

## DATA TYPES AS OBJECTS

Adi Shamir  
Mathematics Dept.  
MIT  
Cambridge, Mass. USA

William W. Wadge  
Computer Science Dept.  
University of Warwick  
Coventry, England

### Abstract

In this paper we present a new approach to the semantics of data types, in which the types themselves are incorporated as elements in the domain of data objects. The approach allows types to have subtypes, allows genuinely polymorphic functions, and gives a precise semantics for recursive type definitions (including definitions with parameters). In addition, the approach yields simple and straight forward methods for proving type properties of recursive definitions. These methods include a new fixedpoint rule which permits case analysis.

### 0. Introduction

Most current type systems (e.g. that of LCF, as in [1]) are based on [2], Church's 1940 paper on a simple theory of types. In these systems there are a number of unrelated ground data types and an arrow operation for forming function types of finite level.

A serious drawback of such systems is the requirement that each data object must have a unique type. Thus integers, for example, cannot be considered as a special case of real numbers, and in general one type cannot be a subtype of another.

A second problem, related to the first, is that polymorphism is not possible; each argument of a function must be of some specified type. For example, we cannot have a single addition operation capable of adding both integers and reals; instead, we need four functions of various types for the various possible combinations. Especially serious is the lack of a general if-then-else conditional whose domains are left unspecified.

The third problem is that in a Church system there are in fact two 'meta'-types, namely types and data objects, which cannot be mixed (LCF

has variables of both meta-types). For example, the natural identity

$$\underline{\text{eveninteger}} + 1 = \underline{\text{oddinteger}}$$

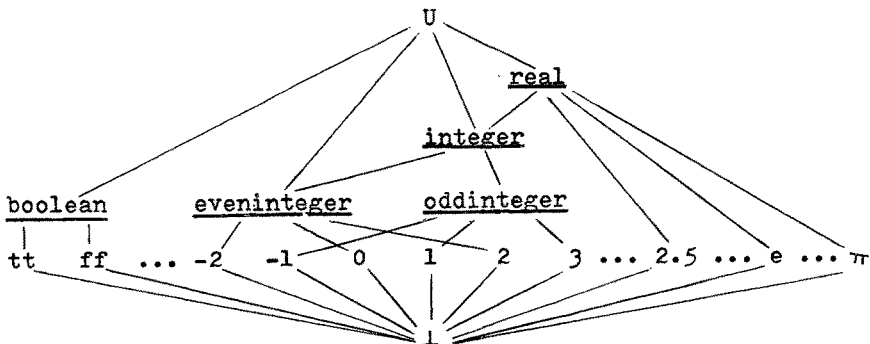
relating the types eveninteger and oddinteger and the data object 1 makes no sense in a Church system (except, perhaps, informally). Even unmixed recursive type equations are problematic because the types do not constitute a domain.

Our approach is to incorporate all the objects and data types together in a single unified domain. Any element of the resulting domain then has two roles:

- (i) it is a data object, at which functions can be defined, including of course functions which are least fixedpoints of recursive definitions;
- (ii) it is the type of all objects which approximate it; the assertions "  $x \sqsubseteq y$  ", "  $x$  is of type  $y$  " and "any object of type  $x$  is of type  $y$  " are therefore all equivalent.

Given an initial domain  $D$  of data objects, and a collection of intuitive data types, we form a new domain  $\hat{D}$  (which we will call a type extension of  $D$ ) by adding the types as new data objects in the sense of (i). The extended  $\sqsubseteq$  relation on  $\hat{D}$  is determined by the equivalence in part (ii): a type object is placed above the data objects of the type in question, and above all its subtypes (so that  $\sqsubseteq$  simultaneously orders types by inclusion and data objects by approximation).

For example, if our original (untyped) domain contains the real numbers, truth values and the undefined element  $\perp$  (ordered as a flat domain), and if we are interested in the types integer, eveninteger, oddinteger, real, boolean and the universal type  $U$ , we get the following extended domain:



Our system, unlike most others, does not associate particular types with elements of  $\hat{D}$ , and so a question like "what is the type of 2?" is meaningless in our system. The relation  $\sqsubseteq$  allows us to answer only questions of the form "is x of type y?". For example, the assertions "2 is of type integer", "2 is of type real", "integer is of type real" and "1 is of type 2" are all true in the  $\hat{D}$  diagrammed above. Note that the least element 1 is of all types, and that everything is of the universal type U.

