to be instantiated by the type system, therefore variables of polymorphic types may force the type inference algorithm to instantiate impredicatively. In combination with the let-generalization rule, programmers may also encode explicit generalization and explicit instantiation. This work provides a nice means for programmers to control how type inference algorithm deals with impredicative instantiations.

## 6.2 TYPE INFERENCE ALGORITHMS WITH SUBTYPING

Type systems in presence of subtyping usually encounter constraints that are not simply equalities as in HM. Therefore constraint solvers used in HM, where unifications are based on equality, cannot be easily extended to support subtyping. Instead, constraints are usually collected as subtyping relations and may delay resolving as the constraints accumulate. Eifrig et al. [1995a,c] proposed systems that are based on *constraint types*, i.e. types expressed together with a set of constraints  $\tau \mid \{\overline{\tau_1 \leq \tau_2}\}$ . Their type checking algorithm checks at each step whether each constraint in the closure of constraint set is *consistent*, which is a set of rules that prevent obvious contradictions appear, such as  $\text{Int} \leq \text{Bool}$ . Our attempt in Section 5.5.2 is similar to this idea, where we exhaustively check if every subtyping relation derived from the bounds is valid. However, our algorithm tries to solve the bounds into concrete types, and it turns out that it never terminates in some complex cases.

Constraint types improve the expressiveness of the type system, yet the type inference algorithm can be quite slow. The size of the constraint is linear in the program size, and the closure can grow to cubic size. Pottier [1998] proposed three methods to simplify constraints in his Ph.D. thesis, aiming at improving the efficiency of type inference algorithms and improving the readability of the resulting types. By *canonization*, constraints are converted to canonical forms with the introduction of new meta-variables; *garbage collection* eliminates constraints that do not affect the type; finally, *minimization* shares nodes within the constraint graph.

Inspired by the simplification strategies of Pottier’s, MLsub [Dolan and Mycroft 2017] suggest that the data flow on the constraint graph can reflect directly on types extended by a richer type system, where lattices operations are used to represent constraints imposed on type variables. Polar types distinguish between input types and output types, and pose different restrictions on them. As a result, constraints can easily be transformed into canonical forms, and the bi-unification algorithm can solve them by simple substitutions. One can also view the MLsub system as a different way to encode the constraint types: instead of a set of

constraints that is stated along with the type, types themselves now contain subtyping constraints with the help of lattice operations. Similarly, the size of constraints may grow with the size of program and affect the readability. Therefore, various simplification algorithms should be used in real applications. A more recent work, the Simple-sub [Parreaux 2020], further simplifies the algorithm of MLsub and is implemented in 500 lines of code. While being equivalent to MLsub, it is a more efficient variant.

## 6.3 TECHNIQUES USED IN TYPE INFERENCE ALGORITHMS

### 6.3.1 ORDERED CONTEXTS IN TYPE INFERENCE

Gundry et al. [2010] revisit algorithm  $\mathcal{W}$  and propose a new unification algorithm with the help of ordered contexts. Similar to DK's algorithm, information of meta-variables flows from input contexts to output contexts. Not surprisingly, its information increase relation has a similar role to DK's context extension. Our algorithm, in contrast, eliminates output contexts and solution records ( $\hat{\alpha} = \tau$ ), simplifying the information propagation process through immediate substitution by collecting all the judgments in a single worklist.

### 6.3.2 THE ESSENCE OF ML TYPE INFERENCE

Constraint-based type inference is adopted by Pottier and Rémy [2005] for ML type systems, which do not employ higher-ranked polymorphism. An interesting feature of their algorithm is that it keeps precise scoping of variables, similarly to our approach. Their algorithm is divided into constraint generation and solving phases (which are typical of constraint-based algorithms). Furthermore an intermediate language is used to describe constraints and their constraint solver utilizes a stack to track the state of the solving process. In contrast, our algorithm has a single phase, where the judgment chains themselves act as constraints, thus no separate constraint language is needed.

### 6.3.3 LISTS OF JUDGMENTS IN UNIFICATION

Some work [Abel and Pientka 2011; Reed 2009] adopts a similar idea to this paper in work on unification for dependently typed languages. Similarly to our work the algorithms need to be very careful about scoping, since the order of variable declarations is fundamental in a dependently typed setting. Their algorithms simplify a collection of unification constraints progressively in a single-step style. In comparison, our algorithm mixes variable declarations with judgments, resulting in a simpler judgment form, while processing them in a similar way. One important difference is that contexts are duplicated in their unification judgments,

which complicates the unification process, since the information of each local context needs to be synchronized. Instead we make use of the nature of ordered context to control the scope of unification variables. While their algorithms focus only on unification, our algorithm also deals with other types of judgments like synthesis. A detailed discussion is given in Section 4.1.3.

