

# A theory of qualified types

Mark P. Jones\*

*Yale University, Department of Computer Science, P.O. Box 208285, New Haven, CT 06520-8285,  
USA*

Received July 1992; revised January 1993

---

## Abstract

This paper describes a general theory of overloading based on a system of qualified types. The central idea is the use of predicates in the type of a term, restricting the scope of universal quantification. A corresponding semantic notion of evidence is introduced and provides a uniform framework for implementing applications of this system, including Haskell style type classes, extensible records and subtyping.

Working with qualified types in a simple, implicitly typed, functional language, we extend the Damas–Milner approach to type inference. As a result, we show that the set of all possible typings for a given term can be characterized by a principal type scheme, calculated by a type inference algorithm.

---

## 1. Introduction

In a language with a polymorphic type system, a term of type  $\forall t.f(t)$  can be treated (possibly after suitable instantiation) as having any of the types in the set:

$$\{f(t) \mid t \text{ is a type}\}.$$

It is natural to consider a more restricted form of polymorphism in which the value taken by  $t$  may be constrained to a particular subset of types. In this situation, we write  $\forall t.\pi(t) \Rightarrow f(t)$ , where  $\pi(t)$  is a predicate on types, for the type of an object that can be treated (after suitable instantiation) as having any of the types in the set:

$$\{f(t) \mid t \text{ is a type such that } \pi(t) \text{ holds}\}.$$

---

\* E-mail: jones-mark@cs.yale.edu.

A term with a restricted polymorphic type of this kind is often said to be *overloaded*, having different interpretations for different argument types.

This paper presents a general theory of overloading based on the use of *qualified types*, which are types of the form  $\pi \Rightarrow \sigma$  denoting those instances of type  $\sigma$  that satisfy the predicate  $\pi$ . The main benefits of using qualified types are:

- A general approach that includes a range of familiar type systems as special cases. Results and tools developed for the general system are immediately applicable to each particular application.
- A precise treatment of the relationship between implicit and explicit overloading. This is particularly useful for describing the implementation of systems supporting qualified types.
- The ability to include local constraints as part of the type of an object. This enables the definition and use of polymorphic overloaded values within a program.

## 2. Outline of the paper

Each of the type systems considered in this paper is parameterized by the choice of a system of predicates on type expressions, whose basic properties are described in Section 3. A number of examples are included to illustrate the use of this framework to describe a range of type systems including Haskell type classes, extensible records and subtyping. Section 4 describes the use of qualified types in the context of polymorphic  $\lambda$ -calculus with explicit typing. This is extended in Section 5 using a general notion of *evidence* to explore the relationship between implicit and explicit overloading. An alternative approach, suitable for use in an implicitly typed language, is introduced in Section 6 using an extension of the ML type system [14] to support qualified types. Although substantially less powerful than polymorphic  $\lambda$ -calculus, we show that the resulting system is suitable for use in a language based on type inference, that allows the type of a term to be determined without explicit type annotations. The development of a suitable type inference algorithm is described in Sections 7 and 8. Finally, Section 9 surveys some areas for further work.

Detailed proofs for many of the results described in this paper may be found in [12,13]; for reasons of space, they cannot be included here.

## 3. Predicates

Each of the type systems considered in this paper is parameterized by the choice of a language of *predicates*  $\pi$  whose properties are described by an entailment relation  $\vdash$  between (finite) sets of predicates. Individual predicates may be written using expressions of the form  $\pi = p \tau_1 \dots \tau_n$  where  $p$  is

a predicate symbol corresponding to an  $n$ -place relation between types; the predicate  $\pi$  represents the assertion that the types  $\tau_1, \dots, \tau_n$  are in this relation. The definition of  $\vdash$  varies from one application to another. The only properties that we will assume are:

- *Monotonicity.*  $P \vdash P'$  whenever  $P \supseteq P'$ .
- *Transitivity.* If  $P \vdash Q$  and  $Q \vdash R$ , then  $P \vdash R$ .
- *Closure property.* If  $P \vdash Q$ , then  $SP \vdash SQ$  for any substitution  $S$  mapping type variables (and hence type expressions) to type expressions.

If  $P$  is a set of predicates and  $\pi$  is a predicate, then we write  $P \vdash \pi$  and  $P, \pi$  as abbreviations for  $P \vdash \{\pi\}$  and  $P \cup \{\pi\}$  respectively.

The following subsections illustrate the languages of predicates that might be used in three applications of qualified types. Only the basic ideas are sketched here; further details are given in [11,13].

### 3.1. Example: type classes

Introduced in [23] and adopted as part of the standard for the programming language Haskell [10], type classes are particularly useful for describing the implementation of standard polymorphic operators such as computable equality. Much of the original motivation for qualified types came from the study of type classes.

Broadly speaking, a *type class* is a family of types (the *instances* of the class) on which a number of values (the *member functions*) are defined. Each predicate symbol corresponds to a user-defined class and a predicate of the form  $C \tau$  represents the assertion that  $\tau$  is an instance of the class named  $C$ . The class  $Eq$  is a standard example whose instances are those types whose elements can be tested for equality using the operator  $(==) :: \forall a. Eq a \Rightarrow a \rightarrow a \rightarrow Bool$ . As a further example, one possible type for a function to test for membership of a value in a list is  $\forall a. Eq a \Rightarrow a \rightarrow [a] \rightarrow Bool$  where  $[a]$  denotes the type of lists of values of type  $a$ .

Differences in the basic approach to type classes are reflected in the properties of the  $\vdash$  relation. In a standard Haskell system we have axioms such as  $\emptyset \vdash EqInt$  and  $Eq a \vdash Eq[a]$ . The same framework can also be used to describe the use of Haskell *superclasses*, and to support the extension to classes with multiple parameters.

Type classes are best suited to systems with a type inference algorithm such as that described in Section 8 where the appropriate instances of each overloaded operator can be determined automatically as part of the type inference process.

### 3.2. Example: extensible records

A record is a set of values labeled by the elements  $l$  of a specified set of *labels*. There has been considerable interest in the use of record types to

model inheritance in object oriented programming languages and a number of different approaches have been considered. We can construct a system of extensible records, strongly reminiscent of [9] using predicates of the form:

$r$  has  $l:t$

indicating that a record of type  $r$  has a field labeled  $l$  of type  $t$ , and

$r$  lacks  $l$

indicating that a record of type  $r$  does not have a field labeled  $l$ . This also requires an extension of the language of type expressions to allow types of the form  $\langle \rangle$  (the empty record, which lacks any fields),  $r \setminus l$  (the type of a record obtained by removing a field labeled  $l$  from a record of type  $r$ ) and  $\langle r \mid l:t \rangle$  (the type of a record obtained by extending a record of type  $r$  with a new field of type  $t$  labeled  $l$ ). The definition of the entailment relation includes axioms such as  $\emptyset \Vdash (\langle \rangle \text{ lacks } l)$  and  $r \text{ lacks } l \Vdash (\langle r \mid l:t \rangle \text{ has } l:t)$ .

The primitive operations of record *restriction*, *extension* and *selection* can then be represented by families of functions (indexed by labels) of type:

$$\begin{aligned} (- \setminus l) &:: \forall r. \forall t. (r \text{ has } l:t) \Rightarrow r \rightarrow r \setminus l \\ (- \mid l = -) &:: \forall r. \forall t. (r \text{ lacks } l) \Rightarrow r \rightarrow t \rightarrow \langle r \mid l:t \rangle \\ (-.l) &:: \forall r. \forall t. (r \text{ has } l:t) \Rightarrow r \rightarrow t \end{aligned}$$

Details of the relationship between this approach and those of [5,19] are given in [9]. The approach described here leads to a promising new implementation of extensible records, outlined in [13]. A similar system, dealing only with the use of records rather than our more general qualified types, has been proposed by Ohori [17].

### 3.3. Example: subtyping

Languages with subtyping can be described using predicates of the form  $\sigma \subseteq \sigma'$ , representing the assertion that  $\sigma$  is a subtype of  $\sigma'$ . Many proposals, including those of [8,15], allow the use of implicit coercions from one type to another. The type systems described in this paper are less flexible, but still of interest as a result of the interaction between subtyping and polymorphism. The extensions required to support implicit coercions are discussed in Section 9.4.

