

Local Type Argument Synthesis with Bounded Quantification

Benjamin C. Pierce
Computer Science Department
Indiana University
Lindley Hall 215
Bloomington, IN 47405, USA
pierce@cs.indiana.edu

David N. Turner
An Teallach Limited
Technology Transfer Center
King's Buildings
Edinburgh, EH9 3JL, UK
dnt@an-teallach.com

Indiana University
CSCI Technical Report #495

November 12, 1997

Abstract

In a companion paper [PT98], we introduced a *local type inference* method for inferring type arguments to polymorphic functions. We show here how our method can be extended to handle Cardelli and Wegner's Kernel Fun variant of F_{\leq} .

1 Introduction

Most statically typed programming languages offer some form of *type inference*, allowing programmers to omit type annotations that can be recovered from context. Such a facility can eliminate a great deal of needless verbosity, making programs easier both to read and to write. Unfortunately, type inference technology has not kept pace with developments in type systems. In particular, the combination of subtyping and parametric polymorphism has been intensively studied for more than a decade in calculi such as System F_{\leq} [CW85, CG92, CMMS94, etc.], but these features have not yet been satisfactorily integrated with practical type inference methods. Part of the reason for this gap is that most work on type inference for this class of languages has concentrated on the difficult problem of developing *complete* methods, which are guaranteed to infer types, whenever possible, for entirely unannotated programs. In a companion paper [PT98], we proposed a simple alternative, refining the idea of *partial* type inference with the additional simplifying principle that missing annotations should be recovered using only types propagated *locally*, from adjacent nodes in the syntax tree.

In [PT98], we presented our type inference scheme for a calculus incorporating both impredicative polymorphism and subtyping. For simplicity, we did not treat *bounded quantification*, where subtyping and quantification interact via upper bounds on type variables. This paper rectifies that shortcoming, showing how our type inference method can be extended to handle Cardelli and Wegner's Kernel Fun variant of F_{\leq} .

There is one caveat: we make some restrictions on the kinds of polymorphic functions we automatically infer type arguments for. In particular, we have so far been unable to deal with inter-dependent bounds: we do not know of a complete algorithm which can synthesize, for example, the type arguments for a function of type $\text{All } (X < \text{Top}, Y < X) S \rightarrow T$. Rather than introduce a potentially unimplementable rule in the specification of type inference, we explicitly disallow this case in our specification: the user must always write explicit type arguments on applications of such functions.

We concentrate here on the technical details of how to extend our local type inference scheme to allow bounded quantification. We refer the reader to [PT98] for a fuller explanation of the basic local type inference scheme. In the next section, we define (a mild extension of) Cardelli and Wegner’s Kernel Fun. Section 3 then shows how we can specify and implement a local type argument inference scheme for this calculus.

2 Kernel Fun

For our explicitly-typed internal language—the target for the type inference methods described in Section 3—we extend Cardelli and Wegner’s Kernel Fun calculus [CW85] of subtyping and impredicative polymorphism. We only give definitions here; the meta-theory of the system is developed in detail in a companion paper [Pie97].

2.1 Syntax

We extend the original Kernel Fun system [CW85] in a few significant ways.

Firstly, we add a minimal type **Bot**. Our type inference algorithm keeps track of various type constraints by calculating the least upper bound and greatest lower bound of pairs of types. The **Bot** type plays a crucial role in these calculations, since without it we could not guarantee that least upper-bounds and greatest lower-bounds always exist. The properties of Kernel Fun with **Bot** are similar to those of pure Kernel Fun, but there are a number of significant differences in details. The properties of the internal language are developed in detail in [Pie97].

Secondly, we extend abstraction and application so that several arguments (including both types and terms) may be passed at the same time. In other words, we favor a “fully uncurried” style of function definition and application (though currying is, of course, still available). This bias will play an important role in our scheme for inferring type arguments in Section 3.

The syntax of types, terms, and typing contexts in the internal language is as follows:

T	::=	X	type variable
		Top	maximal type
		Bot	minimal type
		$\text{All}(\bar{X} <: \bar{T}) \bar{T} \rightarrow T$	function type
e	::=	x	variable
		$\text{fun}[\bar{X} <: \bar{T}] (\bar{x} : \bar{T}) e$	abstraction
		$e[\bar{T}] (\bar{e})$	application
$,$::=	\bullet	empty context
		$, , x : T$	variable binding
		$, , X <: T$	type variable binding

We use the meta-variables A, B, R, S, T and U to range over types; e and f range over terms. We use the notation \bar{X} to denote the sequence X_1, \dots, X_n , and similarly $\bar{x} : \bar{T}$ to denote $x_1 : T_1, \dots, x_n : T_n$. We write $, (X)$ for the bound of X in $,$ and $, (x)$ for the type of x in $,$.

We write $\bar{S} \rightarrow T$ as an abbreviation for the monomorphic function type $\text{All}() \bar{S} \rightarrow T$. Similarly, we write $\text{fun}(\bar{x} : \bar{T}) e$ as an abbreviation for the monomorphic function $\text{fun}[] (\bar{x} : \bar{T}) e$.

Types, terms, contexts, and judgments that differ only in the names of bound variables are regarded as identical. Binders in contexts are assumed to have distinct names; when a new binding is added to a context, we assume that it has been renamed so as to maintain this invariant. The rules for scoping of bound variables are as usual (in $\text{All}(\bar{X} <: \bar{S}) \bar{T} \rightarrow U$, the scope of each X_i is everything to its right: the bounds S_{i+1} through S_n , all of \bar{T} , and U). $FV(T)$, the set of type variables free in T , is defined in the usual way.

2.2 Subtyping

The subtyping relation is the evident extension of Kernel Fun’s subtyping relation with the `Bot` type. We write $\cdot, \bar{S} <: \bar{T}$ to mean “ $|\bar{S}| = |\bar{T}|$ and $\cdot, \vdash S_i <: T_i$ for all $1 \leq i \leq |S|$.”

$$\cdot, \vdash T <: \text{Top} \quad (\text{S-TOP})$$

$$\cdot, \vdash \text{Bot} <: T \quad (\text{S-BOT})$$

$$\cdot, \vdash X <: X \quad (\text{S-REFL})$$

$$\frac{\cdot, \vdash \cdot, (X) <: T}{\cdot, \vdash X <: T} \quad (\text{S-VAR})$$

$$\frac{\cdot, \bar{X} <: \bar{B} \vdash \bar{T} <: \bar{R} \quad \cdot, \bar{X} <: \bar{B} \vdash S <: U}{\cdot, \vdash \text{All}(\bar{X} <: \bar{B}) \bar{R} \rightarrow S <: \text{All}(\bar{X} <: \bar{B}) \bar{T} \rightarrow U} \quad (\text{S-FUN})$$

Two points are worth noting. First, we use the original “Kernel Fun” rule for comparing quantifiers [CW85], in which the upper bounds \bar{B} in the subtyping rule for polymorphic functions are required to be identical, rather than the more powerful but less tractable variant of Curien and Ghelli [CG92, CMMS94].¹ The principal reason for this restriction is that it allows us to define meets and joins of all pairs of types, which may fail to exist in “Full F_{\leq} ” [Ghe90]. Second, for simplicity, we use an algorithmic presentation of subtyping, in which the rules of transitivity and general reflexivity are omitted and recovered as properties of the definition (cf. [Pie97, Section 3.1]).