## 6.4 MECHANICAL FORMALIZATION OF POLYMORPHIC TYPE SYSTEMS

Since the publication of the POPLMARK challenge [Aydemir et al. 2005], many theorem provers and packages provide new methods for dealing with variable binding [Aydemir et al. 2008; Chlipala 2008; Urban 2008]. More and more type systems are formalized with these tools. However, mechanizing certain algorithmic aspects, like unification and constraint solving, has received very little attention and is still challenging. Moreover, while most tools support local (input) contexts in a neat way, many practical type-inference algorithms require more complex binding structures with output contexts or various forms of constraint solving procedures.

Algorithm  $\mathcal{W}$ , as one of the classic type inference algorithms for polymorphic type systems, has been manually proven to be sound and complete with respect to the Hindley-Milner type system [Damas and Milner 1982; Hindley 1969; Milner 1978]. After around 15 years, the algorithm was formally verified by Naraschewski and Nipkow [1999] in Isabelle/HOL [Nipkow et al. 2002]. The treatment of new variables was tricky at that time, while the overall structure follows the structure of Damas’s manual proof closely. Later on, other researchers [Dubois 2000; Dubois and Menissier-Morain 1999] formalized algorithm  $\mathcal{W}$  in Coq [The Coq development team 2017]. Nominal techniques [Urban 2008] in Isabelle/HOL have been developed to help programming language formalizations, and are used for a similar verification [Urban and Nipkow 2008]. Moreover, Garrigue [Garrigue 2015] mechanized a type inference algorithm, with the help of locally nameless [Charguéraud 2012], for Core ML extended with structural polymorphism and recursion.

## PART IV

## EPILOGUE



# 7 CONCLUSION AND FUTURE WORK

## 7.1 CONCLUSION

In this thesis, we proposed new bidirectional type inference algorithms for predicative higher-ranked implicit parametric polymorphic systems. We showed how worklists ease the design of algorithms and also mechanical formalizations. By collecting all judgments in a single worklist, the algorithm performs unification with a bird’s-eye view of all the judgments, therefore propagation between judgments is as simple as a global substitution. Compared with the classical HM unification procedure, our algorithm does not need separated relations (or fixpoints) for unification. From the formalization point of view, eager substitutions are also easier to state and reason in a proof assistant. Therefore, we obtained fully formalized properties for all our developments relatively easily. Overall, we developed the following systems and/or type inference algorithms:

- We developed a worklist algorithm for OL’s higher-ranked subtyping system. The algorithm operates on a worklist of subtyping judgments and a single context where a variable declaration is shared across the worklist. Compared with DK’s algorithm, our approach avoids the use of output contexts, which complicates the scoping of variables and is hard to formalize in a proof assistant. Eager substitutions that solve existential variables are directly applied to the worklist, therefore passing the partial information to the rest judgments. We proved *soundness*, *completeness* and *decidability* of the algorithm in the Abella theorem prover.
- We developed a worklist algorithm for DK’s higher-ranked bidirectional type system. In order to properly encode judgments that output types, such as the type inference judgment, continuation-passing-style *judgment chains* are developed. Compared with the previous work that uses a single context, we further unify the worklist with variable declarations. Such unification results in a much more accurate track of variable scopings, and the algorithm will garbage-collect variables as soon as they are not referred to anymore. Unlike using output contexts as DK’s algorithm, designing rules for worklist context is less likely to contain bugs in terms of variable scoping. Once

again, based on eager substitutions, the algorithm is easy to formalize. We showed *soundness*, *completeness* and *decidability* in the Abella theorem prover.

- We developed a backtracking-based algorithm for a higher-ranked bidirectional type system with object-oriented subtyping. With the introduction of the top and bottom types and relevant subtyping relations, meta-variable instantiations are no longer deterministic like the HM system. The backtracking-based algorithm preserves most characteristics of the previous worklist context, and it “tries” obvious solutions in parallel with detailed analysis as previous work. We proved that the algorithm is always *sound*, and subtyping is complete under the rank-1 restriction.

## 7.2 FUTURE WORK

In this section, we discuss several interesting possibilities to explore in the future.

**TERMINATION THEOREM FOR TYPE INFERENCE WITH SUBTYPING** The algorithm we presented and formalized in Chapter 5 enjoys several important machine-verified properties, including soundness and completeness under a rank-1 restriction. However, we have not yet found a proper termination measure for the algorithmic subtyping rules. What complicates the problem is that common measures are not decreasing all the time. For example, subtyping between function types may increase the number of judgments; and instantiation judgments that split existential variables may increase the number of unsolved existential variables. We are mostly convinced by the fact that a large set of tests performed on randomly generated subtyping judgments do terminate. And we hope that a clever termination argument may be proposed in the near future.

**PRACTICAL IMPREDICATIVE TYPE INFERENCE** Impredicative type inference algorithms are sometimes helpful when dealing with higher-ranked types. Simple forms of impredicativity can be unambiguous and easy to understand by programmers, such as

`head ids`

where  $\text{head} :: \forall p. [p] \rightarrow p$  and  $\text{ids} :: [\forall a. a \rightarrow a]$ . The application  $(\text{head } \text{ids})$  should have type  $\forall a. a \rightarrow a$ , because  $p$  is instantiated to  $\forall a. a \rightarrow a$  and that is the only choice we can make.