### 3.4. Example: combined systems

Although we have dealt with each of the previous three applications individually, it would also be possible to combine elements of each in a single predicate system. This in turn leads to the design of programming languages which support all of these different features within a uniform framework.

#### 4. Polymorphic $\lambda$ -calculus with qualified types

##### 4.1. Basic definitions

In this section, we work with a variant of the polymorphic  $\lambda$ -calculus that includes qualified types using type expressions of the form:

$$\sigma ::= t \mid \sigma \rightarrow \sigma \mid \forall t. \sigma \mid \pi \Rightarrow \sigma$$

where  $t$  ranges over a given set of type variables. The  $\rightarrow$  and  $\Rightarrow$  symbols are treated as right associative infix binary operators with  $\rightarrow$  binding more tightly than  $\Rightarrow$ . Additional type constructors such as those for integers, lists and record types will be used as required. The set of type variables appearing (free) in an expression  $X$  is denoted  $TV(X)$ .

To begin with we use an unmodified form of the (unchecked) terms of polymorphic  $\lambda$ -calculus, given by expressions of the form:

$$M ::= x \mid MN \mid \lambda x:\sigma. M \mid M\sigma \mid \lambda t. M$$

where  $x$  ranges over a given set of term variables. The set of free (term) variables appearing in a term  $M$  will be denoted  $FV(M)$ . Note that we do not provide constructs for the introduction of new overloadings such as **inst** and **over** in [23]. If none of the free variables for a given term have qualified (i.e. overloaded) types, then no overloading will be used in the expression.

##### 4.2. Typing rules

A *type assignment* is a (finite) set of *typing statements* of the form  $x:\sigma$  in which no term variable  $x$  appears more than once. If  $A$  is a type assignment, then we write  $\text{dom } A = \{x \mid (x:\sigma) \in A\}$ , and if  $x$  is a term variable with  $x \notin \text{dom } A$ , then we write  $A, x:\sigma$  as an abbreviation for the type assignment  $A \cup \{x:\sigma\}$ . The type assignment obtained from  $A$  by removing any typing statement for the variable  $x$  is denoted  $A_x$ . A type assignment  $A$  can be interpreted as a function mapping each element of  $\text{dom } A$  to a type scheme. In particular, if  $(x:\sigma) \in A$ , then we write  $A(x) = \sigma$ .

An expression of the form  $P \mid A \vdash M : \sigma$  represents the assertion that the term  $M$  has type  $\sigma$  when the predicates in  $P$  are satisfied and the types of free variables in  $M$  are as specified in the type assignment  $A$ . The typing rules for this system are given in Fig. 1. Most of these are similar to the rules for explicit typing of polymorphic  $\lambda$ -calculus and do not involve the predicate set.

By an abuse of notation, we will also use  $P \mid A \vdash M : \sigma$  as a proposition asserting the existence of a derivation of  $P \mid A \vdash M : \sigma$ .

<b>Standard rules:</b>	(var)	$\frac{(x:\sigma) \in A}{P \mid A \vdash x : \sigma}$
	( $\rightarrow$ E)	$\frac{P \mid A \vdash M : \sigma' \rightarrow \sigma \quad P \mid A \vdash N : \sigma'}{P \mid A \vdash MN : \sigma}$
	( $\rightarrow$ I)	$\frac{P \mid A, x:\sigma' \vdash M : \sigma}{P \mid A \vdash \lambda x:\sigma'.M : \sigma' \rightarrow \sigma}$
<b>Qualified types:</b>	( $\Rightarrow$ E)	$\frac{P \mid A \vdash M : \pi \Rightarrow \sigma \quad P \Vdash \pi}{P \mid A \vdash M : \sigma}$
	( $\Rightarrow$ I)	$\frac{P, \pi \mid A \vdash M : \sigma}{P \mid A \vdash M : \pi \Rightarrow \sigma}$
<b>Polymorphism:</b>	( $\forall$ E)	$\frac{P \mid A \vdash M : \forall t.\sigma}{P \mid A \vdash M\tau : [\tau/t]\sigma}$
	( $\forall$ I)	$\frac{P \mid A \vdash M : \sigma \quad t \notin TV(A) \cup TV(P)}{P \mid A \vdash \lambda t.M : \forall t.\sigma}$

Fig. 1. Typing rules for polymorphic  $\lambda$ -calculus with qualified types.

## 5. Evidence

Although the system of qualified types described in the previous sections is suitable for reasoning about the types of overloaded terms, it cannot be used to describe their evaluation. For example, the knowledge that *Int* is an instance of the class *Eq* is not sufficient to determine the value of the expression  $2 == 3$ ; we must also be provided with the value of the equality operator that makes *Int* an instance of *Eq*. In general, we can only use a term of type  $\pi \Rightarrow \sigma$  if we are also supplied with suitable *evidence* that the predicate  $\pi$  does indeed hold.

This leads us to consider an extension of the term language that makes the role of evidence explicit, using:

- *Evidence expressions*: A language of *evidence expressions*  $e$  denoting evidence values, including a set of *evidence variables*  $v$ .
- *Evidence construction*: A *predicate assignment* is a set of elements of the form  $(v:\pi)$  in which no evidence variable appears more than once. The  $\Vdash$  relation is extended to a three place relation  $P \Vdash e:\pi$ , indicating that it is possible to construct evidence  $e$  for the predicate  $\pi$  in any environment binding the variables in the predicate assignment  $P$  to appropriate evidence values. Thus predicates play a similar role for evidence expressions as types for simple  $\lambda$ -calculus terms.

- *Evidence abstraction*: A term  $M$  of type  $\pi \Rightarrow \rho$  is implemented by a term of the form  $\lambda v:\pi.M'$  where  $v$  is an evidence variable and  $M'$  is a term of type  $\rho$  corresponding to  $M$  using  $v$  in each place where evidence for  $\pi$  is needed.
- *Evidence application*: Each use of an overloaded expression  $N$  of type  $\pi \Rightarrow \rho$  is replaced by a term of the form  $N'e$  where  $N'$  is a term corresponding to  $N$  and  $e$  is an evidence expression for  $\pi$ .
- *Evidence reduction*: The standard rules of computation are augmented by a variant of  $\beta$ -reduction for evidence abstraction and application:

$$(\lambda v.M)e \triangleright_{\beta_e} [e/v]M.$$

Most of the typing rules given in Fig. 1 can be used with the extended system without modification. The only exceptions are the rules for dealing with qualified types; suitably modified versions of these are given in Fig. 2.

$  \begin{array}{c}  (\Rightarrow I) \quad \frac{P, v:\pi \mid A \vdash M : \sigma}{P \mid A \vdash \lambda v:\pi.M : \pi \Rightarrow \sigma} \\  (\Rightarrow E) \quad \frac{P \mid A \vdash M : \pi \Rightarrow \sigma \quad P \vdash e : \pi}{P \mid A \vdash Me : \sigma}  \end{array}  $
---

Fig. 2. Modified rules for qualified types with evidence.

Notice that extending the term language to make the use of evidence explicit gives unicity of type; each well-typed term has a uniquely determined type. This approach is very similar to the techniques used to make polymorphism explicit in the translation from implicit to explicit typed  $\lambda$ -calculus using abstraction and application over types [16]. As in that situation, there is a simple correspondence between derivations in the two systems, described by means of a function *Erase* mapping explicitly overloaded terms to their implicitly overloaded counterparts:

$$\begin{aligned}
 \text{Erase}(x) &= x \\
 \text{Erase}(MN) &= (\text{Erase}(M))(\text{Erase}(N)) \\
 &\vdots \\
 \text{Erase}(\lambda v:\pi.M) &= \text{Erase}(M) \\
 \text{Erase}(Me) &= \text{Erase}(M)
 \end{aligned}$$

The correspondence between the two systems can now be described by:

**Theorem 1.**  $P \mid A \vdash M : \sigma$  using the original typing rules if and only if  $P' \mid A \vdash M' : \sigma$  by a derivation of the same structure in the extended system such that  $P = \{ \pi \mid (v:\pi) \in P' \}$  and  $\text{Erase}(M') = M$ .