It is also important to note that some of the usual properties of presentations of Kernel Fun without `Bot` do not hold here. For instance, $\cdot, \vdash S <: T$ and $\cdot, \vdash T <: S$ do not imply $S = T$ (consider, for example, $X <: \text{Bot} \vdash X <: \text{Bot}$ and $X <: \text{Bot} \vdash \text{Bot} <: X$). This fact is the result of an unfortunate interaction between bounded quantification and `Bot`, and it substantially complicates the proofs of the properties in the remainder of this section as well as the development in Section 3.

We write $\cdot, \vdash S \uparrow T$ for the *least non-variable supertype* of S , defined by repeated promotion of variables:

$$\frac{S \text{ is not a variable}}{\cdot, \vdash S \uparrow S}$$

$$\frac{\cdot, \vdash \cdot, (X) \uparrow T}{\cdot, \vdash X \uparrow T}$$

We write $\cdot, \vdash S \wedge T = M$ for “ M is the meet of S and T in context \cdot ,” and $\cdot, \vdash S \vee T = J$ for “ J is the join of S and T in \cdot .” The definitions of these relations can be found in [Pie97], Section 3.3.

2.3 Typing Rules

The typing relation $\cdot, \vdash e \in T$ is essentially the standard one, except that, as in the definition of subtyping, we use an algorithmic presentation, omitting the usual rule of subsumption (“if $e \in S$ and $S <: T$, then $e \in T$ ”); instead, the rules below calculate for each typable term a single *manifest type*, corresponding to its minimal type in the system with subsumption. For subtyping, the choice of algorithmic presentation was made for the sake of simplicity. Here, it is actually crucial: our type inference methods depend on the fact that a typable term has a unique type, and that this type can easily be predicted by the programmer. (Note that this stylistic choice does not change the set of typable terms.)

In the rules, we write $\cdot, \vdash e \uparrow T$ to abbreviate $\cdot, \vdash e \in S$ and $\cdot, \vdash S \uparrow T$; similarly, $\cdot, \vdash e \in S <: T$ abbreviates $\cdot, \vdash e \in S$ and $\cdot, \vdash S <: T$.

¹A variant on the rule used here would require that the upper bounds be *equivalent*—i.e., each a subtype of the other. Choosing this variant appears to make some of the following development simpler and other parts more complex, sometimes substantially so. It is not clear to us which is globally better.

The typing rule for variables is standard.

$$, \vdash x \in , (x) \quad (\text{T-VAR})$$

The rule for (multi-)abstractions combines the usual rules for term and type abstractions.

$$\frac{, , \bar{x} <: \bar{B}, \bar{x} : \bar{S} \vdash e \in T}{, \vdash \text{fun } [\bar{x} <: \bar{B}] (\bar{x} : \bar{S}) e \in \text{All } (\bar{x} <: \bar{B}) \bar{S} \rightarrow T} \quad (\text{T-ABS})$$

Similarly, the rule for (multi-)applications combines the usual application and polymorphic application rules of Kernel Fun. We calculate the type of the function f and promote it (if it is a variable), to an explicit function type $\text{All } (\bar{x} <: \bar{B}) \bar{S} \rightarrow R$. We then check that the provided type arguments are subtypes of the upper bounds \bar{B} and that the provided term arguments have the expected types. If all these constraints are satisfied, then the result type of the application is found by substituting the actual type arguments into R .

$$\frac{, \vdash f \uparrow \text{All } (\bar{x} <: \bar{B}) \bar{S} \rightarrow R \quad , \vdash T_i <: [T_1/X_1 \dots T_{i-1}/X_{i-1}] B_i \quad , \vdash \bar{e} \in \bar{U} <: [\bar{T}/\bar{X}] \bar{S}}{, \vdash f [\bar{T}] (\bar{e}) \in [\bar{T}/\bar{X}] R} \quad (\text{T-APP})$$

To finish the definition of the typing relation, another rule is required. To see why, note that $\text{Bot} <: \text{All } (\bar{x} <: \bar{B}) \bar{S} \rightarrow T$ for any $\bar{x}, \bar{B}, \bar{S}$, and T . This means that any expression of type Bot should be applicable to any set of well-formed type and expression arguments (if we did not allow for this behavior, we would lose the type soundness property):

$$\frac{, \vdash f \uparrow \text{Bot} \quad , \vdash \bar{e} \in \bar{S}}{, \vdash f [\bar{T}] (\bar{e}) \in \text{Bot}} \quad (\text{T-BOT})$$

Note that this rule gives the expression $f [\bar{T}] (\bar{e})$ result type Bot , the most informative result type for the expression.

2.3.1 Theorem [Uniqueness of manifest types]: Each typable term has a unique type, up to equivalence: if $, \vdash e \in S$ and $, \vdash e \in T$, then $, \vdash S <: T$ and $, \vdash T <: S$. (Indeed, given the definitions we have seen so far, it is even the case that $S = T$; we state the proposition in this weaker version because later we will want to guarantee uniqueness of inferred types only up to equivalence in the subtype relation.)

The definitions of operational and denotational semantics for the above calculus are standard, as are proofs of properties such as subject reduction and absence of runtime errors. Evaluation order may be chosen either call-by-name or call-by-value; function spaces may be interpreted as either total or partial. The only slightly unusual case is the type Bot , which can be interpreted as an empty type (in a total-function semantics) or a type containing only divergent terms (in a partial function semantics).

3 Local Type Argument Synthesis

Our measurements of ML programs in [PT98] showed that type arguments to polymorphic functions are inferred by the ML typechecker on at least one line in every three, in typical programs. Moreover, explicit type arguments rarely have any useful documentation value. As an example, consider the polymorphic identity function id with type $\text{All } (X <: \text{Top}) X \rightarrow X$. Our goal is to allow the programmer to apply the id function without explicitly supplying any type arguments: $\text{id}(3)$ rather than $\text{id}[\text{Int}](3)$.

The main problem we have to address is the fact that, in general, there may be a number of different type arguments that we can pick for a particular function application. For example, both $\text{id}[\text{Int}](x)$ and $\text{id}[\text{Real}](x)$ are valid completions of the term $\text{id}(x)$, where $x \in \text{Int}$ and Int is a subtype of Real . Fortunately, there is usually a good way to choose between all the alternatives: we pick the type arguments

that yield the best (smallest) type for the result. In the case of $\text{id}(x)$, we choose $\text{id}[\text{Int}](x)$ since this has result type Int , which is more informative type than the result type Real of $\text{id}[\text{Real}](x)$.

Sadly, there are cases where there is no best result type. Suppose, for example, that f has type $\text{All}(X<:\text{Top})() \rightarrow (X \rightarrow X)$ (a function which takes a single type argument X and returns a function from X to X). Two possible completions of the term $f()$ are $f[\text{Int}]()$ and $f[\text{Real}]()$, which have result types $\text{Int} \rightarrow \text{Int}$ and $\text{Real} \rightarrow \text{Real}$. These two result types are incomparable in the subtyping relation, so there is no “best” result type available. In this case type argument synthesis will fail, since it is not possible to locally determine the missing type arguments for f (in [PT98] we show how propagating additional contextual information sometimes allows us to avoid this situation).

3.1 Type Inference (Specification)

We define a three-place *type inference* relation:

$$, \vdash e \in T \Rightarrow e'$$

Intuitively, this relation can be read “In context $,$, type annotations can be added to the term e to yield the fully-typed term e' , which has type T .”