Moreover, even though there are programs that have incomparable System F types, impredicative system may offer means like type annotation to manually pick the right one. For example,

`single id`

where  $\text{single} :: \forall a. p \rightarrow [p]$  and  $\text{id} :: \forall a. a \rightarrow a$ . This term has two incomparable types  $[\forall a. a \rightarrow a]$  and  $\forall a. [a \rightarrow a]$ . One might expect that a type annotation might help pick the right case, like  $((\text{single id}) :: [\forall a. a \rightarrow a])$ , which requires a system that allows impredicative instantiation to do so. There are many impredicative systems developed out of various techniques [Emrich et al. 2020; Le Botlan and Rémy 2003; Leijen 2008, 2009; Serrano et al. 2020, 2018; Vytiniotis et al. 2008], and we leave the extension of impredicativity to our algorithms as future work.

**OPTIMIZATION** All the three type inference algorithms described in this thesis are based on eager global substitution on the worklists. The benefit is to reduce the difficulty for mechanical formalizations in theorem provers, especially in Abella. Additionally, eager substitutions represent simple equivalent transformations on the state of worklist, alleviating the complication of variable scoping and reasoning for correctness. However, a naive implementation of our algorithm will be very inefficient, since the worklist is iterated so often whenever an existential variable is solved or partially solved.

Compared with mature algorithms like the algorithm  $\mathcal{W}$  [Milner 1978] and the  $\text{OutsideIn}(X)$  [Vytiniotis et al. 2011] type inference algorithm, which produces substitutions once all the constraints are collected and reduced, our algorithms manipulate the judgments themselves on-the-fly. Efficient implementations will model meta-variables as mutable references and apply the results if they are solved by the constraint solver. DK’s algorithm [Dunfield and Krishnaswami 2013] requires output contexts, which are likely to be implemented in a state monad or other stateful approaches. Solved existential variables encoded in their algorithmic context are different from ours. Instead of removing the declaration of the variable and propagate its solution to all the rest judgments, they extend the entry in the context with its solution,  $\hat{\alpha} = \tau$ , and pass it to the output context. Such entry is easy to implement with an efficient approach using mutable references. We can also adopt such an idea and reformulate our eager substitution with lazy ones by keeping track of solutions to existential variables, and only perform substitution when the judgment is just about to be reduced.

Unfortunately, lazy substitution only solves part of the problem. By using ordered contexts, the scoping of variables is precisely captured by the relative positions in which they are declared in the context. Some algorithmic rules are designed to use the position informa-

tion, where swapping positions might lead to unsound implementation. For example, the algorithmic system of Chapter 5 includes these rules:

$$\begin{aligned}\Gamma[a][\widehat{\beta}] \Vdash a \leq \widehat{\beta} &\longrightarrow_{16} \{\widehat{\beta} := a\} \Gamma[a][\widehat{\beta}] \\ \Gamma[a][\widehat{\beta}] \Vdash \widehat{\beta} \leq a &\longrightarrow_{17} \{\widehat{\beta} := a\} \Gamma[a][\widehat{\beta}]\end{aligned}$$

The hole notation  $\Gamma[a][\widehat{\beta}]$  is used to ensure that  $a$  is declared before  $\widehat{\beta}$ , therefore it is fine for the monotype represented by  $\widehat{\beta}$  to be  $a$ . In a naive implementation, frequent iteration is required to look up the relative positions of variable declarations. We may address the problem for a more efficient algorithm following the idea of Remy et al. [1992].

Moreover, there are also reduction rules that *remove* and/or *insert* variables from(to) the middle of the ordered context, as in the following rules:

$$\begin{aligned}\Gamma[\widehat{\alpha}] \Vdash \widehat{\alpha} \leq A \rightarrow B &\longrightarrow_{14} \{\widehat{\alpha} := \widehat{\alpha}_1 \rightarrow \widehat{\alpha}_2\} (\Gamma[\widehat{\alpha}] \Vdash \widehat{\alpha} \leq A \rightarrow B) \\ &\quad \text{when } \widehat{\alpha} \notin FV(A \rightarrow B) \\ \Gamma[\widehat{\alpha}] \Vdash A \rightarrow B \leq \widehat{\alpha} &\longrightarrow_{15} \{\widehat{\alpha} := \widehat{\alpha}_1 \rightarrow \widehat{\alpha}_2\} (\Gamma[\widehat{\alpha}] \Vdash A \rightarrow B \leq \widehat{\alpha}) \\ &\quad \text{when } \widehat{\alpha} \notin FV(A \rightarrow B)\end{aligned}$$

To sum up, we plan to propose an implementation-friendly version of the algorithm that does not require many eager substitutions, and we also hope that there will be suitable data structures or algorithms which optimizes the complexity of those time-consuming operations on the ordered context.

