

Quantifier elimination and parametric polymorphism in programming languages

HARRY G. MAIRSON*

Department of Computer Science, Brandeis University, Waltham, MA 02254, USA

Abstract

We present a simple and easy-to-understand explanation of ML type inference and parametric polymorphism within the framework of type monomorphism, as in the first order typed lambda calculus. We prove the equivalence of this system with the standard interpretation using type polymorphism, and extend the equivalence to include polymorphic fixpoints. The monomorphic interpretation gives a purely combinatorial understanding of the type inference problem, and is a classic instance of quantifier elimination, as well as an example of Gentzen-style cut elimination in the framework of the Curry–Howard propositions-as-types analogy.

Capsule review

The core type system of ML, Haskell, and other functional languages, often called the Hindley–Milner type system, is the fundamental source of parametric polymorphism in these languages. The seminal papers on this topic explain the core type system essentially ‘from scratch’. Mairson’s paper, however, gives an alternative and provably equivalent explanation, based on the monomorphic first-order typed lambda calculus. That polymorphism can be explained in terms of a monomorphic type system is at first surprising, until one realises that the essential ‘trick’ is to type each instance of let-bound identifiers uniquely. This also explains why this kind of parametric polymorphism is often called ‘let-bound’ polymorphism.

1 Introduction

In his influential paper ‘A theory of type polymorphism in programming’, Robin Milner proposed an extension to the first-order typed λ -calculus which has become known as the core of the ML programming language (Milner, 1978; Harper *et al.* 1990). The extension augmented the monomorphic type language of the first-order typed λ -calculus with *polytypes* (also known as *type schemes*) allowing a limited form of quantification over type variables. The expression language was similarly expanded by introducing the construct *let* $x = E$ in B where, by typing E with a polytype, the free occurrences of x in B could be typed *differently* (i.e. *polymorphically*) by varied instantiations of the quantified variables in the polytype. The added expressiveness of the type language then allowed *let* $x = E$ in B to be typable where the λ -calculus

* Supported by NSF Grant CCR-9017125, and grants from Texas Instruments and from the Tyson Foundation.

'equivalent' $(\lambda x. B)E$ might not be; the classic example of this facility is that *let* $I = \lambda z. z$ in II is typable in ML, while $(\lambda I. II)(\lambda z. z)$ is not first-order typable.

Type polymorphism has since been incorporated into a variety of functional programming languages (Turner, 1985; Hudak and Wadler, 1988). Among its virtues are *static typing*, so that all typing is done at compile time, with the guarantee that typechecked programs will not 'go wrong' at run time; *parametric polymorphism*, so that polymorphically typed code can be reused (via *let*) on abstract data types; and *decidable type inference*, where the compiler can automatically infer the most general type information (the so-called *principal type*) for an expression, so that any typing for the expression is a *substitution instance* of the principal type.

To what extent is Milner's proposal of type polymorphism necessary to achieve this degree of parametric polymorphism? Surprisingly, the type language of the first-order typed λ -calculus is sufficient to support ML-style parametric polymorphism, as long as we use the following inference rule for typing *let*-expressions

$$(let) \quad \frac{\Gamma \triangleright E: \tau_0 \quad \Gamma \triangleright [E/x]B: \tau_1}{\Gamma \triangleright let\ x = E\ in\ B: \tau_1}$$

Any ML program without free variables that is typable in the standard Milner–Damas inference system (Damas and Milner, 1982) is also typable using the classical Curry inference system (Curry and Feys, 1958) augmented with the above rule. Hence, parametric polymorphism as realized in ML may be achieved within the framework of type monomorphism.

Observe that the above inference rule realizes parametric polymorphism ('code reuse') explicitly through the expression $[E/x]B$: namely, each free occurrence of x gets replaced with a separate *copy* of the program E . The example of typing *let* $I = \lambda z. z$ in II , for instance, is reduced to typing $(\lambda z. z)(\lambda z. z)$, so each $\lambda z. z$ may be typed *differently*. The effect is the same as considering the expression to be a *marked redex* $(\lambda I. II)(\lambda z. z)$ ¹ in the theory of labelled reductions (Barendregt, 1984).

The monomorphic realization of ML's parametric polymorphism is not new. A recent survey of type systems in programming (Mitchell, 1990) attributes the observation to Albert Meyer. An earlier appearance of the idea is found in the dissertation of Luis Damas (1985), and in fact a question about it is found in the 1985 postgraduate examination in computing at Edinburgh University (Sannella, 1988).

In this paper we present a simple and easy-to-understand explanation of ML type inference in the framework of type monomorphism, where we prove its equivalence to the standard interpretation using type polymorphism. In addition, we analyze an extension of the ML inference system proposed by Alan Mycroft (1984) allowing fixpoints where the variable appearing in a recursion equation may have a polymorphic type. While type inference for this system is not computable (Kfoury *et al.*, 1990), we show that the inference system nonetheless has a purely monomorphic interpretation.

We believe that the monomorphic interpretation is important because it gives a *purely combinatorial* understanding of a significant fragment of the Girard/Reynolds second-order polymorphic typed λ -calculus (Girard, 1972; Reynolds, 1974). It also provides a classic example of *quantifier elimination*, which in the context of the

Curry–Howard propositions-as-types analogy serves as a sort of Gentzen-style cut elimination. The simple combinatorics of the monomorphic interpretation, which reduces the problem of type inference to first-order unification (Robinson, 1965), has played a central role in a complete analysis of the computational complexity of ML type inference (Kanellakis and Mitchell, 1989; Mairson, 1990; Kanellakis *et al.*, 1991) as well as providing insight into the first significant lower bounds on type inference for higher-order typed λ -calculi (Henglein and Mairson, 1991).