The specification of the type inference relation is quite simple. For each typing rule with conclusion $, \vdash e \in T$ in the explicitly-typed language, the type inference relation contains an analogous rule with conclusion $, \vdash e \in T \Rightarrow e'$, where e' is derived in the obvious way from the fully typed subexpressions yielded by subderivations. To these rules is added one additional rule, handling the case where type arguments are omitted:

$$\frac{\begin{array}{l} , \vdash f \uparrow \text{All}(\bar{X}<:\bar{S})\bar{T} \rightarrow R \Rightarrow f' \quad , \vdash \bar{e} \in \bar{U} \Rightarrow \bar{e}' \\ |\bar{X}| > 0 \quad \bar{X} \cap FV(\bar{S}) = \emptyset \quad , \vdash \bar{A} <: \bar{S} \quad , \vdash \bar{U} <: [\bar{A}/\bar{X}]\bar{T} \\ \forall \bar{B}. (, \vdash \bar{B} <: \bar{S} \text{ and } , \vdash \bar{U} <: [\bar{B}/\bar{X}]\bar{T} \text{ imply } , \vdash [\bar{A}/\bar{X}]R <: [\bar{B}/\bar{X}]R) \end{array}}{\quad , \vdash f(\bar{e}) \in [\bar{A}/\bar{X}]R \Rightarrow f'[\bar{A}](\bar{e}')} \quad (\text{APP-INF SPEC})$$

The condition $|\bar{X}| > 0$ says that type argument synthesis is only required in the case where the function f is polymorphic but there are no explicit type arguments. (For simplicity, we don’t synthesize type arguments in the case where an application node provides some, but not all, of its required type arguments explicitly. This would be easy to do, but does not seem very useful in practice.) The condition $\bar{X} \cap FV(\bar{S}) = \emptyset$ explicitly disallows type argument synthesis in the case where the bounds \bar{S} are inter-dependent (since, at this time, we do not know of a complete solution to this problem).

The type arguments \bar{A} that we pick as the result of our synthesis rule must satisfy a number of conditions. Firstly, they must be legal type arguments for f . (The condition $, \vdash \bar{A} <: \bar{S}$ ensures that the arguments are subtypes of the required bounds \bar{S} , while the condition $, \vdash \bar{U} <: [\bar{A}/\bar{X}]\bar{T}$ ensures that the types of the argument expressions match the types of the function parameters.) Secondly, the final line of the rule asserts that the arguments \bar{A} must be chosen in such a way that any other choice of arguments \bar{B} satisfying the above conditions will yield a less informative result type, i.e., a supertype of $[\bar{A}/\bar{X}]R$.

To state the formal properties of this technique, we need to relate terms in the internal language to terms in the external language. We say that a term e is a *partial erasure* of e' if e can be obtained from e' by erasing some type annotations (i.e., deleting type arguments from one or more applications).

3.1.1 Theorem [Soundness]:

If $, \vdash e \in T \Rightarrow e'$, then e is a partial erasure of e' and $, \vdash e' \in T$.

Proof: Straightforward from the definition. ■

Since we are dealing with a partial type inference technique, we cannot expect a completeness property at this point. However, we can show that the type inference relation is “*locally complete*” in the sense that its specification guarantees that it will find the best values for missing type arguments in a single application, whenever these exist.

3.1.2 Theorem [Partial Completeness]: If $\Gamma, \vdash e \in \mathsf{T}$ (i.e., e is fully typed), then $\Gamma, \vdash e \in \mathsf{T} \Rightarrow e$.

Proof: Straightforward since, in inferring types for a fully typed program, the local type argument synthesis rule can never be used. \blacksquare

It should be emphasized that the APP-INFSPEC rule, together with the other typing rules from the explicitly-typed language, constitutes a complete specification of the type inference relation: it is all that a programmer needs to understand in order to use the language. Only the compiler writer needs to go further into the development in the rest of the section, whose job is to show how the rule we have given can be implemented.

3.2 Variable Elimination

In the constraint-generation algorithm that we present in the next section, it will sometimes be necessary to eliminate all occurrences of a certain set of variables from a given type by promoting or demoting the type until we reach a type in which these variables do not occur. (This property is another crucial reason for choosing the “Kernel Fun” variant of F_{\leq} rather than the “full F_{\leq} ” variant where two polymorphic function types with different upper bounds for their type components are allowed to stand in the subtype relation under appropriate conditions; in the latter system, it can be shown that variables cannot always be eliminated in a most general way [GP97].)

Formally, we write $\Gamma, \vdash S \uparrow^V T$ for the relation “ T is the least supertype of S such that $FV(T) \cap V = \emptyset$ ” and $\Gamma, \vdash S \downarrow_V T$ for the dual relation “ T is the greatest subtype of S such that $FV(T) \cap V = \emptyset$.” The variable-elimination-by-promotion relation can be computed as follows:

$$\Gamma, \vdash \mathsf{Top} \uparrow^V \mathsf{Top} \quad (\text{VU-TOP})$$

$$\Gamma, \vdash \mathsf{Bot} \uparrow^V \mathsf{Bot} \quad (\text{VU-BOT})$$

$$\frac{x \in V \quad \Gamma, \vdash (x) \uparrow^V T}{\Gamma, \vdash x \uparrow^V T} \quad (\text{VU-VAR-1})$$

$$\frac{x \notin V}{\Gamma, \vdash x \uparrow^V x} \quad (\text{VU-VAR-2})$$

$$\frac{FV(\bar{A}) \cap V = \emptyset \quad \Gamma, \bar{x} <: \bar{A} \vdash \bar{S} \downarrow_V \bar{S}' \quad \Gamma, \bar{x} <: \bar{A} \vdash T \uparrow^V T'}{\Gamma, \vdash \mathsf{All}(\bar{x} <: \bar{A}) \bar{S} \rightarrow T \uparrow^V \mathsf{All}(\bar{x} <: \bar{A}) \bar{S}' \rightarrow T'} \quad (\text{VU-FUN-1})$$

$$\frac{FV(\bar{A}) \cap V \neq \emptyset}{\Gamma, \vdash \mathsf{All}(\bar{x} <: \bar{A}) \bar{S} \rightarrow T \uparrow^V \mathsf{Top}} \quad (\text{VU-FUN-2})$$

The definition of $\Gamma, \vdash S \downarrow_V T$ is similar to $\Gamma, \vdash S \uparrow^V T$:

$$\Gamma, \vdash \mathsf{Top} \downarrow_V \mathsf{Top} \quad (\text{VD-TOP})$$

$$\Gamma, \vdash \mathsf{Bot} \downarrow_V \mathsf{Bot} \quad (\text{VD-BOT})$$

$$\frac{x \in V}{\Gamma, \vdash x \downarrow_V \mathsf{Bot}} \quad (\text{VD-VAR-1})$$

$$\frac{x \notin V}{\Gamma, \vdash x \downarrow_V x} \quad (\text{VD-VAR-2})$$

$$\frac{FV(\bar{A}) \cap V = \emptyset \quad , \bar{X} <: \bar{A} \vdash \bar{S} \uparrow^V \bar{S}' \quad , \bar{X} <: \bar{A} \vdash T \downarrow_V T'}{\quad , \vdash \text{All}(\bar{X} <: \bar{A}) \bar{S} \rightarrow T \downarrow_V \text{All}(\bar{X} <: \bar{A}) \bar{S}' \rightarrow T'} \quad (\text{VD-FUN-1})$$

$$\frac{FV(\bar{A}) \cap V \neq \emptyset}{\quad , \vdash \text{All}(\bar{X} <: \bar{A}) \bar{S} \rightarrow T \downarrow_V \text{Bot}} \quad (\text{VD-FUN-2})$$

It is easy to check that, for each variable set V , \uparrow^V and \downarrow_V are total functions. (These functions are similar to the ones used in [GP97], but somewhat simpler because of the presence of **Bot** in our type system.)