Monotonic functions and operations over the original domain can be extended to monotonic functions and operations in many ways. Of particular interest are the tightly extended functions (i.e. the least monotonic extensions) whose definitions are usually clear. For example, the tight extension of addition (with  $\hat{D}$  as above) yields:

$$\begin{aligned} 2 + 3 &= 5 \\ \text{integer} + 3 &= \text{integer} \\ \text{integer} + \text{real} &= \text{real} \\ \text{integer} + \text{ } &= \text{real} \\ \text{oddinteger} + \text{oddinteger} &= \text{eveninteger} \end{aligned}$$

Another important example is the tightly extended if-then-else conditional. For any two objects  $x_1$  and  $x_2$  in D we have:

$$\begin{aligned} \text{if } \perp \text{ then } x_1 \text{ else } x_2 &= \perp \\ \text{if } \text{tt} \text{ then } x_1 \text{ else } x_2 &= x_1 \\ \text{if } \text{ff} \text{ then } x_1 \text{ else } x_2 &= x_2 \\ \text{if } \text{boolean} \text{ then } x_1 \text{ else } x_2 &= x_1 \sqcup x_2 \end{aligned}$$

where  $x_1 \sqcup x_2$  is the least upper bound of  $x_1$  and  $x_2$  in  $\hat{D}$ . For example,

$$\text{if } \text{boolean} \text{ then } 3 \text{ else } \text{eveninteger} = \text{integer}$$

in the  $\hat{D}$  diagrammed above.

Certain functions in  $[\hat{D} \rightarrow \hat{D}]$  (the space of monotonic functions from  $\hat{D}$  to  $\hat{D}$  with the standard ordering) play the same role as the type objects added to D, i.e. they embody intuitive types. Given x and y in  $\hat{D}$ , we define the arrow function  $x \rightarrow y$  to be the following element of  $[\hat{D} \rightarrow \hat{D}]$ :

$$\lambda z \text{ if } z \sqsubseteq x \text{ then } y \text{ else } U .$$

This function maps objects of type x into the object y itself, and all other objects into U. It is easily verified that a function h in  $[\hat{D} \rightarrow \hat{D}]$  lies below  $x \rightarrow y$  iff h applied to an object below x yields a result below y. Thus  $x \rightarrow y$  represents the set of functions which, given an argument of type x, yield a value of type y.

For example, integer→real represents those functions which yield a real number when given an integer. It lies (in  $[\hat{D} \rightarrow \hat{D}]$ ) above the function integer→integer which in turn lies above real→integer.

One of the more important properties of the proposed system is that type properties of least fixedpoints of recursive definitions can be deduced from the values of least fixedpoints of appropriately extended definitions. More precisely, let  $t$  be a term, let  $b$  be a sequence of operations over  $D$  and let  $t_b$  be the functional mapping  $[D \rightarrow D]$  into itself defined by using  $b$  to interpret the base function symbols (such as "+") occurring in  $t$ . Then the least fixedpoint  $f$  of  $t_b$  can be considered as the 'meaning' of the recursive program  $F = t[F]$ . Now let  $\hat{b}$  be formed by extending the operations of  $b$  to monotonic operations over  $\hat{D}$  (e.g. as the addition operation was extended in the above example) and let  $t_{\hat{b}}$  be the corresponding functional mapping  $[\hat{D} \rightarrow \hat{D}]$  into itself, with least fixedpoint  $\hat{f}$ . Then under certain conditions it can be shown that  $\hat{f}$  is a monotonic extension of  $f$ , and thus information about the behaviour of  $f$  on the possibly infinite number of objects of type  $x$  can be deduced from the single value  $\hat{f}(x)$ .

For example, if  $\hat{f}(\text{integer}) = \text{real}$ , and if  $i$  is an integer, then in  $\hat{D}$  we have  $i \sqsubseteq \text{integer}$  and so  $\hat{f}(i) \sqsubseteq \text{integer}$  by monotonicity. But since  $\hat{f}$  is an extension of  $f$ ,  $\hat{f}(i)$  is equal to  $f(i)$  and so the latter must be a real number (or  $\perp$ ).

In our system, all type properties of functions are expressed as inequalities (inclusions). In order to show that the least fixedpoint  $Yt_b$  of the functional  $t_b$  has certain type properties, we show that  $Yt_{\hat{b}} \sqsubseteq g$  for a suitable combination  $g$  of arrow functions. Such inclusions can often be proved using Park's fixedpoint method, i.e. by showing that  $t_{\hat{b}}(g) \sqsubseteq g$ . This check can be performed by direct evaluation of  $t_{\hat{b}}(g)$  and does not require an inductive proof.

An additional advantage of having a single unified domain is that recursive definitions of objects and types are handled in the same way. Some interesting possibilities can be realised by adding parameters to recursive definitions that mix objects and types. Consider, for example, the following type generating recursive definition:

$$S(x,n) = \text{if } n < 0 \text{ then nil else } x * S(x,n-1)$$

(where  $*$  is string concatenation). Then over an appropriate domain,  $S(\text{boolean},3)$  represents the type consisting of strings of booleans of length 3,  $S(U,5)$  the type of strings of arbitrary elements of length 5, and  $S(0,\text{integer})$  strings of 0's of arbitrary length.

Another interesting possibility is the use of type objects to represent errors and exceptional situations. The new objects added are essentially error messages, and they may have an internal structure giving various degrees of information about the nature of the error. For example, the general type error may be above the subtypes overflow, divisionerror, subscripterror, and domainerror. The type subscripterror can itself be above outofrange and notinteger. The last would also be located below the general domainerror, which would also be the result of meaningless combinations such as  $3 + tt$ .