Given a term  $M$  in the original system, each corresponding term using explicit overloading is called a *translation* of  $M$  and can be used to give a semantics for the term. We write  $P' \mid A \vdash M \rightsquigarrow M' : \sigma$  to refer to the translation of a term in a specific context. Note that the translation of a given term may not be uniquely defined (with distinct translations corresponding to distinct derivations of  $P \mid A \vdash M : \sigma$ ). This is discussed in more detail in Section 9.1.

The form of evidence required will vary from one application to another. Suitable choices for each of the examples described in Section 3 are as follows:

- *Type classes*: The evidence for a type class predicate of the form  $C \tau$  is a *dictionary* containing the values of the members of  $C$  at the instance  $\tau$ . For example, in the simplest case, the evidence for a predicate  $Eq \tau$  might be an equality test function for values of type  $\tau$ .
- *Extensible records*: The evidence for a predicate of the form  $(r \text{ lacks } l)$  is the function:

$$(- \mid l = \_) :: \forall t. r \rightarrow t \rightarrow (r \mid l : t)$$

The evidence for a predicate of the form  $(r \text{ has } l : t)$  is the pair of functions:

$$\begin{aligned} (- \setminus l) &:: r \rightarrow r \setminus l \\ (-.l) &:: r \rightarrow t \end{aligned}$$

In practice, a concrete implementation of extensible records is likely to use offsets into a table of values used to store a record as evidence, passing these values to generic functions for updating or selecting from a record as necessary.

- *Subtypes*: The evidence for a predicate  $\sigma \subseteq \sigma'$  is a coercion function that maps values of type  $\sigma$  to values of type  $\sigma'$ .

Further details of the use of evidence in these applications is included in [13].

## 6. An extension of ML using qualified types

Polymorphic  $\lambda$ -calculus is not a suitable language to describe an implicitly typed language in which the need for explicit type annotations is replaced by the existence of a type inference algorithm. In practice, the benefits of type inference are often considered to outweigh the disadvantages of a less powerful type system. The ML type system [7,14] is a well-known example in which the price of type inference is the inability to define functions with polymorphic arguments. Nevertheless, it has proved to be very useful in practice and has subsequently been adopted by a number of other languages.



### 6.1. Basic definitions

Following the definition of types and type schemes in ML we consider a structured language of types, with the principal restriction being the inability to support functions with either polymorphic or overloaded arguments:

$$\begin{aligned}\tau &::= t \mid \tau \rightarrow \tau && \text{types} \\ \rho &::= P \Rightarrow \tau && \text{qualified types} \\ \sigma &::= \forall T. \rho && \text{type schemes}\end{aligned}$$

( $P$  and  $T$  range over finite sets of predicates and finite sets of type variables respectively).

It is convenient to introduce some abbreviations for qualified type and type scheme expressions. In particular, if  $\rho = (P \Rightarrow \tau)$  and  $\sigma = \forall T. \rho$ , then we write:

Abbreviation	Qualified type	Abbreviation	Type scheme
$\tau$	$\emptyset \Rightarrow \tau$	$\rho$	$\forall \emptyset. \rho$
$\pi \Rightarrow \rho$	$P, \pi \Rightarrow \tau$	$\forall t. \sigma$	$\forall (T \cup \{t\}). \rho$
$P' \Rightarrow \rho$	$P \cup P' \Rightarrow \tau$	$\forall T'. \sigma$	$\forall (T \cup T'). \rho$

In addition, if  $\{\alpha_i\}$  is an indexed set of variables, we write  $\forall \alpha_i. \rho$  as an abbreviation for  $\forall \{\alpha_i\}. \rho$ . As usual, type schemes are regarded as equal if they are equivalent up to renaming of bound variables.

Using this notation, any type scheme can be written in the form  $\forall \alpha_i. P \Rightarrow \tau$ , representing the set of qualified types  $\{[\tau_i/\alpha_i]P \Rightarrow [\tau_i/\alpha_i]\tau \mid \tau_i \in \text{Type}\}$  where  $[\tau_i/\alpha_i]$  is the substitution mapping each of the variables  $\alpha_i$  to the corresponding type  $\tau_i$  and  $\text{Type}$  is the set of all simple type expressions (represented by  $\tau$  in the grammar above).

As in [4,7,14], we use a term language based on simple untyped  $\lambda$ -calculus with the addition of a **let** construct to enable the definition and use of polymorphic (and in this case, overloaded) terms.

$$M ::= x \mid MN \mid \lambda x. M \mid \text{let } x = M \text{ in } N$$

A suitable set of typing rules for this system is given in Fig. 3. Note the use of the symbols  $\tau$ ,  $\rho$  and  $\sigma$  to restrict the application of certain rules to specific sets of type expressions.

### 6.2. Constrained type schemes

A typing judgement  $P \mid A \vdash M : \sigma$  assigns a type scheme  $\sigma$  to a term  $M$ , but also constrains uses of this typing to environments satisfying the predicates in  $P$ . This observation motivates the use of *constrained type schemes*, written as pairs of the form  $(P \mid \sigma)$  where  $P$  is a set of predicates and  $\sigma$  is a type scheme. Following the development of type inference in [7], we will define

<b>Standard rules:</b>	(var)	$\frac{(x:\sigma) \in A}{P \mid A \vdash x : \sigma}$
	( $\rightarrow E$ )	$\frac{P \mid A \vdash M : \tau' \rightarrow \tau \quad P \mid A \vdash N : \tau'}{P \mid A \vdash MN : \tau}$
	( $\rightarrow I$ )	$\frac{P \mid A_x, x:\tau' \vdash M : \tau}{P \mid A \vdash \lambda x.M : \tau' \rightarrow \tau}$
<b>Qualified types:</b>	( $\Rightarrow E$ )	$\frac{P \mid A \vdash M : \pi \Rightarrow \rho \quad P \Vdash \pi}{P \mid A \vdash M : \rho}$
	( $\Rightarrow I$ )	$\frac{P, \pi \mid A \vdash M : \rho}{P \mid A \vdash M : \pi \Rightarrow \rho}$
<b>Polymorphism:</b>	( $\forall E$ )	$\frac{P \mid A \vdash M : \forall t.\sigma}{P \mid A \vdash M : [\tau/t]\sigma}$
	( $\forall I$ )	$\frac{P \mid A \vdash M : \sigma \quad t \notin TV(A) \cup TV(P)}{P \mid A \vdash M : \forall t.\sigma}$
<b>Local Definition:</b>	(let)	$\frac{P \mid A \vdash M : \sigma \quad Q \mid A_x, x:\sigma \vdash N : \tau}{P \cup Q \mid A \vdash (\text{let } x = M \text{ in } N) : \tau}$

Fig. 3. ML-like typing rules for qualified types.

an ordering that can be used to describe when one constrained type scheme is more general than another. As a first step, we introduce the concept of generic instances:

**Definition 2.** A qualified type  $R \Rightarrow \mu$  is said to be a *generic instance* of the constrained type scheme  $(P \mid \forall \alpha_i. Q \Rightarrow \tau)$  if there are types  $\tau_i$  such that  $R \Vdash P \cup [\tau_i/\alpha_i]Q$  and  $\mu = [\tau_i/\alpha_i]\tau$ .

The principal motivation for the definition of the ordering ( $\leq$ ) between type schemes is that a statement of the form  $\sigma' \leq \sigma$  should mean that it is possible to use an object of type  $\sigma$  wherever an object of type  $\sigma'$  is required.

**Definition 3.** The constrained type scheme  $(Q \mid \eta)$  is said to be *more general* than a constrained type scheme  $(P \mid \sigma)$ , written  $(P \mid \sigma) \leq (Q \mid \eta)$ , if every generic instance of  $(P \mid \sigma)$  is also a generic instance of  $(Q \mid \eta)$ .