3.2.1 Lemma [Soundness of variable elimination]:

1. If $\vdash S \uparrow^V T$ then $FV(T) \cap V = \emptyset$ and $\vdash S <: T$.
2. If $\vdash S \downarrow_V T$ then $FV(T) \cap V = \emptyset$ and $\vdash T <: S$.

Proof: By a straightforward simultaneous induction on variable-elimination derivations. ■

3.2.2 Lemma [Completeness of variable elimination]:

1. If $\vdash S <: T$ and $FV(T) \cap V = \emptyset$, then $\vdash S \uparrow^V R$ with $\vdash R <: T$.
2. If $\vdash T <: S$ and $FV(T) \cap V = \emptyset$, then $\vdash S \downarrow_V R$ with $\vdash T <: R$.

Proof: See [Pie97]. ■

3.3 Constraints

Next, we introduce the constraints that will be manipulated by our algorithm. We need constraints of two forms, one for recording the fact that a type variable X must be exactly equal to some type T (for example, X must be exactly equal to **Bot** in order to make $\text{All}(Y <: X) Y \rightarrow Y$ a subtype of $\text{All}(Y <: \text{Bot}) Y \rightarrow Y$), and the other for recording the fact that a variable X must lie between two types S and T (for example, X must lie between A and B in order to make $X \rightarrow X$ a subtype of $A \rightarrow B$).

Formally, a V -constraint has one of the forms

$$\begin{array}{ll} [T] & \text{equality constraint} \\ [S, T] & \text{subtyping constraint} \end{array}$$

with the additional constraint that all the free variables of S and T are distinct from V . A type R *satisfies* a constraint c , written $\vdash R \in c$, if

$$\begin{array}{ll} c = [S] & \text{and } R = S \\ \text{or } c = [S, T] & \text{and } \vdash S <: R \text{ and } \vdash R <: T. \end{array}$$

The *maximal* and *minimal* types satisfying a given constraint are defined in the obvious way:

$$\begin{array}{ll} \max([S, T]) = T & \max([S]) = S \\ \min([S, T]) = S & \min([S]) = S \end{array}$$

An \bar{X}/V -constraint set C is a finite map from \bar{X} to V -constraints. The empty \bar{X}/V -constraint set, written \emptyset , maps each variable X_i to the constraint $[\text{Bot}, \text{Top}]$. The singleton \bar{X}/V -constraint set $\{X_i \mapsto c\}$ maps X_i to the constraint c and every other X_j to $[\text{Bot}, \text{Top}]$. The *meet* of two V -constraints is defined as follows (for all cases other than those specified below, the meet is undefined):

$$\begin{array}{ll} [S] \wedge [S] & = [S] \\ [S] \wedge [U, V] & = [S] \quad \text{if } \vdash U <: S <: V \\ [S, T] \wedge [U] & = [U] \quad \text{if } \vdash S <: U <: T \\ [S, T] \wedge [U, V] & = [J, M] \quad \text{if } \vdash S \vee U = J \text{ and } \vdash T \wedge V = M \end{array}$$

The meet operation is extended pointwise to constraint sets.

$$(C \wedge D)(X_i) = C(X_i) \wedge D(X_i)$$

We write $\bar{C} \wedge \bar{D}$ to abbreviate $C_1 \wedge \dots \wedge C_m \wedge D_1 \wedge \dots \wedge D_n$.

3.4 Constraint Generation

Our constraint generation rules have the form

$$, \vdash_{\bar{X}}^V S <: T \Rightarrow C$$

and define a partial function that, given a typing context $, ,$ a set of type variables V , a set of unknowns \bar{X} , and two types S and T , calculates the minimal \bar{X}/V -constraint set C guaranteeing that $, \vdash S <: T$.

The set V allows us to avoid generating nonsensical constraint sets in which bound variables are mentioned outside their scopes (this part of the constraint generation problem is similar to *mixed-prefix unification* [Mil92]). For example, if we are interested in constraining X so that $\text{All}(Y <: \text{Top}) () \rightarrow (Y \rightarrow Y)$ is a subtype of $\text{All}(Y <: \text{Top}) () \rightarrow X$, we should not return the constraint set $\{X \mapsto [Y \rightarrow Y, \text{Top}]\}$, since Y would be out of scope. Instead, we should return the constraint set $\{X \mapsto [\text{Bot} \rightarrow \text{Top}, \text{Top}]\}$, which is in fact the weakest constraint on X guaranteeing that $\text{All}(Y <: \text{Top}) () \rightarrow (Y \rightarrow Y)$ is a subtype of $\text{All}(Y <: \text{Top}) () \rightarrow X$.

Our constraint generation algorithm is defined by the following collection of rules. In the definition, we suppose that $\bar{X} \cap V = \emptyset$. More importantly, we assume (and recursively maintain) the invariant that only one of S and T mentions the variables \bar{X} (i.e. either $FV(S) \cap \bar{X} = \emptyset$ or $FV(T) \cap \bar{X} = \emptyset$). This is crucial to the completeness of the algorithm, since it ensures we only have to solve a matching problem (modulo subtyping) rather than a unification problem.

$$, \vdash_{\bar{X}}^V T <: \text{Top} \Rightarrow \emptyset$$

$$, \vdash_{\bar{X}}^V \text{Bot} <: T \Rightarrow \emptyset$$

$$\frac{Y \in \bar{X} \quad T \Downarrow_V R}{, \vdash_{\bar{X}}^V Y <: T \Rightarrow \{Y \mapsto [\text{Bot}, R]\}}$$

$$\frac{Y \in \bar{X} \quad T \Uparrow^V R}{, \vdash_{\bar{X}}^V T <: Y \Rightarrow \{Y \mapsto [R, \text{Top}]\}}$$

$$, \vdash_{\bar{X}}^V Y <: Y \Rightarrow \emptyset$$

$$\frac{, \vdash_{\bar{X}}^V (Y) <: T \Rightarrow C}{, \vdash_{\bar{X}}^V Y <: T \Rightarrow C}$$

$$\frac{\begin{array}{l} , \vdash_{\bar{X}}^V \bar{A} \equiv \bar{B} \Rightarrow \bar{K}, \bar{D} \quad V' = V \cup \{\bar{Y}\} \\ \bar{Y} \cap V = \emptyset \quad \bar{Y} \cap \bar{X} = \emptyset \quad FV(\sigma) \cap \bar{Y} = \emptyset \\ , \bar{Y} <: \bar{K} \vdash_{\bar{X}}^{V'} \bar{T} <: \bar{R} \Rightarrow \bar{C} \quad , \bar{Y} <: \bar{K} \vdash_{\bar{X}}^{V'} S <: U \Rightarrow D \end{array}}{, \vdash_{\bar{X}}^V \text{All}(\bar{Y} <: \bar{A}) \bar{R} \rightarrow S <: \text{All}(\bar{Y} <: \bar{B}) \bar{T} \rightarrow U \Rightarrow D \wedge \bar{D} \wedge \bar{C}}$$