2 Preliminaries

2.1 Expressions

We consider ML expressions defined by the grammar

$$\mathcal{E} ::= x \mid \mathcal{E} \mathcal{E} \mid \lambda x. \mathcal{E} \mid \text{let } x = \mathcal{E} \text{ in } \mathcal{E} \mid \text{fix } x. \mathcal{E}$$

where x ranges over a set \mathcal{V} of *expression variables*. Excluding expressions of the form $\text{fix } x. E$ where $E \in \mathcal{E}$, the language considered is known as *Core ML* (see, for example, Mitchell and Harper, 1988). The syntax of Core ML is just that of the λ -calculus augmented with the polymorphic *let* construct. We write $FV(E)$ to denote the *free variables* of E . We allow α -renaming and β -reduction as in the λ -calculus, as well as reduction of *let*-expressions following the rule

$$\text{let } x = E \text{ in } B \rightarrow_{\text{let}} [E/x]B$$

For more details concerning reductions in the λ -calculus and ML, see Barendregt (1984), Hindley and Seldin (1987), and Harper *et al.* (1990).

2.2 Types

The syntax of *types* is given by the grammar

$$\begin{aligned} \mathcal{T}_0 & ::= t \mid \mathcal{T}_0 \rightarrow \mathcal{T}_0 \\ \mathcal{T} & ::= \mathcal{T}_0 \mid \forall t. \mathcal{T} \end{aligned}$$

where t ranges over a set \mathcal{TV} of *type variables*. We refer to $\tau \in \mathcal{T}_0$ as *monotypes*, and $\sigma \in \mathcal{T}$ as *polytypes* (sometimes also called *type schemes*).

We define a partial order \leq on \mathcal{T}_0 as $\tau_1 \leq \tau_2$ iff there exists a substitution $\Sigma: \mathcal{TV} \rightarrow \mathcal{T}_0$ such that $\Sigma\tau_1 \equiv \tau_2$, where \equiv denotes syntactic equivalence (overloaded for use on expressions as well). We interpret polytypes as sets of monotypes, using the following interpretation:

$$\begin{aligned} \langle \alpha \rangle & = \{ \alpha \} \quad \text{where } \alpha \text{ is a monotype} \\ \langle \forall t. \alpha \rangle & = \bigcup_{\tau \in \mathcal{T}_0} \langle [\tau/t] \alpha \rangle, \end{aligned}$$

where $[\tau/t] \alpha$ denotes the substitution of τ for free occurrences of t in α .

The interpretation of polytypes as sets of monotypes allows the definition of a partial order \sqsubseteq on \mathcal{T} : we write $\sigma_1 \sqsubseteq \sigma_2$ iff $\langle \sigma_2 \rangle \subseteq \langle \sigma_1 \rangle$. It is easy to see, for instance, that $\forall t. \alpha \sqsubseteq [\tau/t] \alpha$ for any polytype α and monotype τ . Note the minimal element of \mathcal{T} is $\forall t. t$, since $\langle \forall t. t \rangle = \mathcal{T}_0$. We further define an equivalence relation on \mathcal{T} as σ_1

$\cong \sigma_2$ iff $\langle \sigma_2 \rangle = \langle \sigma_1 \rangle$, and write $[\sigma]$ to denote the equivalence class $\{\alpha : \alpha \cong \sigma\}$. (This equivalence class definition will be used when we wish to argue that the names and order of bound variables in a polytype are not significant.) When σ is a polytype, we write $\bar{\sigma}$ to denote the monotype derived by removing all quantifiers from σ ; when τ is a monotype, we write $\check{\tau}$ to denote the polytype derived by quantifying over some subset of type variables in τ .

2.3 Inference rules

Expressions are associated with types using a fixed set of *inference rules*. We describe two such systems of rules: the first being the standard one given by Damas and Milner (1982), which we call the *polytype system*, and the second one a variant called the *monotype system*. As its name suggests, the monotype system associates expressions with monotypes only. The major point of this paper is to show simply why this limitation is not truly a restriction.

The inference rules manipulate an expression called a *type judgement*, written $\Gamma \triangleright E : \sigma$, where $E \in \mathcal{E}$, $\sigma \in \mathcal{T}$, and $\Gamma : FV(E) \rightarrow \mathcal{T}$. The type judgement is read as ‘with environment (context) Γ , expression E has type σ .’ In the λ -calculus, environments associate values to free variables in an expression, while in this case the environment is used to associate types with the free variables. We write $t_p \Gamma \triangleright E : \sigma$ (respectively, $t_M \Gamma \triangleright E : \sigma$) to mean that $\Gamma \triangleright E : \sigma$ is a derivable judgement in the polytype (respectively, monotype) system.

We give below the inference rules for the polytype and monotype systems. The polytype system is due to Damas and Milner (1982), and the monotype system is essentially due to Curry and Feys (1958) augmented with the rule for *let*. Observe the use in rule (*let_p*) of types with quantifiers (namely, the binding for x), requiring the rules (*gen_p*) and (*inst_p*) for quantifier introduction and elimination. For more details on type inference rules, we recommend Milner (1978), Cardelli (1984), Hancock (1987), and Wand (1987).

