



Efficient data representation in polymorphic languages

Xavier Leroy

► **To cite this version:**

Xavier Leroy. Efficient data representation in polymorphic languages. RR-1264, INRIA. 1990. <inria-00075295>

HAL Id: inria-00075295

<https://hal.inria.fr/inria-00075295>

Submitted on 24 May 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



UNITÉ DE RECHERCHE
INRIA-ROCQUENCOURT

Institut National
de Recherche
en Informatique
et en Automatique

Domaine de Voluceau
Rocquencourt
B.P. 105
78153 Le Chesnay Cedex
France
Tél.:(1)39 63 55 11

Rapports de Recherche

N°1264

Programme 2

*Calcul symbolique, Programmation
et Génie logiciel*

EFFICIENT DATA REPRESENTATION IN POLYMORPHIC LANGUAGES

Xavier LEROY

Août 1990

Représentations de données efficaces dans les langages polymorphes

Xavier Leroy¹

Résumé

Les langages avec typage polymorphe (p.ex. ML) utilisent généralement des représentations de données à la Lisp pour leur implémentation—tout doit tenir dans un mot, si nécessaire en étant alloué dans le tas et manipulé via un pointeur. Ceci s’explique par le fait que, contrairement aux langages classiques avec typage statique (Pascal), on ne peut pas donner à la compilation un et un seul type à chaque expression, condition *sine qua non* pour pouvoir utiliser des représentations plus efficaces (p.ex. des valeurs non allouées occupant plusieurs mots). Dans ce rapport, on montre comment tirer parti des informations de typage statique pour mélanger de manière correcte deux styles de représentations de données dans l’implémentation d’un langage polymorphe: on utilise des représentations spécialisées, particulièrement efficaces, quand les types sont complètement déterminés à la compilation, et des représentations uniformes, à la Lisp, dans les autres cas.

Efficient data representation in polymorphic languages

Xavier Leroy

Abstract

Languages with polymorphic types (e.g. ML) have traditionally been implemented using Lisp-like data representations—everything has to fit in one word, if necessary by being heap-allocated and handled through a pointer. The reason is that, in contrast with conventional statically-typed languages such as Pascal, it is not possible to assign one unique type to each expression at compile-time, an absolute requirement for using more efficient representations (e.g. unallocated multi-word values). In this paper, we show how to take advantage of the static polymorphic typing to mix correctly two styles of data representation in the implementation of a polymorphic language: specialized, efficient representations are used when types are fully known at compile-time; uniform, Lisp-like representations are used otherwise.

¹Ecole Normale Supérieure et INRIA Rocquencourt, projet Formel

1 Introduction

Most programming languages include some kind of type system. Among the numerous motivations for using a type system, I shall focus on two main goals: 1- to make programs safer, and 2- to allow for better compilation.

The first concern is to ensure data integrity in programs. Many operations are meaningless when they are performed on values of the wrong kind, such as, for instance, applying a boolean as if it was a function. In these cases, the results are either meaningless, unpredictable values, or a hardware fault. One of the aims of a type system is to prevent such run-time type errors. From this standpoint, typing can be either static (performed at compile-time), or dynamic (performed at run-time, just before type-constrained operations). But in any case, typing must be strong: there should be no way to claim that a value has a given type when in fact it does not.

Another aim of a type system is to support efficient compilation. Most hardware architectures are somehow typed, in the sense that some resources are dedicated to operate on certain kinds of data. For instance, many processors have two sets of registers, one set to hold integers and pointers, and the other to hold floating-point values. On integer registers, no instructions are provided to perform floating-point computations, and vice-versa. In addition, floating-point registers are usually wider than integer registers, therefore a floating-point number cannot fit in an integer register. When mapping a programming language on these architectures, it is important to know which values are floating-point numbers and which ones are not, otherwise the correctness of the compilation would be compromised. For instance, to compile an assignment `a := b`, it is crucial to know the types of variables `a` and `b`, in order to determine their exact size, and copy the right number of bytes. Of course, this type information must be available at compile-time, hence the emphasis is on static typing, here. From this standpoint, strong and weak typing are equally acceptable, provided that typing violations are explicitly mentioned in the source code.

Given a strong type system with static checking, both approaches can be combined in a single language, as in Algol-68, Pascal or Modula-2 for instance. This gives the best of both worlds: type safety in a well-compiled language. However, it seems difficult to achieve the same results with more powerful type systems than the one of Pascal. Using types to determine the size (and other relevant information) of a value works fine when every value has exactly one type. This condition is too restrictive in practice: one cannot write “generic” functions once for all and then apply them to data of several types, provided that it makes sense. For instance, a sorting function should be able to operate on many kinds of arrays (e.g. arrays of integers and arrays of strings), assuming a suitable comparison function is provided.

Therefore, many advanced type systems lift the restriction that every value has a unique, statically-known type, through the introduction of concepts such as type abstraction, subtyping, inheritance, and polymorphism. See Cardelli and Wegner’s [5] for a uniform presentation of these features. In the following, I concentrate on polymorphism, as in ML [12]. Polymorphic type systems allow type expressions to contain universally quantified type variables. For instance, a sorting function should, for all types T , take an array of elements of type T and return another array of elements of type T , given in addition an ordering predicate over T , that is a function taking a pair of elements of type T and returning a boolean. Its polymorphic type is therefore $\forall T. (T \times T \rightarrow \text{Bool}) \rightarrow \text{Array}(T) \rightarrow \text{Array}(T)$. As the quantification implies, the type variable T can be substituted by any actual type, such as `Int` or `String`. In a sense, the type formula given above summarizes all the possible types for the sorting function.

When a value can safely belong to several types, we cannot always determine statically all its characteristics. In the previous example, the sorting function is likely to copy elements of the array given as parameter, but the type of these elements can be any instance of the type variable T , that is any type. Therefore, we do not know at compile-time how many bytes to move to perform the copy. Consider two solutions to this problem.

One is to defer the compilation of polymorphic function until they are actually applied to values of known types. Then, we can deduce the type which the polymorphic function is used with, and therefore get all the information (e.g. sizes) we need to compile. Of course, we will have to compile several specialized versions of a polymorphic functions, if it is used with several different types. This technique is often used for Ada’s “generics”. Its strength is that it allows the use of efficient representations, and the production of good code (as efficient as if the function was monomorphic). However, it results in code duplication and loss of separate compilation. In addition, it is hard to maintain the illusion that polymorphic functions are still first-class objects. For instance, compilation of polymorphic functions built on top of other polymorphic functions, or functions taking polymorphic objects as arguments must be deferred in turn. This may lead to an explosion in the number of specializations of a polymorphic function we have to compile—this number could even be infinite, as in the following example, written in ML syntax:

```
datatype 'a chain = Empty | Cell of 'a * 'a list chain;
fun length_of_chain Empty = 0
  | length_of_chain (Cell(x,rest)) = 1 + length_of_chain rest;
```

The other approach to polymorphic compilation is to revert to *uniform* data representations. That is, data representation for all types share a common format, allowing “generic” operations such as parameter passing, function return, assignment, ... to be compiled without knowing the type of the values they manipulate. In particular, all representations must have the same size, usually one machine word, and any data which does not fit in one word has to be heap-allocated and represented by a pointer (assuming pointers do fit in one word). In addition, the calling conventions must be exactly the same for all functions. With these representation constraints, type information is no longer needed to generate correct code, and compiling polymorphic functions is no longer a problem. This approach is used in most implementations of ML to date.

The main drawback of the second approach is that uniform representations are not as efficient as the kind of *specialized* representations used in monomorphic languages. First, all data which does not fit naturally in one word must be heap-allocated. This is much more expensive than carrying around multi-word representations in several registers, resulting in high heap consumption and frequent garbage collection. Second, communications across function calls are not very efficient: a function returning a floating-point number, just computed in a floating-point register, has to copy it in the heap and return a pointer to it, only to have the callee dereference that pointer, and reload the number in a floating-point register before using it.

This use of uniform data representations is not the main reason why current ML compilers do not produce as efficient code as, say, Modula compilers. However, as modern compiling techniques are applied to the ML language, I think this representation problem will show up as a serious bottleneck for efficient compilation of ML, and similarly for other polymorphic languages.

In this paper, I set out to reconcile the safety offered by strong static typing, the convenience and conceptual cleanliness of having polymorphic functions compiled only once, just like regular functions, and the efficiency of specialized data representations. Since the first two requirements

imply polymorphic functions must work on uniform representations anyway, the third requirement has to be relaxed slightly. In this paper, I attempt to mix uniform and specialized representations in the implementation of a polymorphic language. It is intended that monomorphic functions work only on specialized representations, the uniform representations being used only to communicate with polymorphic functions. The main question is, then, How can one infer where to insert the necessary conversions between representations?