**CLASS HIERARCHY OF NOMINAL SUBTYPING** Nowadays, most mainstream object-oriented programming languages support nominal subtyping, which means that the subtyping relation is defined by the programmers explicitly through language constructs such as inheritance. For example, we may define a `Pos` class to represent positive integers, which inherits the `Int` class. The subtyping relation  $\text{Pos} \leq \text{Int}$  automatically holds. An algorithmic subtyping judgment  $\widehat{\alpha} \leq \text{Int}$  now might have several solutions: `Int`, `Pos`, or  $\perp$ . In order to calculate a *range* of solutions instead of just simple types, lazy algorithms that operate on bounds fit the task better. We plan to explore extensions of the algorithm described in Section 5.5.2. The algorithm progressively collects upper and lower bounds of existential variables, and finally checks if valid solutions are satisfying all the bounds.

Nonetheless, we encounter termination problems when trying to adopt the lazy approach. Experiments indicate that loop dependency in combination with the instantiation rule ( $\widehat{\alpha} \leq A \rightarrow B$ ) might cause infinite loops. As a practical compromise, we may impose restrictions on existential variables so that they are not instantiated to function types when they are likely

to be a class type. We also hope that better treatment on loop dependencies can be developed, especially under assumptions of real-world programming tasks. For future work, we aim at a terminating and sound algorithm for practical object-oriented type inference.

**TYPE INFERENCE WITH RECURSIVE TYPES** Among all the algorithms we proposed in this thesis, we reject subtyping judgments like

$$\hat{\alpha} \leq A \rightarrow B \text{ where } \hat{\alpha} \in \text{FV}(A) \cup \text{FV}(B)$$

because no HM monotype satisfies a subtyping relation like

$$\tau \leq \tau \rightarrow \text{Int}$$

and unfolding such algorithmic judgment might cause an infinite loop. However, recursive types may satisfy such judgment, for example,

$$\mu x. x \rightarrow \text{Int} \leq (\mu x. x \rightarrow \text{Int}) \rightarrow \text{Int}$$

holds since the unfolding of  $\mu x. x \rightarrow \text{Int}$  is  $(\mu x. x \rightarrow \text{Int}) \rightarrow \text{Int}$ . In other words, if we extend our type system by including recursive types, we would accept such judgment with a principal recursive type, instead of rejecting it because we cannot find such a type. The MLsub [Dolan and Mycroft 2017] already accepts recursive types, and it deals with recursive constraints in the way we described above. It is an interesting question if adding recursive types may improve the expressiveness of our system and remain complete at the same time.

**BOUNDED QUANTIFICATION AND F-BOUNDED QUANTIFICATION** Bounded quantification [Cardelli and Wegner 1985] and F-bounded quantification [Canning et al. 1989] are techniques commonly seen in object-oriented programming languages. Bounded quantification allows constraints imposed on universally quantified variables:

$$\forall(a \leq \text{Int}). a \rightarrow a \rightarrow \text{String}$$

F-Bounded quantification further extends the constraints so that recursive references of the variables are also allowed. Although the subtyping is proven to be undecidable in presence of both F-bounded quantification and variance on generics, Kennedy and Pierce [2007] proposed three forms of restrictions that we can explore further. Combining bidirectional type checking with bounded quantifications may further improve the experience for object-oriented programmers.

## BIBLIOGRAPHY

[Citing pages are listed after each reference.]

Andreas Abel, Guillaume Allais, Aliya Hameer, Brigitte Pientka, Alberto Momigliano, Steven Schäfer, and Kathrin Stark. 2018. POPLMark Reloaded: Mechanizing Proofs by Logical Relations. *Submitted to the Journal of functional programming* (2018). [cited on page 9]

Andreas Abel and Brigitte Pientka. 2011. Higher-order dynamic pattern unification for dependent types and records. In *International Conference on Typed Lambda Calculi and Applications*. Springer, 10–26. [cited on pages 44, 50, and 105]

Brian Aydemir, Arthur Charguéraud, Benjamin C. Pierce, Randy Pollack, and Stephanie Weirich. 2008. Engineering Formal Metatheory. In *Proceedings of the 35th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '08)*. [cited on pages 9 and 106]

Brian E Aydemir, Aaron Bohannon, Matthew Fairbairn, J Nathan Foster, Benjamin C Pierce, Peter Sewell, Dimitrios Vytiniotis, Geoffrey Washburn, Stephanie Weirich, and Steve Zdancewic. 2005. Mechanized metatheory for the masses: The POPLmark challenge. In *The 18th International Conference on Theorem Proving in Higher Order Logics*. [cited on pages 9 and 106]

Yves Bertot, Benjamin Grégoire, and Xavier Leroy. 2006. A Structured Approach to Proving Compiler Optimizations Based on Dataflow Analysis. In *Proceedings of the 2004 International Conference on Types for Proofs and Programs (TYPES'04)*. [cited on page 9]