2.3.1 Core ML inference rules for the polytype system

$$\begin{array}{l}
 (\text{var}_p) \quad \frac{}{\Gamma \cup \{x : \sigma\} \triangleright x : \sigma} \\
 (\text{gen}_p) \quad \frac{\Gamma \triangleright E : \sigma(t) \quad [t \notin FV(\Gamma)]}{\Gamma \triangleright E : \forall t. \sigma(t)} \\
 (\text{inst}_p) \quad \frac{\Gamma \triangleright E : \forall t. \sigma(t)}{\Gamma \triangleright E : [\tau/t] \sigma} \\
 (\text{abs}_p) \quad \frac{\Gamma \cup \{x : \tau_0\} \triangleright E : \tau_1}{\Gamma \triangleright \lambda x. E : \tau_0 \rightarrow \tau_1} \\
 (\text{app}_p) \quad \frac{\Gamma \triangleright M : \tau_0 \rightarrow \tau_1 \quad \Gamma \triangleright N : \tau_0}{\Gamma \triangleright MN : \tau_1} \\
 (\text{let}_p) \quad \frac{\Gamma \triangleright E : \sigma \quad \Gamma \cup \{x : \sigma\} \triangleright B : \tau}{\Gamma \triangleright \text{let } x = E \text{ in } B : \tau}
 \end{array}$$

2.3.2 Core ML inference rules for the monotype system

$$\begin{array}{l}
 (var_M) \quad \frac{}{\Gamma \cup \{x: \tau\} \triangleright x: \tau} \\
 (app_M) \quad \frac{\Gamma \triangleright M: \tau_0 \rightarrow \tau_1 \quad \Gamma \triangleright N: \tau_0}{\Gamma \triangleright MN: \tau_1} \\
 (abs_M) \quad \frac{\Gamma \cup \{x: \tau_0\} \triangleright E: \tau_1}{\Gamma \triangleright \lambda x. E: \tau_0 \rightarrow \tau_1} \\
 (let_M) \quad \frac{\Gamma \triangleright E: \tau_0 \quad \Gamma \triangleright [E/x]B: \tau_1}{\Gamma \triangleright let\ x = E\ in\ B: \tau_1}
 \end{array}$$

3 Equivalence of the polytype and monotype systems

What should it mean when we say that the monotype system is ‘equivalent’ to the polytype system? A first guess might be that for any expression E , monotype τ , and context Γ , $\vdash_P \Gamma \triangleright E: \tau$ iff $\vdash_M \Gamma \triangleright E: \tau$. However, if Γ contains polytypes, it is not a valid context for a monotype judgement. Unfortunately, if we try $\vdash_P \Gamma \triangleright E: \tau$ iff $\vdash_M \Gamma_0 \triangleright E: \tau$, where Γ_0 are the monotype bindings of Γ , the statement is not true: consider $\vdash_P \{I: \forall t. t \rightarrow t\} \triangleright II: t \rightarrow t$ iff $\vdash_M \emptyset \triangleright II: t \rightarrow t$ – the monotype judgement is clearly false.

A second guess might be to insist that E be a closed term. A proof along these lines can indeed be given, where we proceed by a double induction on the structure of E and the maximum number of *let*-reductions needed to reduce E to *let-normal* form (see the Appendix of Kanellakis *et al.*, 1991). However, the proof is overly tedious and technical, and requires an understanding of minimal complete developments in the λ -calculus (Barendregt, 1984; Hindley and Seldin, 1987). But most of all, it contradicts an overwhelming sentiment that the equivalence we want is something very *simple* which should be *easy* to prove. What, then, should the equivalent of $\vdash_P \{I: \forall t. t \rightarrow t\} \triangleright II: t \rightarrow t$ be in the monotype system? (Note II is not closed.) We propose the following: the monotype equivalent should be $\vdash_M \emptyset \triangleright [E/I]II: t \rightarrow t$, where $[E/I]II$ is a closed term, and $\vdash_P \emptyset \triangleright E: \forall t. t \rightarrow t$. Of course, we are thinking in this case of $E \equiv \lambda z. z$.

What justifies this specification of equivalence? Polytypes are a kind of *shorthand* in the spirit of ‘code reuse’ and the Gentzen ‘cut’. In identifying an expression variable with a polytype, there is an implicit assumption that a piece of code exists with that type; what we have done in this example is simply to *insert* the code in place of its variable representative.

Generalizing from the example, we propose the following as a reasonable definition of ‘equivalence’.

Definition 3.1

Let $\Gamma = \{w_1: \alpha_1, w_2: \alpha_2, \dots, w_m: \alpha_m\}$ be any context of monotype bindings, and $\Gamma' = \{y_1: \beta_1, y_2: \beta_2, \dots, y_n: \beta_n\}$ be any context of polytype bindings. Let $\mathcal{F} = \{F_1, F_2, \dots, F_n\}$

be a set of terms where β_j is the principal type of F_j ; specifically, we insist that for $1 \leq j \leq n$,

$$\begin{aligned} &\vdash_P \Gamma \cup \{y_1 : \beta_1, y_2 : \beta_2, \dots, y_{j-1} : \beta_{j-1}\} \triangleright F_j : \beta_j \\ &\vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_{j-1}/y_{j-1}] F_j : \bar{\beta}_j \end{aligned}$$

where $\bar{\beta}_j$ is β_j with all quantifiers removed.

Given this framework, we can say precisely what is meant by equivalence. Let E be any term and τ be any monotype; then

$$\vdash_P \Gamma \cup \Gamma' \triangleright E : \tau \text{ if and only if } \vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_n/y_n] E : \tau \quad (1)$$

It is clear, and indeed natural, that the equivalent monotype judgement should be contingent on the explicit substitution of code represented at the type level by polytypes. In the case of a closed term E with empty contexts, we have $\vdash_P \emptyset \triangleright E : \tau$ iff $\vdash_M \emptyset \triangleright E : \tau$, as in Kanellakis *et al.* (1991). However, inspired by the example of Tait’s (1967) strong normalization theorem for the first-order typed λ -calculus, we have facilitated the proof by strengthening the induction hypothesis of what is to be a syntax-directed induction on E .

Before continuing with the proof, we introduce a standard structural lemma allowing us to ‘normalize’ derivations in the polytype system for use in a syntax-directed proof.