The remainder of the paper is organized as follows: section 2 presents the problem of data representation in the case of a monomorphic calculus. Uniform and specialized representations are contrasted by giving two compilation schemes for a simple stack-based abstract machine. This is mostly implementor’s wisdom, but presented in a uniform setting. Using the same approach, section 3 tackles the representation problem in the case of a polymorphic calculus. The main novelty of this paper—combining polymorphism and specialized representations—is informally presented here, then formalized using an intermediate calculus, with a restricted notion of polymorphism. Section 4 aims at showing that specialized representations can be profitably used in the ML language. To this end, the main features of the type system of ML are recalled, and their compatibility with specialized representations checked. Section 5 presents the interactions between data representation and garbage collection, which is required by most modern programming languages. We show that specialized representations do not preclude garbage collection, and that they alleviate the need for tagging—a costly operation. Finally, we give a few concluding remarks in section 6

2 The monomorphic case

In this section, we present the problem of data representation in the simple case where every term has exactly one type, and this type is known at compile-time.

2.1 A simply-typed language

We consider a small language based on the simply-typed λ -calculus with constants. The only data structures are pairs. The constants include integer and floating-point numbers, of base types `Int` and `Float`, as well as primitive operations such as `succ_int` and `add_float`.

The syntax of this calculus is as follows. We write i for an integer, f for a floating-point number, c for a constant, x or y for a variable, a or b for terms, and A, B for types.

$$\begin{aligned} c & ::= i \mid f \mid \text{succ_int} \mid \text{add_float} \mid \dots \\ a & ::= c \mid x \mid \lambda x : A. b \mid b(a) \mid (a, b) \mid a.\text{fst} \mid a.\text{snd} \\ A & ::= \text{Int} \mid \text{Float} \mid A \rightarrow B \mid A \times B \end{aligned}$$

Typing rules are classical. They are written in structural operational semantics style [15], as a set of axioms and inference rules defining the judgement “under assumptions E , term a has type A ”, written $E \vdash a : A$. The typing environment E consists in a sequence of assumptions of the form $x : A$, meaning that variable x is assumed to have type A . For each constant c , we write $T(c)$ for its associated type; in particular, $T(i) = \text{Int}$, $T(f) = \text{Float}$, $T(\text{add_float}) = \text{Float} \times \text{Float} \rightarrow \text{Float}$, and so on.

$$E \vdash c : T(c) \qquad E_1, x : A, E_2 \vdash x : A$$

$$\begin{aligned}
\mathcal{U}_E(i) &= \text{Const}(i) \\
\mathcal{U}_E(f) &= \text{Const}(\langle f_{low}, f_{high} \rangle) \\
\mathcal{U}_E(x) &= \text{Access}(\text{pos}_x(E)) \\
\mathcal{U}_E(\lambda x : A. b) &= \text{Abstr}(\mathcal{U}_{x:A, E}(b), \text{Return}) \\
\mathcal{U}_E(b(a)) &= \mathcal{U}_E(b), \mathcal{U}_E(a), \text{Apply} \\
\mathcal{U}_E(a, b) &= \mathcal{U}_E(a), \mathcal{U}_E(b), \text{Pair} \\
\mathcal{U}_E(a.\text{fst}) &= \mathcal{U}_E(a), \text{First} \\
\mathcal{U}_E(a.\text{snd}) &= \mathcal{U}_E(a), \text{Second} \\
\mathcal{U}_E(\text{add_float}(a, b)) &= \mathcal{U}_E(a), \mathcal{U}_E(b), \text{AddFloat} \\
\mathcal{U}_E(\text{add_float}) &= \text{Abstr}(\text{Access}(0), \text{First}, \text{Access}(0), \text{Second},) \\
&\quad \text{AddFloat}, \text{Return}
\end{aligned}$$

Figure 1: Compilation scheme with uniform representations

$$\begin{array}{c}
\frac{x : A, E \vdash b : B}{E \vdash \lambda x : A. b : A \rightarrow B} \\
\frac{E \vdash a : A \quad E \vdash b : B}{E \vdash (a, b) : A \times B} \\
\frac{E \vdash b : A \rightarrow B \quad E \vdash a : A}{E \vdash b(a) : B} \\
\frac{E \vdash a : A \quad E \vdash b : B}{E \vdash a.\text{fst} : A \quad E \vdash a.\text{snd} : B}
\end{array}$$

2.2 Evaluation using uniform representations

We now give a compilation scheme for this small language. The target machine is a stack-based machine with environments, where functional values are represented by closures. Call-by-value is assumed, to be consistent with the strict semantics of ML. This machine is very close to Landin's SECD [10] and Cardelli's FAM [3]. The compilation scheme is given in figure 1. It is straightforward, except maybe for the treatment of variables. The value of a variable is to be found at run-time in the environment (a tuple of values). To access this variable, we need to know its position inside this tuple. That's the *raison d'être* of the compilation environment E , similar in structure to typing environments: it records the name of all free variables, in the order they will appear in the run-time environment. Then, the position of a variable x in an environment described by E is simply:

$$\begin{aligned}
\text{pos}_x(\emptyset) &\text{ is undefined} \\
\text{pos}_x(x : A, E) &= 0 \\
\text{pos}_x(y : A, E) &= 1 + \text{pos}_x(E)
\end{aligned}$$

This simple compilation scheme does not make use of typing information. This fact has deep consequences on the way data are represented in the machine, namely that *all data must fit in one word*. Indeed, the functions $\lambda x : A. x$, for all types A , have exactly the same code, and are applied in exactly the same way. This means that the instructions to apply a function to an argument, fetch

	Code	Stack	Environment
Before	<code>Const(v), C</code>	S	γ
After	C	v, S	γ
Before	<code>Access(n), C</code>	S	$\gamma = \langle v_1, \dots, v_n, \dots \rangle$
After	C	v_n, S	γ
Before	<code>Abstr(C₁), C</code>	S	γ
After	C	$\langle \langle C_1 \rangle, \gamma \rangle, S$	γ
Before	<code>Apply, C</code>	$v, \langle \langle C_1 \rangle, \langle \Gamma_1 \rangle \rangle, S$	γ
After	C_1	$\langle C \rangle, \gamma, S$	$\langle v, \Gamma_1 \rangle$
Before	<code>Return, C</code>	$v, \langle C_0 \rangle, \gamma_0, S$	γ
After	C_0	v, S	γ_0
Before	<code>Pair, C</code>	v_2, v_1, S	γ
After	C	$\langle v_1, v_2 \rangle, S$	γ
Before	<code>First, C</code>	$\langle v_1, v_2 \rangle, S$	γ
After	C	v_1, S	γ
Before	<code>Second, C</code>	$\langle v_1, v_2 \rangle, S$	γ
After	C	v_2, S	γ
Before	<code>AddFloat, C</code>	$\langle f'_{low}, f'_{high} \rangle, \langle f_{low}, f_{high} \rangle, S$	γ
After	C	$\langle (f + f')_{low}, (f + f')_{high} \rangle, S$	γ

Figure 2: An abstract machine with uniform representations

a value from the environment, and returning a value, must operate uniformly on data of any type. This implies that all data representations have the same size, and in case of register machines with several register classes, that they all use the same register class.

As a consequence, data which do not fit in one word are allocated in the heap, and handled through a pointer. (We assume that any pointer fits in one word). We write $\langle S \rangle$ for a pointer to the sequence of words S , located in the heap. For instance, the pair of values v_1 and v_2 is represented by $\langle v_1, v_2 \rangle$, and similarly for closures. Regarding constants, we assume that integers may fit in one word, but that high-precision floating-point numbers f require two words, written f_{low} and f_{high} , therefore f is represented by $\langle f_{low}, f_{high} \rangle$.

These representations lead to the transition function given in figure 2.

2.3 Inefficiencies of uniform representations

The evaluation mechanism presented above is simple, but not very efficient. As an example, let us consider the function $f(x, y) = x + 2y$, where x and y are reals, represented in floating-point. Since our calculus does not directly support functions with several arguments, some transformation is required. We can make it into a function taking a pair:

$$f_{pair} = \lambda z : \text{Float} \times \text{Float}. \text{add_float}(z.\text{fst}, \text{add_float}(z.\text{snd}, z.\text{snd}))$$

or into a function taking x and returning another function (this technique is known as *currying*):

$$f_{cur} = \lambda x : \text{Float}. \lambda y : \text{Float}. \text{add_float}(x, \text{add_float}(y, y)).$$

$$\begin{aligned}
\mathcal{S}_E(i^{\text{Int}}) &= \text{Const}(i) \\
\mathcal{S}_E(f^{\text{Float}}) &= \text{Const}(f_{\text{high}}), \text{Const}(f_{\text{low}}) \\
\mathcal{S}_E(x^A) &= \text{Access}(p + s - 1), \dots, \text{Access}(p) \\
&\quad \text{where } p = \text{pos}_x(E) \text{ and } s = \text{size}(A) \\
\mathcal{S}_E(\lambda x : A. b^B) &= \text{Abstr}(\mathcal{S}_{x:A,E}(b^B), \text{Return}_{\text{size}(B)}) \\
\mathcal{S}_E(b(a^A)) &= \mathcal{S}_E(b), \mathcal{S}_E(a), \text{Apply}_{\text{size}(A)} \\
\mathcal{S}_E(a, b) &= \mathcal{S}_E(a), \mathcal{S}_E(b) \\
\mathcal{S}_E(a^{A \times B}.fst) &= \mathcal{S}_E(a), \text{First}_{\text{size}(A), \text{size}(B)} \\
\mathcal{S}_E(a^{A \times B}.snd) &= \mathcal{S}_E(a), \text{Second}_{\text{size}(A), \text{size}(B)} \\
\mathcal{S}_E(\text{add_float}(a, b)) &= \mathcal{S}_E(a), \mathcal{S}_E(b), \text{AddFloat} \\
\mathcal{S}_E(\text{add_float}) &= \text{Abstr}(\text{Access}(0), \text{Access}(1), \text{Access}(2), \text{Access}(3), \\
&\quad \text{AddFloat}, \text{Return}_2)
\end{aligned}$$

Figure 3: Compilation scheme with specialized representations

Both versions are inefficient in terms of heap allocation and memory accesses. First, each floating-point addition must allocate two words in the heap to store its result, and perform three memory accesses. This is especially absurd in the case of the innermost addition, whose result is used only once, by the next instruction. Admittedly, a simple analysis of the code could detect that, and avoid allocating the intermediate result. But the final result must be allocated anyway, as required by the calling convention.