In the clause for quantifiers (whose bounds must match exactly rather than modulo subtyping), we need an auxiliary “matching relation” $, \vdash_{\bar{X}}^V S \equiv T \Rightarrow U, C$, which yields both a constraint set C whose solutions make S and T identical and a type U that is equal to whichever of S and T is concrete (recall that the variables \bar{X} do not occur in one of S or T). The definition of this relation follows the same lines as the main constraint generator:

$$, \vdash_{\bar{X}}^V \text{Top} \equiv \text{Top} \Rightarrow \text{Top}, \emptyset$$

$$, \vdash_{\bar{X}}^V \text{Bot} \equiv \text{Bot} \Rightarrow \text{Bot}, \emptyset$$

$$\frac{Y \in \bar{X} \quad V \cap FV(T) = \emptyset}{, \vdash_{\bar{X}}^V Y \equiv T \Rightarrow T, \{Y \mapsto [T]\}}$$

$$\frac{Y \in \bar{X} \quad V \cap FV(T) = \emptyset}{, \vdash_{\bar{X}}^V T \equiv Y \Rightarrow T, \{Y \mapsto [T]\}}$$

$$\begin{array}{c}
\frac{Y \notin \bar{X}}{\text{, } \vdash_{\bar{X}} Y \equiv Y \Rightarrow Y, \emptyset} \\
\text{, } \vdash_{\bar{X}} \bar{A} \equiv \bar{B} \Rightarrow \bar{K}, \bar{D} \quad V' = V \cup \{\bar{Y}\} \\
\bar{Y} \cap V = \emptyset \quad \bar{Y} \cap \bar{X} = \emptyset \quad FV(\sigma) \cap \bar{Y} = \emptyset \\
\text{, } \bar{Y} \langle: \bar{K} \vdash_{\bar{X}'} \bar{T} \equiv \bar{S} \Rightarrow \bar{L}, \bar{C} \quad \text{, } \bar{Y} \langle: \bar{K} \vdash_{\bar{X}'} U \equiv V \Rightarrow M, D \\
\hline
\text{, } \vdash_{\bar{X}} \text{All}(\bar{Y} \langle: \bar{A}) \bar{R} \rightarrow S \equiv \text{All}(\bar{Y} \langle: \bar{B}) \bar{T} \rightarrow U \Rightarrow \text{All}(\bar{Y} \langle: \bar{K}) \bar{L} \rightarrow M, D \wedge \bar{D} \wedge \bar{C}
\end{array}$$

Note that the C returned by these algorithms is always an \bar{X}/V -constraint set. Also, if $\text{, } \vdash_{\bar{X}} S \langle: T \Rightarrow C$ and the variables \bar{X} do not appear in S or T , then the constraint set C is always empty. The constraint generator in this case is effectively just the subtyping relation.

3.5 Soundness and Completeness of Constraint Generation

Each constraint set returned by the constraint generator characterizes a collection of substitutions associating concrete types with the names of the missing type parameters. Formally, a *substitution* is a finite map from type variables to types. More precisely, an \bar{X}/V -substitution is one whose domain is \bar{X} and whose codomain contains types whose free variables are distinct from V . We write $\sigma[X_i \mapsto T]$ for the substitution that behaves like σ everywhere except at X_i , where its value is T .

We say that an \bar{X}/V -substitution σ *satisfies* an \bar{X}/V -constraint set C , written $\text{, } \vdash \sigma \in C$, iff $\text{, } \vdash \sigma(X_i) \in C(X_i)$ for each variable X_i . A constraint set is called *satisfiable* if there is some substitution that satisfies it. (Note that this condition can be checked very easily, by verifying that $\text{, } \vdash \min(C(X_i)) \langle: \max(C(X_i))$ for each X_i .)

If C and D are two \bar{X}/V -constraint sets such that $\text{, } \vdash \sigma \in C$ implies $\text{, } \vdash \sigma \in D$ for all σ , we say that C is *more demanding than* D . Note that the meet of constraint sets defined above yields a greatest lower bound in this ordering and that the empty constraint set is maximal (i.e., least demanding).

Before we can prove soundness and completeness for the constraint generator, we need analogous lemmas for the auxiliary “matching constraint” generator.

3.5.1 Lemma [Soundness of matching constraint generation]: Suppose that either $FV(S) \cap \bar{X} = \emptyset$ or $FV(T) \cap \bar{X} = \emptyset$. If $\text{, } \vdash_{\bar{X}} S \equiv T \Rightarrow U, C$ and $\sigma \in C$, then $\sigma S = \sigma T = U$.

Proof: Straightforward induction. ■

3.5.2 Lemma [Completeness of matching constraint generation]: Let σ be an \bar{X}/V -substitution with $\bar{X} \cap V = \emptyset$, and let S and T be types such that either $FV(S) \cap \bar{X} = \emptyset$ or $FV(T) \cap \bar{X} = \emptyset$. If $\sigma S = \sigma T$, then $\text{, } \vdash_{\bar{X}} S \equiv T \Rightarrow \sigma S, C$ for some C such that $\text{, } \vdash \sigma \in C$.

Proof: By induction on the structure of $\sigma S (= \sigma T)$. We only give the most interesting cases of the proof. The remaining cases follow easily using the induction hypothesis and, for the function case, the fact that $\text{, } \vdash \sigma \in C$ and $\text{, } \vdash \sigma \in D$ implies $\text{, } \vdash \sigma \in C \wedge D$.

Case: $S = Y$ where $Y \in \bar{X}$

It must be the case that $FV(T) \cap \bar{X} = \emptyset$, since $Y \in \bar{X}$. We therefore have $\sigma Y = \sigma T = T$. It must also be the case that $FV(T) \cap V = \emptyset$, since T occurs in the codomain of σ and σ is a \bar{X}/V -substitution. We therefore have $\text{, } \vdash_{\bar{X}} S \equiv T \Rightarrow T, \{Y \mapsto [T]\}$ and $\text{, } \vdash \sigma \in \{Y \mapsto [T]\}$ as required.

Case: $T = Y$ where $Y \in \bar{X}$

Similar. ■

3.5.3 Proposition [Soundness of constraint generation]: Suppose that $FV(\text{, }) \cap \bar{X} = \emptyset$, that $\text{dom}(\text{, }) \cap \bar{X} = \emptyset$, and that either $FV(S) \cap \bar{X} = \emptyset$ or $FV(T) \cap \bar{X} = \emptyset$. If $\text{, } \vdash_{\bar{X}} S \langle: T \Rightarrow C$ and $\text{, } \vdash \sigma \in C$, then $\text{, } \vdash \sigma S \langle: \sigma T$.

Proof: By induction on the derivation of $\text{, } \vdash_{\bar{X}} S \langle: T \Rightarrow C$.

Case: $\text{, } \vdash_{\bar{X}} S \langle: \text{Top} \Rightarrow \emptyset$

Immediate, since $\sigma\text{Top} = \text{Top}$ and $\cdot, \vdash \sigma S <: \text{Top}$, no matter what σS is.

Case: $\cdot, \vdash_{\bar{X}}^V \text{Bot} <: T \Rightarrow \emptyset$

Immediate, since $\sigma\text{Bot} = \text{Bot}$ and $\cdot, \vdash \text{Bot} <: \sigma T$ no matter what σT is.

