# Simple Type Inference for Structural Polymorphism<sup>†</sup>

Jacques Garrigue Research Institute for Mathematical Sciences Kyoto University, 606-8502 Kyoto, JAPAN garrigue@kurims.kyoto-u.ac.jp

ABSTRACT

We propose a new way to mix constrained types and type inference, where the interaction between the two is minimal. By using local constraints embedded in types, rather than the other way round, we obtain a system which keeps the usual structure of an Hindley-Milner type system. In practice, this means that it is easy to introduce local constraints in existing type inference algorithms.

Eventhough our system is notably weaker than general constraintbased type systems, making it unable to handle subtyping for instance, it is powerful enough to accomodate many features, from simple polymorphic records à la Ohori to Objective Caml's polymorphic variants, and accurate typing of pattern matching (*i.e.* polymorphic message dispatch), all these through tiny variations in the constraint part of the system.

# 1. INTRODUCTION

Type inference for *structural polymorphism* has been a long standing area of research. By structural polymorphism, we mean a form of parametric polymorphism in which some types may have a partially known structure, often denoted by a notion of row variable or kinded variable. Applications range from polymorphic records [14, 17] to polymorphic variants [16, 5] and objects with first-class (dynamic) messages [12]. All these systems were proposed as extensions of the Hindley-Milner type system [10], with growing degrees of complexity.

A recent trend is to simplify these type systems by moving to a constraint-based framework [13, 11, 18, 15]. This indeed greatly improved the understanding of the various systems, as they can all be described as particular instances of the HM(X) framework, on various constraint domains. However, the switch to a constraint system also means that

<sup>†</sup>A preliminiary version of this paper has been presented at the 9<sup>th</sup> FOOL Workshop, Portland, OR, January 2002. previous work on Hindley-Milner does not directly apply. This has practical consequences: it is hard to extend an existing algorithm to work with constraints, as all types must now be embedded in constraints. This has also theoretical consequences: take for instance Hindley-Milner extended with first-class polymorphism [7]; the extension intimately depends on the way polymorphism is handled in the original type system, and as HM(X) handles it in a radically different way (a type variable can be both constrained and polymorphic), it is unclear how it could handle such an extension.

This is all the more disturbing as HM(X) is somehow "too powerful" for structural polymorphism. It was designed to accomodate subtyping, which is carefully avoided in structural polymorphism.

A natural way to go is then to search for a weaker system, not handling subtyping, but closer to the original Damas-Milner formulation [3]. This is what we present here. We do it by restricting the role of constraints to individual type nodes: each constraint lives inside its type node, and can only influence outside types by requiring two children of the node to be equal. Since the embedding is done the other way round, existing type inference algorithms can easily be extended with such constrained types. For instance, the handling of polymorphic variants in Objective Caml 3.00[9] can be seen as an instance of this framework. We keep the original definition of free variable, making this system compatible with first-class polymorphism for instance. Locality is also beneficial for producing readable type errors: rather than producing an unsolvable set of arbitrary constraints, we can localize the problem, and translate it back to a meaningful form.

Contributions of this paper can be found at the application level: to our knowledge, both the problem of discarding irrelevant type information in shrinking types, and of accurately typing pattern matching, were still open for type systems without subtyping-like mechanisms. The solutions we propose here are simple instances of our framework.

Our system is very close to Ohori's polymorphic records [14], reusing the concept of kind as constraint. It is made modular by completely separating the descriptions of types and constraints. New typed constructs can be added by merely specifying a new constraint domain, and new typed constants, without needing to introduce new type inference rules. First, we will present by examples various forms of structural polymorphism. Then we introduce the notion of *constraint domain*, which is used as parameter to our framework. All forms of structural polymorphism we present here appear to have very similar constraint domains. In a 4th section we introduce the type system itself, with its specific notion of substitution. Section 5 introduces terms and typing rules. Section 6 and 7 are devoted to unification and type reconstruction. Section 8 studies in more detail the differences with HM(X). We conclude discussing the expressive power of this system. Proofs are given in appendix.

# 2. STRUCTURAL POLYMORPHISM

We present here various forms of structural polymorphism, on a gradual scale of increasing complexity. In order to uniformize the presentation, we write  $K \triangleright \tau$  for types, where K is a kinding environment, containing constraints for individual type variables, under which the open type  $\tau$  is to be understood. We use a few built-in functions, to obtain interesting types: integer addition +, string concatenation  $\uparrow$ , conversions  $string\_of\_int : string \rightarrow int$  and  $float : int \rightarrow float$ .

## 2.1 Records à la Ohori

Arguably this is the simplest form of structural polymorphism considered in the literature. Record values have monomorphic types, and polymorphism is only used for typing field access.

 $\begin{array}{l} \{name = \texttt{"Jacques"}, age = 30\} : \{name : string, age : int\} \\ \texttt{fun } x \to x.age + 1 \qquad : \alpha :: \{age : int\} \triangleright \alpha \to int \\ \texttt{fun } x \to x.name \uparrow \texttt{" is "} \uparrow string\_of\_int x.age \\ : \alpha :: \{name : string, age : int\} \triangleright \alpha \to string \end{array}$ 

Intuitively, the kinding  $\alpha :: \{age : int\}$  means that  $\alpha$  should at least have a field age, and this field should have type *int*. Accessing several fields result in a kind containing all of them. Two kinds are compatible if they agree on the types of their common fields. A type satisfies a kind if all the required fields are provided, with the correct types.

Note that the above formulation is a direct adaptation of Ohori's. Fitting it inside the local constraint framework will require a small change in presentation, without changing expressive power.

#### 2.2 **Records and variants with masking**

Ohori's records are weak in that record values have only structurally monomorphic types. We might want to allow building lists of records containing different fields, with only part of them common to all members. This requires giving polymorphic types not only to field access, but also to record values. A type system allowing it was first proposed by Rémy [16], but here we use a more intuitive formalism [5], which is also closer to Ohori's.

$$\begin{cases} name = \texttt{"Jacques"}, age = 30 \\ : \alpha :: (\{name : string, age : int\}, \emptyset, \{name, age\}) \triangleright \alpha \\ \texttt{fun } x \rightarrow x.age + 1 \quad : \alpha :: (\{age : int\}, \{age\}, \mathcal{L}) \triangleright \alpha \rightarrow int \\ \texttt{let } l_1 = [\{name = \texttt{"Jacques"}, age = 30\}, \\ \{name = \texttt{"Serge"}, weight = 13\}]$$

 $\begin{array}{ll} l_1 &: \ \alpha :: (\{name : string, age : int, weight : int\}, \emptyset, \{name\}) \\ & \triangleright \alpha \ list \end{array}$ 

Kinds are now represented by a triple (T, L, U). Along with the type of each field, we have two sets of labels. Required labels L (union of all field accesses, a lower bound) form a subset of available labels U (intersection of the possible record values, an upper bound).  $\mathcal{L}$  is the set of all labels. You can see how fields are masked in the third example: while we have types for the *name*, *age* and *weight* fields, only the *name* field is accessible.

The combined type  $\alpha :: (\{name : string, age : int\}, \{age\}, \{name, age\})$  is an acceptable description for both  $\{name = "Jacques", age = 30\}$ , which makes name and age available, and fun  $x \to x.age + 1$ , which only requires age. The need for two distinct sets of labels stems from the natural appearence of such combined types during type inference. Since attempting to access an unavailable label would be a type error, this also explains why, for a constraint kind to be meaningful (to have a solution), the set of required labels should be included in the set of available labels.

By duality, the same types can be used to describe polymorphic variants. The basic idea is that case-analysis of a variant can receive the same type as a record, while the variant itself would get the type of a field accessor.

 $Number(5) : \alpha :: ({Number : int}, {Number}, \mathcal{L}) \triangleright \alpha$ let  $l_2 = [Number(5), Face("King")]$ 

- $l_{2} : \alpha :: (\{Number : int, Face : string\}, \{Number, Face\}, \mathcal{L}) \\ \triangleright \alpha \ list$
- let  $f_1 =$ function  $Number(n) \rightarrow string\_of\_int n$ |  $Face(name) \rightarrow name$
- $\begin{array}{l} f_1 \ : \ \alpha :: (\{Number : int, Face : string\}, \emptyset, \{Number, Face\}) \\ \ \triangleright \ \alpha \rightarrow string \end{array}$

Our two sets of labels have now different meanings. The first one is the set of present constructors, or required cases, which must be handled by case-analysis; the second one is the set of handled constructors, or available cases, which is accepted by all case-analyses.

An advantage of not using predefined sum types, is that we can make case analysis modular. Consider the following function  $f_2$ , which uses a special syntax for dispatch. g(respectively h) will only receive A or B (respectively C or D). This can be reflected at the type level by requiring them to handle only relevant cases.

$$\begin{array}{l} \text{let } f_2 = \text{fun } g \to \text{fun } h \to \text{function } (A|B) \ as \ x \to g \ x \\ & \mid (C|D) \ as \ x \to h \ x \end{array} \\ f_2 \ : \ \alpha_1 :: (\{A : \alpha_A, B : \alpha_B\}, \{A, B\}, \mathcal{L}), \\ \alpha_2 :: (\{C : \alpha_C, D : \alpha_D\}, \{C, D\}, \mathcal{L}), \\ \alpha_3 :: (\{A : \alpha_A, B : \alpha_B, C : \alpha_C, D : \alpha_D\}, \emptyset, \{A, B, C, D\}) \\ & \triangleright \ (\alpha_1 \to \alpha) \to (\alpha_2 \to \alpha) \to \alpha_3 \to \alpha \end{array}$$

Our choice of keeping all constraints local (constraining only one variable) makes our system slightly weaker than Rémy's, which has row and presence variables. For instance we have no way to relate the set of required cases in  $\alpha_1$  and  $\alpha_3$ , which forces us to make the safe assumption  $\{A, B\}$  in  $\alpha_1$ . Rémy's system handles requirement for each constructor as an independent presence variable, which may be shared between two different variant types. This allows more precise typing, but at the cost of harder to understand types, with lots of variables, of different sorts. Experience suggests that kinded variables, by reducing the number of variables to one by record or variant type, and keeping only one sort of type variables, make reading types much easier.