The passing of the two parameters is also inefficient. In the case of the uncurried form, the caller has to build a pair of the two arguments, which means allocating two words in the heap and performing two memory writes, only to have the callee discard the pair and solely use its components, at the cost of one memory access for each use of a parameter. In the case of the curried form, the main flaw is the building of an intermediate closure between the passing of the first and the second argument. (This closure corresponds to the partial application of the function to its first argument.)

To be more efficient, it is clear now that we have to lift the restriction that any value must either fit in one word or be allocated, and be able to handle unallocated multi-word values. To do so, we need to statically keep track of the size of all values and results. (In case of a register machine with several classes of registers, we would have to record the suitable register class for each value). Obviously, all this information is already contained in the typing of the program; what we shall present now is a compilation scheme taking advantage of the types.

2.4 Evaluation using specialized representations

The new compilation function is given in figure 3. It corresponds to the case where data representations are as “flat” as possible: floating-point numbers are not allocated, pairs are simple concatenations of the sequences of words representing their components, and for closures, only the

	Code	Stack	Environment
Before	$\text{Const}(v); C$	S	γ
After	C	v, S	γ
Before	$\text{Access}(n), C$	S	$\gamma = \langle v_1, \dots, v_n, \dots \rangle$
After	C	v_n, S	γ
Before	$\text{Abstr}(C_1), C$	S	γ
After	C	$\langle C_1 \rangle, \gamma, S$	γ
Before	Apply_i, C	$v_i, \dots, v_1, \langle C_1 \rangle, \langle E_1 \rangle, S$	γ
After	C_1	$\langle C \rangle, \gamma, S$	$\langle v_1, \dots, v_i, E_1 \rangle$
Before	Return_i, C	$v_i, \dots, v_1, \langle C_0 \rangle, \gamma_0, S$	γ
After	C_0	v_i, \dots, v_1, S	γ_0
Before	$\text{First}_{i,j}, C$	$w_j, \dots, w_1, v_i, \dots, v_1, S$	γ
After	C	v_i, \dots, v_1, S	γ
Before	$\text{Second}_{i,j}, C$	$w_j, \dots, w_1, v_i, \dots, v_1, S$	γ
After	C	w_j, \dots, w_1, S	γ
Before	$\text{AddFloat}, C$	$f'_{low}, f'_{high}, f_{low}, f_{high}, S$	γ
After	C	$(f + f')_{low}, (f + f')_{high}, S$	γ

Figure 4: An abstract machine with specialized representations

environment tuple is allocated, but the pair of the code pointer and the environment pointer is unallocated.

We use the convention that all terms are subexpressions of a given closed term a_0 , the whole program. We write a^A to indicate that a was given the type A in the (unique) typing derivation of a_0 , in the empty environment. This annotation is used to determine the size (number of words used in the representation) of data of type A , which determines in turn the position of the first word of a variable x in the run-time environment:

$$\begin{array}{ll}
\text{size}(\text{Int}) &= 1 & \text{pos}_x(\emptyset) &\text{is undefined} \\
\text{size}(\text{Float}) &= 2 & \text{pos}_x(x : A, E) &= 0 \\
\text{size}(A \rightarrow B) &= 2 & \text{pos}_x(y : A, E) &= \text{size}(A) + \text{pos}_x(E) \\
\text{size}(A \times B) &= \text{size}(A) + \text{size}(B) & &
\end{array}$$

The transition function for the corresponding machine is given in figure 4.

With this new evaluation mechanism, the previous example function ($f(x, y) = x + 2y$) executes much more efficiently. The uncurried version takes as argument an unallocated pair of unallocated floating-point numbers, that is, four words on the stack. Intermediate results are held in the stack, without any heap allocation or heap accesses. The final result, an unallocated float, is returned to the caller as two words on top of the stack. The curried version benefits similarly from unallocated floats. In addition, the intermediate closure returned to the caller between the passing of the first and second arguments is not allocated either, but left as two words on the stack (a code pointer, an environment pointer), ready to be applied to the second argument.

2.5 Performance comparisons

Specialized representations lead to less heap allocation and less memory accesses than uniform representations. In the case of uniform representations, multi-word values are always allocated in the heap when created, and reloaded when used, while this is not true in the case of specialized representations.

On the other hand, specialized representations generate more stack or register moves. The savings in heap accesses far outweighs them, except in extreme cases where some data are discarded. For instance, applying a function taking a 10-word argument and returning a constant requires ten stack moves using specialized representations, and only one using uniform representations. This does not happen frequently in actual programming, however.

3 The polymorphic case

3.1 A polymorphic language

We now consider a polymorphic language based on the second-order λ -calculus, as introduced by Girard [7] and independently by Reynolds [16]. At the level of types, we introduce universal quantification, with the intent that a term of type $\forall X. A[X]$ can be used with types $A[B]$ for all types B . At the level of terms, the corresponding elimination construct is application of a term a to a type B , written $a(B)$. The introduction construct is abstraction over a type variable X , written $\Lambda X. A$.

$$\begin{aligned} a & ::= c \mid x \mid \Lambda X. a \mid a(B) \mid \lambda x : A. b \mid b(a) \mid (a, b) \mid a.\mathbf{fst} \mid a.\mathbf{snd} \\ A & ::= X \mid \forall X. A \mid \mathbf{Int} \mid \mathbf{Float} \mid A \rightarrow B \mid A \times B \end{aligned}$$

Second-order λ -calculus is one of the purest and most general approaches to polymorphism, but very few programming languages implement it in its full generality (Poly [11], Quest [4]). The ML language proposes a restricted version of it: it requires that universal quantifiers be in prenex position: that all type expressions are of the form $\forall X_1 \dots \forall X_n. A$, where A does not contain quantifiers. This makes type inference possible, using the well-known Damas-Milner algorithm [6], while type inference for second-order λ -calculus is still an open problem. In the following, we do not need the prenex quantification hypothesis, and therefore consider arbitrary quantification.

Informally, typechecking rules are those of the simply-typed language extended by the following two rules:

$$\frac{E \vdash a : A}{E \vdash \Lambda X. a : \forall X. A} \qquad \frac{E \vdash a : \forall X. A}{E \vdash a(B) : A\{X \leftarrow B\}}$$

The actual rules are slightly more complex, since we must take care of the scope of type variables. This means that not all well-formed type expressions are valid types in a given context, and similarly for environments. We use two auxiliary predicates, $E \vdash A$ type, meaning that A is a valid type in environment E , and $\vdash E$ env, meaning that E is a valid environment.

$$\frac{\vdash E \text{ env}}{E \vdash c : T(c)} \qquad \frac{\vdash E_1, x : A, E_2 \text{ env}}{E_1, x : A, E_2 \vdash x : A}$$

$$\begin{array}{c}
\frac{E \vdash A \text{ type} \quad E, x : A \vdash b : B}{E \vdash \lambda x : A. b : A \rightarrow B} \\
\frac{E, X \text{ type} \vdash a : A}{E \vdash \Lambda X. a : \forall X. A} \\
\frac{E \vdash a : A \quad E \vdash b : B}{E \vdash (a, b) : A \times B} \\
\frac{E \vdash b : A \rightarrow B \quad E \vdash a : A}{E \vdash b(a) : B} \\
\frac{E \vdash a : \forall X. A \quad E \vdash B \text{ type}}{E \vdash a(B) : A\{X \leftarrow B\}} \\
\frac{E \vdash a : A \times B}{E \vdash a.\text{fst} : A \quad E \vdash a.\text{snd} : B}
\end{array}$$