Case: $\cdot, \vdash_{\bar{X}}^V Y <: T \Rightarrow C$ where $Y \in \bar{X}$, $T \Downarrow_V R$ and $C = \{Y \mapsto [\text{Bot}, R]\}$

Since $\cdot, \vdash \sigma \in C$ we have $\cdot, \vdash \sigma Y <: R$. Since $Y \in \bar{X}$, we know that $\sigma T = T$; also, Lemma 3.2.1(2) tells us that $\cdot, \vdash R <: T$. We therefore have $\cdot, \vdash \sigma Y <: R <: T = \sigma T$, as required.

Case: $\cdot, \vdash_{\bar{X}}^V S <: Y \Rightarrow C$ where $Y \in \bar{X}$, $S \Uparrow^V R$ and $C = \{Y \mapsto [R, \text{Top}]\}$

Since $\cdot, \vdash \sigma \in C$ we have $\cdot, \vdash R <: \sigma Y$. Since $Y \in \bar{X}$, we know that $\sigma S = S$; also, Lemma 3.2.1(1) tells us that $\cdot, \vdash S <: R$. We therefore have $\cdot, \vdash \sigma S = S <: R <: \sigma Y$, as required.

Case: $\cdot, \vdash_{\bar{X}}^V Y <: Y \Rightarrow \emptyset$

Since, by assumption, the variables \bar{X} do not appear free in both S and T , it must be the case that $Y \notin \bar{X}$. Thus, $\sigma Y = Y$ and $\cdot, \vdash Y <: Y$, as required.

Case: $\cdot, \vdash_{\bar{X}}^V Y <: T \Rightarrow C$ where $\cdot, \vdash_{\bar{X}}^V (Y) <: T \Rightarrow C$

Using the induction hypothesis, we obtain $\cdot, \vdash \sigma((Y)) <: \sigma T$. Since $Y \in \text{dom}(\cdot, \vdash)$, we know that $Y \notin \bar{X}$. The fact that $FV((Y)) \cap \bar{X} = \emptyset$ follows from our assumption that $FV(\cdot, \vdash) \cap \bar{X} = \emptyset$. We therefore have $\cdot, \vdash (Y) <: \sigma T$. Using the S-VAR rule, we obtain $\cdot, \vdash Y <: \sigma T$, which is what we need since $\sigma Y = Y$.

Case: $\cdot, \vdash_{\bar{X}}^V \text{All}(\bar{Y} <: \bar{A}) \bar{R} \rightarrow S <: \text{All}(\bar{Y} <: \bar{B}) \bar{T} \rightarrow U \Rightarrow D \wedge \bar{D} \wedge \bar{C}$ where $\cdot, \vdash_{\bar{X}}^V \bar{A} \equiv \bar{B} \Rightarrow \bar{K}, \bar{D}$ and $\cdot, \bar{Y} <: \bar{K} \vdash_{\bar{X}}^{V'} \bar{T} <: \bar{R} \Rightarrow \bar{C}$ and $\cdot, \bar{Y} <: \bar{K} \vdash_{\bar{X}}^{V'} S <: U \Rightarrow D$ and $V' = V \cup \{\bar{Y}\}$ and $\bar{Y} \cap V = \emptyset$ and $\bar{Y} \cap \bar{X} = \emptyset$ and $FV(\sigma) \cap \bar{Y} = \emptyset$

The side conditions on \bar{Y} imply that σ is a valid \bar{X}/V' -substitution. Therefore, our assumption that $\cdot, \vdash \sigma \in D \wedge \bar{D} \wedge \bar{C}$, plus the fact that $FV(\bar{K}) \cap \bar{X} = \emptyset$, is sufficient to allow us to use the induction hypothesis to prove that $\cdot, \bar{Y} <: \bar{K} \vdash \sigma S <: \sigma U$ and $\cdot, \bar{Y} <: \bar{K} \vdash \sigma \bar{T} <: \sigma \bar{R}$. Moreover, Lemma 3.5.1 tells us that $\sigma \bar{A} = \sigma \bar{B} = \bar{K}$. By the subtyping rule for functions, we conclude that $\cdot, \vdash \text{All}(\bar{Y} <: \sigma \bar{A}) \sigma \bar{R} \rightarrow \sigma S <: \text{All}(\bar{Y} <: \sigma \bar{B}) \sigma \bar{T} \rightarrow \sigma U$. The result now follows, since $FV(\sigma) = \emptyset$ ensures that $\text{All}(\bar{Y} <: \sigma \bar{A}) \sigma \bar{R} \rightarrow \sigma S = \sigma(\text{All}(\bar{Y} <: \bar{A}) \bar{R} \rightarrow S)$ and $\text{All}(\bar{Y} <: \sigma \bar{B}) \sigma \bar{T} \rightarrow \sigma U = \sigma(\text{All}(\bar{Y} <: \bar{B}) \bar{T} \rightarrow U)$. ■

3.5.4 Proposition [Completeness of constraint generation]: Let σ be an \bar{X}/V -substitution with $\bar{X} \cap V = \emptyset$, and let S and T be types such that either $FV(S) \cap \bar{X} = \emptyset$ or $FV(T) \cap \bar{X} = \emptyset$. Let \cdot, \vdash be a context such that $\bar{X} \cap \text{dom}(\cdot, \vdash) = \emptyset$ and $FV(\cdot, \vdash) \cap \bar{X} = \emptyset$. If $\cdot, \vdash \sigma S <: \sigma T$, then $\cdot, \vdash_{\bar{X}}^V S <: T \Rightarrow C$ and $\cdot, \vdash \sigma \in C$.

Proof: By induction on the depth of a derivation of $\cdot, \vdash \sigma S <: \sigma T$.

Case: $S = Y$ where $Y \in \bar{X}$

We have $\cdot, \vdash_{\bar{X}}^V Y <: T \Rightarrow C$ where $C = \{Y \mapsto [\text{Bot}, R]\}$ and $T \Downarrow_V R$. Now, since σ is a \bar{X}/V -substitution, we know that $FV(\sigma Y) \cap V = \emptyset$, and therefore, using Lemma 3.2.2, we have $\cdot, \vdash \sigma Y <: R$. This ensures that $\cdot, \vdash \sigma \in C$, as required.

Case: $T = Y$ where $Y \in \bar{X}$

We have $\cdot, \vdash_{\bar{X}}^V S <: Y \Rightarrow C$ where $C = \{Y \mapsto [R, \text{Top}]\}$ and $S \Uparrow^V R$. Now, since σ is a \bar{X}/V -substitution, we know that $FV(\sigma Y) \cap V = \emptyset$, and therefore, using Lemma 3.2.2, we have $\cdot, \vdash R <: \sigma Y$. This ensures that $\cdot, \vdash \sigma \in C$, as required.

Case: $T = \text{Top}$

Immediate, since $\cdot, \vdash_{\bar{X}}^V S <: \text{Top} \Rightarrow \emptyset$ and $\cdot, \vdash \sigma \in \emptyset$.

Case: $S = \text{Bot}$

Immediate, since $\cdot, \vdash_{\bar{X}}^V \text{Bot} <: T \Rightarrow \emptyset$ and $\cdot, \vdash \sigma \in \emptyset$.

Case: $S = Y$ and $T = Y$ where $Y \notin \bar{X}$

Immediate, since $\vdash_{\bar{X}}^V Y < Y \Rightarrow \emptyset$ and $\sigma \in \emptyset$.