## 2.3 Discarding masked types

One may wonder about why one should keep all field types in the type of  $l_1$ , a list of records with some masked fields. If fields *age* and *weight* are actually unavailable, why should their types matter? Clearly, this is not the case with subtyping, which would allow to discard not only the fields, but also their types. However, our requirement of principal type inference makes impossible to simply forget the type.

For instance, let us consider the list:

let 
$$l_3 = [\{name = "Jacques", age = 30\},$$
  
 $\{name = "Rachel", age = 6.5\}]$ 

The two records disagree on the type for *age*, either *int* or *float*. This should trigger a type error. Yet, if we choose to discard the type for *age* in  $l_1$ , we would be able to type  $l_1@[\{name = "Rachel", age = 6.5\}]$ , which contains the untypable  $l_3$  as a sublist.

A solution to this problem is to only trigger an error when a field with conflicting types is used. This amounts to allowing conjunctive types in the types of non-required fields, only forcing them to be equal when they are required. Then we can give the following type to  $l_3$ .

 $l_3 : \alpha :: (\{name : string, age : int \land float\}, \emptyset, \{name, age\}) \triangleright \alpha$ 

Since the two types for the field *age* are incompatible, it cannot be accessed. And if it disappears from the list of available fields, we can just drop these types as useless information.

While the absence of error may seem strange in the above case (even though it is necessary for coherence), it is a natural way to allow the construction of heterogeneous collections of objects in a system without subtyping.

The same mechanism also applies to polymorphic variants, and is actually used in Objective Caml 3.00. This allows principal type inference of  $f_5$  for the following case, solving the long lasting problem of "masked but not discarded" types.

$$\begin{array}{l} \mathsf{let} \ f_3 = \mathsf{function} \ Number(n) \to n \mid Face() \to 15\\ f_3 \ : \ \alpha :: (\{Number : int, Face : unit\}, \emptyset, \{Number, Face\})\\ \ \triangleright \ \alpha \to int\\ \mathsf{let} \ f_4 = \mathsf{function} \ Number(n) \to n/2\\ f_4 \ : \ \alpha :: (\{Number : int\}, \emptyset, \{Number\})\\ \ \triangleright \ \alpha \to int\\ \mathsf{let} \ f_5 = \mathsf{fun} \ x \to (f_1(x), f_3(x), f_4(x))\\ f_5 \ : \ \alpha :: (\{Number : int\}, \emptyset, \{Number\})\\ \ \triangleright \ \alpha \to (string \times int \times int)\end{array}$$

#### 2.4 Accurate matching and dynamic messages

Polymorphic variants provide already more precise typing than classical predefined variant types, since a specific variant type is computed for every case-analysis, even if the same constructors are partly reused.

We can go further, and allow the type of the result to depend on which case was used; this happens for instance when a message is dispatched between methods whose result types differ. Such a mechanism was first proposed by Aiken *et al.* in the context of soft typing for dynamically typed languages [1]. Pottier reformulated it as an instance of the HM(X)framework, calling it *accurate analysis of pattern matching* [15]. In both cases, the type of the result is described by a conditional type, depending on the input. Moreover, they use global inclusion constraints to deal with the problem, while we use only local ones.

Concretely, we add a new **cases** construct, which has the same syntax as **function**, but allows each case to return a value of different type. Here is an example of the expected behaviour:

let  $f_6 = \text{cases Number}() \rightarrow 10 \mid Face() \rightarrow "King"$  $f_6 \ (Number()) : int \hookrightarrow 10$  $f_6 \ (Face()) : string \hookrightarrow "King"$ 

Our basic idea to accommodate different return types, is to add another set of type bindings, corresponding to the results of each branch of the pattern matching. The special label **return** is used to indicate the result type of the whole pattern matching. If a constructor is required, then its associated type is unified with the type of the result.

$$f_{6} : \alpha :: (\{Number : unit, Face : unit\}, \\ \{Number : int, Face : string, return : \alpha_{1}\}, \\ \emptyset, \{Number, Face\}) \\ \triangleright \alpha \to \alpha_{1}$$

$$\begin{array}{l} \text{let } f_7 = \text{fun } (x: \alpha_2 :: (\emptyset, \emptyset, \{Number\}, \mathcal{L}) \triangleright \alpha_2) \to f_6 \ x \\ f_7 : \alpha :: (\{Number : unit, Face : unit\}, \\ \{Number : int, Face : string, \texttt{return} : int\}, \\ \{Number\}, \{Number, Face\}) \\ \triangleright \alpha \to int \end{array}$$

You can see how inference proceeds in  $f_7$ . The type annotation on x is only there to simulate application on a monomorphic constructor. Note how the type of the variant itself ends up including the type of the result, avoiding the need for a global constraint.

This mechanism is then extended to deal with the case when the same variant is passed to several pattern-matching functions, with different return types: fun  $x \to (f x, g x)$ . If we have only one return type set in our variant type, then the monomorphism of x would force return types in f and g to be identical. But we provide several sets, one by function, by a mechanism we call *generativity*. This way the type of x may contain different return types for different functions, without interference between them. We leave details for the formal development.

A simple and useful application of this extension is to provide polymorphic decomposing functions for variants, without overloading.

let  $arg = cases Number(n : int) \rightarrow n \mid Face(s : string) \rightarrow s$ arg :  $\alpha :: (\{Number : int, Face : string\},$ {*Number* : *int*, *Face* : *string*, return :  $\alpha_1$ },  $\emptyset$ , {*Number*, *Face*})  $\triangleright \alpha \rightarrow \alpha_1$ let  $tag = cases Number(n : int) \rightarrow 1 \mid Face(s : string) \rightarrow 2$  $tag : \alpha :: ({Number : int, Face : string}),$ {*Number* : *int*, *Face* : *int*, return :  $\alpha_1$ },  $\emptyset$ , {*Number*, *Face*})  $\triangleright \alpha \rightarrow \alpha_1$ let  $decomp = fun \ x \rightarrow (tag \ x, arg \ x)$ decomp :  $\alpha :: (\{Number : int, Face : string\},$  $\{\{Number : int, Face : int, return : \alpha_1\},\$  $\{Number : int, Face : string, return : \alpha_2\}\},\$  $\emptyset$ , {*Number*, *Face*})  $\triangleright \alpha \to \alpha_1 \times \alpha_2$ decomp (Face "King") :  $int \times string \hookrightarrow (2, "King")$ decomp (Number 10) :  $int \times int \hookrightarrow (1, 10)$ 

In a more object-oriented context, the same kind of mechanism was called dynamic messages by Nishimura [12]; later it was reformulated using different instances of HM(X) [11, 15]. Dynamic messages can be encoded using a small amount of syntactic sugar, but the way we rely on polymorphism in pattern-matching functions makes our system weaker.

This can be seen on the following example. Objects  $\{|...|\}$  are case-analyzing functions, and message passing  $o \leftarrow m$  is the application o m.

$$\begin{array}{l} \mathsf{let} \ o_1 = \{|Number(x) = float \ x; Face(s : string) = s|\}\\ o_1 : \ \alpha :: (\{Number : int, Face : string\}, \\ \{Number : float, Face : string, \mathsf{return} : \alpha_1\}, \\ \emptyset, \{Number, Face\})\\ \vartriangleright \alpha \to \alpha_1\\ o_1 \leftarrow Number \ 6 : \ float \ \hookrightarrow 6.0\\ o_1 \leftarrow Face \ "King" : \ string \ \hookrightarrow "King"\\ \mathsf{let} \ f_{11} = \mathsf{fun} \ o \to (o \leftarrow Number \ 6, o \leftarrow Face \ "King")\\ f_8 : \ \alpha_1 :: (\{Number : int, Face : string\}, \emptyset, \\ \{Number, Face\}, \mathcal{L})\\ \vartriangleright (\alpha_1 \to \alpha) \to \alpha \times \alpha\\ f_8 \ o_1 : \ untypable \end{array}$$

As you can see, the encoding is sound, and different methods can have different result types. However, if we apply the same monomorphic object to two different messages, they are merged in one, and later we are unable to distinguish their result types. For this reason, the application  $f_8 o_1$  cannot be typed: constraints involving both *Number* and *Face* cannot be solved simultaneously. Nishimura's system did not require polymorphism, since message and object types were not merged, rather a global subtyping constraint was generated.

While complex applications of dynamic messages cannot be encoded directly in our system, simpler cases where objects are only used as second class values can be handled. Moreover, if one needs real first class objects, first class polymorphism [7] is available to solve the problem. This only requires to write object types when they appear as function parameters.

# 3. CONSTRAINTS DOMAINS

Contrary to usual constraint-based type systems, our local constraints are grafted on top of the Hindley-Milner type system, rather than mixing constraints and types at the same level. This means that we can reduce requirements on constraints to a minimum: their interaction with the rest of the type system will be minimal anyway. In particular, we do not introduce any syntax for constraints: they are black boxes, and need just be able to answer some questions. While not essential, we think this freedom is important, as the choice of how to represent constraints is relevant to their understanding. Except for this extra freedom, nothing would prevent one to formalize our constraint domains as special kinds of cylindric algebras, which can be plugged into HM(X).

A constraint domain describes a class of constraints, and how they interact with the type system. A particular instance of the type system may contain several constraint domains, as long as all their operations and values are clearly distinguished. For simplicity, we will only consider type systems operating on a single constraint domain.

DEFINITION 1. A constraint domain C is composed of the following items.