The auxiliary predicates are defined as follows:

$$\begin{array}{c}
\vdash \emptyset \text{ env} \\
\frac{\vdash E \text{ env} \quad X \notin \text{Dom}(E)}{\vdash E, X \text{ type env}} \\
\frac{\vdash E \text{ env}}{E \vdash \text{Int type} \quad E \vdash \text{Float type}} \\
\frac{E, X \text{ type} \vdash A \text{ type}}{E \vdash \forall X. A \text{ type}} \\
\vdash \emptyset \text{ env} \\
\frac{\vdash E \text{ env} \quad x \notin \text{Dom}(E) \quad E \vdash A \text{ type}}{\vdash E, x : A \text{ env}} \\
\frac{\vdash E_1, X \text{ type}, E_2 \text{ env}}{E_1, X \text{ type}, E_2 \vdash X \text{ type}} \\
\frac{E \vdash A \text{ type} \quad E \vdash B \text{ type}}{E \vdash A \rightarrow B \text{ type} \quad E \vdash A \times B \text{ type}}
\end{array}$$

3.2 Evaluation using uniform representations

In the case where all data have uniform, single-word representations, the abstract machine needs no special provision to accommodate polymorphic programs. The evaluation mechanism of section 2.2 already implements polymorphism in some sense, since, for instance, the term $\lambda x:\text{Int}. x$ is compiled in such a way that it can be applied to data of any type, not just `Int`, and return it unmodified. Therefore, we can use the abstract machine of figure 2 as is, along with the compilation scheme of figure 1. We just have to add the two following cases, stating that abstraction on a type variable and application to a type have no computational contents (in this case, they are mere typechecking annotations):

$$\begin{aligned}
\mathcal{U}_E(\Lambda X. a) &= \mathcal{U}_E(a) \\
\mathcal{U}_E(a(B)) &= \mathcal{U}_E(a)
\end{aligned}$$

3.3 Evaluation using specialized representations

Let us try now to implement a polymorphic language with non-uniform data representations, such as multi-word values. Things are not as easy as in the monomorphic case. We already dismissed the approach of compiling several specializations of a polymorphic term on demand. Therefore, when compiling a polymorphic term, we are left with no alternative but to assume that all values of unknown types (that is, whose type is a type variable) are represented in a uniform way, as in the previous section. However, when we have more information on the type of a value, we would like to use specialized representations, for the sake of efficiency. In particular, we hope to compile fully monomorphic terms as efficiently as in the case of the monomorphic calculus.

This requirement implies that every value has two representations, a uniform one, to be used for communication with polymorphic functions, and a specialized, efficient one, possibly spanning several words and taking advantage of special hardware, to be used the rest of the time. We refer to the former as “the wrapped representation”, and to the latter as “the unwrapped representation.” When we apply a polymorphic function to a value of known type, it is likely that the value will be unwrapped, while the function expects it wrapped. Therefore, the compiler will sometimes have to insert explicit coercions between the two representations; we write `wrapA` and `unwrapA` for the coercions operating on values of type *A*. (We mention the type *A* explicitly to emphasize that those coercions are not polymorphic functions operating uniformly on all data, but rather functions defined in an *ad-hoc* way for each type *A*.)

Consider the example of the `reverse_pair` function:

$$\text{reverse_pair} = \Lambda X. \Lambda Y. \lambda z : X \times Y. (z.\text{snd}, z.\text{fst}).$$

Since the types of `z.fst` and `z.snd` are unknown, these values must be wrapped. However, `z` itself is known to be a pair, so there is no need to wrap it. Therefore, the calling conventions of `reverse_pair` is as follows: it expects two words on the stack, which are wrapped representations of `z.fst` and `z.snd`, and returns two words on the stack. Now, let us consider the following application:

$$\text{reverse_pair} (\text{Float}) (\text{Int} \times \text{Float}) (3.14, (7, 2.718)).$$

We assume that the two floating-point constants are allocated “flat”, as well as the two pairs. The argument is therefore represented by the five words:

$$(3.14)_{\text{low}}, (3.14)_{\text{high}}, 7, (2.718)_{\text{low}}, (2.718)_{\text{high}}.$$

Before passing it to `reverse_pair`, some transformations are required: wrap the first component of the pair (the first two words,) resulting in the single word $\langle (3.14)_{\text{low}}, (3.14)_{\text{high}} \rangle$; similarly for the second component, leading to $\langle 7, (2.718)_{\text{low}}, (2.718)_{\text{high}} \rangle$. The resulting two words are a suitable argument for `reverse_pair`. On return, the stack holds the two words:

$$\langle 7, (2.718)_{\text{low}}, (2.718)_{\text{high}} \rangle, \langle (3.14)_{\text{low}}, (3.14)_{\text{high}} \rangle.$$

Two steps of unwrapping lead to the 5-tuple $7, (2.718)_{\text{low}}, (2.718)_{\text{high}}, (3.14)_{\text{low}}, (3.14)_{\text{high}}$, which is the unwrapped representation of $((7, 2.718), 3.14)$, as expected.

Let us now consider an example involving higher-order functions:

$$\begin{aligned} \text{map_pair} &= \Lambda X. \Lambda Y. \lambda f : X \rightarrow Y. \lambda z : X \times X. (f z.\text{fst}, f z.\text{snd}) \\ \text{int_of_float} &: \text{Float} \rightarrow \text{Int} \\ \text{map_pair} (\text{Float}) (\text{Int}) (\text{int_of_float}) &(3.14, 2.718) \end{aligned}$$

According to our principles, the functional `map_pair` expects its parameter *f* to be a function taking one word (a wrapped representation) as argument, and returning one word (another wrapped representation). In addition, the parameter *z* should be an unwrapped pair of two wrapped values, as previously. However, the `int_of_float` primitive function expects an unwrapped floating-point argument (two words), and produces an unwrapped integer. Therefore, `map_pair` cannot be applied directly to `int_of_float`; it must be given a version of `int_of_float` which takes a wrapped

floating-point number as argument, and returns a wrapped integer as result, that is, with obvious notations:

$$\lambda x : \text{Wrapped}(\text{Float}). \text{wrap}_{\text{Float}}(\text{int_of_float}(\text{unwrap}_{\text{Float}}(x))).$$

The rest of this example proceeds as above. The important point is that higher-order functions may require their functional arguments to be transformed in order to accommodate wrapped arguments or results instead of unwrapped ones, and vice-versa. This transformation does not require recompilation of the function. It merely puts some “stub code” around it, performing the right `wrap` and `unwrap` operations.

The rest of this section formalizes a compilation scheme based on the ideas above. This is a two-step process: first, a translation into another polymorphic calculus, where the duality of wrapped/unwrapped representations is taken into account; then, a code generation phase, combining cases from section 2.4 for the unwrapped values and section 2.2 for the wrapped, uniformly represented values.

3.3.1 A restricted polymorphic calculus

First, the distinction between wrapped and unwrapped representations is made explicit in the types, through the introduction of a new type operator, `Wrapped`. Informally, for all types A , the type `Wrapped(A)` contains all wrapped representations of values of type A . At the level of terms, we add the operators `wrapA` and `unwrapA`, which map A to `Wrapped(A)` and conversely. Then, we restrict polymorphism by requesting that type variables range over the class of wrapped types, that is all `Wrapped(A)` where A is a type, instead of the full class of types. By analogy with bounded quantification [5], we use the notation $\forall X \leq \text{Wrapped}. A$ for this restricted universal quantification. Conversely, for type application $a(B)$, we require that B is a wrapped type, that is either `Wrapped(B')`, or a type variable Y . The syntax of the restricted calculus is therefore as follows:

$$\begin{aligned} a &::= c \mid x \mid \text{wrap}_A(a) \mid \text{unwrap}_A(a) \mid \Lambda X \leq \text{Wrapped}. a \\ &\quad \mid a(B) \mid \lambda x : A. b \mid b(a) \mid (a, b) \mid a.\text{fst} \mid a.\text{snd} \\ A &::= X \mid \text{Wrapped}(A) \mid \forall X \leq \text{Wrapped}. A \mid \text{Int} \mid \text{Float} \mid A \rightarrow B \mid A \times B \end{aligned}$$

The typing rules are almost the same as those of section 3.1, with additional rules for the `Wrapped`, `wrap` and `unwrap` operators, and a different treatment of type application. To help distinguishing both calculi, we use \Vdash instead of \vdash for the typing judgements of this calculus. Here are the rules that differ from the one of the general calculus:

$$\begin{array}{c} \frac{E \Vdash a : A}{E \Vdash \text{wrap}_A(a) : \text{Wrapped}(A)} \qquad \frac{E \Vdash a : \text{Wrapped}(A)}{E \Vdash \text{unwrap}_A(a) : A} \\ \\ \frac{E, X \text{ type } \Vdash a : A}{E \Vdash \Lambda X \leq \text{Wrapped}. a : \forall X \leq \text{Wrapped}. A} \\ \\ \frac{E \Vdash a : \forall X \leq \text{Wrapped}. A \quad E \Vdash Y \text{ type}}{E \Vdash a(Y) : A\{X \leftarrow Y\}} \qquad \frac{E \Vdash a : \forall X \leq \text{Wrapped}. A \quad E \Vdash B \text{ type}}{E \Vdash a(\text{Wrapped}(B)) : A\{X \leftarrow B\}} \end{array}$$

$$\frac{E, X \text{ type} \quad E \Vdash A \text{ type}}{E \Vdash \forall X \leq \text{Wrapped}. A \text{ type}} \qquad \frac{E \Vdash A \text{ type}}{E \Vdash \text{Wrapped}(A) \text{ type}}$$