Case: $\vdash Y < \sigma T$ where $\vdash (Y) < \sigma T$

Since $FV(\cdot) \cap \bar{X} = \emptyset$ we have $\sigma(\cdot(Y)) = \cdot(Y)$, so we can use the induction hypothesis to prove that $\vdash_{\bar{X}}^V (Y) < T \Rightarrow C$ and $\vdash \sigma \in C$. The result follows directly, since $\vdash_{\bar{X}}^V Y < T \Rightarrow C$.

Case: $\vdash \text{All}(\bar{Y} < \sigma \bar{A}) \sigma \bar{R} \rightarrow \sigma S < \text{All}(\bar{Y} < \sigma \bar{B}) \sigma \bar{T} \rightarrow \sigma U$, where $\sigma \bar{A} = \sigma \bar{B}$ and $\vdash \bar{Y} < \sigma \bar{B} \vdash \sigma \bar{T} < \sigma \bar{R}$ and $\vdash \bar{Y} < \sigma \bar{B} \vdash \sigma S < \sigma U$

Since we identify type expressions up to alpha-conversion, we may suppose (wlog) that the \bar{Y} are chosen so that $\bar{Y} \cap V = \emptyset$, $\bar{Y} \cap \bar{X} = \emptyset$, and $FV(\sigma) \cap \bar{Y} = \emptyset$. Thus, σ is a valid \bar{X}/V' -substitution and we can use the induction hypothesis to prove that $\vdash \bar{Y} < \sigma \bar{B} \vdash_{\bar{X}'}^V \bar{T} < \bar{R} \Rightarrow \bar{C}$ and $\vdash \sigma \in \bar{C}$, and similarly, $\vdash \bar{Y} < \sigma \bar{B} \vdash_{\bar{X}'}^V S < U \Rightarrow D$ and $\vdash \sigma \in D$. Using Lemma 3.5.2, we have that $\vdash_{\bar{X}}^V \bar{A} \equiv \bar{B} \Rightarrow \sigma \bar{B}, \bar{D}$ and $\vdash \sigma \in \bar{D}$. So, by the constraint generation rule for function types, we have $\vdash_{\bar{X}}^V \text{All}(\bar{Y} < \bar{A}) \bar{S} \rightarrow P < \text{All}(\bar{Y} < \bar{B}) \bar{T} \rightarrow Q \Rightarrow D \wedge \bar{D} \wedge \bar{C}$. Finally, by the fact that $D \wedge \bar{D} \wedge \bar{C}$ is a greatest lower bound, we have $\vdash \sigma \in D \wedge \bar{D} \wedge \bar{C}$, as required. ■

3.6 Calculating Type Arguments

Having generated a set of constraints for the missing type parameters \bar{X} , the final job of the local constraint solver is to choose values for \bar{X} that make the result type of the whole application as small as possible. Depending on where the variables \bar{X} occur in R , this may involve choosing the smallest possible values for some variables and the largest for others. For example, if $R = X \rightarrow Y$ and we have generated the constraint set $\{X \mapsto [S, T], Y \mapsto [U, V]\}$, the smallest possible value for R is found by maximizing X and minimizing Y —i.e., by taking the substitution $\sigma = [X \mapsto T, Y \mapsto U]$. It may also be the case that no substitution for the variables yields a minimal result type; for example, if $R = X \rightarrow X$ and we have the constraint set $\{X \mapsto [\text{Int}, \text{Top}]\}$, then both $\text{Int} \rightarrow \text{Int}$ and $\text{Top} \rightarrow \text{Top}$ are solutions but neither is a subtype of the other. Local type argument synthesis fails in this case (as required by the specification in Section 3.1).

We begin by formalizing the ways in which maximizing or minimizing X affect the final result type.

1. We say that R is *covariant in X* if $\vdash [S/X]R < [T/X]R$ whenever $\vdash S < T$.
2. We say that R is *contravariant in X* if $\vdash [T/X]R < [S/X]R$ whenever $\vdash S < T$.
3. We say that R is *invariant in X* if $\vdash [S/X]R < [T/X]R$ only when both $\vdash S < T$ and $\vdash T < S$.
4. We say that R is *rigid in X* if $\vdash [S/X]R < [T/X]R$ only when $S = T$.

It is easy to check whether R is covariant, contravariant, invariant, or rigid in a given variable X by examining where X occurs in R (to the right or left of arrows, in the bounds of type binders, etc.).

Next, we need a technical definition characterizing types whose equivalence classes in the subtype relation are singletons. For example, if $X < \text{Bot}$, then $X \rightarrow X$ is equivalent, but not identical, to $\text{Bot} \rightarrow \text{Bot}$: indeed its equivalence class has several members. On the other hand, Top is only equivalent to itself. Formally, we call a type variable a *bottom variable* if it is bounded by Bot or by another bottom variable. Now, let \cdot be a context and S a type whose free variables are in $\text{dom}(\cdot)$. We say that S is *rigid under* \cdot if

- $S = \text{Top}$;
- $S = \text{Bot}$ and no variable in \cdot is bounded by Bot ;
- $S = X$ and X is not a bottom variable;
- $S = \text{All}(\bar{X}) \bar{S} \rightarrow T$ with each S_i rigid under \cdot , $X_1 : S_1, \dots, X_{i-1} : S_{i-1}$ and T rigid under \cdot , $\bar{X} : \bar{S}$;

Extending the notion of rigidity from types to constraints, we say that a \cdot -constraint c is *rigid* if it admits only one solution—i.e., if either $c = [S]$ or else $c = [S, S]$, where S is rigid under \cdot . Similarly, c is said to be *tight* if it admits only one solution, up to equivalence—i.e., if either $c = [S]$ or else $c = [S, T]$ with both $\vdash S < T$ and $\vdash T < S$.

3.6.1 Lemma: If S is rigid under \cdot , then every type equivalent to S is syntactically equal to S —i.e., $\cdot \vdash S <: T$ and $\cdot \vdash T <: S$ together imply that S and T are identical.

Proof: See [Pie97, Lemma 4.1.2]. ■

3.6.2 Corollary: If c is rigid under \cdot , and $\cdot \vdash S \in c$ and $\cdot \vdash T \in c$, then S and T are identical.

With the forgoing definitions in hand, we can now show how to choose values for the variables \bar{X} that will minimize R (or else determine that this is not possible). Let C be a satisfiable \bar{X}/V -constraint set and R a type whose free variables are in $dom(\cdot) \cup \{\bar{X}\}$. Let the substitution σ_{CR} be defined (when it exists) as follows:

For each $X_i \dots$

- if R is covariant in X_i ,
then $\sigma_{CR}(X_i) = \min(C(X_i))$
- else if R is contravariant in X_i ,
then $\sigma_{CR}(X_i) = \max(C(X_i))$
- else if R is invariant in X_i
and $C(X_i)$ is tight,
then $\sigma_{CR}(X_i) = \min(C(X_i))$
- else if R is rigid in X_i ,
and $C(X_i)$ is rigid,
then $\sigma_{CR}(X_i) = \min(C(X_i))$
- else σ_{CR} is undefined.

It remains to verify that the substitution σ_{CR} chosen in this way is indeed the best possible. Let C be an \bar{X}/V -constraint set, and σ be a \bar{X}/V -substitution. We say that σ is *minimal* for C and R , written $\cdot \vdash \sigma \in C \Downarrow R$, if $\cdot \vdash \sigma \in C$ and, for all \bar{X}/V -substitutions σ' with $\cdot \vdash \sigma' \in C$, we have $\cdot \vdash \sigma R <: \sigma' R$.