It is straightforward to show that this defines a preorder on the set of constrained type schemes, such that a qualified type  $\rho$  is a generic instance

of the type scheme  $\sigma$  if and only if  $\rho \leq \sigma$ . We will write  $(P \mid \sigma) \simeq (Q \mid \eta)$  to indicate when two constrained type schemes are equivalent with respect to  $(\leq)$ , i.e. when each is more general than the other. The following properties are easily established:

- $\sigma \simeq (\emptyset \mid \sigma)$  for any type scheme  $\sigma$ .
- If  $\rho$  is a qualified type and  $P$  is a set of predicates, then  $(P \mid \rho) \simeq P \Rightarrow \rho$ .
- If  $\sigma$  is a type scheme and  $P$  is a set of predicates, then  $(P \mid \sigma) \leq \sigma$ .
- If  $\sigma' \leq \sigma$  and  $P' \vdash P$ , then  $(P' \mid \sigma') \leq (P \mid \sigma)$ .
- If none of the variables  $\alpha_i$  appear in  $P$ , then the constrained type scheme  $(P \mid \forall \alpha_i. \rho)$  is equivalent to the type scheme  $\forall \alpha_i. P \Rightarrow \rho$ . Thus every constrained type scheme can be represented by a simple type scheme using a renaming of bound variables.

The application of a substitution  $S$  to a constrained type scheme  $(P \mid \sigma)$  is defined by  $S(P \mid \sigma) = (SP \mid S\sigma)$ . The next proposition describes an important property of the ordering on constrained type schemes.

**Proposition 4.** *For any substitution  $S$  and constrained type schemes  $(P \mid \sigma)$  and  $(Q \mid \eta)$ :*

$$(P \mid \sigma) \leq (Q \mid \eta) \Rightarrow S(P \mid \sigma) \leq S(Q \mid \eta).$$

### 6.3. Ordering of type assignments

The definition of constrained type schemes and the ordering  $(\leq)$  extends naturally to an ordering on (constrained) type assignments.

**Definition 5.** If  $A$  and  $A'$  are type assignments and  $P$  and  $P'$  are sets of predicates, then we say that  $(P \mid A)$  is *more general* than  $(P' \mid A')$ , written  $(P' \mid A') \leq (P \mid A)$ , if  $\text{dom } A = \text{dom } A'$  and  $(P' \mid A'(x)) \leq (P \mid A(x))$  for each  $x \in \text{dom } A$ .

The results of the previous section can be used to prove that this ordering on type assignments is reflexive, transitive and preserved by substitutions. In this paper, we will only use the special case where  $P = \emptyset$  in which case we write  $(P' \mid A') \leq A$ . This can be interpreted as indicating that each of the types assigned to a variable in  $A$  is more general than the type assigned in  $A'$  in any environment that satisfies the predicates in  $P'$ .

### 6.4. Generalization

Given a derivation  $P \mid A \vdash M : \tau$ , it is useful to have a notation for the most general type scheme that can be obtained for  $M$  from this derivation using the rules  $(\Rightarrow I)$  and  $(\forall I)$  given in Fig. 3.

**Definition 6.** The *generalization* of a qualified type  $\rho$  with respect to a type assignment  $A$  is written  $Gen(A, \rho)$  and defined by:

$$Gen(A, \rho) = \forall (TV(\rho) \setminus TV(A)).\rho.$$

In other words, if  $\{\alpha_i\} = TV(\rho) \setminus TV(A)$ , then  $Gen(A, \rho) = \forall \alpha_i.\rho$ . The following propositions describe the interaction of generalization with predicate entailment and substitution.

**Proposition 7.** Suppose that  $A$  is a type assignment,  $P$  and  $P'$  are sets of predicates and  $\tau$  is a type. Then  $Gen(A, P' \Rightarrow \tau) \leq Gen(A, P \Rightarrow \tau)$  whenever  $P' \vdash P$ .

**Proposition 8.** If  $A$  is a type assignment,  $\rho$  is a qualified type and  $S$  is a substitution, then:

$$Gen(SA, S\rho) \leq S(Gen(A, \rho)).$$

Furthermore, there is a substitution  $R$  such that:

$$RA = SA, \quad SGen(A, \rho) = Gen(RA, R\rho).$$

## 7. A syntax-directed approach

The typing rules in Fig. 3 provide clear descriptions of the treatment of each of the syntactic constructs of the term and type languages. Unfortunately, they are not suitable for use in a type inference algorithm where it should be possible to determine an appropriate order in which to apply the typing rules by a simple analysis of the syntactic structure of the term whose type is required.

In this section, we introduce an alternative set of typing rules with a single rule for each syntactic construct in the term language. We refer to this as the *syntax-directed* system because it has the following important property: *all typing derivations for a given term  $M$  (if there are any) have the same structure, uniquely determined by the syntactic structure of  $M$ .* We regard the syntax-directed system as a tool for exploring the type system of Section 6 and we establish a congruence between the two systems so that results about one can be translated into results about the other. The advantages of working with the syntax-directed system are:

- The rules are better suited to use in a type inference algorithm; having found types for each of the subterms of a given term  $M$ , there is at most one rule that can be used to obtain a type for the term  $M$  itself.
- Only type expressions are involved in the matching process. Type schemes and qualified types can only appear in type assignments.

- There are fewer rules and hence fewer cases to be considered in formal proofs.

A similar approach is described in [4] which gives a deterministic set of typing rules for ML and outlines their equivalence to the rules in [7].

### 7.1. Syntax-directed typing rules

The typing rules for the syntax-directed system are given in Fig. 4. Typings in this system are written in form  $P \mid A \vdash^s M : \tau$ , where  $\tau$  ranges over the set of type expressions rather than the set of type schemes as in the typing judgements of Section 6. Other than this, the principal differences between the two systems are in the rules  $(var)^s$  and  $(let)^s$  which use the operations of instantiation and generalization introduced in Sections 6.2 and 6.4.

$(var)^s$	$\frac{(x:\sigma) \in A}{P \mid A \vdash^s x : \tau} \quad (P \Rightarrow \tau) \leq \sigma$
$(\rightarrow E)^s$	$\frac{P \mid A \vdash^s M : \tau' \rightarrow \tau \quad P \mid A \vdash^s N : \tau'}{P \mid A \vdash^s MN : \tau}$
$(\rightarrow I)^s$	$\frac{P \mid A_x, x:\tau' \vdash^s M : \tau}{P \mid A \vdash^s \lambda x. M : \tau' \rightarrow \tau}$
$(let)^s$	$\frac{P \mid A \vdash^s M : \tau \quad P' \mid A_x, x:\sigma \vdash^s N : \tau'}{P' \mid A \vdash^s (\text{let } x = M \text{ in } N) : \tau'} \quad \sigma = \text{Gen}(A, P \Rightarrow \tau)$

Fig. 4. Syntax-directed inference system.

### 7.2. Properties of the syntax-directed system

The following proposition illustrates the parametric polymorphism present in the syntax-directed system; instantiating the free type variables in a derivable typing with arbitrary types produces another derivable typing.

**Proposition 9.** *If  $P \mid A \vdash^s M : \tau$  and  $S$  is a substitution, then  $SP \mid SA \vdash^s M : S\tau$ .*

A similar result is established in [6] where it is shown that for any derivation  $A \vdash M : \tau$  in the usual (non-deterministic) ML type system and any substitution  $S$ , there is a derivation  $SA \vdash M : S\tau$  which can be chosen in such a way that the height of the latter is bounded by the height of the former. This additional condition is needed to ensure the validity of proofs by induction on the size of a derivation. This complication is avoided by the syntax-directed

system; the derivations in Proposition 9 are guaranteed to have the same structure because the term  $M$  is common to both.