Lemma 3.2

Let $\Pi \vdash_P \Gamma \triangleright E : \vec{v}. \tau$ where Γ is a context, τ is a monotype, and \vec{v} denotes a (possibly empty) list of quantified variables. Then if E is not a variable, there exists a proof $\Pi' \vdash_P \Gamma \triangleright E : \tau$ where the last rule used in Π' is either (var_P) , (abs_P) , (app_P) , or (let_P) .

Proof

Observe that $\Pi \vdash_P \Gamma \triangleright E : \vec{v}. \tau$ is a syntax-directed proof, except for the use of (gen_P) and $(inst_P)$. The lemma states that the final uses of (gen_P) and $(inst_P)$ can be removed. The proof proceeds by induction on the number of such uses; in the basic case, clearly $\vec{v}. \tau \equiv \tau$.

For the inductive step, we must consider only two cases:

Case 1: The proof Π ends using the rule (gen_P)

$$(gen_P) \quad \frac{\Gamma \triangleright E : \vec{v}. \tau(t)}{\Gamma \triangleright E : \forall t. \vec{v}. \tau(t)}$$

Simply remove the last step of the proof to remove one quantifier from the type, and apply the inductive hypothesis.

Case 2: The proof Π ends using the rule $(inst_P)$

$$(inst_P) \quad \frac{\Gamma \triangleright E : \forall t. \vec{v}. \tau(t)}{\Gamma \triangleright E : \vec{v}. \tau(\alpha)}$$

where α is a monotype. Observe that the last rules appearing in the proof are a series of uses of (gen_p) and $(inst_p)$, where the former adds a quantifier, and the latter removes a quantifier. As such, they act like a stack. Identify the point (I) in the proof where t is universally quantified

$$(gen_p) \quad \frac{\Gamma \triangleright E: \vec{\forall}. \tau'(t)}{\Gamma \triangleright E: \forall t. \forall. \tau'(t)}$$

We now proceed as follows:

1. In the subproof rooted at (I), replace all free occurrences of t by α .
2. In the deductions from (I) until the end of the proof, remove the binding $\forall t$, replacing newly free occurrences of t by α .
3. Remove the conclusions of (I) and the final inference. The proof now has two fewer uses of (gen_p) and $(inst_p)$, so we can apply the inductive hypothesis. \square

Given the lemma and the stated assumptions on Γ , Γ' , and \mathcal{F} , we prove the above statement (1) in Definition 3.1 via structural induction on E , proceeding by case analysis. We assume by renaming that no variable is bound or quantified more than once.

Case $E \equiv w_i$

Then necessarily $\tau \equiv \alpha_i$ and $\vdash_p \Gamma \cup \Gamma' \triangleright w_i: \alpha_i$; since $[F_1/y_1][F_2/y_2] \dots [F_n/y_n] w_i \equiv w_i$, the result is immediate.

Case $E \equiv y_j$

In the forward direction, assume $\vdash_p \Gamma \cup \Gamma' \triangleright y_j: \tau$. Since $\vdash_p \Gamma \cup \Gamma' \triangleright y_j: \beta_j$ is a principal typing, we know $\beta_j \sqsubseteq \tau$; since $\tau \in \mathcal{T}_0$, we know there exists a substitution Σ such that $\Sigma \beta_j = \tau$. But since $[F_1/y_1][F_2/y_2] \dots [F_n/y_n] y_j \equiv [F_1/y_1][F_2/y_2] \dots [F_{j-1}/y_{j-1}] F_j$, and

$$\Pi \vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_{j-1}/y_{j-1}] F_j: \vec{\beta}_j$$

we know

$$\vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_{j-1}/y_{j-1}] F_j: \tau$$

by applying Σ to all types appearing in the proof Π .

In the reverse direction, suppose $\vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_{j-1}/y_{j-1}] F_j: \tau$; then by principality we know $\Sigma \vec{\beta}_j = \tau$. Given $\vdash_p \Gamma \cup \Gamma' \triangleright y_j: \beta_j$, we derive $\vdash_p \Gamma \cup \Gamma' \triangleright y_j: \tau$ by instantiating the \forall -bound variables of β_j according to Σ .

Case $E \equiv GH$

To prove ‘only if’, if $\vdash_p \Gamma \cup \Gamma' \triangleright GH: \tau$, then $\vdash_p \Gamma \cup \Gamma' \triangleright G: \tau' \rightarrow \tau$ and $\vdash_p \Gamma \cup \Gamma' \triangleright H: \tau'$ for some monotype τ' , by Lemma 3.2. From induction on G and H , we know $\vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_n/y_n] G: \tau' \rightarrow \tau$ and $\vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_n/y_n] H: \tau'$, so the result follows by use of (app_M) . Note the implications are all reversible, except that Lemma 3.2 is not needed.

Case $E \equiv \lambda x. G$

To prove ‘only if’, if $\vdash_P \Gamma \cup \Gamma' \triangleright \lambda x. G : \tau' \rightarrow \tau''$ where $\tau \equiv \tau' \rightarrow \tau''$, then by Lemma 3.2 we know $\vdash_P \Gamma \cup \{x : \tau'\} \cup \Gamma' \triangleright G : \tau''$. By induction on G (with a larger monotype context, since the binding for x is added), we have $\vdash_M \Gamma \cup \{x : \tau'\} \triangleright [F_1/y_1][F_2/y_2] \dots [F_n/y_n] G : \tau''$, and by (abs_M) , $\vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_n/y_n] \lambda x. G : \tau' \rightarrow \tau''$. Again, all implications are reversible.