3.6.3 Proposition:

1. If the substitution σ_{CR} exists, then it is a minimal substitution for C and R .
2. If σ_{CR} is undefined, then C and R have no minimal substitution.

Proof:

1. Suppose σ_{CR} exists, and that σ' is another substitution such that $\cdot \vdash \sigma' \in C$. We must show that $\cdot \vdash \sigma_{CR} R <: \sigma' R$.

Let $n = |\bar{X}|$, and construct a sequence of substitutions $\sigma_0, \dots, \sigma_n$ as follows:

$$\begin{aligned} \sigma_0 &= \sigma_{CR} \\ \sigma_i &= \sigma_{i-1}[X_i \mapsto \sigma'(X_i)] \quad \text{if } i \geq 1. \end{aligned}$$

Note that $\sigma_n = \sigma'$. We now argue that $\cdot \vdash \sigma_{i-1} R <: \sigma_i R$ for each $i \geq 1$.

- If R is covariant in X_i , then, by definition, $\sigma_{i-1} X_i = \sigma_{CR}(X_i) = \min(C(X_i))$, and thus $\cdot \vdash \sigma_{i-1}(X_i) <: \sigma_i(X_i)$. But this implies that $\cdot \vdash \sigma_{i-1} R <: \sigma_i R$, by the definition of covariance.
- Similarly, if R is contravariant in X_i , then $\sigma_{i-1} X_i = \sigma_{CR}(X_i) = \max(C(X_i))$, and thus $\cdot \vdash \sigma_i(X_i) <: \sigma_{i-1}(X_i)$, which implies that $\cdot \vdash \sigma_{i-1} R <: \sigma_i R$, by the definition of contravariance.
- If R is invariant in X_i , then $\sigma_{i-1} X_i = \sigma_{CR}(X_i) = \max(C(X_i))$, and we also know that $\cdot \vdash \min(C(X_i)) <: \max(C(X_i)) <: \min(C(X_i))$. But since $\cdot \vdash \min(C(X_i)) <: \sigma_i(X_i) <: \min(C(X_i))$, we have by transitivity that $\cdot \vdash \sigma_i(X_i) <: \sigma_{i-1}(X_i) <: \sigma_i(X_i)$, which, by the definition of invariance, yields $\cdot \vdash \sigma_{i-1} R <: \sigma_i R$.
- Finally, if R is rigid in X_i , then $\sigma_{i-1}(X_i) = \sigma_i(X_i)$, and so $\cdot \vdash \sigma_{i-1} R <: \sigma_i R$ by reflexivity of subtyping.

We have thus shown that $\cdot, \vdash \sigma_{C\mathbf{R}}\mathbf{R} = \sigma_0\mathbf{R} <: \sigma_1\mathbf{R} <: \dots <: \sigma_n\mathbf{R} = \sigma'\mathbf{R}$, and the desired result follows by transitivity of subtyping.

2. If $\sigma_{C\mathbf{R}}$ is undefined, then either C is unsatisfiable (in which case the result holds trivially) or else C is satisfiable and we must show that no substitution that satisfies it is minimal. So suppose, for a contradiction, that σ is minimal for C and \mathbf{R} . There are two cases to consider, depending on why $\sigma_{C\mathbf{R}}$ failed to be defined:
 - (a) For some X_i , \mathbf{R} is invariant in X_i but $C(X_i)$ is not a tight constraint. In this case, we know that there must be some T such that $\cdot, \vdash T \in C(X_i)$ but such that either $\cdot, \vdash \sigma(X_i) \not<: T$ or $\cdot, \vdash T \not<: \sigma(X_i)$. We can then construct a substitution $\sigma' = [X_i \mapsto T]$ such that $\cdot, \vdash \sigma' \in C$ and, since X_i is invariant in \mathbf{R} , such that $\cdot, \vdash \sigma\mathbf{R} \not<: \sigma'\mathbf{R}$, contradicting our assumption that σ is minimal for C and \mathbf{R} .
 - (b) For some X_i , \mathbf{R} is rigid in X_i but $C(X_i)$ is not a rigid constraint. In this case, we know that there must be some T different from $\sigma(X_i)$ such that $\cdot, \vdash T \in C(X_i)$. We can then construct a substitution $\sigma' = \sigma[X_i \mapsto T]$ such that $\cdot, \vdash \sigma' \in C$ and, since X_i is rigid in \mathbf{R} , such that $\cdot, \vdash \sigma\mathbf{R} \not<: \sigma'\mathbf{R}$, contradicting our assumption that σ is minimal for C and \mathbf{R} . \blacksquare

3.6.4 Corollary: The algorithmic rule

$$\frac{\begin{array}{c} \cdot, \vdash f \in F \quad \cdot, \vdash F \uparrow \text{All } (\overline{X} <: \overline{S}) \overline{T} \rightarrow \mathbf{R} \\ \cdot, \vdash \overline{e} \in \overline{U} \quad |\overline{X}| > 0 \quad \overline{X} \cap FV(\overline{S}) = \emptyset \\ \cdot, \vdash_{\overline{X}} \overline{A} <: \overline{S} \Rightarrow \overline{C} \quad \cdot, \vdash_{\overline{X}} \overline{U} <: \overline{T} \Rightarrow \overline{D} \quad \cdot, \vdash \sigma \in (\overline{C} \wedge \overline{D}) \Downarrow \mathbf{R} \end{array}}{\cdot, \vdash f(\overline{e}) \in \sigma\mathbf{R} \Rightarrow f[\sigma\overline{X}](\overline{e})}$$

is equivalent to the declarative rule given in Section 3.1.

Acknowledgements

The paper was mostly written while Turner was visiting Indiana University in the summer of '97. Pierce was partially supported by NSF grant CCR-9701826, *Principled Foundations for Programming with Objects*.

References

- [CG92] Pierre-Louis Curien and Giorgio Ghelli. Coherence of subsumption: Minimum typing and type-checking in F_{\leq} . *Mathematical Structures in Computer Science*, 2:55–91, 1992. Also in Carl A. Gunter and John C. Mitchell, editors, *Theoretical Aspects of Object-Oriented Programming: Types, Semantics, and Language Design* (MIT Press, 1994).
- [CMMS94] Luca Cardelli, Simone Martini, John C. Mitchell, and Andre Scedrov. An extension of system F with subtyping. *Information and Computation*, 109(1–2):4–56, 1994. A preliminary version appeared in TACS '91 (Sendai, Japan, pp. 750–770).
- [CW85] Luca Cardelli and Peter Wegner. On understanding types, data abstraction, and polymorphism. *Computing Surveys*, 17(4), December 1985.
- [Ghe90] Giorgio Ghelli. *Proof Theoretic Studies about a Minimal Type System Integrating Inclusion and Parametric Polymorphism*. PhD thesis, Università di Pisa, March 1990. Technical report TD-6/90, Dipartimento di Informatica, Università di Pisa.
- [GP97] Giorgio Ghelli and Benjamin Pierce. Bounded existentials and minimal typing. *Theoretical Computer Science*, 1997. To appear.
- [Mil92] Dale Miller. Unification under a mixed prefix. *Journal of Symbolic Computation*, 14(4):321–358, October 1992.
- [Pie97] Benjamin C. Pierce. Bounded quantification with bottom. Technical Report 492, Computer Science Department, Indiana University, 1997.
- [PT98] Benjamin C. Pierce and David N. Turner. Local type inference. In *Conference Record of POPL '98: the 25th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, 1998. To appear. Full version available as Indiana University CSCI technical report #493.