- 1. A theory  $\mathcal{T}_{\mathcal{C}}$  with an entailment relation  $\models$  satisfying the following properties
  - (a) There is a constraint  $\bot$ , such that for any C we have  $\bot \models C$ .
  - (b) A constraint C such that  $C \models \perp$  is invalid. Validity is decidable.
  - (c) Entailment is reflexive and transitive:  $C \models C$ ; if  $C \models C'$  and  $C' \models C''$  then  $C \models C''$ .
  - (d) For any two constraints C and C', there is a constraint  $C \wedge C'$  such that  $C \wedge C' \models C$ ,  $C \wedge C' \models C'$ , and for all C'' such that  $C'' \models C$  and  $C'' \models C'$ , we have  $C'' \models C \wedge C'$ .
- 2. An observation relation  $\vdash$  checking some atomic properties of a constraint:  $C \vdash p(a)$  where p and a are respectively a predicate and a symbol for the domain. Observation should be compatible with entailment:

If  $C \models C'$  and  $C' \vdash p(a)$  then  $C \vdash p(a)$ .

- 3. A set of relating predicates of the form  $r(a, \tau)$ , which relate symbols and types. Some of these predicates are said to be generative, that is they have an extra index  $\epsilon$ , which can be existentially quantified:  $\exists \epsilon. r^{\epsilon}(a, \tau)$ .
- 4. A set of propagation rules  $\mathcal{E}_{\mathcal{C}}$ , of the form (depending on whether r is generative or not)

$$\forall x. (r(x, \alpha_1) \land r(A, \alpha_2) \land p(x) \Rightarrow \alpha_1 = \alpha_2) \\ \forall x \forall \epsilon. (r^\epsilon(x, \alpha_1) \land r^\epsilon(A, \alpha_2) \land p(x) \Rightarrow \alpha_1 = \alpha_2)$$

where A is either the same variable x or a symbol.

Our requirements on the entailment relation  $\models$  are as free as one can get. Basically, we only need a way to distinguish valid constraints from invalid ones, and build the intersection of two constraints.

The observation relation  $\vdash$  is a consequence of the representation independency of our constraints: we need an explicit way to relate them to the rest of the world. It interfaces constraints, which are semantic, with propagation rules, which are syntactic.

Relating predicates are similar in flavor to features in feature algebras [2]. Each constraint will be coupled with a set R of instances of these predicates (*e.g.*  $R = \{r(a, int), r(b, int), r(c, \alpha list)\}$ ), in order to relate its semantics with types. R need not always describe functions: coherence will be enforced by the propagation rules. Generative predicates are used for type information depending on the constraint, but which should not be identified throughout all its occurences. The idea is that merging two quantified sets of predicates will not mix them:  $\exists \epsilon. \{r^{\epsilon}(a, \tau)\} \cup \exists \epsilon. \{r^{\epsilon}(a, \tau')\} =$  $\exists \epsilon \epsilon'. \{r^{\epsilon}(a, \tau), r^{\epsilon'}(a, \tau')\}$ . Among our examples, only the typing of accurate matching uses generative predicates. Other cases can be handled in a system without  $\epsilon$ 's.

Propagation rules specify how C and R constrain other types. They come in two forms: either A is x, and r only has to be functional for every symbol a such that  $C \vdash p(a)$ ; or A is some symbol b, and additionally the image of a by r should be the same as the image of b (*i.e.* the second case includes the first one by transitivity). There is an additional variation for generative predicates, expressing that the functionality condition only applies when identical  $\epsilon$ 's are involved.

The syntax may seem restrictive, but relaxing the definition would not give more expressive power: arbitrarily complex conditions can be encoded in the constraint itself. Keeping the syntax of the rules simple makes checking the applicability of rules easy: one just has to iterate on the (finitary) relating predicates, and check observation for ground atoms. On the other hand, the restriction to the form  $\alpha_1 = \alpha_2$  on the right hand side is crucial: this is what allows our unification algorithm to terminate.

We say that a constraint is *exact* if it cannot be further refined, *i.e.* if it is only entailed by itself and  $\perp$ . This is the equivalent of a ground type in the world of constraints. Whether a constraint is exact or not does not impact anything in the theory, but knowing it may be helpful when reading types.

We now consider the constraint domains associated to examples in the previous sections.

#### Records à la Ohori

The original formulation did mix required fields and their types in a single kind. We have to distinguish the two, to adapt to our framework. Moreover, since we do not extend monomorphic types themselves, monomorphic records should also be described by exact constraints.

A constraint is a pair (L, x) of a finite set of labels L, together with a mark x distinguishing exact types (1) from refinable ones (0). Entailment on refinable types is containment:  $(L, x) \models (L', 0)$  iff  $L \supset L'$ , which makes set union the conjunction operation. For exact types entailment is only reflexive. All such constraints are valid, so we must add a  $\perp$  element. We do not need to observe anything particular, so we just use  $C \vdash true(l)$  for any l. The type associated to each label is described by a predicate  $r(l, \tau)$ , and the fact a label has only one type is captured by the unique propagation rule:

$$r(l, \alpha_1) \wedge r(l, \alpha_2) \wedge true(l) \Rightarrow \alpha_1 = \alpha_2.$$

#### Records with maskable fields

Constraints are represented by a pair of sets (L, U), L a finite set of accessed labels, and U either the set of all labels  $\mathcal{L}$ , or a finite set of available labels. Entailment is defined by  $(L, U) \models (L', U')$  iff  $L \supset L'$  and  $U \subset U'$ . We can choose  $(\mathcal{L}, \emptyset)$  as  $\bot$ . The validity check is  $(L, U) \models \bot$  when  $L \not\subset U$ . One can observe required labels:  $(L, U) \vdash req(l)$  iff  $l \in L$ . We can either keep the previous propagation rule, and obtain records without discarding, or switch to the following new rule, which allows conjunctive typing, by delaying the equality constraint until the field is accessed.

 $r(l, \alpha_1) \wedge r(l, \alpha_2) \wedge req(l) \Rightarrow \alpha_1 = \alpha_2.$ 

#### Accurate matching

As for records with maskable fields, constraints are pairs (L, U), with the same entailment relation. We reuse also the conditional propagation rule, but we also add a second generative relation predicate  $v^{\epsilon}(l, \tau)$ , and a second propagation rule, acting on result types:

 $v^{\epsilon}(l, \alpha_1) \wedge v^{\epsilon}(\mathsf{return}, \alpha_2) \wedge req(l) \Rightarrow \alpha_1 = \alpha_2.$ 

Both argument and return types will only be unified if the case corresponding to l is required, that is if the method l is called.

### 4. TYPES AND KINDS

Simple types  $\tau$  are defined as usual. Polytypes  $\sigma$  are extended with a kinding environment K which restricts possible instances for constrained variables.

au	::=	$\alpha$	type variable
		u	base type
	Í	au  ightarrow  au	function type
Κ	::=	$\emptyset \mid \mathbf{K}, \alpha :: (C, \exists \bar{\epsilon}.R)$	kinding environment
$\sigma$	::=	$\tau \mid \forall \alpha \dots \alpha . K \triangleright \tau$	polytypes

K is a set of bindings  $\alpha :: (C, \exists \bar{\epsilon}.R)$ , C a constraint and R a set of relations to types, describing the possible values admitted for the type  $\alpha$ . Note that recursive types can be defined using a mutually recursive kinding environment, *i.e.* where kinds are related to each other. The set of relations  $\exists \bar{\epsilon}.R$  is to be understood modulo  $\alpha$ -conversion of the  $\epsilon$ 's, and polytypes modulo  $\alpha$ -conversion of quantified type variables.

DEFINITION 2. A kind  $k = (C, \exists \bar{\epsilon}.R)$  is well formed if

- 1. the constraint C is satisfiable.
- 2. for each  $t(x, \alpha_1) \wedge t(A, \alpha_2) \wedge p(x) \Rightarrow \alpha_1 = \alpha_2$  in  $\mathcal{E}$ (where t = r or  $r^{\epsilon}$ ) and each  $t(a, \tau_1)$  and  $t([a/x]A, \tau_2)$ in R such that  $C \vdash p(a)$ , we have  $\tau_1 = \tau_2$ .

We define kinding environments as containing only well formed kinds. Notice that all type variables do not necessarily have a kind, only those that represent constrained types do.

Free variables  $FV_K(\sigma)$  of a polytype  $\sigma$  under a kinding environment K are defined as the minimum set satisfying the following equations. We write  $FV(\sigma)$  when K is clear from the context.

$$\begin{array}{lll} \mathrm{FV}_{\mathrm{K}}(\forall \alpha_{1} \ldots \alpha_{n}, \mathrm{K}' \triangleright \tau) &=& \mathrm{FV}_{\mathrm{K}, \mathrm{K}'}(\tau) \setminus \{\alpha_{1}, \ldots, \alpha_{n}\} \\ \mathrm{FV}_{\mathrm{K}, \alpha::(C, \exists \bar{\epsilon}, R)}(\alpha) &=& \{\alpha\} \cup \mathrm{FV}_{\mathrm{K}}(R) \\ & & \mathrm{FV}_{\mathrm{K}}(u) &=& \emptyset \\ \mathrm{FV}_{\mathrm{K}}(\tau_{1} \to \tau_{2}) &=& \mathrm{FV}_{\mathrm{K}}(\tau_{1}) \cup \mathrm{FV}_{\mathrm{K}}(\tau_{2}) \end{array}$$

Our notion of free variables is an important difference with systems like HM(X). We explicitly force all free variables inside a kind to be free inside any type containing a variable of this kind. We obtain thus the property that possible instances of a type only depend of its free variables, like in traditional Hindler-Milner type systems, whereas in HM(X) constraints may include unrelated variables.

DEFINITION 3. A type substitution  $\theta$ , extended as usual on monotypes and polytypes, is admissible between the kinding environments K and K', written  $K \vdash \theta : K'$ , if for all  $\alpha :: (C, \exists \epsilon. R)$  in K,  $\theta(\alpha)$  is a type variable  $\alpha'$  and it satisfies the following properties.