Case $E \equiv \text{let } y = G \text{ in } H$

If $\vdash_P \Gamma \cup \Gamma' \triangleright \text{let } y = G \text{ in } H : \tau$, then by Lemma 3.2 and (let_P) , $\vdash_P \Gamma \cup \Gamma' \triangleright G : \sigma$ for (principal) polytype σ , and $\vdash_P \Gamma \cup \Gamma' \cup \{y : \sigma\} \triangleright H : \tau$. By $(inst_P)$, $\vdash_P \Gamma \cup \Gamma' \triangleright G : \bar{\sigma}$, so by induction on G we have a proof

$$\Pi \vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_n/y_n] G : \bar{\sigma}.$$

Hence by induction on H with polytype context $\{y_1 : \beta_1, y_2 : \beta_2, \dots, y_n : \beta_n, y : \sigma\}$ and associated code $\{F_1, F_2, \dots, F_n, G\}$ we have $\vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_n/y_n][G/y] H : \tau$, which we rewrite as $\vdash_M \Gamma \triangleright [[F_1/y_1][F_2/y_2] \dots [F_n/y_n] G/y][F_1/y_1][F_2/y_2] \dots [F_n/y_n] H : \tau$. Using proof Π above and (let_M) , we then have a proof of

$$\vdash_M \Gamma \triangleright \text{let } y = [F_1/y_1][F_2/y_2] \dots [F_n/y_n] G \text{ in } [F_1/y_1][F_2/y_2] \dots [F_n/y_n] H : \tau,$$

which is syntactically identical to

$$\vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_n/y_n] \text{let } y = G \text{ in } H : \tau.$$

Once again, the argument is reversible.

Given our above definition of equivalence, we then have:

Theorem 3.3

The polytype and monotype inference systems for Core ML are equivalent.

Corollary 3.4

Let E be a closed term and τ be a monotype. Then $\vdash_P \emptyset \triangleright E : \tau$ if and only if $\vdash_M \emptyset \triangleright E : \tau$.

3.1 Parametric polymorphism, cut elimination, and proof theory

In the well-known Curry–Howard propositions-as-types analogy, we read $E : \sigma$ not as ‘expression E has type σ ’, but ‘expression E is a proof of proposition σ ’. In this case, the function \rightarrow in types is read as logical implication. An environment Γ then serves as a set of labelled assumptions, so that a type judgement $\Gamma \triangleright E : \sigma$ elaborates a logical *sequent* $\tilde{\Gamma} \vdash \sigma$.

Proofs in sequent calculus, like type derivations, can be written in the form of trees, where the leaves form propositional hypotheses. The logical formalism of *cancelling hypotheses* via \rightarrow -introduction is reflected in removing type assumptions and introducing λ -abstraction. (For a further detailed but elementary discussion, see van Daalen, 1979.)

The process of β -reduction in the simply typed λ -calculus can be interpreted as a transformation on proof trees: if $\Pi_M \vdash \Gamma \triangleright \lambda x. M : \tau_1 \rightarrow \tau_2$ and $\Pi_N \vdash \Gamma \triangleright N : \tau_1$,

by *modus ponens* we get $\Pi_{MN} \vdash \Gamma \triangleright (\lambda x. M) N : \tau_2$. In this case we also know $\vdash \Gamma \triangleright [N/x]M : \tau_2$, since $\Pi'_M \vdash \Gamma \cup \{x : \tau_1\} \triangleright M : \tau_2$; we replace the *assumption* $x : \tau_1$ (appearing as a *leaf* in the proof tree Π'_M) by the *subtree (proof)* Π_N . Interpreted at the proof theory level, this transformation is an example of what is called, after Gentzen (1969), *cut elimination*, since $(\lambda x. M) N$ represents a proof where τ_1 is proved *once*, a ‘shortcut’ over $[N/x]M$, which may require many proofs of τ_1 .

The parametric polymorphism in Core ML introduced by the *let* construct can be viewed as a more powerful form of cut-elimination. The cut-elimination via β -reduction allows one proof of a \forall -free proposition to be used several times, while cut-elimination via *let*-reduction allows one proof of a proposition to be used several times, provided that the ‘use’ is always monomorphic (\forall -free). Rather than prove $P = \{\alpha \rightarrow \alpha, \beta \rightarrow \beta, (\alpha \rightarrow \beta) \rightarrow (\alpha \rightarrow \beta)\}$, for example, we construct one proof of $\forall t. t \rightarrow t$, and instantiate t appropriately. The propositions in P have a *most general unifier*, namely a proposition π such that $\pi \leq p$ for each $p \in P$. We make the related observation that the monomorphic inference rules for Core ML show that the *principal type property* proved in Milner (1978) is a straightforward consequence of the existence of most general unifiers in the first-order domain.

Similar to β -reduction, *let*-reduction can be viewed as a proof transformation. Since the expression *let* $x = E$ in B may have a polytype assigned to x , each use of E in $[N/x]B$ can instantiate the quantifiers differently. The ‘same’ proof is recycled to generate structurally similar propositions.

4 A characterization of polymorphic recursion by monotypes

The inference rules we have described thus far for typing ML programs do not include a rule for typing fixpoints, and hence do not allow recursion. In ML, fixpoints are constrained to be monomorphic; as such, the polytype system is usually extended with

$$(fix_{ML}) \quad \frac{\Gamma \cup \{x : \tau\} \triangleright E : \tau}{\Gamma \triangleright \mathbf{fix} \ x. E : \tau}$$

It is not difficult to show that when this rule is added to the polytype system, $\vdash_p \Gamma \triangleright \mathbf{fix} \ x. E : \tau$ iff $\vdash_p \Gamma \triangleright \lambda x. Eq \ E \ x : \tau$, where $Eq \equiv \lambda p. \lambda q. Kp(\lambda r. K(r \ p)(r \ q))$, given the usual definition $K \equiv \lambda x. \lambda y. x$, since Eq has principal type $\forall t. t \rightarrow t \rightarrow t$. As a consequence, adding monomorphic fixpoint does not make type inference particularly more complex.