3.3.2 Translation into the restricted calculus

Here, we provide translations for terms and types of the original, polymorphic calculus into the restricted calculus given above. The translation function is written $[\]$. On types, it simply consists in restricting all quantifications, as follows:

$$[\forall X. A] = \forall X \leq \text{Wrapped}. [A]$$

and then it is extended as a congruence over all types. On terms, the translation transforms application to a type, so as to generate the kind of “stub” code needed in the example above. We have:

$$\begin{aligned} [\Lambda X. a] &= \Lambda X \leq \text{Wrapped}. [a] \\ [a^{\forall X. A}(B)] &= \mathcal{T}_{X \leftarrow B}([a](\text{Wrapped}[B])^A) \end{aligned}$$

Again, it is extended as a congruence over all terms. The hard work is done by the twin auxiliary functions $\mathcal{T}_{X \leftarrow B}(a^A)$ and $\overline{\mathcal{T}}_{X \leftarrow B}(a^A)$, defined inductively on A in figure 5. They are responsible for specializing X to B in any term a of type A . The function \mathcal{T} deals with positive occurrences of the coercion, and $\overline{\mathcal{T}}$ with negative occurrences. When $A = X$, we simply unwrap or wrap a . When X is not free in A , and especially when A is a base type, nothing needs to be done. When A is a product type, we specialize recursively its two components, and pair them together. When A is a function type, we build a function that takes an argument x , specializes it recursively, applies a to x , and specializes the result. To specialize the argument x , we switch to the other transformation ($\overline{\mathcal{T}}$ instead of \mathcal{T} and conversely), because of the contravariance of the arrow. The case of a universal type is similar, but easier.

An example of translation is given in figure 6. The translation usually creates many β -redexes, which we reduced on the fly for the sake of readability.

It remains to show that this translation is sensible, in particular, that it preserves semantics. First, we can map terms of the restricted calculus back to the full calculus by erasing the `wrap` and `unwrap` nodes, and the bounded quantifications. It is easy to see that for all a , the translation $[a]$ is mapped back to a term which reduces to a . This means that if we identify the wrapped and unwrapped representations, a and $[a]$ produce the same results. It remains to show that in $[a]$ we do not use a wrapped value when an unwrapped one is expected and conversely. However, such an error is caught by the type system of the restricted calculus. Hence, we just have to prove that the translation of a well-typed term is well-typed. The following lemma expresses the correctness of the auxiliary translation functions with respect to types.

Lemma 1 *Let a be a term, E be an environment of the restricted calculus, and A, B be two types of the full polymorphic calculus.*

- *If $E \Vdash a : [A]\{X \leftarrow \text{Wrapped}[B]\}$, then $E \Vdash \mathcal{T}_{X \leftarrow B}(a^A) : [A]\{X \leftarrow B\}$.*
- *If $E \Vdash a : [A]\{X \leftarrow B\}$, then $E \Vdash \overline{\mathcal{T}}_{X \leftarrow B}(a^A) : [A]\{X \leftarrow \text{Wrapped}[B]\}$.*

$$\begin{aligned}
\mathcal{T}(a^X) &= \text{unwrap}_{[B]}(a) \\
\mathcal{T}(a^Y) &= a \text{ if } X \neq Y \\
\mathcal{T}(a^{\text{Int}}) &= a \\
\mathcal{T}(a^{\text{Float}}) &= a \\
\mathcal{T}(a^{A_1 \rightarrow A_2}) &= \lambda x : [A_1] \{X \leftarrow B\}. \mathcal{T}(a(\overline{\mathcal{T}}(x^{A_1}))^{A_2}) \text{ where } x \text{ not free in } a \\
\mathcal{T}(a^{A_1 \times A_2}) &= (\mathcal{T}(a.\text{fst}^{A_1}), \mathcal{T}(a.\text{snd}^{A_2})) \\
\mathcal{T}(a^{\forall X. A}) &= a \\
\mathcal{T}(a^{\forall Y. A}) &= \Lambda Y \leq \text{Wrapped}. \mathcal{T}(a(Y)^A) \\
\\
\overline{\mathcal{T}}(a^X) &= \text{wrap}_{[B]}(a) \\
\overline{\mathcal{T}}(a^Y) &= a \text{ if } X \neq Y \\
\overline{\mathcal{T}}(a^{\text{Int}}) &= a \\
\overline{\mathcal{T}}(a^{\text{Float}}) &= a \\
\overline{\mathcal{T}}(a^{A_1 \rightarrow A_2}) &= \lambda x : [A_1] \{X \leftarrow \text{Wrapped}[B]\}. \overline{\mathcal{T}}(a(\mathcal{T}(x^{A_1}))^{A_2}) \text{ where } x \text{ not free in } a \\
\overline{\mathcal{T}}(a^{A_1 \times A_2}) &= (\overline{\mathcal{T}}(a.\text{fst}^{A_1}), \overline{\mathcal{T}}(a.\text{snd}^{A_2})) \\
\overline{\mathcal{T}}(a^{\forall X. A}) &= a \\
\overline{\mathcal{T}}(a^{\forall Y. A}) &= \Lambda Y \leq \text{Wrapped}. \overline{\mathcal{T}}(a(Y)^A)
\end{aligned}$$

Figure 5: Generation of stub code to accommodate restricted polymorphism. ($\mathcal{T}(a^A)$ and $\overline{\mathcal{T}}(a^A)$ abbreviate $\mathcal{T}_{X \leftarrow B}(a^A)$ and $\overline{\mathcal{T}}_{X \leftarrow B}(a^A)$, respectively.)

Proof: The proof is by induction on A . Since X and B are fixed, we abbreviate $\mathcal{T}_{X \leftarrow B}(a^A)$ and $\overline{\mathcal{T}}_{X \leftarrow B}(a^A)$ as $\mathcal{T}(a^A)$ and $\overline{\mathcal{T}}(a^A)$, respectively. The cases $A = \text{Int}$, $A = \text{Float}$, $A = Y \neq X$ are trivial. In the case $A = X$, the propositions reduce to

$$\frac{E \Vdash a : \text{Wrapped}([B])}{E \Vdash \text{unwrap}_{[B]}(a) : [B]} \qquad \frac{E \Vdash a : [B]}{E \Vdash \text{wrap}_{[B]}(a) : \text{Wrapped}([B])}$$

which are special cases of the typing rules for `wrap` and `unwrap`. The cases $A = A_1 \times A_2$ and $A = \forall Y. A'$ are straightforward by induction hypothesis. Finally, in the case $A = A_1 \rightarrow A_2$, we can

$$\begin{aligned}
& \mathcal{T}_{X \leftarrow \text{Float}}(\text{double}^{(X \rightarrow X) \rightarrow (X \rightarrow X)}) \\
&= \lambda f : \text{Float} \rightarrow \text{Float}. \mathcal{T}(\text{double}(\overline{\mathcal{T}}(f^{X \rightarrow X}))^{X \rightarrow X}) \\
&= \lambda f : \text{Float} \rightarrow \text{Float}. \lambda x : \text{Float}. \\
&\quad \text{unwrap}_{\text{Float}}(\text{double}(\lambda y : \text{Wrapped}(\text{Float}). \text{wrap}_{\text{Float}}(f(\text{unwrap}_{\text{Float}}(y)))) \\
&\quad\quad\quad (\text{wrap}_{\text{Float}}(x)))
\end{aligned}$$

Hence:

$$\begin{aligned}
& [\text{double}^{\forall X. (X \rightarrow X) \rightarrow (X \rightarrow X)}(\text{Float})(\lambda x : \text{Float}. \text{add_float}(x)(1))(3.14)] \\
& \xrightarrow{*} \text{unwrap}_{\text{Float}}(\text{double}(\lambda y : \text{Wrapped}(\text{Float}). \text{wrap}_{\text{Float}}(\text{add_float}(\text{unwrap}_{\text{Float}}(y))(1))) \\
&\quad\quad\quad (\text{wrap}_{\text{Float}}(3.14)))
\end{aligned}$$

Figure 6: An example of translation

build the following derivation:

$$\begin{array}{c}
E, x : [A_1\{X \leftarrow B\}] \Vdash x : [A_1\{X \leftarrow B\}] \\
\hline\hline
\begin{array}{c}
E, x : [A_1\{X \leftarrow B\}] \Vdash \overline{\mathcal{T}}(:^{A_1})[A_1]\{X \leftarrow \text{Wrapped}[B]\} \\
a : [A_1 \rightarrow A_2]\{X \leftarrow \text{Wrapped}[B]\}
\end{array} \\
\hline
E, x : [A_1\{X \leftarrow B\}] \Vdash a(\overline{\mathcal{T}}(x^{A_1})) : [A_2]\{X \leftarrow \text{Wrapped}[B]\} \\
\hline\hline
E, x : [A_1\{X \leftarrow B\}] \Vdash \mathcal{T}(a(\overline{\mathcal{T}}(x^{A_1}))^{A_2}) : [A_2\{X \leftarrow B\}] \\
\hline
E \Vdash \lambda x : [A_1\{X \leftarrow B\}]. \mathcal{T}(a(\overline{\mathcal{T}}(x^{A_1}))^{A_2}) : [(A_1 \rightarrow A_2)\{X \leftarrow B\}]
\end{array}$$

that leads to the first expected result $E \Vdash \mathcal{T}(a^A) : [A\{X \leftarrow B\}]$. Here, double horizontal bars correspond to the use of the induction hypothesis with $A = A_1$ or $A = A_2$. A symmetric derivation, where \mathcal{T} and $\overline{\mathcal{T}}$ are exchanged, leads to the second expected result. \square