1. 
$$\alpha' :: (C', \exists \overline{\epsilon}'.R') \in \mathbf{K}'$$

2. 
$$C' \models C$$

3. there is an index substitution  $\eta : \overline{\epsilon} \to \overline{\epsilon}'$  such that  $\eta(\theta(R)) \subseteq R'$ .

Condition 1 ensures that all kinded variables are mapped to kinded variables. Condition 2 ensures that constraints are instantiated correctly (according to entailment). Condition 3 ensures that all type constraints are kept, while making provision for changes in relation indices.

An immediate consequence of this definition is that one cannot substitute a ground type for a kinded variable. This does not mean that kinded variables always stay polymorphic: the constraint may be an exact one, which cannot be refined further. This only reflects our design choice of only changing the constraint domain, without extending the types themselves.

For both kinding environments and substitutions, we will write  $K|_D$  (resp.  $\theta|_D$ ) for their restriction to variables in D, and  $K|_{\overline{D}}$  (resp.  $\theta|_{\overline{D}}$ ) for their restriction to variables outside of D.

#### Discarding useless types

The above definition of substitution requires one to keep eventually useless relational information (that cannot be matched by any propagation rule anymore). In practice, it may be useful to allow forgetting some relations, to lighten kind descriptions. However we must be careful of not dropping free type variables, as they may have an impact on which types can be made polymorphic. Let us consider the relation  $t(a, \tau)$  in R, where t = r or  $r^{\epsilon}$ . If for all propagation rule  $t(x, \alpha_1) \wedge t(A, \alpha_2) \wedge p(x) \Rightarrow \alpha_1 = \alpha_2$ in  $\mathcal{E}$ , we have<sup>1</sup>  $C \vdash \neg p(x)$ , and moreover if A = a, we have  $\forall y.C \vdash \neg p(y)$ , then  $t(a, \tau)$  will never cause propagation. We call such a  $t(a, \tau)$  inert under propagation rules  $\mathcal{E}$  and constraint C. Such an inert atom can be safely replaced by a dummy  $r_{fv}(\tau)$ , which just makes sure that free variables are kept. Moreover,  $r_{fv}(\tau)$  can be replaced by  $\{r_{fv}(\alpha) \mid \alpha \in$  $FV(\tau)\}$  at any point. For this reason we can simplify Rthrough the reflectively and transitively closed relation  $\cong_{\mathcal{E}}$ :

$$R \cup \{r(a,\tau)\} \cong_{C\mathcal{E}} R \cup \{r_{fv}(\tau') \mid \tau' \in T\}$$
  
when  $r(a,\tau)$  inert and  $\mathrm{FV}_{\emptyset}(\tau) = \mathrm{FV}_{\emptyset}(T)$ 

 $R \cong_{C\mathcal{E}} R' \Rightarrow (C, \exists \bar{\epsilon}. R) \cong_{\mathcal{E}} (C, \exists \bar{\epsilon}. R')$ 

Note that  $\cong_{\mathcal{E}}$  is not symmetric: this is a refinement from left to right.  $\cong_{\mathcal{E}}$  commutes with substitution.

PROPOSITION 1 (DISCARDING POSTPONEMENT). If K,  $\alpha :: k \vdash \theta : K', \theta(\alpha) :: k' \text{ and } k_0 \cong_{\mathcal{E}} k$ , then there is a kind  $k'_0$  such  $k'_0 \cong_{\mathcal{E}} k'$  and  $K, \alpha :: k_0 \vdash \theta : K, \theta(\alpha) :: k'_0$ .

The above proposition means that one can arbitrarily postpone discarding of useless types. For this reason we do not consider discarding in the rest of this paper, and assume it will be done before presenting types to the programmer.

# 5. TERMS AND TYPING RULES

## Expressions are the standard ones, $e := x \mid fun \ x \rightarrow e \mid e e$ , core lambda

Type judgments are extended with a kinding environment,

 $K; \Gamma \vdash e : \tau$ 

where K is a well-formed kinding environment and  $\Gamma$  is a set of bindings  $x : \sigma$  from term variables to polytypes.

Typing rules appear in figure 1. They are in the syntaxdirected style, instantiating and generalizing in one step, which simplifies the handling of eventual mutual recursion in the kinding environment. GENERALIZE<sup>2</sup> and LET are two rules, but they are always used together.

The extra rule CONSTANT is a variation on VARIABLE, intended to allow easy definition of extra operations through polymorphically typed constants. As in the informal presentation, the pair  $K_0 \triangleright \tau$  is to be understood as the polytype  $\forall FV_{K_0}(\tau).K_0 \triangleright \tau$ .

The following lemma is paramount in proving the soundness of type inference.

<sup>&</sup>lt;sup>1</sup>Here  $\neg p$  is some observable predicate in the complement of p, *i.e.*  $\forall C.(C \vdash p(x) \land C \vdash \neg p(x)) \Rightarrow C \models \bot$ . As an observable, it is preserved by entailment. Since constraints are to be understood under the open world hypothesis,  $C \vdash$  $\neg p(a)$  is not equivalent to  $\neg C \vdash p(a)$ .

<sup>&</sup>lt;sup>P</sup>( $\alpha$ ) a since the same case of GENERALIZE happens if K contains a variable  $\alpha$  such that  $FV_{K}(\alpha) \cap B \neq \emptyset$  but  $\alpha \notin B$ . The meaning of  $\alpha$  in  $K|_{\overline{B}}$  is not the same as in K, but since then  $\alpha$  would not be in  $FV_{K}(\Gamma)$  either (if it were, then we would have  $FV_{K}(\alpha) \subset FV_{K}(\Gamma)$ , and  $FV_{K}(\Gamma) \cap B = \emptyset$ ), this does not matter anyway.

VARIABLE	
$K, K_0 \vdash \theta : K$	$\operatorname{Dom}(\theta) \subset B$
$\overline{\mathbf{K};\Gamma,x:\forall B.K_0}$	$\triangleright  \tau \vdash x : \theta(\tau)$
Abstraction $\mathbf{K}; \Gamma, x : \tau \vdash e :$	au'
$\overline{\mathrm{K};\Gamma\vdashfun\ x}\to$	$e:\tau\to\tau'$
$\begin{array}{l} \text{Application} \\ \text{K}; \Gamma \vdash e_1 : \tau \rightarrow \end{array}$	$ au'  \mathbf{K}; \Gamma \vdash e_2 :  au$
$\overline{\mathrm{K};\Gamma\vdash e_1\ e_2:\tau}$	,

```
 \begin{array}{l} \underset{K; \Gamma \vdash e: \tau}{\operatorname{Generalize}} & B = \operatorname{FV}_{\mathrm{K}}(\tau) \setminus \operatorname{FV}_{\mathrm{K}}(\Gamma) \\ \overline{\mathrm{K}|_{\overline{B}}; \Gamma \vdash e: \forall B.\mathrm{K}|_{B} \triangleright \tau} \\ \underset{K; \Gamma \vdash e_{1}: \sigma \quad \mathrm{K}; \Gamma, x: \sigma \vdash e_{2}: \tau}{\operatorname{K}; \Gamma \vdash \operatorname{let} x = e_{1} \text{ in } e_{2}: \tau} \\ \end{array}
```

Figure 1: Typing rules

Records à la Ohori				
$\operatorname{Tconst}(\operatorname{record}_{l_1\dots l_n})$	=	$\alpha_0 :: ((\{l_1,, l_n\}, 1), \{r(l_1, \alpha_1),, r(l_n, \alpha_n)\})$		
		$\triangleright \alpha_1 \to \dots \to \alpha_n \to \alpha_0$		
$\mathrm{Tconst}(get_l)$	=	$\alpha_0 :: ((\{l\}, 0), \{r(l, \alpha_1)\}) \triangleright \alpha_0 \to \alpha_1$		
Records with maskable fields				
$\operatorname{Tconst}(\operatorname{record}_{l_1\dots l_n})$	=	$lpha_0 :: ((\emptyset, \{l_1,, l_n\}), \{r(l_1, lpha_1),, r(l_n, lpha_n)\})$		
		$\triangleright \alpha_1 \to \dots \to \alpha_n \to \alpha_0$		
$Tconst(get_l)$	=	$\alpha_0 :: ((\{l\}, \mathcal{L}), \{r(l, \alpha_1)\}) \triangleright \alpha_0 \to \alpha_1$		
Variants with maskable const	truct	tors		
$Tconst(tag_l)$	=	$\alpha_0 :: ((\{l\}, \mathcal{L}), \{r(l, \alpha_1)\}) \triangleright \alpha_1 \to \alpha_0$		
$Tconst(match_{l_1l_n})$	=	$\alpha_0 :: ((\emptyset, \{l_1,, l_n\}), \{r(l_1, \alpha_1),, r(l_n, \alpha_n)\})$		
		$\triangleright (\alpha_1 \to \alpha) \to \dots \to (\alpha_n \to \alpha) \to \alpha_0 \to \alpha$		
$\operatorname{Tconst}(\operatorname{split}_{l_1 \dots l_k, l_{k+1} \dots l_n})$	=	$\alpha' ::: ((\{l_1,, l_k\}, \mathcal{L}), \{r(l_1, \alpha_1),, r(l_k, \alpha_k)\}),$		
· · · · · · · · · · · · · · · · · · ·		$\alpha'' :: ((\{l_{k+1},, l_n\}, \mathcal{L}), \{r(l_{k+1}, \alpha_{k+1}),, r(l_n, \alpha_n)\}),$		
		$\alpha_0 :: ((\emptyset, \{l_1,, l_n\}), \{r(l_1, \alpha_1),, r(l_n, \alpha_n)\})$		
		$\triangleright (\alpha' \to \alpha) \to (\alpha'' \to \alpha) \to \alpha_0 \to \alpha$		
Accurate matching				
$Tconst(cases l_1l_n)$	=	$\alpha_0 :: ((\emptyset, \{l_1,, l_n\}), \exists \epsilon. \{r(l_1, \alpha_1),, r(l_n, \alpha_n),$		
		$v^{\epsilon}(l_1, \alpha'_1), \dots, v^{\epsilon}(l_n, \alpha'_n), v^{\epsilon}(return, \alpha)\})$		
		$\triangleright (\alpha_1 \to \alpha'_1) \to \cdots \to (\alpha_n \to \alpha'_n) \to \alpha_0 \to \alpha$		

Figure 2: Examples of constructs with their associated constraints

LEMMA 2 (TYPE SUBSTITUTION). If  $K, \Gamma \vdash e : \tau$  and  $K \vdash \theta : K'$ , then  $K', \theta(\Gamma) \vdash e : \theta(\tau)$ .