Alan Mycroft (1984) proposed a more powerful variant to the above rule, whereby *fix*-bound variables could occur *polymorphically*

$$(fix_p) \quad \frac{\Gamma \cup \{x : \sigma\} \triangleright E : \sigma}{\Gamma \triangleright \mathbf{fix} \ x. E : \sigma}$$

In this rule, σ is a *polytype*. It has recently been shown by Kfoury *et al.* (1990) that type inference in the presence of such a polymorphic fixpoint is undecidable.

In this section, we show that polymorphic fixpoint can also be described using type monomorphism only. As the polymorphic inference system has been augmented with the rule (fix_p) , we add the following rules to the monotype system

$$\begin{array}{l} (\perp_M) \quad \Gamma \triangleright \mathbf{fix} \ x. x : \tau \\ (fix_M) \quad \frac{\Gamma \triangleright E_k : \tau \quad \Gamma \triangleright E_{k+1} : \tau}{\Gamma \triangleright \mathbf{fix} \ x. E : \tau} \end{array}$$

where for any term E with free variable x , we define

$$\begin{aligned} E_0 &\equiv \perp \equiv \mathbf{fix} \ x. x \\ E_1 &\equiv \mathbf{let} \ x_0 = E_0 \text{ in } [x_0/x] E \\ &\dots \\ E_{k+1} &\equiv \mathbf{let} \ x_k = E_k \text{ in } [x_k/x] E \end{aligned}$$

Henceforth, we refer to the polytype system as augmented with rule (fix_p) , and the monotype inference systems as augmented with rules (\perp_M) and (fix_M) . We observe that, properly speaking, (fix_M) is actually a rule *schema*, since its syntax varies with the integer k . However, it should be noted that *all* the inference rules are actually schemas! The monomorphic rules for typing polymorphic fixpoint have a simple explanation. An initial approximation $\perp : \forall t. t$ is made for the fixpoint, and the principal types of the E_k are repeatedly computed to better approximate the fixpoint until (possible) convergence.

To carry out this approach, we must show that for a given term E_k with principal type μ_k approximating the least fixpoint, and a known (type) fixpoint σ of $\mathbf{fix} \ x. E$ in the polytype system, it follows that the principal type μ_{k+1} of $E_{k+1} \equiv \mathbf{let} \ x_k = E_k \text{ in } [x_k/x] E$ always satisfies $\mu_k \sqsubseteq \mu_{k+1} \sqsubseteq \sigma$. Since there are (up to renaming of \forall -bound variables) only a finite number of types σ' satisfying $\sigma' \sqsubseteq \sigma$, we know by a pigeonhole argument that the sequence must converge. (We could instead show that the types form a complete partial order, from which convergence of the sequence is assured, but we prefer to proceed using a more combinatorial approach.)

We begin by indicating how polytype inferences can be derived from monotype inferences.

Proposition 4.1

Let $\Gamma, \Gamma', \mathcal{F}$ be defined as in Definition 3.1. If

$$\vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_n/y_n] E_j : \bar{\mu}$$

is a *principal* typing for $j \in \{k, k+1\}$, and E is \mathbf{fix} -free, then $\vdash_p \Gamma \cup \Gamma' \triangleright \mathbf{fix} \ x. E : \mu$.

Proof

By Theorem 3.3, we know that

$$\vdash_p \Gamma \cup \Gamma' \triangleright \mathbf{let} \ x_k = E_k \text{ in } [x_k/x] E : \mu,$$

so that

$$\vdash_p \Gamma \cup \Gamma' \cup \{x : \mu\} \triangleright E : \mu,$$

and hence by rule (fix_p)

$$\vdash_p \Gamma \cup \Gamma' \triangleright \mathbf{fix} \ x. E : \mu. \quad \square$$

Observe that in the above Proposition, μ must be the *principal* type of E_k and E_{k+1} , as principality is required for the proof of Theorem 3.3. To prove the converse, namely that monotype inferences can be derived from polytype inferences, a bit more detail is required. We begin with the following simple observation:

Proposition 4.2

If $\vdash_p \Gamma \cup \Gamma' \triangleright \mathbf{fix} \ x. E : \sigma$, then $\vdash_p \Gamma \cup \Gamma' \triangleright E_k : \sigma$ for all $k \geq 0$, and the *principal* type μ_k of E_k in environment $\Gamma \cup \Gamma'$ satisfies $\mu_k \sqsubseteq \sigma$.

Proof

By induction on k . The case $k = 0$ is trivial. For $k \geq 0$, recall by (fix_p) that $\vdash_p \Gamma \cup \Gamma' \cup \{x : \sigma\} \triangleright E : \sigma$. To show the same judgement holds of E_{k+1} , observe that as $E_{k+1} \equiv \mathbf{let} \ x_k = E_k \ \text{in} \ [x_k/x]E$, we must have by (let_p) $\vdash_p \Gamma \cup \Gamma' \triangleright E_k : \sigma'$ and $\vdash_p \Gamma \cup \Gamma' \cup \{x_k : \sigma'\} \triangleright [x_k/x]E : \sigma$; by inductive hypothesis, take $\sigma' \equiv \sigma$. Since $\mathbf{fix} \ x. E$ is typable, then all the E_k are also typable. As such, each E_k must have a principal type $\mu_k \sqsubseteq \sigma$. □

Our goal is to now show that some successive E_k, E_{k+1} must have the same *principal* type.