Peter Canning, William Cook, Walter Hill, Walter Olthoff, and John C. Mitchell. 1989. F-Bounded Polymorphism for Object-Oriented Programming. In *Proceedings of the Fourth International Conference on Functional Programming Languages and Computer Architecture* (Imperial College, London, United Kingdom) (*FPCA '89*). Association for Computing Machinery, New York, NY, USA, 273–280. <https://doi.org/10.1145/99370.99392> [cited on page 112]

- Luca Cardelli. 1988. A semantics of multiple inheritance. *Information and Computation* 76, 2 (1988), 138–164. [https://doi.org/10.1016/0890-5401\(88\)90007-7](https://doi.org/10.1016/0890-5401(88)90007-7) [cited on page 8]
- Luca Cardelli and Peter Wegner. 1985. On Understanding Types, Data Abstraction, and Polymorphism. *ACM Comput. Surv.* 17, 4 (Dec. 1985), 471–523. <https://doi.org/10.1145/6041.6042> [cited on page 112]
- Bor-Yuh Evan Chang, Adam Chlipala, and George C. Necula. 2006. A Framework for Certified Program Analysis and Its Applications to Mobile-code Safety. In *Proceedings of the 7th International Conference on Verification, Model Checking, and Abstract Interpretation (VMCAI’06)*. [cited on page 9]
- Arthur Charguéraud. 2012. The Locally Nameless Representation. *Journal of Automated Reasoning* 49, 3 (01 Oct 2012), 363–408. [cited on page 106]
- Paul Chiusano and Runar Bjarnason. 2015. Unison. <http://unisonweb.org> [cited on page 7]
- Adam Chlipala. 2008. Parametric Higher-order Abstract Syntax for Mechanized Semantics. In *Proceedings of the 13th ACM SIGPLAN International Conference on Functional Programming (ICFP ’08)*. [cited on page 106]
- Jacek Chrząszcz. 1998. Polymorphic subtyping without distributivity. In *Mathematical Foundations of Computer Science 1998*, Luboš Brim, Jozef Gruska, and Jiří Zlatuška (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 346–355. [cited on pages 19 and 102]
- Alonzo Church. 1932. A Set of Postulates for the Foundation of Logic. *Annals of Mathematics* 33, 2 (1932), 346–366. <http://www.jstor.org/stable/1968337> [cited on page 4]
- Alonzo Church. 1941. *The calculi of lambda-conversion*. Number 6. Princeton University Press. [cited on page 4]
- 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)*. [cited on pages 5, 13, 30, and 106]
- N.G de Bruijn. 1972. Lambda calculus notation with nameless dummies, a tool for automatic formula manipulation, with application to the Church-Rosser theorem. *Indagationes Mathematicae (Proceedings)* 75, 5 (1972), 381–392. [https://doi.org/10.1016/1385-7258\(72\)90034-0](https://doi.org/10.1016/1385-7258(72)90034-0) [cited on page 10]

- Disciple Development Team. 2017. The Disciplined Disciple Compiler. <http://disciple.ouroborus.net/> [cited on page 8]
- Stephen Dolan and Alan Mycroft. 2017. Polymorphism, Subtyping, and Type Inference in MLsub. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages* (Paris, France) (POPL 2017). Association for Computing Machinery, New York, NY, USA, 60–72. <https://doi.org/10.1145/3009837.3009882> [cited on pages 23, 24, 79, 104, and 112]
- Catherine Dubois. 2000. Proving ML type soundness within Coq. *Theorem Proving in Higher Order Logics* (2000), 126–144. [cited on pages 9 and 106]
- Catherine Dubois and Valerie Menissier-Morain. 1999. Certification of a type inference tool for ML: Damas–Milner within Coq. *Journal of Automated Reasoning* 23, 3 (1999), 319–346. [cited on pages 9 and 106]
- 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)*. [cited on pages 7, 9, 21, 28, 30, 32, 43, 44, 48, 65, 101, and 110]
- Joshua Dunfield and Neelakantan R. Krishnaswami. 2019. Sound and Complete Bidirectional Typechecking for Higher-rank Polymorphism with Existentials and Indexed Types. *Proc. ACM Program. Lang.* 3, POPL, Article 9 (Jan. 2019), 28 pages. [cited on pages 48, 49, and 101]
- Jonathan Eifrig, Scott Smith, and Valery Trifonov. 1995a. Sound Polymorphic Type Inference for Objects (OOPSLA ’95). Association for Computing Machinery, New York, NY, USA, 169–184. <https://doi.org/10.1145/217838.217858> [cited on page 104]
- Jonathan Eifrig, Scott Smith, and Valery Trifonov. 1995b. Type Inference for Recursively Constrained Types and its Application to OOP. *Electronic Notes in Theoretical Computer Science* 1 (1995), 132 – 153. [https://doi.org/10.1016/S1571-0661\(04\)80008-2](https://doi.org/10.1016/S1571-0661(04)80008-2) MFPS XI, Mathematical Foundations of Programming Semantics, Eleventh Annual Conference. [cited on page 79]
- Jonathan Eifrig, Scott Smith, and Valery Trifonov. 1995c. Type Inference for Recursively Constrained Types and its Application to OOP. *Electronic Notes in Theoretical Computer Science* 1 (1995), 132–153. [https://doi.org/10.1016/S1571-0661\(04\)80008-2](https://doi.org/10.1016/S1571-0661(04)80008-2) MFPS XI, Mathematical Foundations of Programming Semantics, Eleventh Annual Conference. [cited on page 104]