For dynamic semantics, constants providing extra operations are coupled with  $\delta$ -rules, of the form  $A[e_1, \ldots, e_n] \rightarrow A'[e_1, \ldots, e_n]$ , where A and A' are application contexts.

$$A[x_1, ..., x_n] ::= c \mid x_i \mid (A[x_1, ..., x_n] \mid A[x_1, ..., x_n])$$

Each  $\delta$ -rule should be type safe: if  $K_0$ ;  $\Gamma \vdash A[x_1, \ldots, x_n] : \tau$ , where  $\Gamma = \{x_1 : \tau_1, \ldots, x_n : \tau_n\}$ , then  $K_0$ ;  $\Gamma \vdash A'[x_1, \ldots, x_n] : \tau$  must be provable.

If all the  $\delta$ -rules are type safe, we can prove subject reduction in the usual way. E[] denotes an arbitrary single hole context, not restricted to application contexts.

THEOREM 4 (SUBJECT REDUCTION). If  $K; \Gamma \vdash E[e] : \sigma$ and  $e \to e'$  by a  $\delta$ -rule, BETA :  $((fun \ x \to e_1) \ e_2) \to [e_2/x]e_1$ or Let : let  $x = e_1$  in  $e_2 \to [e_2/x]e_1$ , then  $K; \Gamma \vdash E[e'] : \sigma$ 

Subject reduction is only a step on the way to type soundness (*i.e.* well-typed programs do not go wrong), the next one being to show that normal forms of programs typable in the empty environment are values [20]. While easy, this part is not modular (reduction rules may interact). For this reason we leave it to specific instances.

#### Examples

Some typed constructs are given in figure 2.

For record calculi the  $\delta$ -rule is

$$\operatorname{get}_{l_i} (\operatorname{record}_{l_1 \dots l_n} e_1 \dots e_n) \to e_i.$$

For variants, which share the same constraint as records with maskable fields, we have three  $\delta$ -rules. The first one is the dual of field extraction, while the two others implement polymorphic dispatch. Notice that with dispatch the functions' types are not directly connected to the input type.

For variants with accurate matching, we do not need to change the typing of the tag constructor; it is enough to add a new cases constructor for polymorphic functions. The  $\delta$ -rule is the same as for match, only the type differs.

cases 
$$l_1 \dots l_n f_1 \dots f_n$$
 (tag $l_i e$ )  $\rightarrow f_i e$ 

## 6. UNIFICATION

We define unification on monotypes through rules of the form

$$\frac{\varphi}{\varphi'}$$

where  $\varphi$  are preconditions and  $\varphi'$  conclusions.  $\varphi$  is a *unification problem* composed both of *kinding constraints*  $\alpha :: (C, \exists \bar{\epsilon}.R)$  and *equality constraints*  $\tau_1 \doteq \tau_2$  connected by  $\wedge$ .

Both  $\wedge$  and  $\doteq$  are commutative, and  $\wedge$  is associative. A variable may not be kinded by more than one kinding constraint. Note that the "constraints" here, along with  $\wedge$ , have nothing to do with the previously defined constraint domains; we just reuse a natural terminology in a different context, where constraints are purely syntactic.

An equality constraint is said to be *solved* in  $\varphi$  if it is of the form  $\alpha \doteq \tau$  and  $\alpha$  only appears as strict subterm in other constraints of  $\varphi$ . A kinding constraint is *solved* if its kind is well-formed. A unification problem  $\varphi$  is solved if all its constraints are solved.

A solution to a unification problem  $\varphi$ , whose kinding constraints form a kinding environment K (its *basis*), is a kinding environment K' and an admissible substitution  $K \vdash \theta$ : K' such that for all equality constraint  $\tau_1 \doteq \tau_2$  in  $\varphi$ , we have  $\theta(\tau_1) = \theta(\tau_2)$ .

Rules for unification are in figure 3.  $sort_{\varphi}$  is a partial function from types to  $\{u_1, \ldots, u_n, \rightarrow, k\}$ , left undefined on non-kinded variables:

$sort_{\varphi}(u)$	=	u	u a base type
$sort_{\varphi}(\tau_1 \to \tau_2)$	=	$\rightarrow$	
$sort_{\varphi}(\alpha)$	=	k	$\alpha :: (C, \exists \bar{\epsilon}.R) \in \varphi$

PROPOSITION 5. Rewriting a unification problem leads either to  $\perp$  or to a solved problem.

From a solved unification problem  $\varphi$  one can directly read a pair of a kinding environment K (the kinding constraints) and a substitution  $\theta$  (the solved equality constraints), such that  $K \vdash \theta : K$ .

A substitution  $K \vdash \theta_1 : K_1$  is said to be more general than  $K \vdash \theta_2 : K_2$  if there is a substitution  $K_1 \vdash \theta : K_2$  such that  $\theta_2 = \theta \circ \theta_1$ .

PROPOSITION 6. The solution  $K' \vdash \theta : K'$  read from the solved form obtained by rewriting a unification problem  $\varphi$  of basis K is also  $K \vdash \theta : K'$ , and it is the most general unifier for this problem.

### 7. TYPE RECONSTRUCTION

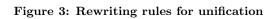
Type reconstruction is done by translating a typing problem  $K; \Gamma \triangleright e : \tau$  into a unification problem  $\varphi$ , and solving it to obtain a substitution  $K \vdash \theta : K'$ .

The algorithm in figure 4 does not depend on a specific evaluation order, equality and kinding constraints only needing to be solved eagerly for let nodes. The function  $solve(\varphi)$ extracts an admissible substitution from a normal form of  $\varphi$ .

A solution to a typing problem  $K; \Gamma \triangleright e : \tau$  is a substitution  $K \vdash \theta : K'$  such that  $K'; \theta(\Gamma) \vdash e : \theta(\tau)$  is derivable.

THEOREM 7. If  $K; \Gamma \triangleright e : \tau$  can be reduced to  $K \vdash \theta : K'$ by the type reconstruction algorithm,  $K'; \theta(\Gamma) \vdash e : \theta(\tau)$  is derivable, and  $\theta$  is the most general solution; otherwise it reduces to  $\bot$  and there is no solution.

INCOMPATIBLE $\frac{\varphi \wedge \tau_1 \doteq \tau_2}{ } \text{ when } sort_{\varphi}(\tau_1) \neq sort_{\varphi}(\tau_2)$	$\frac{\text{Redundancy}}{\varphi \land \tau \doteq \tau}$
$\frac{\varphi \land \alpha \doteq \tau}{\bot} \text{ when } \alpha \neq \tau \text{ and } \alpha \in \mathrm{FV}_{\emptyset}(\tau)$	Function $\frac{\varphi \wedge \tau_1 \rightarrow \tau_2 \doteq \tau_1' \rightarrow \tau_2'}{\varphi \wedge \tau_1 \doteq \tau_1' \wedge \tau_2 \doteq \tau_2'}$
SUBSTITUTION $\frac{\varphi \land \alpha \doteq \tau}{\varphi[\tau/\alpha] \land \alpha \doteq \tau}  \text{when } \alpha :: (C, \exists \bar{\epsilon}.R) \notin \varphi \text{ and } \alpha \notin \mathrm{FV}_{\emptyset}(\tau)$ and $\alpha \in \mathrm{FV}(\varphi)$ and $\tau \neq \beta \lor \beta \in \mathrm{FV}(\varphi)$	
$\frac{\substack{\varphi \land \alpha :: (C, \exists \epsilon. R)}{\bot} \text{ when } C \models \bot$	
$\frac{C_{\text{ONSTRAINT}}}{\varphi \land \alpha_1 :: (C_1, \exists \bar{\epsilon}_1.R_1) \land \alpha_2 :: (C_2, \exists \bar{\epsilon}_2.R_2) \land \alpha_1 \doteq \alpha_2}{\varphi \land \alpha :: (C_1 \land C_2, \exists \bar{\epsilon}_1 \bar{\epsilon}_2.R_1 \cup R_2) \land \alpha_1 \doteq \alpha \land \alpha_2 \doteq \alpha} \alpha \text{ fresh, } \bar{\epsilon}_1 \cap \bar{\epsilon}_2$	$= \emptyset$
$\begin{array}{l} \begin{array}{l} \text{Propagation} \\ t(x,\alpha_1) \wedge t(A,\alpha_2) \wedge p(x) \Rightarrow \alpha_1 = \alpha_2 \in \mathcal{E} \\ \hline \varphi \wedge \alpha :: (C, \exists \overline{\epsilon}.R)  t(a,\tau_1) \in R  r([a/x]A,\tau_2) \in R  C \vdash p(a) \\ \hline \varphi \wedge \alpha :: (C, \exists \overline{\epsilon}.R) \wedge \tau_1 \doteq \tau_2 \end{array} \text{ when } \tau$	$\tau_1 \neq \tau_2$