There is also a form of polymorphism over the sets of environments in which a particular typing can be used, as described by the following proposition:

**Proposition 10.** *If  $P \mid A \models M : \tau$  and  $Q \vdash P$ , then  $Q \mid A \models M : \tau$ .*

Recall that an ordering  $\sigma' \leq \sigma$  is intended to mean that, at least for the purposes of type inference, it is possible to use an object of type  $\sigma$  whenever with an object of type  $\sigma'$  is required. In much the same way, given two type assignments such that  $A' \leq A$  (so that the type assigned to each variable in  $A$  is more general than the corresponding type in  $A'$ ), we would expect that any typing that can be derived using  $A'$  could also be derived from  $A$ . The following proposition establishes a slightly more general form of this result:

**Proposition 11.** *If  $P \mid A' \models M : \tau$  and  $(P \mid A') \leq A$ , then  $P \mid A \models M : \tau$ .*

The hypothesis  $(P \mid A') \leq A$  means that the types assigned to variables in  $A$  are more general than those given by  $A'$  in any environment that satisfies the predicates in  $P$ . For example:

$$\begin{aligned} & (EqInt \mid \{(=): Int \rightarrow Int \rightarrow Bool\}) \\ & \leq \{(=): \forall a. Eq a \Rightarrow a \rightarrow a \rightarrow Bool\} \end{aligned}$$

and hence, by the proposition above, it is possible to replace an integer equality function with a generic equality function of type  $\forall a. Eq a \Rightarrow a \rightarrow a \rightarrow Bool$  in any environment that satisfies  $EqInt$ .

### 7.3. Relationship with original type system

In order to use the syntax-directed system as a tool for reasoning about the type system described in Section 6, we need to investigate the way in which the existence of a derivation in one system determines the existence of derivations in the other.

Our first result establishes the soundness of the syntax-directed system with respect to the original typing rules, showing that any derivable typing in the former system is also derivable in the latter.

**Theorem 12.** *If  $P \mid A \models M : \tau$ , then  $P \mid A \vdash M : \tau$ .*

The translation of derivations in the original type system to those of the syntax-directed system is less obvious. For example, if  $P \mid A \vdash M : \sigma$ , then it will not in general be possible to derive the same typing in the syntax-directed system because  $\sigma$  is a type scheme, not a simple type. However, for any

derivation  $P' \mid A \vdash^s M : \tau$ , Theorem 12 guarantees the existence of a derivation  $P' \mid A \vdash M : \tau$  and hence  $\emptyset \mid A \vdash M : \text{Gen}(A, P' \Rightarrow \tau')$  by Definition 6. The following theorem shows that it is always possible to find a derivation in this way such that the inferred type scheme  $\text{Gen}(A, P' \Rightarrow \tau')$  is more general than the constrained type scheme  $(P \mid \sigma)$  determined by the original derivation.

**Theorem 13.** *If  $P \mid A \vdash M : \sigma$ , then  $P' \mid A \vdash^s M : \tau$  for some set of predicates  $P'$  and type  $\tau$  such that  $(P \mid \sigma) \leq \text{Gen}(A, P' \Rightarrow \tau)$ .*

## 8. Type inference

In this section, we give an algorithm for calculating a typing for a given term, using an extension of Milner's algorithm W [14] to support qualified types. We show that the typings produced by this algorithm are derivable in the syntax-directed system and that they are, in a certain sense, the most general typings possible. Combining this with the results of the previous section, the algorithm can be used to reason about the type system in Section 6.

### 8.1. Unification

This section describes the unification algorithm which is a central component of the type inference algorithm. A substitution  $S$  is called a *unifier* for the type expressions  $\tau$  and  $\tau'$  if  $S\tau = S\tau'$ . The following theorem is due to Robinson [20].

**Theorem 14** (Unification algorithm). *There is an algorithm whose input is a pair of type expressions  $\tau$  and  $\tau'$  such that:*

- *Either the algorithm succeeds with a substitution  $U$  as its result and the unifiers of  $\tau$  and  $\tau'$  are precisely those substitutions of the form  $RU$  for any substitution  $R$ . The substitution  $U$  is called a most general unifier for  $\tau$  and  $\tau'$ , and is denoted  $\text{mgu}(\tau, \tau')$ .*
- *Or the algorithm fails and there are no unifiers for  $\tau$  and  $\tau'$ .*

In the following, we write  $\tau \stackrel{U}{\sim} \tau'$  for the assertion that the unification algorithm succeeds by finding a most general unifier  $U$  for  $\tau$  and  $\tau'$ .

### 8.2. A type inference algorithm

Following the presentation of [19], we describe the type inference algorithm using the inference rules in Fig. 5. These rules use typings of the form

$$P \mid TA \vdash^w M : \tau$$

$(var)^w$	$\frac{(x:\forall\alpha_i. P \Rightarrow \tau) \in A}{[\beta_i/\alpha_i]P \mid A \Vdash^w x : [\beta_i/\alpha_i]\tau}$	$\beta_i \text{ new}$
$(\rightarrow E)^w$	$\frac{P \mid TA \Vdash^w M : \tau \quad Q \mid T'TA \Vdash^w N : \tau' \quad T'\tau \stackrel{U}{\sim} \tau' \rightarrow \alpha}{U(T'P \cup Q) \mid UT'TA \Vdash^w MN : U\alpha}$	$\alpha \text{ new}$
$(\rightarrow I)^w$	$\frac{P \mid T(A_x, x:\alpha) \Vdash^w M : \tau}{P \mid TA \Vdash^w \lambda x. M : T\alpha \rightarrow \tau}$	$\alpha \text{ new}$
$(let)^w$	$\frac{P \mid TA \Vdash^w M : \tau \quad P' \mid T'(TA_x, x:\sigma) \Vdash^w N : \tau'}{P' \mid T'TA \Vdash^w (\text{let } x = M \text{ in } N) : \tau'}$	
$\sigma = Gen(TA, P \Rightarrow \tau)$		

Fig. 5. Type inference algorithm W.

where  $P$  is a set of predicates,  $T$  is a substitution,  $A$  is a type assignment,  $M$  is a term and  $\tau$  is a simple type expression. The typing rules can be interpreted as an attribute grammar in which  $A$  are  $M$  inherited attributes, while  $P$ ,  $T$  and  $\tau$  are synthesized. One of the advantages of this choice of notation is that it highlights the relationship between W and the syntax-directed system. This point is illustrated by the following theorem.

**Theorem 15.** *If  $P \mid TA \Vdash^w M : \tau$ , then  $P \mid TA \Vdash^s M : \tau$ .*

Combining this with the result of Theorem 12 gives the following important corollary.

**Corollary 16** (Soundness of W). *If  $P \mid TA \Vdash^w M : \tau$ , then  $P \mid TA \vdash M : \tau$ .*

With the exception of  $(let)^w$ , each of the rules in Fig. 5 introduces “new” variables; i.e. variables that do not appear in the hypotheses of the rule nor in any other distinct branches of the complete derivation. Note that it is always possible to choose type variables in this way because the set of type variables is assumed to be countably infinite. In the presence of new variables, it is convenient to work with a weaker form of equality on substitutions, writing  $S \approx R$  to indicate that  $St = Rt$  for all but a finite number of new variables  $t$ . In most cases, we can treat  $S \approx R$  as  $S = R$ , since the only differences between the substitutions occur at variables that are not used elsewhere in the algorithm.

This notation enables us to give an accurate statement of the following result which shows that the typings obtained by W are, in a precise sense, the most general derivable typings for a given term.



**Theorem 17.** *If  $P \mid SA \vdash^S M : \tau$ , then  $Q \mid TA \vdash^W M : \nu$  for some  $Q, T$  and  $\nu$  and there is a substitution  $R$  such that  $S \approx RT$ ,  $\tau = R\nu$  and  $P \vdash RQ$ .*

Combining the result of Theorem 17 with that of Theorem 13 we obtain a similar completeness result for  $W$  with respect to the type system of Section 6.

**Corollary 18** (Completeness of  $W$ ). *If  $P \mid SA \vdash M : \sigma$ , then  $Q \mid TA \vdash^W M : \nu$  for some  $Q, T$  and  $\nu$  and there is a substitution  $R$  such that  $(P \mid \sigma) \leq RGen(TA, Q \Rightarrow \nu)$  and  $S \approx RT$ .*

### 8.3. Principal type schemes

A term  $M$  is *well-typed* under a type assignment  $A$  if  $P \mid A \vdash M : \sigma$  for some  $P$  and  $\sigma$ . It is natural to try to characterize the set of constrained type schemes  $(P \mid \sigma)$  for which such a derivation can be found. This can be described using the concept of a *principal type scheme*:

**Definition 19.** A *principal type scheme* for a term  $M$  under a type assignment  $A$  is a constrained type scheme  $(P \mid \sigma)$  such that  $P \mid A \vdash M : \sigma$ , and  $(P' \mid \sigma') \leq (P \mid \sigma)$  whenever  $P' \mid A \vdash M : \sigma'$ .