Proposition 4.3

For all types σ , $\mathcal{S}_\sigma = \{\{\alpha\} : \alpha \sqsubseteq \sigma\}$ is a finite set.

Proof

Define the *length* of a monotype as $|t| = 1, |\tau_0 \rightarrow \tau_1| = |\tau_0| + |\tau_1|$. If $\alpha \sqsubseteq \sigma$, then $\langle \alpha \rangle \supseteq \langle \sigma \rangle$, hence $\bar{\sigma} \in \langle \alpha \rangle$. Since $|\llbracket \tau/t \rrbracket \beta| \geq |\beta|$ (substitution cannot decrease length), we know $|\bar{\alpha}| \leq |\bar{\sigma}|$, and without loss of generality, the number of quantifiers preceding $\bar{\alpha}$ in α is bounded by $|\bar{\alpha}|$. □

Proposition 4.4

Let $\vdash_p \Gamma \triangleright E_i : \mu_i$ be a principal typing, where E is *fix*-free. Then for all $i \geq 0$, $\mu_i \sqsubseteq \mu_{i+1}$.

Proof

By induction on i . The basis is when $i = 0$: in this case, $\mu_0 \equiv \forall t. t \sqsubseteq \mu_1$. For the inductive step, assume by inductive hypothesis that $\mu_i \sqsubseteq \mu_{i+1}$. Then $\vdash_p \Gamma \cup \{x : \mu_i\} \triangleright E : \mu_{i+2}$, since we can take the proof $\vdash_p \Gamma \cup \{x : \mu_{i+1}\} \triangleright E : \mu_{i+2}$, and note that any instantiation of $x : \mu_{i+1}$ to a monotype can also be carried out if $x : \mu_i$. Since $\vdash_p \Gamma \triangleright E_i : \mu_i$, by (let_p), it is clear that $\vdash_p \Gamma \triangleright E_{i+1} : \mu_{i+2}$. By the principality of $\vdash_p \Gamma \triangleright E_{i+1} : \mu_{i+1}$, we then know that $\mu_{i+1} \sqsubseteq \mu_{i+2}$. □

Lemma 4.5

Let E be *fix*-free. Then $\vdash_p \Gamma \cup \{x : \sigma\} \triangleright E : \sigma$ if and only if there exists $k \geq 0$ and type $\sigma' \sqsubseteq \sigma$ where $\vdash_p \Gamma \triangleright E_k : \sigma'$ and $\vdash_p \Gamma \triangleright E_{k+1} : \sigma'$.

Proof

Proposition 4.1 proves the ‘if’ direction. As for ‘only if’, recall that $\mathcal{S}_\sigma = \{\alpha : \alpha \sqsubseteq \sigma\}$. By Proposition 4.2, $\mu_i \sqsubseteq \sigma$ for all $i \geq 0$, hence $[\mu_i] \in \mathcal{S}_\sigma$. As \mathcal{S}_σ is finite by Proposition 4.3, there exists by the pigeonhole principle $0 \leq k < \ell \leq |\mathcal{S}_\sigma|$ where $[\mu_k] = [\mu_\ell]$. But by Proposition 4.4, $\mu_k \sqsubseteq \mu_{k+1} \sqsubseteq \mu_\ell$, hence $[\mu_k] = [\mu_{k+1}]$; take $\sigma' \equiv \mu_k$. \square

Finally, we can state the equivalence theorem for the polytype and monotype inference systems with polymorphic fixpoint:

Theorem 4.6

Let $\Gamma, \Gamma', \mathcal{F}$ be as in Definition 3.1. Then

$$\vdash_P \Gamma \cup \Gamma' \triangleright E' : \sigma \quad \text{if and only if} \quad \vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_n/y_n] E' : \bar{\sigma}.$$

Proof

We augment the induction proof of Theorem 3.3 with the case $E' \equiv \mathbf{fix} \ x. E$. Without loss of generality, assume σ is a principal typing. If $\vdash_P \Gamma \cup \Gamma' \triangleright \mathbf{fix} \ x. E : \sigma$, then by the argument of Lemma 3.2 $\vdash_P \Gamma \cup \Gamma' \cup \{x : \sigma\} \triangleright E : \sigma$.

If E is **fix**-free, by the previous Theorem there exist $k \geq 0$ and $\sigma' \sqsubseteq \sigma$ such that $\vdash_P \Gamma \cup \Gamma' \triangleright E_j : \bar{\sigma}'$ for $j \in \{k, k+1\}$. Since the E_j are **fix**-free, by Theorem 3.3 we have $\vdash_M \Gamma \triangleright [F_1/y_1][F_2/y_2] \dots [F_n/y_n] E_j : \bar{\sigma}'$. Let $\tilde{E} \equiv [F_1/y_1][F_2/y_2] \dots [F_n/y_n] E$, so that $\vdash_M \Gamma \triangleright \tilde{E}_j : \bar{\sigma}'$ for $j \in \{k, k+1\}$; by (\mathbf{fix}_M) we have $\vdash_M \Gamma \triangleright \mathbf{fix} \ x. \tilde{E} : \bar{\sigma}'$. However, note that $\mathbf{fix} \ x. \tilde{E} \equiv [F_1/y_1][F_2/y_2] \dots [F_n/y_n] \mathbf{fix} \ x. E$.

In the case that E is not **fix**-free, we observe that using the inductive hypothesis, the above claims about **fix**-free expressions also hold with such stipulation. We then repeat the argument. \square

5 Final remarks

It seems obvious that the polytype and monotype inference systems should be equivalent in their expressive power. When we use an expression defined with *let* and having a polytype, we instantiate the quantified type variables to be in accordance with the type context. Had we the code instead, we could type the code differently in each instance. In the ML module system, identifiers are bound to types without code, so that type inference can still take place; an obvious use for this facility is in incremental compilation. Of course, the module could instead give the code, but in practice the type is shorter. There are, however, examples where the type is much larger than the code, and these pathological examples provide the foundation for lower bounds on type inference (Kanellakis and Mitchell, 1989; Mairson, 1990; Kanellakis *et al.*, 1991). In short: most general *specifications* (i.e. types) can be considerably longer than the programs implementing the specifications when the specification language is rich enough.