$\mathbf{K};\Gamma \triangleright e:\tau$	=	$(\Gamma \triangleright e : \tau) \land \mathcal{K}$
$\Gamma, x: (\forall B. \mathbf{K} \triangleright \tau_1) \triangleright x: \tau$	=	$\tau \doteq \theta(\tau_1) \land \bigwedge_{\alpha_i :: (C_i, \exists \bar{\epsilon}_i . R_i) \in \mathbf{K}} \beta_i :: (C_i, \exists \bar{\epsilon}_i . \theta(R_i))$
		where $B = \{\alpha_1, \dots, \alpha_n\}$ and $\theta = \{\alpha_1 \mapsto \beta_1, \dots, \alpha_n \mapsto \beta_n\}$
$\Gamma \triangleright fun \ x \to e : \tau$	=	$(\Gamma, x : \alpha_1 \triangleright e : \alpha_2) \land \tau \doteq \alpha_1 \to \alpha_2$
$\Gamma \triangleright e_1 \ e_2 : \tau$	=	$(\Gamma \triangleright e_1 : \alpha \to \tau) \land (\Gamma \triangleright e_2 : \alpha)$
$\Gamma \triangleright let \ x = e_1 \ in \ e_2 : \tau$	=	$\varphi _{\overline{B}} \wedge (\Gamma, x : (\forall B.\mathbf{K} _B \triangleright \theta(\alpha)) \triangleright e_2 : \tau)$
		where $\varphi = (\Gamma \triangleright e_1 : \alpha)$ and $K \vdash \theta : K = \text{solve}(\varphi)$ and $B = \text{FV}_{K}(\theta(\alpha)) \setminus \text{FV}_{K}(\theta(\Gamma))$
$\Gamma \triangleright c : \tau$	=	$\tau \doteq \theta(\tau_0) \land \bigwedge_{\alpha_i :: (C_i, \exists \bar{\epsilon}_i . R_i) \in \mathcal{K}_0} \beta_i ::: (C_i, \exists \bar{\epsilon}_i . \theta(R_i))$
		where $\operatorname{Tconst}(c) = \operatorname{K}_0 \triangleright \tau_0$ , $\operatorname{FV}_{\operatorname{K}_0}(\tau_0) = \{\alpha_1,, \alpha_n\}$ and $\theta = \{\alpha_1 \mapsto \beta_1, \ldots, \alpha_n \mapsto \beta_n\}$

Figure 4: Reconstruction algorithm

# 8. RELATION TO HM(X)

While we have hinted at subtle differences with HM(X) in previous sections, one may still wonder why we cannot present this system as a subframework of HM(X). The subframework approach not only would relieve us from part of the proof burden (which, luckily, was not so big), it would also allow for a direct comparison with other systems developped inside HM(X).

In this section we explain in detail why this system does not fit inside the HM(X) framework, and why attempting to make it fit into would loose some theoretical and practical properties. In order to do so, we will first define an HM(X) version of our framework, and show on an example how it differs from our version. Then we will formalize some distinctive properties.

## 8.1 Generalization

Obtaining an HM(X) version of our framework is very easy. One just needs to change the definition of the GENERALIZE typing rule.

$$\begin{array}{l} \text{Generalize-HM}(\mathbf{X}) \\ \underline{\mathbf{K}_1, \mathbf{K}_2; \Gamma \vdash e : \tau} \quad B \cap (\mathbf{FV}_{\emptyset}(\mathbf{K}_1) \cup \mathbf{FV}_{\emptyset}(\Gamma)) = \emptyset \\ \overline{\mathbf{K}_1, \exists B.\mathbf{K}_2; \Gamma \vdash e : \forall B.\mathbf{K}_2 \triangleright \tau} \end{array}$$

To make the comparison easier, we consider a syntax-directed version of this rule.

$$\begin{array}{l} \text{GENERALIZE-HM}(\mathbf{X}) \\ \text{K}; \Gamma \vdash e : \tau \quad B = (\mathrm{FV}_{\emptyset}(\tau) \cup \mathrm{FV}_{\emptyset}(\mathbf{K}_B)) \setminus \mathrm{FV}_{\emptyset}(\Gamma) \\ \hline \exists B. \mathbf{K}; \Gamma \vdash e : \forall B. \mathbf{K}_B \triangleright \tau \end{array}$$

where  $K_B$  is the subset of K containing all kinds related to B (either as name of the kind, or inside R). We solve the mutual recursion between B and  $K_B$  by choosing the smallest B statisfying the equation.

Here is a concrete example where the HM(X) version would infer a different type. Let's consider the function

fun 
$$x \to$$
let  $y = ($ function  $A(z) \to z \mid B(z) \to z) x$  in  $y$ 

If we write K for the kinding environment  $\alpha :: (\emptyset, \{A, B\}, \{A : \beta, B : \beta\})$ , the typing derivation for the function application becomes (independently of the framework used)

$$K; x : \alpha \vdash (function \ A(z) \rightarrow z \mid B(z) \rightarrow z) \ x : \beta$$

According to our framework,  $\beta$  is a free variable of  $\alpha$ , so we cannot generalize this type. However, using rather GENERA-LIZE-HM(X), one can reach the following conclusion.

$$\mathbf{K} ; x : \alpha \vdash y : \forall \beta. \alpha :: (\emptyset, \{A, B\}, \{A:\beta, B:\beta\}) \triangleright \beta$$

You can see here the peculiar form of the polymorphism obtained: the kinding part in  $\forall B.K \triangleright \tau$  is no longer restricted to variables of *B* (as requires our type reconstruction algorithm). And a polymorphic variable ( $\beta$ ) may actually depend on a non-polymorphic one ( $\alpha$ ).

In the above example, the difference in typing is only superficial, making the extra prolymorphism spurious: actually one gets exactly the same polymorphism when considering the whole term. While there exist examples where GENERALIZE-HM(X) provides really more polymorphism, they are rare, and there are workarounds.

## 8.2 Relevant variables

A formal way to distinguish the two systems is by considering relevant type variables. Intuitively, a variable is relevant if the meaning of a type depends on it proper (*i.e.* its whole meaning is included in the type).

More specifically, we consider the meaning of a type as the set of terms typable by one of its instances, and we say that a type variable  $\alpha$  is *relevant* to a type  $K \triangleright \tau$ , if for some instance  $K' \triangleright \tau'$  of this type, there is a term e with a derivation  $K'; \emptyset \vdash e : \tau'$ , such that any substitution  $K \vdash \theta : K''$ , verifying  $K''; \emptyset \vdash e : \theta(\tau)$ , has either to substitute  $\alpha$  or have it bound in K'' to a kind with a different constraint (a kind with different related types does not imply relevance).

Our system has the following property: all relevant variables of  $K \triangleright \tau$  are included in  $FV_K(\tau)$ . This is not the case with HM(X) based systems. For instance, in the above example,  $\beta$  is clearly relevant to  $K \triangleright \alpha$  (*e.g.* A(1) binds  $\beta$  to *int*), but it is not in  $FV_{\emptyset}(\alpha) = \{\alpha\}$ .

This property gives a clean way to print a type, with all relevant type information, by inserting a kind where its variable occurs. For instance, in Objective Caml, we would write

$$[ < A \text{ of } \beta \mid B \text{ of } \beta ] \rightarrow \beta$$

for the above term's type.

This property is also needed for tracing sharing in types, as we do for semi-explicit first-class polymorphism [7]. Sharing is easily obtained by making each explicit polytype a kind  $\alpha :: [\sigma]$  and using  $\alpha$  for each shared occurence of this polytype, following the idea in the section 2.6 of [7]. Since HM(X) fails to capture relevance in the type system, our first-class polymorphism would not have principal inference for it.

A possibly simpler way to characterize the relation to relevance is at the type inference level. In our system, if  $\alpha \notin FV_K(\tau_1) \cup FV_K(\tau_2)$ , then solving  $K \wedge \tau_1 = \tau_2$  will not bind  $\alpha$ . Of course, this is not true for HM(X).

This has technical consequences. For instance this property is needed to do "in place generalization", that is marking directly generic variables in a type. This is only possible when generalizable variables are guaranteed not to be shared with types in the environment. While "in place generalization" is not a goal in itself, detecting the extent of the enclosing type which should be copied, as expected by GENERALIZE-HM(X), would require a fair amount of extra machinery.

This is also an important property for programmers, as they can easily identify the consequences of a type constraint: to prove that a constraint does not affect a type, one only needs to check that they do not share free variables. With HM(X), the baseline is to consider the transitive closure of all constraints sharing some variables, which potentially contains lots of irrelevant variables.

# 9. CONCLUSION

We proposed a type system where constraints are local to type nodes, and showed how it can be applied to many type inference problems, which fit between usual parametric polymorphism and systems requiring the power of subtyping. Its simplicity, and the unperspicuous way in which it extends the original type inference algorithm, make it a good candidate for practical implementations and further extensions.

While we think that this system has a wide range of applicability, as is shown by a new structural encoding for accurate matching, it does not cover the full range of structural polymorphism. For instance, Rémy's row and presence variables provide for slightly more precise typing than the encoding we give. Similarly Kennedy's dimension types use a multi-sorted algebra [8]. Nonetheless, if we restrict sharing in types to one single sort of type variables, as is done in Objective Caml's type system (which is itself based on Rémy's work), then we come back into the expressive area of local constraints.

An interesting conjecture is whether local constraints can handle all one-sorted ML-like type systems without the "power of subtyping". This would confirm the intuition that there is a qualitative jump when going from systems without subtyping to systems which allow it. In particular, local constraints fail to encode simply typed dynamic messages or record concatenation [19], correctly showing that, while this is not immediately explicit, all these systems require some form of "subtyping" or multi-directional constraints.