- Frank Emrich, Sam Lindley, Jan Stolarek, James Cheney, and Jonathan Coates. 2020. FreezeML: Complete and Easy Type Inference for First-Class Polymorphism. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation* (London, UK) (*PLDI 2020*). Association for Computing Machinery, New York, NY, USA, 423–437. <https://doi.org/10.1145/3385412.3386003> [cited on pages 6, 104, and 110]
- Phil Freeman. 2017. PureScript. <http://www.purescript.org/> [cited on pages 7 and 8]
- Andrew Gacek. 2008. The Abella Interactive Theorem Prover (System Description). In *Proceedings of IJCAR 2008 (Lecture Notes in Artificial Intelligence)*. [cited on pages 10, 40, 43, 44, 61, 71, and 95]
- Jacques Garrigue. 2015. A certified implementation of ML with structural polymorphism and recursive types. *Mathematical Structures in Computer Science* 25, 4 (2015), 867–891. [cited on pages 9 and 106]
- Andrew Gill, John Launchbury, and Simon L. Peyton Jones. 1993. A Short Cut to Deforestation. In *Proceedings of the Conference on Functional Programming Languages and Computer Architecture (FPCA '93)*. [cited on page 6]
- Jean-Yves Girard. 1971. Une Extension De L'Interpretation De Gödel a L'Analyse, Et Son Application a L'Elimination Des Coupures Dans L'Analyse Et La Theorie Des Types. In *Proceedings of the Second Scandinavian Logic Symposium*, J.E. Fenstad (Ed.). Studies in Logic and the Foundations of Mathematics, Vol. 63. Elsevier, 63–92. [https://doi.org/10.1016/S0049-237X\(08\)70843-7](https://doi.org/10.1016/S0049-237X(08)70843-7) [cited on page 4]
- Jean-Yves Girard. 1972. *Interprétation fonctionnelle et élimination des coupures de l'arithmétique d'ordre supérieur*. Ph.D. Dissertation. Éditeur inconnu. [cited on page 4]
- Adam Gundry, Conor McBride, and James McKinna. 2010. Type Inference in Context. In *Proceedings of the Third ACM SIGPLAN Workshop on Mathematically Structured Functional Programming (MSFP '10)*. [cited on pages 44 and 105]
- Roger Hindley. 1969. The principal type-scheme of an object in combinatory logic. *Transactions of the american mathematical society* 146 (1969), 29–60. [cited on pages 5, 13, and 106]
- Mark P. Jones. 1995. Functional Programming with Overloading and Higher-Order Polymorphism. In *Advanced Functional Programming (Lecture Notes in Computer Science 925)*. [cited on page 6]

- Andrew Kennedy and Benjamin C. Pierce. 2007. On Decidability of Nominal Subtyping with Variance. In *International Workshop on Foundations and Developments of Object-Oriented Languages (FOOL/WOOD)* (international workshop on foundations and developments of object-oriented languages (fool/wood) ed.). <https://www.microsoft.com/en-us/research/publication/on-decidability-of-nominal-subtyping-with-variance/> [cited on page 112]
- A. J. Kfoury and J. Tiuryn. 1992. Type Reconstruction in Finite Rank Fragments of the Second-Order  $\lambda$ -Calculus. *Inf. Comput.* 98, 2 (June 1992), 228–257. [https://doi.org/10.1016/0890-5401\(92\)90020-G](https://doi.org/10.1016/0890-5401(92)90020-G) [cited on page 17]
- Casey Klein, John Clements, Christos Dimoulas, Carl Eastlund, Matthias Felleisen, Matthew Flatt, Jay A. McCarthy, Jon Rafkind, Sam Tobin-Hochstadt, and Robert Bruce Findler. 2012. Run Your Research: On the Effectiveness of Lightweight Mechanization. In *Proceedings of the 39th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Philadelphia, PA, USA) (POPL '12). 285–296. [cited on pages 9 and 44]
- Ralf Lämmel and Simon Peyton Jones. 2003. Scrap Your Boilerplate: A Practical Design Pattern for Generic Programming. In *Proceedings of the 2003 ACM SIGPLAN International Workshop on Types in Languages Design and Implementation (TLDI '03)*. [cited on page 6]
- John Launchbury and Simon L. Peyton Jones. 1994. Lazy Functional State Threads. In *Proceedings of the ACM SIGPLAN 1994 Conference on Programming Language Design and Implementation* (Orlando, Florida, USA) (PLDI '94). Association for Computing Machinery, New York, NY, USA, 24–35. <https://doi.org/10.1145/178243.178246> [cited on page 6]
- John Launchbury and Simon L. Peyton Jones. 1995. State in Haskell. *LISP and Symbolic Computation* 8, 4 (1995), 293–341. [cited on page 6]
- Didier Le Botlan and Didier Rémy. 2003. MLF: Raising ML to the Power of System F. In *Proceedings of the Eighth ACM SIGPLAN International Conference on Functional Programming (ICFP '03)*. [cited on pages 6, 7, 102, and 110]
- Didier Le Botlan and Didier Rémy. 2009. Recasting MLF. *Information and Computation* 207, 6 (2009), 726–785. <https://doi.org/10.1016/j.ic.2008.12.006> [cited on page 102]
- Daan Leijen. 2008. HMF: Simple Type Inference for First-class Polymorphism. In *Proceedings of the 13th ACM SIGPLAN International Conference on Functional Programming (ICFP '08)*. [cited on pages 6, 7, 103, and 110]