The following result gives a sufficient condition for the existence of principal type schemes, by showing how they can be constructed from typings produced by  $W$ .

**Corollary 20.** *Suppose that  $M$  is a term,  $A$  is a type assignment and*

$$Q \mid TA \vdash^W M : \nu \text{ for some } Q, T \text{ and } \nu.$$

*Then  $Gen(TA, Q \Rightarrow \nu)$  is a principal type scheme for  $M$  under  $TA$ .*

Combining this with Corollary 18 gives a necessary condition for the existence of principal type schemes: a term is well-typed if and only if it has a principal type scheme. Furthermore, if it exists, a suitable principal type can be calculated using the type inference algorithm  $W$ .

**Corollary 21** (Principal Type Theorem). *Let  $M$  be a term and  $A$  an arbitrary type assignment. The following conditions are equivalent:*

- $M$  is well-typed under  $A$ .
- $Q \mid TA \vdash^W M : \nu$  for some  $Q$  and  $\nu$  and there is a substitution  $R$  such that  $RTA = A$ .
- $M$  has a principal typing under  $A$ .

## 9. Extensions and topics for further work

### 9.1. The coherence problem

It is important to point out that the type systems described by the rules in the previous sections are not *coherent* (in the sense of [1]). In other words, it is possible to construct translations  $P \mid A \vdash M \rightsquigarrow M'_1 : \sigma$  and  $P \mid A \vdash M \rightsquigarrow M'_2 : \sigma$  in which the terms  $M'_1$  and  $M'_2$  are not equivalent, and hence the semantics of  $M$  are not well-defined.

For an example in which the coherence problem arises, consider the term  $out\ (in\ x)$  under the predicate assignment  $P = \{u : C\ Int, v : C\ Bool\}$  and the type assignment:

$$A = \{x : Int, in : \forall a. C\ a \Rightarrow Int \rightarrow a, out : \forall a. C\ a \Rightarrow a \rightarrow Int\}$$

for some unary predicate symbol  $C$ . Instantiating the quantified type variable in the type of  $in$  (and hence also in that of  $out$ ) with the types  $Int$  and  $Bool$  leads to distinct derivations  $P \mid A \vdash out\ (in\ x) : Int$  in which the corresponding translations,  $out\ u\ (in\ u\ x)$  and  $out\ v\ (in\ v\ x)$  are clearly not equal.

Note that the principal type scheme of  $out\ (in\ x)$  in this example is  $\forall a. C\ a \Rightarrow Int$  and that the type variable  $a$  (the source of the lack of coherence in the derivations above) appears only in the predicate qualifying the type of the term, not in the type itself. Motivated by the functional programming language Haskell [10], we say that a type of the form  $\forall \alpha_i. P \Rightarrow \tau$  is *unambiguous* if  $\{\alpha_i\} \cap TV(P) \subseteq TV(\tau)$ . Extending the results of this paper to describe the use of translations in the syntax-directed system and the type inference algorithm, we have established the following coherence result:

**Theorem 22.** *If  $P \mid A \vdash M \rightsquigarrow M'_1 : \sigma$  and  $P \mid A \vdash M \rightsquigarrow M'_2 : \sigma$  and the principal type scheme of  $M$  in  $A$  is unambiguous, then the translations  $M'_1$  and  $M'_2$  are equivalent.*

This generalizes an earlier result by Blott [2] for the special case of [23]. Full details are included in [13] and we expect to describe this work more fully in a forthcoming paper.

### 9.2. Eliminating evidence parameters

Using translations as described in Section 5, a term  $M$  of type  $\forall \alpha_i. P \Rightarrow \tau$  will be implemented by a term of the form  $\lambda v_1. \dots \lambda v_n. M'$ , where  $P = \{\pi_1, \dots, \pi_n\}$  and each  $v_i$  is an evidence variable for the corresponding predicate  $\pi_i$ . The following subsections outline a number of situations in which it is useful to reduce or even eliminate the use of evidence parameters, either to obtain a more efficient implementation or to avoid unnecessary repeated calculations.

### 9.2.1. Simplification

The translation of a term whose type is qualified by a set of predicates  $P$  requires one evidence abstraction for each element of  $P$ . Thus the number of evidence parameters that are required can be reduced by finding a smaller set of predicates  $Q$ , equivalent to  $P$  in the sense that  $P \vdash Q$  and  $Q \vdash P$  (and hence the type of the new term is equivalent to that of the original term). In this situation, we have a compromise between reducing the number of evidence parameters required and the cost of constructing evidence for  $P$  from evidence for  $Q$ . The process of simplification can be formalized by allowing the rule:

$$\frac{P \mid A \Vdash^w M : \nu \quad P \vdash Q \quad Q \vdash P}{Q \mid A \Vdash^w M : \nu}$$

to be used at any stage during type inference to simplify the inferred predicate set. It is relatively straightforward to show that this rule is sound and that the extended algorithm still calculates principal (but potentially simplified) type schemes.

In general, the task of finding an optimal set of predicates with which to replace  $P$  is likely to be intractable. One potentially useful approach would be to determine a minimal subset  $Q \subseteq P$  such that  $Q \vdash P$ . To see that this is likely to be a good choice, note that:

- $P \vdash Q$ , by monotonicity of  $\vdash$  and hence  $Q$  is equivalent to  $P$  as required.
- Since  $Q \subseteq P$ , the number of evidence abstractions required using  $Q$  is less than or equal to the number required when using  $P$ .
- The construction of evidence for a predicate in  $P$  using evidence for  $Q$  is trivial for each predicate that is already in  $Q$ .

### 9.2.2. Evidence parameters considered harmful

The principal motivation for including the **let** construct in the term language was to enable the definition and use of polymorphic and overloaded values. In practice, the same construct is also used for a number of other purposes:

- To avoid repeated evaluation of a value that is used at a number of points in an expression.
- To create *cyclic data structures* using recursive bindings [3].
- To enable the use of identifiers as abbreviations for the subexpressions of a large expression.

Unfortunately, the addition of evidence parameters to the value defined in a **let** expression may mean that the evaluation of an overloaded term will not behave as intended. For example, if  $f : \forall a. C a \Rightarrow Int \rightarrow a$ , then we have a translation:

$$\text{let } x = f \ 0 \text{ in } (x, x) \rightsquigarrow \lambda v. \text{let } x = (\lambda v. f \ v \ 0) \text{ in } (x \ v, x \ v)$$

and the evaluation of  $xv$  in the translation is no longer shared. There are a number of potential solutions to this problem. In the example above, one method would be to rewrite the translation as:

$$\lambda v. \text{let } x = (f v 0) \text{ in } (x, x).$$

This is the kind of translation which will be obtained using a *monomorphism restriction* such as that proposed for Haskell [10] which restricts the amount of overloading that can be used in particular syntactic forms of binding. Note that this approach is only useful when the variable defined in the **let** expression is not required to have a polymorphic type in the scope of that definition.

### 9.2.3. Constant and locally-constant overloading

Consider the typing of local definitions in the type system of Section 6 using the rule:

$$\frac{P \mid A \vdash M : \sigma \quad Q \mid A_x, x : \sigma \vdash N : \tau}{P \cup Q \mid A \vdash (\text{let } x = M \text{ in } N) : \tau}$$

Notice that this allows some of the predicates constraining the typing of  $M$  (i.e. those in  $P$ ) to be retained as a constraint on the environment in the conclusion of the rule rather than being included in the type scheme  $\sigma$ . However, in the corresponding rule  $(\text{let})^s$  for the syntax-directed system, all of the predicates constraining the typing of  $M$  are included in the inferred type scheme  $\text{Gen}(A, P \Rightarrow \tau)$ :

$$\frac{P \mid A \vdash^s M : \tau \quad P' \mid A_x, x : \text{Gen}(A, P \Rightarrow \tau) \vdash^s N : \tau'}{P' \mid A \vdash^s (\text{let } x = M \text{ in } N) : \tau'}$$

As a consequence, evidence parameters are needed for all of the predicates in  $P$ , even if some of the corresponding evidence values are the same for each occurrence of  $x$  in  $N$ . In particular, this includes *constant* evidence (for predicates with no free type variables) and *locally-constant* evidence (for predicates, each of whose free variables also appears free in  $A$ ).