Ironically the last counter-example standing against a characterization of structural polymorphism by local constraints seems to be the label-selective lambda-calculus and its type system [6]. Function types are of the form  $l: \tau \to \tau$ , module the congruence  $l_1:\tau_1 \to l_2:\tau_2 \to \tau = l_2:\tau_2 \to l_1:\tau_1 \to \tau$ . Clearly, this bears no relation to subtyping, as shown by the existence of a typing algorithm without constraints, yet one has to separate  $l_1:\tau_1 \to \tau$  from  $\tau$ , which is not possible with the local constraints presented in this paper. The intuitive encoding would be to describe function types as kinds containing all applicable labels, and add a propagation rule of the form  $l_1 \neq l_2 \wedge r(l_1, \alpha_1) \wedge r(l_2, \alpha_2) \Rightarrow \alpha_1 =$  $\alpha' \wedge \alpha' :: (\{l_2\}, \{r(l_2, \alpha_2)\})$ , introducing equalities with newly built types. Not only does it break our proof of termination for unification, but in general this form of propagation rules could be used to encode subtyping. This suggests that labelselective lambda-calculus requires some form of non-local relations, which is coherent with the presence of a fundamentally non-local requirement in the COMPLETION rule of its unification algorithm [4], that the return type variable after an arbitrary number of arrows should differ<sup>3</sup>. Actually finding what are these non-local relations may be an interesting topic for further work, as the same phenomenon appears in systems with explicit row variables [16, 15].

## **10. REFERENCES**

 A. Aiken, E. L. Wimmers, and T. K. Lakshman. Soft typing with conditional types. In *Proc. ACM* Symposium on Principles of Programming Languages, pages 163–173, Portland, Oregon, Jan. 1994.

- [2] H. Aït-Kaci, A. Podelski, and G. Smolka. A feature-based constraint system for logic programming with entailment. *Theoretical Computer Science*, 122(1-2):263-283, Jan. 1994.
- [3] L. Damas and R. Milner. Principal type-schemes for functional programs. In Proc. ACM Symposium on Principles of Programming Languages, pages 207–212, 1982.
- [4] J. P. Furuse and J. Garrigue. A label-selective lambda-calculus with optional arguments and its compilation method. RIMS Preprint 1041, Research Institute for Mathematical Sciences, Kyoto University, Oct. 1995.
- [5] J. Garrigue. Programming with polymorphic variants. In *ML Workshop*, Baltimore, Sept. 1998.
- [6] J. Garrigue and H. Aït-Kaci. The typed polymorphic label-selective λ-calculus. In Proc. ACM Symposium on Principles of Programming Languages, pages 35–47, 1994.
- [7] J. Garrigue and D. Rémy. Extending ML with semi-explicit higher order polymorphism. *Information* and Computation, 155:134–171, Dec. 1999.
- [8] A. J. Kennedy. Dimension types. In Proceedings of the 5th European Symposium on Programming, volume 788 of Lecture Notes in Computer Science, pages 348–362. Springer, 1994.
- [9] X. Leroy, D. Doligez, J. Garrigue, D. Rémy, and J. Vouillon. *The Objective Caml system release 3.00, Documentation and user's manual.* Projet Cristal, INRIA, Apr. 2000.
- [10] R. Milner. A theory of type polymorphism in programming. Journal of Computer and System Sciences, 17:348–375, 1978.
- [11] M. Müller and S. Nishimura. Type inference for first-class messages with feature constraints. *International Journal of Foundations of Computer Science*, 11(1):29–63, 2000.
- [12] S. Nishimura. Static typing for dynamic messages. In Proc. ACM Symposium on Principles of Programming Languages, San Diego, California, Jan. 1998.
- [13] M. Odersky, M. Sulzmann, and M. Wehr. Type inference with constrained types. *Theory and Practice* of Object Systems, 5(1):35–55, 1999.
- [14] A. Ohori. A polymorphic record calculus and its compilation. ACM Trans. Prog. Lang. Syst., 17(6):844–895, Nov. 1995.
- [15] F. Pottier. A versatile constraint-based type inference system. Nordic Journal of Computing, 7(4):312–347, Nov. 2000.
- [16] D. Rémy. Typechecking records and variants in a natural extension of ML. In Proc. ACM Symposium on Principles of Programming Languages, pages 77–87, 1989.

<sup>&</sup>lt;sup>3</sup>This requirement was missing in [6].

- [17] D. Rémy. Type inference for records in a natural extension of ML. In C. A. Gunter and J. C. Mitchell, editors, *Theoretical Aspects Of Object-Oriented Programming. Types, Semantics and Language Design.* MIT Press, 1993.
- [18] M. Sulzmann. A General Framework for Hindley/Milner Type Systems with Constraints. PhD thesis, Yale University, Department of Computer Science, May 2000.
- [19] M. Wand. Type inference for record concatenation and multiple inheritance. *Information and Computation*, 93:1–15, 1991.
- [20] A. K. Wright and M. Felleisen. A syntactic approach to type soundness. *Information and Computation*, 115(1):38–94, Nov. 1994.

# APPENDIX

## A. PROOFS OF THEOREMS

PROPOSITION 1 (DISCARDING POSTPONEMENT). If K,  $\alpha :: k \vdash \theta : K', \theta(\alpha) :: k' \text{ and } k_0 \cong_{\mathcal{E}} k$ , then there is a kind  $k'_0$  such  $k'_0 \cong_{\mathcal{E}} k'$  and  $K, \alpha :: k_0 \vdash \theta : K, \theta(\alpha) :: k'_0$ .

PROOF. Suppose  $k_0 = (C, \exists \bar{\epsilon}.R \cup R_0)$  and  $k = (C, \exists \bar{\epsilon}.R \cup R_1)$ , with  $R_0$  and  $R_1$  minimal. By definition of  $\cong_{\mathcal{E}}$ , all  $r(a, \tau)$  in  $R_0$  are inert for C. By definition of admissibility,  $k' = (C', \exists \bar{\epsilon}.R')$ , with  $C' \models C$ , so that all elements of  $R_0$  are also inert for C'; and  $\theta(R_1) \subset R'$ , so that  $k'_0 = (C', \exists \bar{\epsilon}.R' \cup \theta(R_0)) \cong_{\mathcal{E}} (C', \exists \bar{\epsilon}.R' \cup \theta(R_1)) = k'$ .  $\Box$ 

LEMMA 2 (TYPE SUBSTITUTION). If  $K, \Gamma \vdash e : \tau$  and  $K \vdash \theta : K'$ , then  $K', \theta(\Gamma) \vdash e : \theta(\tau)$ .

PROOF. Induction on the derivation of  $K, \Gamma \vdash e : \tau$ .

The two important steps are VARIABLE and GENERALIZE. For VARIABLE, we have  $K, K_0 \vdash \theta_0 : K$  and  $Dom(\theta_0) \subset B$ ; we can choose B outside of  $\theta$  and  $K': \forall \alpha \in B, \theta(\alpha) = \alpha$ , and  $FV(K') \cap B = \emptyset$ . By composition, we construct  $\theta_1 = (\theta \circ \theta_0)|_B$ , such that  $K', \theta(K_0) \vdash \theta_1 : K'$ : conditions (1) and (2) are satisfied by transitivity. We check condition (3) for kinded variables:

- If  $\alpha \notin B$ , then  $\theta_1(\alpha) = \alpha$ ,  $\alpha :: (C, \exists \bar{\epsilon}.R) \in K'$  and  $FV_{\emptyset}(R) \cap B = \emptyset$ , so that  $\theta_1(R) = R$ .
- If  $\alpha \in B$ , then  $\theta_1(\alpha) = \theta(\theta_0(\alpha))$ , and  $\alpha :: (C, \exists \bar{e}.R) \in K_0$ . For any variable  $\alpha' \in FV_{\emptyset}(R)$ , we have  $\theta_1(\theta(\alpha')) = \theta(\theta_0(\alpha'))$ : either  $\alpha' \in B$ , and  $\theta(\alpha') = \alpha'$ , so that  $\theta_1(\theta(\alpha')) = \theta_1(\alpha') = \theta(\theta_0(\alpha'))$ , or  $\alpha' \notin B$ , and  $\theta_0(\alpha') = \alpha'$ , giving  $\theta_1(\theta(\alpha')) = \theta(\alpha') = \theta(\theta_0(\alpha'))$ . Combined with  $\theta_1(\alpha) :: (C', \exists \bar{e}.R') \in K'$  implying  $\theta(\theta_0(R)) \subset R'$ , we obtain  $\theta_1(\theta(R)) \subset R'$ , giving condition (3) since  $\alpha :: (C, \exists \bar{e}.\theta(R)) \in \theta(K_0)$ .

Note also that  $\theta_1 \circ \theta|_{\overline{B}} = \theta \circ \theta_0$ . Then we can deduce  $K'; \theta(\Gamma), x: \forall B, \theta(K_0) \triangleright \theta(\tau) \vdash x: \theta(\theta_0(\tau)).$ 

For GENERALIZE, again we choose B outside of  $\theta$  and K'. Then  $K|_{\overline{B}}, \theta(K|_B) \vdash \theta : K', \theta(K|_B)$ . By induction hypothesis we have  $\mathbf{K}', \theta(\mathbf{K}|_B); \theta(\Gamma) \vdash e : \theta(\tau)$ . Then we can deduce  $\mathbf{K}'; \theta(\Gamma) \vdash e : \forall B.\theta(\mathbf{K}|_B) \triangleright \theta(\tau)$ .  $\Box$ 

LEMMA 3 (TERM SUBSTITUTION). If  $K; \Gamma \vdash e : \sigma$  and  $K; \Gamma, x : \sigma \vdash e' : \tau'$ , then  $K; \Gamma \vdash e'[e/x] : \tau'$ .