- Daan Leijen. 2009. Flexible Types: Robust Type Inference for First-Class Polymorphism. In *Proceedings of the 36th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Savannah, GA, USA) (POPL '09). Association for Computing Machinery, New York, NY, USA, 66–77. <https://doi.org/10.1145/1480881.1480891> [cited on pages 103 and 110]
- Xavier Leroy et al. 2012. The CompCert verified compiler. *Documentation and user's manual*. INRIA Paris-Rocquencourt (2012). [cited on page 9]
- Alberto Martelli and Ugo Montanari. 1982. An Efficient Unification Algorithm. *ACM Trans. Program. Lang. Syst.* 4, 2 (April 1982), 258–282. [cited on page 50]
- Dale Miller. 1992. Unification under a mixed prefix. *Journal of Symbolic Computation* 14, 4 (1992), 321–358. [https://doi.org/10.1016/0747-7171\(92\)90011-R](https://doi.org/10.1016/0747-7171(92)90011-R) [cited on page 29]
- Dale Miller. 2000. Abstract Syntax for Variable Binders: An Overview. In *CL 2000: Computational Logic (Lecture Notes in Artificial Intelligence)*. [cited on pages 40 and 71]
- Robin Milner. 1978. A theory of type polymorphism in programming. *Journal of computer and system sciences* 17, 3 (1978), 348–375. [cited on pages 4, 5, 13, 16, 49, 106, and 110]
- John C. Mitchell. 1984. Coercion and Type Inference. In *Proceedings of the 11th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages* (Salt Lake City, Utah, USA) (POPL '84). Association for Computing Machinery, New York, NY, USA, 175–185. <https://doi.org/10.1145/800017.800529> [cited on pages 8 and 102]
- Wolfgang Naraschewski and Tobias Nipkow. 1999. Type inference verified: Algorithm W in Isabelle/HOL. *Journal of Automated Reasoning* 23, 3 (1999), 299–318. [cited on pages 9, 16, and 106]
- Tobias Nipkow, Lawrence C Paulson, and Markus Wenzel. 2002. *Isabelle/HOL: a proof assistant for higher-order logic*. Vol. 2283. Springer Science & Business Media. [cited on page 106]
- 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). [cited on pages 6, 7, 16, 18, 22, 28, 44, 83, and 101]
- Lionel Parreaux. 2020. The Simple Essence of Algebraic Subtyping: Principal Type Inference with Subtyping Made Easy (Functional Pearl). *Proc. ACM Program. Lang.* 4, ICFP, Article 124 (Aug. 2020), 28 pages. <https://doi.org/10.1145/3409006> [cited on page 105]

- Simon Peyton Jones and Mark Shields. 2004. Lexically-scoped type variables. (2004). <http://research.microsoft.com/en-us/um/people/simonpj/papers/scoped-tyvars/> Draft. [cited on pages 71 and 76]
- 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 7, 19, and 101]
- 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 page 7]
- Andrew M. Pitts. 2003. Nominal logic, a first order theory of names and binding. *Information and Computation* 186, 2 (2003), 165–193. [https://doi.org/10.1016/S0890-5401\(03\)00138-X](https://doi.org/10.1016/S0890-5401(03)00138-X) Theoretical Aspects of Computer Software (TACS 2001). [cited on page 10]
- François Pottier. 1998. *Type inference in the presence of subtyping: from theory to practice*. Ph.D. Dissertation. INRIA. [cited on pages 79 and 104]
- François Pottier and Didier Rémy. 2005. *Advanced Topics in Types and Programming Languages*. The MIT Press, Chapter The Essence of ML Type Inference, 387–489. [cited on pages 50 and 105]
- Jason Reed. 2009. Higher-order Constraint Simplification in Dependent Type Theory. In *Proceedings of the Fourth International Workshop on Logical Frameworks and Meta-Languages: Theory and Practice (LFMTP '09)*. [cited on pages 44, 50, and 105]
- Didier Remy, E De Recherche, Et En Automatique, Domaine De Voluceau, and Calcul Symbolique Programmation. 1992. Extension of ML Type System with a Sorted Equational Theory on Types. [cited on page 111]
- Didier Rémy and Boris Yakobowski. 2008. From ML to MLE: Graphic Type Constraints with Efficient Type Inference. *SIGPLAN Not.* 43, 9 (Sept. 2008), 63–74. <https://doi.org/10.1145/1411203.1411216> [cited on page 102]
- John C. Reynolds. 1974. Towards a theory of type structure. In *Programming Symposium*, B. Robinet (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 408–425. [cited on page 4]
- John C. Reynolds. 1983. Types, Abstraction and Parametric Polymorphism. In *Proceedings of the IFIP 9th World Computer Congress*. [cited on page 75]