From the relationship between the type inference algorithm  $W$  and the syntax-directed system, it follows that  $W$  has the same behaviour; indeed, this is essential to ensure that  $W$  calculates principal types: if  $x \notin FV(N)$ , then none of the environment constraints described by  $P$  need be reflected by the constraints on the complete expression in  $P'$ .

However, if  $x \in FV(N)$ , it is possible to find a set  $F \subseteq P$  such that  $P' \Vdash F$  and hence the type scheme assigned to  $x$  can be replaced by  $\text{Gen}(A, (P \setminus F) \Rightarrow \tau)$ , potentially decreasing the number of evidence parameters required by  $x$ . To see this, suppose that  $\text{Gen}(A, P \Rightarrow \tau) = (\forall \alpha_i. P \Rightarrow \tau)$ . A straightforward induction, based on the hypothesis that  $x \in FV(N)$ , shows that  $P' \Vdash [\tau_i / \alpha_i] P$  for some types  $\tau_i$ . If we now define:

$$FP(A, P) = \{ (v : \pi) \in P \mid TV(\pi) \subseteq TV(A) \},$$

then  $F = FP(A, P)$  is the largest subset of  $P$  that is guaranteed to be unchanged by the substitution  $[\tau_i/\alpha_i]$ . These observations suggest that  $(let)^s$  could be replaced by the two rules:

- In the case where  $x \notin FV(N)$ :

$$\frac{P \mid A \vdash^s M : \tau \quad P' \mid A \vdash^s N : \tau'}{P' \mid A \vdash^s (\text{let } x = M \text{ in } N) : \tau'} (let)_f^s$$

The typing judgement involving  $M$  serves only to preserve to property that all subterms of a well-typed term are also well-typed.

- In the case where  $x \in FV(N)$ :

$$\frac{P \mid A \vdash^s M : \tau \quad P' \mid A_x, x : Gen(A, P \setminus F \Rightarrow \tau) \vdash^s N : \tau' \quad P' \vdash F}{P' \mid A \vdash^s (\text{let } x = M \text{ in } N) : \tau'} (let)_b^s$$

where  $F = FP(A, P)$ .

While these rules retain the syntax-directed character necessary for use in a type inference algorithm, they are not suitable for typing top-level definitions (such as those in Haskell or ML) which are treated as **let** expressions in which the scope of the defined variable is not fully determined at compile-time.

A more realistic approach would be to use just  $(let)_b^s$  in place of  $(let)^s$ , with the understanding that type schemes inferred by  $W$  are only guaranteed to be principal in the case where  $x \in FV(N)$  for all subterms of the form **let**  $x = M$  **in**  $N$  in the term whose type is being inferred. Justification for this approach is as follows:

- For a top-level declaration of the identifier  $x$ , we can take the scope of the declaration to be the set of all terms that might reasonably be evaluated in the scope of such a declaration, which of course includes the term  $x$ .
- For **let** expressions in which the scope of the defined variable is known, the local definition in an expression of the form **let**  $x = M$  **in**  $N$  is redundant, and the expression is semantically equivalent to  $N$ . However, expressions of this form are sometimes used in implicitly typed languages to force a less general type than might otherwise be obtained by the type inference mechanism. For example, if  $(=)$  is an integer equality function and 0 is an integer constant, then  $\lambda x. \text{let } y = (x = 0) \text{ in } x$  has principal type scheme  $Int \rightarrow Int$ , whereas the principal type scheme for  $\lambda x. x$  is  $\forall a. a \rightarrow a$ . Such ad-hoc ‘coding-tricks’ become unnecessary if the term language is extended to allow explicit type declarations.

In a practical implementation, it would be useful to arrange for suitable diagnostic messages to be generated whenever an expression of the form **let**  $x = M$  **in**  $N$  with  $x \notin FV(N)$  is encountered; this would serve as a warning to

the programmer that the principal type property may be lost (in addition to catching other potential program errors).

### 9.3. Satisfiability

One of the most important features of the systems of qualified types described in this paper is the ability to move “global” constraints on a typing derivation into the type of an object using  $(\Rightarrow I)$ :

$$\frac{P, \pi \mid A \vdash M : \rho}{P \mid A \vdash M : \pi \Rightarrow \rho}$$

This is essential in many situations where overloading is combined with polymorphism: without the ability to move predicates from the first component of a typing  $P \mid A \vdash M : \rho$  into the type of an object we would not be able to apply  $(\forall I)$  for any type variables appearing in  $TV(P)$ , severely limiting the use of polymorphism.

On the other hand, with the formulation of the typing rules used in the previous sections there is no attempt to guarantee that the predicates introduced into the type of an object using  $(\Rightarrow I)$  are satisfiable. As we have already mentioned, an object of type  $\pi \Rightarrow \rho$  can only be used if we can provide evidence for the predicate  $\pi$ . If no such evidence can be obtained, then any object with this type is useless.

This problem was noted by Volpano and Smith [22] for the special case of the system of type classes described in [23]. With this in mind, they gave a stronger definition of well-typing that includes testing for satisfiability of an inferred type scheme and showed that, in an unrestricted version of the Wadler–Blott system, the process of determining whether a particular term is well-typed is undecidable. The framework used in this paper allows us to separate typability from predicate entailment and to identify the problem as undecidability of the latter. Nevertheless, the difficulty remains.

On the one hand we could simply ignore the problem since it will never be possible to resolve the overloading for an object with an unsatisfiable type scheme and hence any attempt to use it will fail. On the other hand, it would certainly be useful if the type system could be used to identify such objects at the point where they are defined and produce suitable error diagnostics to assist the programmer. One possibility would be to modify the rule for typing **let** expressions with:

$$\frac{P \mid A \vdash M : \sigma \quad Q \mid A_x, x : \sigma \vdash N : \tau \quad P_0 \text{ sat } \sigma}{P, Q \mid A \vdash (\text{let } x = M \text{ in } N) : \tau}$$

to ensure satisfiability with respect to a fixed set of predicates  $P_0$ , where:

$$P_0 \text{ sat } (\forall \alpha_i. P \Rightarrow \tau) \Leftrightarrow \exists \nu_i. P_0 \vdash [\nu_i / \alpha_i] P.$$

The following properties of this relationship between predicate sets and type schemes are easily established and show that this notion of *satisfiability* is well-behaved with respect to our use of polymorphism, entailment and ordering:

- If  $P \text{ sat } \sigma$ , then  $SP \text{ sat } S\sigma$  for any substitution  $S$ .
- If  $P \text{ sat } \sigma$  and  $Q \vdash P$ , then  $Q \text{ sat } \sigma$ .
- If  $P \text{ sat } \sigma'$  and  $(P \mid \sigma') \leq \sigma$ , then  $P \text{ sat } \sigma$ .

We conjecture that, if we restrict our attention to derivations  $P \mid A \vdash M : \sigma$  for which  $P_0 \vdash P$ , then the development of a principal type algorithm and coherence conditions described in the previous sections will extend naturally to deal with this extension. Note however that we will require decidability of  $P_0 \text{ sat } \sigma$  for arbitrary  $P_0$  and  $\sigma$  to ensure decidability of type checking.