PROOF. Easy induction on the derivation tree. When grafting GENERALIZE somewhere, use a type substitution  $K \vdash \theta : K|_{\overline{B}}$ . When substituting under GENERALIZE, *B* may grow. In such case, preserve  $K|_{\overline{B}}$  by duplicating the kind under a different name, and instantiate extra variables with VARIABLE.  $\Box$ 

THEOREM 4 (SUBJECT REDUCTION). If K;  $\Gamma \vdash E[e] : \sigma$ and  $e \rightarrow e'$  by a  $\delta$ -rule, BETA :  $((fun \ x \rightarrow e_1) \ e_2) \rightarrow [e_2/x]e_1$ or Let : let  $x = e_1$  in  $e_2 \rightarrow [e_2/x]e_1$ , then K;  $\Gamma \vdash E[e'] : \sigma$ 

PROOF. First we consider the case when E[e] = e. If  $e \to e'$  is the instance of the  $\delta$ -rule  $A[x_1, \ldots, x_n] \to A'[x_1, \ldots, x_n]$ , then, supposing  $e_1 : \tau_2, \ldots, e_n : \tau_n$ , we have already a derivation of  $K_0; x_1 : \tau_1, \ldots, x_n : \tau_n \vdash e'[x_1, \ldots, x_n] : \tau$ . We conclude by term substitution. Similarly  $\beta$ -reduction and letreduction are immediate consequences of term substitution.

If the redex occurs inside a deeper context, extend the proof by induction on the depth of the context.  $\Box$ 

PROPOSITION 5. Rewriting a unification problem leads either to  $\perp$  or to a solved problem.

PROOF. First, we verify that unification terminates. One is allowed to apply rules in any order, except that (1) failure is eager (*i.e.* INCOMPATIBLE, CYCLIC and BAD CONSTRAINT are given highest priority), and (2) PROPAGATION only applies when no other rule applies (*i.e.* it has lowest priority).

Our measure is a lexicographical ordering on the tuple (number of kind constraints, number of pending propagations, number of unsolved variables, number of arrows, number of equality constraints). CONSTRAINT reduces the number of kind constraints by 1, while possibly increasing the number of pending propagations. PROPAGATION eventually reduces the number of pending propagations (if substitutions make  $\tau_1$  equal to  $\tau_2$ ) or the number of kind constraints (if it causes a constraint unification); no new propagation will occur until one of these happen. SUBSTITUTION reduces the number of unsolved variables. FUNCTION reduces the number of arrows. REDUNDANCY reduces the number of equations.

Then we check that reduction cannot be stuck. If a constraint is stuck, one of its members is not solved. If this is an equality constraint, either SUBSTITUTION applies, or it is  $\alpha \doteq \tau[\alpha]$  and CYCLIC applies, or one member is a kinded variable and either CONSTRAINT OF INCOMPATIBLE applies, or both are not type variables and one of INCOMPATIBLE, REDUNDANCY OF FUNCTION applies. If this is a kinding constraint, then either the constraint is not satisfiable and BAD CONSTRAINT applies, or a propagation equation is not satisfied, and since other rules do not apply, PROPAGATION applies.

PROPOSITION 6. The solution  $K' \vdash \theta : K'$  read from the solved form obtained by rewriting a unification problem  $\varphi$  of basis K is also  $K \vdash \theta : K'$ , and it is the most general unifier for this problem.

PROOF. For each rewriting rule, we verify that it is sound and complete. Together with soundness, we also check changes in kindings were necessary.

Rules INCOMPATIBLE, CYCLIC, REDUNDANCY and FUNCTION are immediate: the upper and lower constraints accept the same solutions.

For SUBSTITUTION, soundness is clear:  $\theta(\alpha) = \theta(\tau)$  then if  $\theta$  is a solution of  $\varphi[\tau/\alpha]$ , it is also a solution of  $\varphi$ . Looking at kindings, let  $K_1$  be the basis for  $\varphi$ , and  $K_2$  the basis for  $\varphi[\tau/\alpha]$ . We have  $K_1 \vdash \{\alpha \mapsto \tau\} : K_2$ , and  $K_2 \vdash \theta$ : K', since  $\theta(\alpha) = \theta(\tau)$ ,  $\theta \circ \{\alpha \mapsto \tau\} = \theta$ , and  $K_1 \vdash \theta$ :  $\tau$ . For completeness, the side-conditions guarantee that  $\alpha$  has no kind (reduction rules only introduce kinds for fresh variables). Then it is safe to replace  $\alpha$  by  $\tau$  in  $\varphi$ , and the same solutions are accepted.

For BAD CONSTRAINT, no substitution can apply to  $\alpha$ , so that the constraint is equivalent to  $\perp$ .

For CONSTRAINT, if we have a solution  $\mathbf{K} \vdash \theta : \mathbf{K}'$  of the upper side, then it has to map both  $\alpha_1$  and  $\alpha_2$  to  $\alpha' :: (C', \exists \bar{\epsilon}'. R')$ , with  $C' \models C_1$  and  $C' \models C_2$ . Since  $C_1 \land C_2$  is the weakest constraint implying both  $C_1$  and  $C_2$ , we can extend  $\theta$  into  $\mathbf{K}, \alpha :: (C_1 \land C_2, \exists \bar{\epsilon}_1 \bar{\epsilon}_2. R_1 \cup R_2) \vdash \theta\{\alpha \mapsto \alpha'\} : \mathbf{K}'$ , and it is a solution to the lower side  $(\theta(\alpha) = \theta(\alpha_1) = \theta(\alpha_2) = \alpha')$ , showing completeness. Reciprocally, with  $\mathbf{K}_1$  and  $\mathbf{K}_2$  bases for the upper and lower sides, if we have a solution  $\mathbf{K}_2 \vdash \theta :$  $\mathbf{K}'$  of the lower side, then  $\mathbf{K}_1 \vdash \{\alpha_1 \mapsto \alpha, \alpha_2 \mapsto \alpha\} : \mathbf{K}_2$ , and since  $\theta(\alpha) = \theta(\alpha_1) = \theta(\alpha_2), \theta \circ \{\alpha_1 \mapsto \alpha, \alpha_2 \mapsto \alpha\} = \theta$ , and  $\mathbf{K}_1 \vdash \theta : \mathbf{K}'$  is a solution to the upper side.

For PROPAGATION, soundness is immediate (the lower side subsumes the upper side). Any solution  $K \vdash \theta : K'$  of the upper side must have  $\theta(\tau_1) = \theta(\tau_2)$  to be admissible, so it is also complete.

Soundness shows that the solution  $\mathbf{K}' \vdash \theta : \mathbf{K}'$  we read from a solved unification problem, is also a solution  $\mathbf{K} \vdash \theta : \mathbf{K}'$ to the original problem. Thanks to completeness, we know that any other solution  $\mathbf{K} \vdash \theta_1 : \mathbf{K}''$  of the original problem can be extended in a solution  $\mathbf{K}' \vdash \theta_2 : \mathbf{K}''$  of the solved problem. Since  $\theta_2$  satisfies all the solved constraints, there is a  $\mathbf{K}' \vdash \theta_3 : \mathbf{K}''$  such that  $\theta_2 = \theta_3 \circ \theta$ , and, when restricted to the original domain,  $\theta_1 = (\theta_3 \circ \theta)|_{\mathbf{FV}(\varphi)}$ . So  $\theta$  is the most general unifier.  $\Box$ 

THEOREM 7. If  $K; \Gamma \triangleright e : \tau$  can be reduced to  $K \vdash \theta : K'$ by the type reconstruction algorithm,  $K'; \theta(\Gamma) \vdash e : \theta(\tau)$  is derivable, and  $\theta$  is the most general solution; otherwise it reduces to  $\bot$  and there is no solution.

PROOF. We verify by induction on the structure of e that an inference problem is translated into an equivalent unification problem. The first case only splits a complete problem in its kinding environment and typing problem.

The second case handles the VARIABLE rule. Kinded  $\beta_i$ 's ensure soundness: if  $\mathbf{K}' = \{\beta_i :: (C_i, \exists \bar{\epsilon}_i.\theta(R_i))\}_{\alpha_i::(C_i, \exists \bar{\epsilon}_i.R_i) \in \mathbf{K},}$  then  $\mathbf{K} \vdash \theta : \mathbf{K}'$ , so that  $\mathbf{K}', \mathbf{K} \vdash \theta : \mathbf{K}'$  (since  $\theta|_{\text{Dom}(\mathbf{K}')} = id$ ), and any solution  $\mathbf{K}' \vdash \theta' : \mathbf{K}''$  of the right hand side is such that  $\mathbf{K}'', \mathbf{K} \vdash (\theta' \circ \theta)|_B : \mathbf{K}''$ ; this proves  $\mathbf{K}''; \theta'(\Gamma), x : \forall B.\theta'(\mathbf{K}) \triangleright \theta'(\tau_1) \vdash x : (\theta' \circ \theta)|_B(\theta'(\tau_1)), \text{ and } (\theta' \circ \theta)|_B(\theta'(\tau_1)) = \theta'(\theta(\tau_1)) = \theta'((\tau_1)) = \theta'(\tau)$  gives the left hand side. For completeness, suppose  $\theta'$  is a solution to the left hand side. Then  $\theta'|_B$  is the instantiation substitution. We can extend  $\theta'$  in  $\theta_1 = \theta' \circ \theta^{-1}$ , since  $\beta_i$ 's are fresh, and  $\theta_1$  is still a solution of the left hand side.

Abstraction and Application are as in ML.

The LET rule requires to solve constraints on  $e_1$  before proceeding. This corresponds to generating a derivation, and adding the GENERALIZE step. The result is then used to type  $e_2$ . Information related to variables of B is discarded from the global constraint as superfluous. This is sound as constraints were already solved once, and cannot be refined as these variables are not accessible through the environment. Completeness requires considering the case when we do not use the GENERALIZE step on  $e_1$ . Generalizable variables are kept as free variables. Then they cannot be generalized when typing  $e_2$ , since they are in the environment. Since this is equivalent to choosing  $\theta = id$  in the corresponding VARIABLE rules, and the VARIABLE case is complete, LET is complete.

CONSTANT is similar to VARIABLE.

Since we already proved rewriting on unification problems to terminate, type inference terminates with either a solution or  $\perp$ .  $\Box$