- J C Reynolds. 1985. Three Approaches to Type Structure. In *Proc. of the International Joint Conference on Theory and Practice of Software Development (TAPSOFT) Berlin, March 25-29, 1985 on Mathematical Foundations of Software Development, Vol. 1: Colloquium on Trees in Algebra and Programming (CAAP'85)* (Berlin, Germany). Springer-Verlag, Berlin, Heidelberg, 97–138. [cited on page 8]
- Alejandro Serrano, Jurriaan Hage, Simon Peyton Jones, and Dimitrios Vytiniotis. 2020. A Quick Look at Impredicativity. *Proc. ACM Program. Lang.* 4, ICFP, Article 89 (Aug. 2020), 29 pages. <https://doi.org/10.1145/3408971> [cited on pages 6, 103, and 110]
- Alejandro Serrano, Jurriaan Hage, Dimitrios Vytiniotis, and Simon Peyton Jones. 2018. Guarded Impredicative Polymorphism. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2018)*. [cited on pages 6, 7, 103, and 110]
- Zhendong Su, Alexander Aiken, Joachim Niehren, Tim Priesnitz, and Ralf Treinen. 2002. The First-Order Theory of Subtyping Constraints. In *Proceedings of the 29th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Portland, Oregon) (POPL '02). Association for Computing Machinery, New York, NY, USA, 203–216. <https://doi.org/10.1145/503272.503292> [cited on page 81]
- The Coq development team. 2017. The Coq proof assistant. <https://coq.inria.fr/> [cited on page 106]
- Jerzy Tiuryn and Pawel Urzyczyn. 1996. The subtyping problem for second-order types is undecidable. In *Proceedings 11th Annual IEEE Symposium on Logic in Computer Science*. [cited on page 5]
- Valery Trifonov and Scott Smith. 1996. Subtyping constrained types. In *Static Analysis*, Radhia Cousot and David A. Schmidt (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 349–365. [cited on page 79]
- Christian Urban. 2008. Nominal techniques in Isabelle/HOL. *Journal of Automated Reasoning* 40, 4 (2008), 327–356. [cited on page 106]
- Christian Urban and Tobias Nipkow. 2008. Nominal verification of algorithm W. *From Semantics to Computer Science. Essays in Honour of Gilles Kahn* (2008), 363–382. [cited on pages 9 and 106]
- Dimitrios Vytiniotis, Simon Peyton Jones, Tom Schrijvers, and Martin Sulzmann. 2011. Outsidein(x) Modular Type Inference with Local Assumptions. *J. Funct. Program.* 21, 4–5

- (Sept. 2011), 333–412. <https://doi.org/10.1017/S0956796811000098> [cited on page 110]
- Dimitrios Vytiniotis, Stephanie Weirich, and Simon Peyton Jones. 2008. FPH: First-class Polymorphism for Haskell. In *Proceedings of the 13th ACM SIGPLAN International Conference on Functional Programming (ICFP '08)*. [cited on pages 6, 7, 103, and 110]
- 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 6, 15, 16, and 101]
- Ningning Xie and Bruno C. d. S. Oliveira. 2018. Let Arguments Go First. In *Programming Languages and Systems*, Amal Ahmed (Ed.). Springer International Publishing, Cham, 272–299. [cited on page 102]
- Jinxu Zhao, Bruno C. d. S. Oliveira, and Tom Schrijvers. 2018. Formalization of a Polymorphic Subtyping Algorithm. In *ITP (Lecture Notes in Computer Science, Vol. 10895)*. Springer, 604–622. [cited on page 12]
- Jinxu Zhao, Bruno C. d. S. Oliveira, and Tom Schrijvers. 2019. A Mechanical Formalization of Higher-Ranked Polymorphic Type Inference. *Proc. ACM Program. Lang.* 3, ICFP, Article 112 (July 2019), 29 pages. <https://doi.org/10.1145/3341716> [cited on page 12]