Proposition 1 *Let a, A, E be a term, a type, and an environment of the full calculus. If $E \vdash a : A$, then in the restricted calculus, $[E] \Vdash [a] : [A]$.*

Proof: A trivial inductive argument on the lengths of the proofs shows that $\vdash E$ env implies $\Vdash [E]$ env, and $E \vdash A$ type implies $[E] \Vdash [A]$ type. Then, the proposition above can be proved by induction on the length of the proof of $E \vdash a : A$. The sole interesting case is the one of application to a type. Then, the original proof is:

$$\frac{\begin{array}{c} \vdots \\ E \vdash a : \forall X. A \end{array} \quad \begin{array}{c} \vdots \\ E \vdash B \text{ type} \end{array}}{E \vdash a(B) : A\{X \leftarrow B\}}$$

	Code	Stack	Environment
Before	Box _{<i>i</i>} , <i>C</i>	v_1, \dots, v_i, S	γ
After	<i>C</i>	$\langle v_1, \dots, v_i \rangle, S$	γ
Before	Unbox _{<i>i</i>} , <i>C</i>	$\langle v_1, \dots, v_i \rangle, S$	γ
After	<i>C</i>	v_1, \dots, v_i, S	γ

Figure 7: Additional instructions for wrapping and unwrapping.

Using the induction hypothesis and the previous lemma, we get the following derivation:

$$\begin{array}{c}
\frac{[E] \Vdash [a] : \forall X. [A] \quad [E] \Vdash [B] \text{ type}}{[E] \Vdash [a](\text{Wrapped}[B]) : [A]\{X \leftarrow \text{Wrapped}[B]\}} \\
\hline
\hline
[E] \Vdash \mathcal{T}_{X \leftarrow B}([a](\text{Wrapped}[B])^A) : [A]\{X \leftarrow B\}
\end{array}$$

which leads to the desired result. \square

3.3.3 Compiling the restricted calculus

To compile the restricted calculus, we use the same abstract machine with specialized representations as in section 2.4 (figure 4). The compilation scheme is almost the same. We simply state that the size of a value of type `Wrapped(A)` is always one, and give two additional rules for the translation of `wrap` and `unwrap`. The only constraints are that `unwrapA` must be the inverse of `wrapA`, and `wrapA` must produce one-word data. A simple solution for `wrapA` is to heap-allocate data occupying more than one word, and keep one-word data unchanged; symmetrically, `unwrapA` performs nothing if `size(A) = 1`, and dereferences the value of `a` otherwise. We add the two corresponding instructions `Box` and `Unbox` in the abstract machine (figure 7). The resulting compilation scheme is given in figure 8, with the `size` function defined as follows:

$$\begin{aligned}
\text{size}(\text{Int}) &= 1 \\
\text{size}(\text{Float}) &= 2 \\
\text{size}(A \rightarrow B) &= 2 \\
\text{size}(A \times B) &= \text{size}(A) + \text{size}(B) \\
\text{size}(\text{Wrapped}(A)) &= 1
\end{aligned}$$

3.4 Performance comparison

As in the case of the monomorphic calculus, specialized representations lead to less heap allocation and pointer dereferencing, at the expense of more stack or register moves. In the case of uniform representations, each primitive performs implicit `unwrap` operations on its arguments, and `wrap` operations on its results. Using specialized representations, we managed to remove some `wrap` and `unwrap` operations when we have enough type information.

$$\begin{aligned}
\mathcal{P}_E(i^{\text{Int}}) &= \text{Const}(i) \\
\mathcal{P}_E(f^{\text{Float}}) &= \text{Const}(f_{\text{high}}), \text{Const}(f_{\text{low}}) \\
\mathcal{P}_E(x^A) &= \text{Access}(p + s - 1), \dots, \text{Access}(p) \\
&\quad \text{where } p = \text{pos}_x(E) \text{ and } s = \text{size}(A) \\
\mathcal{P}_E(\text{wrap}_A(a)) &= \mathcal{P}_E(a) \text{ if } \text{size}(A) = 1 \\
\mathcal{P}_E(\text{wrap}_A(a)) &= \mathcal{P}_E(a), \text{Box}_{\text{size}(A)} \text{ if } \text{size}(A) > 1 \\
\mathcal{P}_E(\text{unwrap}_A(a)) &= \mathcal{P}_E(a) \text{ if } \text{size}(A) = 1 \\
\mathcal{P}_E(\text{unwrap}_A(a)) &= \mathcal{P}_E(a), \text{Unbox}_{\text{size}(A)} \text{ if } \text{size}(A) > 1 \\
\mathcal{P}_E(\lambda x : A. b^B) &= \text{Abstr}(\mathcal{P}_{x:A,E}(b^B), \text{Return}_{\text{size}(B)}) \\
\mathcal{P}_E(b(a^A)) &= \mathcal{P}_E(b), \mathcal{P}_E(a), \text{Apply}_{\text{size}(A)} \\
\mathcal{P}_E(\Lambda X. a) &= \mathcal{P}_E(a) \\
\mathcal{P}_E(a(B)) &= \mathcal{P}_E(a) \\
\mathcal{P}_E(a, b) &= \mathcal{P}_E(a), \mathcal{P}_E(b) \\
\mathcal{P}_E(a^{A \times B}. \text{fst}) &= \mathcal{P}_E(a), \text{First}_{\text{size}(A), \text{size}(B)} \\
\mathcal{P}_E(a^{A \times B}. \text{snd}) &= \mathcal{P}_E(a), \text{Second}_{\text{size}(A), \text{size}(B)} \\
\mathcal{P}_E(\text{add_float}(a, b)) &= \mathcal{P}_E(a), \mathcal{P}_E(b), \text{AddFloat} \\
\mathcal{P}_E(\text{add_float}) &= \text{Abstr}(\text{Access}(0), \text{Access}(1), \text{Access}(2), \text{Access}(3), \\
&\quad \text{AddFloat}, \text{Return}_2)
\end{aligned}$$

Figure 8: Compilation scheme for the restricted polymorphic calculus, mixing uniform and specialized representations.

It is hard to give more precise comparisons, due to the following fact: in the case of uniform representation, wrapping and unwrapping are performed by the primitive operations, therefore unwrapping (on the arguments) is delayed as much as possible, while wrapping (of the results) is performed immediately. We say that with uniform representations, unwrapping is lazy, while wrapping is eager. In case of specialized representations, wrapping is delayed until a value is passed to a polymorphic function, while unwrapping the result of a polymorphic function is performed as soon as possible. Hence, in the case of specialized representations, unwrapping is eager and wrapping is lazy. This leads to a rather different behavior, presumably favorable, since wrapping is much more expensive than unwrapping.

On the negative side, it is true that the stub code inserted to change representations in the case of functions introduces additional function calls. Reductions at compile-time can often eliminate them, but for example if the function being transformed is a parameter of a functional, an additional call will remain.

4 Toward actual programming languages

The highly stylized calculi used in previous sections are not yet close to actual programming languages. They lack many important features: the only data structures are pairs, there are neither variants, nor general records; and recursive types are not allowed. In this section, we shall see how to integrate these features into the calculus of the previous section, and what additional constraints they put on the choice of data representations. Finally, we discuss the application to the ML language.

4.1 Sum types

While the previous calculi have a rudimentary notion of records, they have no notion of the symmetric data structure: variants (tagged unions). In typed calculi, variants are traditionally presented by sum types $A + B$. The elements of that type are either of type A , with tag `left`, or of type B , with tag `right`. The introduction constructs are the two constructors `inleft` and `inright`; the elimination construct is pattern-matching on the tag.

$$\frac{E \vdash a : A \quad E \vdash B \text{ type}}{E \vdash \text{inleft}(a) : A + B} \qquad \frac{E \vdash A \text{ type} \quad E \vdash b : B}{E \vdash \text{inright}(b) : A + B}$$

$$\frac{E \vdash a : A + B \quad E, x : A \vdash c : C \quad E, y : B \vdash d : C}{E \vdash \text{case } a \text{ of } \text{inleft}(x) \rightarrow c \mid \text{inright}(y) \rightarrow d : C}$$

Values of a sum type $A + B$ are usually represented as (tag, value) pairs. For instance, `inleft(a)` is represented as $(0, a)$, and `inright(b)` as $(1, b)$. Pattern-matching is then a simple test on the first component of the pair. Such pairs are not regular pairs, since the type of the second component depends on the value of the first one (A if it is 0, B if it is 1). This is no problem if uniform representations are used. In case of specialized representations, however, it may be the case that an element of type A and an element of type B have incompatible representations (e.g. different sizes); therefore, elements of type $A + B$ would have two incompatible representations, depending on the value of the tag, which is not known at compile-time.