Another, more positive, application of satisfiability is to enable the calculation of more accurate types for given terms. As an example, consider the function  $\lambda r. (r.l, r.l)$  using the record selection operator described in Section 3.2 which has principal type scheme:

$$\forall r. \forall a. \forall b. (r \text{ has } l : a, r \text{ has } l : b) \Rightarrow r \rightarrow (a, b).$$

However, for any given record type  $r$ , the types assigned to the variables  $a$  and  $b$  must be identical since they both correspond to the same field in  $r$ . It would therefore seem quite reasonable to treat  $f$  as having a *principal satisfiable type scheme*:

$$\forall r. \forall a. (r \text{ has } l : a) \Rightarrow r \rightarrow (a, a).$$

To see how this might be dealt with more formally, recall the treatment of the ordering between type schemes in Section 6.2. Writing the set of generic instances of a type scheme as:

$$\llbracket \forall \alpha_i. P \Rightarrow \tau \rrbracket = \{ Q \Rightarrow [\nu_i / \alpha_i] \tau \mid \nu_i \in \text{Type}, Q \vdash [\nu_i / \alpha_i] P \},$$

the ordering on type schemes is described by:

$$\sigma \leq \sigma' \Leftrightarrow \llbracket \sigma \rrbracket \subseteq \llbracket \sigma' \rrbracket$$

In a similar way can define the *generic satisfiable instances* of a type scheme with respect to a predicate set  $P_0$  as:

$$\llbracket \forall \alpha_i. P \Rightarrow \tau \rrbracket_{P_0}^{\text{sat}} = \{ [\nu_i / \alpha_i] \tau \mid \nu_i \in \text{Type}, P_0 \vdash [\nu_i / \alpha_i] P \}$$

and define a satisfiability ordering, again with respect to  $P_0$ , by:

$$\sigma \leq_{P_0}^{\text{sat}} \sigma' \Leftrightarrow \llbracket \sigma \rrbracket_{P_0}^{\text{sat}} \subseteq \llbracket \sigma' \rrbracket_{P_0}^{\text{sat}}$$

We can formalize the notion of principal satisfiable type in the same way as in Section 8.3 using the  $(\leq_{P_0}^{\text{sat}})$  ordering in place of  $(\leq)$ . For the example above,

both of the type schemes given are principal satisfiable type schemes for the term  $\lambda r.(r.l, r.l)$ . The first of these is the type scheme that would be obtained using our type inference algorithm, but it would clearly be preferable if the algorithm could be modified to give the second alternative. Further investigation is needed to discover effective procedures or heuristics for calculating more informative types that can be used to support this extension.

#### 9.4. The use of subsumption

The typing rules in Fig. 1 are only suitable for reasoning about systems with explicit coercions. For example, if  $Int \subseteq Real$ , then we can use an addition function:

$$add :: \forall a. a \subseteq Real \Rightarrow a \rightarrow a \rightarrow Real$$

to add two integers together, obtaining a real number as the result. More sophisticated systems, such as those in [8,15], cannot be described without adding a form of the rule of subsumption:

$$\frac{P \mid A \vdash M : \tau' \quad P \Vdash \tau' \subseteq \tau}{P \mid A \vdash M : \tau}$$

Each use of this rule corresponds to an implicit coercion; the addition of two integers to obtain a real result can be described without explicit overloading using a function:

$$add :: Real \rightarrow Real \rightarrow Real$$

with two implicit coercions from  $Int$  to  $Real$ . As a further example, in the framework of Section 4, the polymorphic identity function  $\lambda t.\lambda x:t.x$  can be treated as having type  $\forall a.\forall b.a \subseteq b \Rightarrow a \rightarrow b$  and hence acts as a generic coercion function.

No attempt has been made to deal with systems including the rule of subsumption in the development of the type inference algorithm in Section 8, which is therefore only suitable for languages using explicit coercions. The results of [8] and [21] are likely to be particularly useful in extending the present system to support the use of implicit coercions.

#### Acknowledgments

I am very grateful for conversations with members of the Department of Computing Science, Glasgow, whose suggestions prompted me to investigate the ideas presented in this paper, and in particular to Phil Wadler for his comments on several earlier versions of this paper. I would also like to thank



the referees for ESOP'92 for their suggestions and advice concerning my original submission to the symposium.

This work was carried out while the author was a member of the Programming Research Group, Oxford University Computing Laboratory, UK, with financial support from the Science and Engineering Research Council of Great Britain.

## References

- [1] V. Breazu-Tannen, T. Coquand, C.A. Gunter and A. Scedrov, Inheritance and coercion, in: *Proceedings Fourth Annual Symposium on Logic in Computer Science* (1989).
- [2] S. Blott, An approach to overloading with polymorphism, Ph.D. Thesis, Department of Computing Science, University of Glasgow, Glasgow, Scotland (1990).
- [3] R.S. Bird and P. Wadler, *Introduction to Functional Programming* (Prentice Hall, 1989).
- [4] D. Clément, J. Despeyroux, T. Despeyroux and G. Kahn, A simple applicative language: Mini-ML, in: *Proceedings ACM Symposium on LISP and Functional Programming* (1986).
- [5] L. Cardelli and J.C. Mitchell, Operations on records, in: *Fifth International Conference on Mathematical Foundations of Programming Language Semantics*, Lecture Notes in Computer Science **442** (Springer, Berlin, 1990).
- [6] L. Damas, Type assignment in programming languages, PhD Thesis, CST-33-85, University of Edinburgh, Edinburgh, Scotland (1985).
- [7] L. Damas and R. Milner, Principal type schemes for functional programs, in: *Proceedings 8th Annual ACM symposium on Principles of Programming Languages*, Albuquerque, NM (1982).
- [8] Y.-C. Fuhr and P. Mishra, *Polymorphic Subtype Inference: Closing the Theory-Practice Gap*, Lecture Notes in Computer Science **352** (Springer, Berlin, 1990).
- [9] R.W. Harper and B.C. Pierce, Extensible records without subsumption, Tech. Report CMU-CS-90-102, Carnegie Mellon University, School of Computer Science, Pittsburgh, PA (1990).
- [10] P. Hudak, S. Peyton Jones and P. Wadler, eds., Report on the programming language Haskell, version 1.2, *ACM SIGPLAN Notices* **27** (5) (1992).
- [11] M.P. Jones, Towards a theory of qualified types, Tech. Report PRG-TR-6-91, Programming Research Group, Oxford University Computing Laboratory, Oxford (1991).
- [12] M.P. Jones, Type inference for qualified types, Tech. Report PRG-TR-10-91, Programming Research Group, Oxford University Computing Laboratory, Oxford (1991).
- [13] M.P. Jones Qualified types: theory and practice, D.Phil. Thesis, Programming Research Group, Oxford University Computing Laboratory, Oxford (1992).
- [14] R. Milner, A theory of type polymorphism in programming, *J. Comput. Syst. Sci.* **17** (3) (1978).
- [15] J.C. Mitchell, Coercion and type inference (summary), in: *Proceedings 11th Annual ACM Symposium on Principles of Programming Languages*, Salt Lake City, UT (1984).
- [16] J.C. Mitchell A type-inference approach to reduction properties and semantics of polymorphic expressions, in: G. Huet, ed., *Logical Foundations of Functional Programming* (Addison Wesley, Reading, MA, 1990).
- [17] A. Ohori, A compilation method for ML-style polymorphic record calculi, in: *Proceedings 19th Annual ACM Symposium on Principles of Programming Languages* (1992).
- [18] S. Peyton Jones and P. Wadler A static semantics for Haskell (draft), Department of Computing Science, University of Glasgow, Glasgow, Scotland (1990).
- [19] D. Rémy, Typechecking records and variants in a natural extension of ML, in: *Proceedings 16th Annual ACM Symposium on Principles of Programming Languages*, Austin, TX (1989).
- [20] J.A. Robinson, A machine-oriented logic based on the resolution principle, *J. ACM* **12** (1965).
- [21] G. Smith, Polymorphic type inference for languages with overloading and subtyping, Ph.D. Thesis, Department of Computer Science, Cornell University, Ithaca, NY (1991).

- [22] D. Volpano and G. Smith, On the complexity of ML typability with overloading, in: *Proceedings 5th ACM Conference on Functional Programming Languages and Computer Architecture*, Lecture Notes in Computer Science **523** (Springer, Berlin, 1991).
- [23] P. Wadler and S. Blott, How to make *ad-hoc* polymorphism less *ad-hoc*, in: *Proceedings 16th Annual ACM Symposium on Principles of Programming Languages*, Austin, TX (1989).