The equivalence proofs we have given are based on a fairly straightforward structural induction. The contribution of this paper, for the most part, is to give a precise definition of the equivalence. The lesson is simple: type polymorphism is not

needed when you do not reuse code, and instead use separate copies of the same code. Our equivalence proofs explain a theory of type monomorphism in programming, where it becomes clear that the type polymorphism found in ML-like languages admits straightforward quantifier elimination procedures.

Acknowledgements

I would like to thank Gerd Hillebrand, Paris Kanellakis and Lincoln Wallen for several stimulating and helpful discussions. In addition, I would like to acknowledge the generous hospitality of the Computer Science Department at UC Santa Barbara, the Music Academy of the West and the Cate School of Carpenteria during my stay in Santa Barbara in the summer of 1990, when I began work on this paper.

References

- Barendregt, H. 1984. *The Lambda Calculus: Its Syntax and Semantics*. North-Holland.
- Cardelli, L. 1984. Basic polymorphic type-checking. *Science of Computer Programming*, 8 (2): 147–172.
- Curry, H. B. and Feys, R. 1958. *Combinatory Logic I*. North-Holland.
- van Daalen, D. 1979. *Logic and Structure*. Springer-Verlag.
- Damas, L. 1985. *Type assignment in programming languages*. PhD dissertation, Computer Science Department, Edinburgh University.
- Damas, L. and Milner, R. 1982. Principal type schemes for functional programs. In *9th ACM Symposium on Principles of Programming Languages*, 207–212, January.
- Gentzen, G. 1969. *The Collected Papers of Gerhard Gentzen* (ed. E. Szabo). North-Holland.
- Girard, J.-Y. 1972. *Interprétation fonctionnelle et élimination des coupures de l'Arithmétique d'ordre supérieur*. Thèse de Doctorat d'Etat, Université de Paris VII.
- Hancock, P. 1987. Polymorphic type-checking. In Peyton-Jones, S., *The Implementation of Functional Programming Languages*. Prentice-Hall.
- Harper, R., Milner, R. and Tofte, M. 1990. *The Definition of Standard ML*. MIT Press.
- Henglein, F. and Mairson, H. 1991. The complexity of type inference for higher-order typed lambda calculi. In *Proc. 18th ACM Symposium on the Principles of Programming Languages*, 119–130, January.
- Hindley, R. and Seldin, J. 1987. *Introduction to Combinators and Lambda Calculus*. Cambridge University Press.
- Hudak, P. and Wadler, P. L., eds. *Report on the functional programming language Haskell*. Yale University Technical Report YALEU/DCS/RR656.
- Kanellakis, P. C. and Mitchell, J. C. 1989. *Polymorphic unification and ML typing*. Brown University Technical Report CS-89-40, August. (Also in *Proc. 16th ACM Symposium on the Principles of Programming Languages*, 105–115, January.)
- Kanellakis, P. C., Mairson, H. G. and Mitchell, J. C. 1991. Unification and ML type reconstruction. In *Computational Logic: Essays in Honor of Alan Robinson*, J.-L. Lassez and G. Plotkin, eds., MIT Press, 1991.
- Kfoury, A. J., Tiuryn, J. and Urzyczyn, P. 1990. Undecidability of the semi-unification problem. In *Proc. 22nd ACM Symposium on Theory of Computing*, May.
- Mairson, H. G. 1990. Deciding ML typability is complete for deterministic exponential time. In *Proc. 17th ACM Symposium on the Principles of Programming Languages*, 382–401, January.
- Milner, R. 1978. A theory of type polymorphism in programming. *J. Computer and System Sciences*, 17: 348–375.

- Mitchell, J. C. 1990. Type systems for programming languages. Volume B, pp. 365–468. In van Leeuwen *et al.*, eds., *Handbook of Theoretical Computer Science*. North-Holland.
- Mitchell, J. C. and Harper, R. 1988. The essence of ML. In *Proc. 15th ACM Symposium on Principles of Programming Languages*, 28–46, January.
- Mycroft, A. 1984. Polymorphic types schemes and recursive definitions. In M. Paul and B. Robinet, eds., *Proc. International Symposium on Programming*. Volume 167 of *Lecture Notes in Computer Science*, Springer-Verlag, 217–228.
- Reynolds, J. C. 1974. Towards a theory of type structure. In *Proc. Paris Colloquium on Programming*, Volume 19 of *Lecture Notes in Computer Science*, Springer-Verlag, 408–425.
- Robinson, J. A. 1965. A machine oriented logic based on the resolution principle. *J. of the ACM*, 12 (1): 23–41.
- Sannella, D., ed. 1988. *Postgraduate Examination Questions in Computation Theory, 1978–1988*. Laboratory for Foundations of Computer Science, Report ECS-LFCS-88-64, Edinburgh University.
- Tait, W. W. 1967. Intensional interpretation of functionals of finite type I. *J. Symbolic Logic*, 32: 198–212.
- Turner, D. A. 1985. Miranda: A non-strict functional language with polymorphic types. In *IFIP International Conference on Functional Programming and Computer Architecture*. Volume 20 of *Lecture Notes in Computer Science*, Springer-Verlag, 1–16.
- Wand, M. 1987. A simple algorithm and proof for type inference. *Fundamenta Informaticae*, 10.