A first solution is to take a representation compatible with both the one of A and the one of B . For instance, regarding the size of $A + B$, we could take

$$\text{size}(A + B) = 1 + \max(\text{size}(A), \text{size}(B))$$

so that the second component of the dependent pair is always large enough to contain an element of A , or an element of B . This is the traditional implementation of variants in Pascal. The main drawback is some waste of space when $\text{size}(A)$ is not equal to $\text{size}(B)$. This becomes serious when the sums are heap-allocated, as parts of a large data structure.

Another solution is to systematically use uniform representations for the argument of a sum constructor. That way, a sum is always represented by one word for the tag, and one word for the argument, hence $\text{size}(A + B) = 2$. The necessary wrapping operations are performed by the constructors `inleft` and `inright`.

	A recursive, B recursive,	A recursive, B not recursive,	A not recursive, B recursive,	A not recursive, B not recursive
$[A \times B]$	$\text{Wrapped}[A] \times \text{Wrapped}[B]$	$\text{Wrapped}[A] \times [B]$	$[A] \times \text{Wrapped}[B]$	$[A] \times [B]$
$[a^A, b^B]$	$(\text{wrap}_A[a], \text{wrap}_B[b])$	$(\text{wrap}_A[a], [b])$	$([a], \text{wrap}_B[b])$	$([a], [b])$
$[a^{A \times B}.\text{fst}]$	$\text{unwrap}_A([a].\text{fst})$	$\text{unwrap}_A([a].\text{fst})$	$[a].\text{fst}$	$[a].\text{fst}$
$[a^{A \times B}.\text{snd}]$	$\text{unwrap}_B([a].\text{snd})$	$[a].\text{snd}$	$\text{unwrap}_B([a].\text{snd})$	$[a].\text{snd}$

Figure 9: Translation from general polymorphism to restricted polymorphism, in the presence of recursive types

Both approaches are compatible with polymorphism. The compilation scheme given in section 3.3 easily extends to sums, with two additional cases for the generation of stub code:

$$\mathcal{T}(a^{A_1+A_2}) = \text{case } a \text{ of } \begin{array}{l} \text{inleft}(x) \rightarrow \text{inleft}(\mathcal{T}(x^{A_1})) \\ | \text{inright}(y) \rightarrow \text{inright}(\mathcal{T}(y^{A_2})) \end{array}$$

and similarly for $\overline{\mathcal{T}}(a^{A_1+A_2})$.

4.2 Recursive types

Recursive types describe recursive data structures. In actual programming languages, they are usually introduced through the ability to name types and to refer to that name in the definition of the associated type. In a typed calculus, they arise when type expressions are not restricted to be finite trees any more, but allowed to be rational trees. We use the notation $\mu X. A$ for cycles; anywhere in a type expression, $\mu X. A$ can be replaced by $A\{X \leftarrow \mu X. A\}$, and conversely. For instance, the type of lists of integers is $\mu L. \text{Unit} + \text{Int} \times L$, where Unit is a special type containing exactly one value.

No other modifications of the type system are needed for recursive types. In particular, they introduce no additional typing rules. The main difficulty with recursive types is that the definitions we gave by induction on types are now possibly ill-founded: they should now read as a set of recursive equations, which may have no solution. For instance, with a naive *size* function, the size s of integer lists, defined as $\mu L. \text{Unit} + \text{Int} \times L$, must verify $s = 1 + \max(0, 1 + s)$, and this equation has no positive solution. This corresponds to the well-known fact that lists cannot be statically allocated “flat”; pointers must be used somewhere. Indeed, if we define integer lists as $\mu L. \text{Unit} + \text{Int} \times \text{Wrapped}(L)$, the size s' of this type is now $s' = 1 + \max(0, 1 + 1) = 3$, as expected.

To guarantee that all recursive types can be represented, it is necessary to introduce additional pointers in the representations. For instance, we may require that a value of a recursive type must be heap-allocated, and represented by a pointer, if it is a component of a structure (product or sum). This treatment can be integrated into the translation from general polymorphism to restricted polymorphism, as shown in figure 9.

Similarly, the stub code generated during the translation from general polymorphism to restricted polymorphism may now involve recursive functions. For instance, the function $f(a) =$

$\mathcal{T}_{X \leftarrow \text{Float}}(a^{\mu Y \cdot \text{Unit} + X \times Y})$, which specializes generic lists to lists of floats, is defined as:

$$f(a) = \text{case } a \text{ of } \begin{array}{l} \text{inleft}(n) \rightarrow \text{inleft}(n) \\ | \text{inright}(p) \rightarrow \text{inright}(\text{unwrap}_{\text{Float}}(p.\text{fst}), f(p.\text{snd})) \end{array}$$

That is, we should copy the list while applying the transformation $\text{unwrap}_{\text{Float}}$ to each of its elements. This is not practical, since this copying takes time and space proportional to the size of the original list, and is not correct when lists can be physically updated. Instead, we restrict further the representations of sums and products, and require that their components be systematically wrapped. This ensures that data structures never need to be recursively copied when converting between uniform and specialized representations. In particular, we may now take:

$$\begin{aligned} \mathcal{T}_{X \leftarrow B}(a^{A_1 \times A_2}) &= \overline{\mathcal{T}}_{X \leftarrow B}(a^{A_1 \times A_2}) = a \\ \mathcal{T}_{X \leftarrow B}(a^{A_1 + A_2}) &= \overline{\mathcal{T}}_{X \leftarrow B}(a^{A_1 + A_2}) = a \end{aligned}$$

We already came to that solution in the case of sums. For products, it may look overly restrictive. For instance, a pair of floating-point numbers cannot be represented “flat” any more, and requires heap-allocating both numbers. We shall see now how the introduction of general records instead of mere pairs alleviates this problem.

4.3 Records and variants

Actual programming languages, such as ML, provide general records and variants as data structures (“concrete types”, in ML terminology), not only binary sums and binary products. Sums and products allow better naming of the components of a data structure, as well as more precise control over their types. For instance, coordinates in the plane could be described by the following record type:

```
type coord2 = {x: Float; y: Float}
```

which is more precise than $\text{Float} \times \text{Float}$. In particular, no polymorphic function can destructure this record. Therefore, it can be represented efficiently by four words allocated “flat”, without ever having to copy it with its components wrapped.

However, it is still possible to define polymorphic data structures by parameterizing a concrete type declaration, using type variables. For instance, binary polymorphic sum and product can be defined in the language as follows:

```
type X * Y = { fst: X; snd: Y };
type X + Y = inleft of X | inright of Y;
```

As we saw previously the components of such structures must use wrapped representations, in order to avoid copying. A simple, natural way to ensure this, while preserving efficient representations for types such as `coord2`, is to say that the representation of a concrete type is chosen once for all when it is defined. Components whose type is a type variable are, as usual, represented in the uniform way. In other terms, parameterized concrete types have one representation, whatever the parameters are: $\text{List}(\text{Int})$ and $\text{List}(\text{Float})$ and $\text{List}(X)$ where X is a type variable all have compatible representations. This ensures that no recursive copying of data structures is ever needed.

The necessary wrapping and unwrapping are performed by the primitives for creation and access. To generate the right conversions, no special rules are needed: it suffices to treat these

primitives as regular functions, possibly polymorphic. For instance, the projection `.x` for the type `coord2` can be given the type `coord2 → Float`. Similarly, `.fst` has type $\forall X. \forall Y. X \times Y \rightarrow X$, and `inleft` has type $\forall X. \forall Y. X \rightarrow X + Y$. When these constructors are specialized to particular values of X and Y , the translation technique of section 3.3 automatically inserts the right coercions.

4.4 Application to ML

We have already covered the main features of ML. They can be integrated into compilation using specialized representations, at the cost of additional constraints on the representations. But it is sometimes necessary to revert to uniform representations. One may fear that these constraints are so restrictive that there is little benefit from specialized representations. This is not the case, however, as specialized representations address important weaknesses of ML implementations.

First of all, specialized representations allow efficient operations on base types which do not fit in a word, and in particular on floating-point numbers. This is an absolute requirement for any general-purpose programming language.

Second, though in ML all functions have exactly one argument (and one result), specialized representations makes it possible to pass a tuple of arguments to a function without heap-allocating the tuple, but simply by putting its components on the stack or in registers. Thus we get the efficiency of functions with several arguments (and several results as well), without having to modify the source language, nor restrict polymorphism in any way.

More generally, specialized representations allow monomorphic programs to be compiled just as efficiently as in traditional monomorphic languages, such as Modula for instance. The price to pay for polymorphism is paid only by those programs that actually use it. This contrasts with a tradition among high-level languages: that advanced features systematically hamper the performance of programs, even if they are not used.

5 Relation with garbage collection

This short section outlines the interaction between data representation and garbage collection. Independently of the algorithm used, garbage collection requires the ability to walk (parts of) the memory graph. The collector must know the roots of the graph, that is, all statically-allocated data containing pointers into the heap (e.g. registers, global variables, and the stack). Further, given an object allocated in the heap, it must be able to find all pointers into the heap contained in that object. Data representations must be designed to permit this memory graph walk.

5.1 Using tagged representations

In the case of dynamically-typed languages (Lisp, Smalltalk), interfacing with a garbage collector is straightforward, since these languages already require all data to carry their own type at run-time, in order to perform dynamic type checking, and to implement primitive operations such as the Lisp predicates `consp` and `integerp`. This is achieved through some kind of tagging scheme. Here is a simple one [8]: a value is a 32-bit word, representing either an integer or a pointer in the heap. The integer i is represented by the bit-field $2i + 1$, hence the low-order bit of an integer is 1. Pointers are word-aligned, hence their low-order bit is 0. They point to a one-word header, containing the

size of the block, and one bit to say whether the block contains tagged values, possibly pointing to other blocks, or unstructured data such as the characters of a string.

Tagging has several drawbacks on a conventional architecture. Tag bits waste space: integers are limited to 31 bits, and one extra word is needed for each block in the heap. Tagging also involves additional computation. Using the scheme above, multiplying two integers means shifting one to the right, decrementing the other, multiplying as usual both intermediate results, and increment the result.

5.2 Using type information and untagged representations

Static monomorphic typing makes tagging unnecessary, since the information whether a value is a pointer or not is already contained in its type, and in case of a pointer, the size of the block pointed to, and where pointers are located in the block as well. To communicate this information to the garbage collector, the compiler can save in the object file the structure of user-defined types, and associate to each global variable (a code for) its type, and to each function the type and location of its local variables. Then, using the stack frames to find active functions, the collector can determine the type of each root of the memory graph, use it to find which are pointers and which are not, find the types of all fields of the blocks pointed to, and recurse.

However, propagating type information at GC-time may be costly, especially if the graph walk is not depth-first. Therefore, in practice, blocks in the heap carry a header containing their type. This header makes it possible to find the size of a block, and all the pointers it contains, without having to propagate types. This header is added when the block is created, as in the tagged model. However, no other tagging is needed on values. For instance, integers are 32-bit wide, and native machine operations are used on them, since they cannot be confused with pointers, thanks to static typing. This technique is used in some implementations of Modula-3 [13].

5.3 The polymorphic case

The addition of polymorphism reintroduces the problem of identifying pointers at GC-time: the argument of a polymorphic function such as the identity function may be an integer as well as a pointer, and therefore it seems necessary again to tag pointers and integers so that they cannot have the same bit pattern. As a consequence, ML implementations usually use Lisp-like tagging, just for the purpose of garbage collection; for instance, the tagging scheme described above is used in Standard ML of New Jersey [2].

As pointed out by Appel [1], this can be avoided: the actual type of a polymorphic datum (such as the parameter of the identity function) can be determined at GC-time by looking at the call graph, and propagating actual types down it, using compile-time type information. But as in the monomorphic case, types must then be propagated down the memory graph, and this is unrealistic (Appel didn't implement it, either). Of course, it is not possible to precompute the types of the blocks, as in the monomorphic case, since these types are not totally known at compile-time. For instance, the function $\Lambda X. \lambda x : X. (x, x)$ builds a pair without knowing whether its components are pointers or unboxed integers.

In the spirit of section 3.3, my approach is to mix tagged and untagged representations, and reserve the tagged ones for values whose types are not statically known. It suffices to require that uniform representations are not only one-word wide, but also tagged so as to distinguish pointers from unallocated values. In this model, if the type of a value is a type variable, then the value

is correctly tagged, and a test on the tag shows whether it is a pointer or not. Otherwise, the type is not a type variable, and then it contains enough information to guarantee that the value is a pointer or not, without having to test the value at GC-time, and hence, to tag it. Conversely, blocks are again adorned with a header containing their (static) type. This type may contain type variables, but the transformations of section 3.3 ensures that the fields whose types are statically unknown are uniformly represented, hence tagged. In any case, it is possible to find all pointers contained in a block, without having to propagate type information.

Tagging and untagging are performed when converting from specialized to uniform representations, and vice-versa, that is, by the `wrap` and `unwrap` operations. Using the tagging scheme described above, we simply have to redefine `wrapInt` and `unwrapInt`, so that they convert from i to $2i + 1$ and conversely:

$$\begin{aligned} \mathcal{P}_E(\text{wrap}_{\text{Int}}(a)) &= \mathcal{P}_E(a, \text{Const}(1), \text{ShiftLeftInt}, \text{Const}(1), \text{AddInt}) \\ \mathcal{P}_E(\text{unwrap}_{\text{Int}}(a)) &= \mathcal{P}_E(a, \text{Const}(1), \text{ShiftRightInt}) \end{aligned}$$

If $\text{size}(A) > 1$, `wrapA` allocates its argument in the heap and returns a pointer to it, already well-tagged. The remaining case is $\text{size}(A) = 1$ and $A \neq \text{Int}$; in this case, `wrapA` does nothing, but this is correct, since the only possibility is $A = \text{Wrapped}(B)$, and by hypothesis, values of type `Wrapped(B)` are always well-tagged.

6 Conclusions and related work

The techniques presented in this paper combine the cleanliness and expressiveness of unrestricted polymorphism with the efficiency of specialized representations. We developed one possible application, getting better implementation of an existing polymorphic language, ML. The dual application is to add full polymorphism to a conventional, Algol-like language while keeping efficient compilation of existing programs. One may hope that both approaches would converge toward languages with powerful type systems and efficient implementations, without putting strong constraints on the type system (for instance, abstract data types, or records with subtyping are easily accommodated), nor on the runtime system (an evidence is that interfacing with a garbage collector is still possible).

Another interesting feature of our technique is to facilitate interfacing with existing code written in another language, such as libraries written in C: adopting unwrapped representations which are compatible with the ones of C (at least for the base types) alleviates the need for coercions between C and ML data formats.

There is comparatively little work on the problem of data representations in high-level languages. The fact that some typing is necessary to implement efficiently floating-point numbers and records is demonstrated in the design of many conventional imperative languages (such as C [9]), but is scarcely ever stated explicitly. In the area of Lisp compiling, tagging is known to be a performance bottleneck, and several attempts were made to avoid it, at least locally. A systematic approach to this problem can be found in a recent paper by Peterson [14]. He considers mixing tagged and untagged representations, in the setting of an untyped language, and focuses on finding an optimal mix of representations, one that minimizes total execution time. However, his analysis is expensive, and does not consider higher-order functions. Relying on type information, as presented herein, is certainly sub-optimal, but much more practical.

References

- [1] Andrew Appel. Run-time tags aren't necessary. *Lisp and Symbolic Computation*, 2(2):153–162, June 1989.
- [2] Andrew Appel. A runtime system. Technical report, Princeton University, 1989.
- [3] Luca Cardelli. The Functional Abstract Machine. *Polymorphism*, 1(1), 1983.
- [4] Luca Cardelli. Typeful programming. Technical Report 45, DEC Systems Research Center, 1989.
- [5] Luca Cardelli and Peter Wegner. On understanding types, data abstraction, and polymorphism. *Computing surveys*, 17(4), 1985.
- [6] Luis Damas and Robin Milner. Principal type-schemas for functional programs. In *Proc. Symp. Principles of Programming Languages*, 1982.
- [7] Jean-Yves Girard. *Interprétation fonctionnelle et élimination des coupures de l'arithmétique d'ordre supérieur*. Thèse d'Etat, Université Paris VII, 1972.
- [8] Adele Goldberg and David Robson. *Smalltalk-80: the language and its implementation*. Addison-Wesley, 1983.
- [9] Brian W. Kernighan and Dennis M. Ritchie. *The C programming language*. Addison-Wesley, second edition 1988.
- [10] P. J. Landin. The mechanical evaluation of expressions. *The Computer Journal*, pages 308–320, 1964.
- [11] David C. J. Matthews. Poly manual. Technical Report 63, Computer Laboratory, University of Cambridge, 1985.
- [12] Robin Milner, Mads Tofte, and Robert Harper. *The definition of Standard ML*. The MIT Press, 1990.
- [13] Greg Nelson, Luca Cardelli, James Donahue, Lucille Glassman, Mick Jordan, and Bill Kalsow. Modula-3 report (revised). Technical Report 52, DEC Systems Research Center, 1989.
- [14] John Peterson. Untagged data in tagged environments: choosing optimal representations at compile-time. In *Functional Programming Languages and Computer Architecture*, 1989.
- [15] Gordon D. Plotkin. A structural approach to operational semantics. Technical Report DAIMI FN-19, Aarhus University, 1981.
- [16] J. C. Reynolds. Toward a theory of type structure. In *Colloquium on Programming*, volume 19 of *Lecture Notes in Computer Science*. Springer-Verlag, 1974.