We can also have unions of standard (error-free) types and error messages, which play the role of warning messages. The distinction between errors, warnings and error-free types is illustrated by the following possible equations

$$\begin{aligned} \text{integer/positiveinteger} &= \text{rational} \\ \text{integer}/0 &= \text{divisionerror} \\ \text{integer/integer} &= \text{rational} \sqcup \text{divisionerror} \end{aligned}$$

The last identity is a warning (which could be issued during type checking) that division of an integer by an integer could lead to an error.

We now proceed to a more formal development of the system just outlined, together with more detailed examples. Because space is limited, some proofs will be omitted, and others only outlined.

### 1. The Construction of D

By a domain we mean a partially ordered set  $D$  such that

- (i)  $D$  has a least element;
- (ii) Any directed set of elements of  $D$  has a least upper bound.

We will use " $\sqsubseteq_D$ ", " $\perp_D$ " and " $\sqcup_D$ " to denote the ordering on  $D$ , the least element of  $D$  and the lub operation over  $D$ , respectively. The subscript  $D$  will be omitted when no confusion is likely.

By a data type over  $D$  we mean a subset  $x$  of the universe of  $D$  such that  $x$  is

- (i) closed downwards, i.e. if  $d_0$  and  $d_1$  are in  $x$ , and if  $d_0 \sqsubseteq d_1$ , then  $d_1 \in x$  implies  $d_0 \in x$ ;
- (ii) closed under lub, i.e. if  $s$  is a subset of  $x$  and is a directed subset of  $D$ , then  $\text{lub } s \in x$ .

Sets with these properties are also called ideals.

Classical examples of types are the set of all integers, of all reals, of all strings, and so on. In addition, for the purpose of type checking, we may consider nonstandard types such as "integers greater than or equal to 91" or "strings of 0's and 1's only".

Not every set is a type; our methods require (i) and (ii) as above. In particular,  $\perp$  must be an element of every type and so the set of defined integers, for example, is not a type.

With every element  $d$  of  $D$  is associated the type  $\tilde{d}$  defined as follows:

$$\tilde{d} = \{d' : d' \sqsubseteq d\},$$

i.e.  $\tilde{d}$  is the set of all elements which approximate  $d$ .

The expanded domain  $\hat{D}$  is formed by adding (in effect) a collection of types to  $D$ . Not any collection will do, however. By a type structure over  $D$  we mean a collection  $T$  of types over  $D$  satisfying the following conditions:

- (i)  $\tilde{d} \in T$  whenever  $d \in D$ ;
- (ii) the universal type  $U$ , the set of all elements of  $D$ , is in  $T$ ;
- (iii) the set intersection of the types in any non-empty subcollection of  $T$  is again in  $T$ .

The domain  $\hat{D}$  is the set  $T$  together with set inclusion as the order.

THEOREM I. For any domain  $D$  and any type structure  $T$  over  $D$ , if  $\hat{D} = (T, \subseteq)$  then

- (i)  $\hat{D}$  is a complete lattice (and therefore a domain);
- (ii)  $\perp$  is the least element of  $\hat{D}$ ;
- (iii) for any  $d_0, d_1$  in  $D$ ,  $d_0 \sqsubseteq_D d_1$  iff  $d_0 \sqsubseteq \hat{D} d_1$ ;
- (iv) if  $s$  is a  $\sqsubseteq_D$ -directed subset of  $D$ ,  $e = \sqcup s$ , then  $\{\tilde{d} : d \in s\}$  is a  $\sqsubseteq_{\hat{D}}$ -directed subset of  $\hat{D}$  and  $\hat{e}$  is its  $\hat{D}$ -lub.

Thus  $\hat{D}$  contains an isomorphic copy of  $D$ , and so can be considered an extension of  $D$  (and we will often treat it as such).

Property (iv) is particularly important when least fixedpoints are discussed because we do not want the structure of lub's in  $D$  to be changed in  $\hat{D}$ . Note, however, that if  $s$  is an arbitrary nondirected set with a lub  $x$  in  $D$ , then the lub in  $\hat{D}$  of the corresponding subset of  $\hat{D}$  may not be  $\tilde{x}$ . Note also that it may be necessary to add to  $T$  some sets whose sole purpose is to fill in gaps in the intersection structure of the elements of  $D$ .

A given domain may be extended in many different ways. One way is to let  $T$  be the collection of all possible data types over the domain in question. In most cases this turns out to be a highly undesirable extension, since its structure is so rich as to become unmanageable (for example, it may be extremely difficult to extend a given operation). In practice, smaller extensions are easier to handle than larger ones, and it seems best to restrict the added types as far as possible to those actually needed.

## 2. Function Domains

The system developed in this paper is in a sense "data type free" but not "function type free" since we preserve the separation between function domains of different orders (e.g. between  $D$  and  $[D \rightarrow D]$ ). For the sake of simplicity we consider only functions of a single argument; the extension to multi-argument functions is straight forward.

If  $D$  and  $E$  are any two domains, we define  $[D \rightarrow E]$  to be the collection of all monotonic functions from  $D$  to  $E$  with the usual (pointwise) ordering. That  $[D \rightarrow E]$  is a domain is easily verified.

A function  $\hat{f}$  in  $[\hat{D} \rightarrow \hat{D}]$  is an extension of a function  $f$  in  $[D \rightarrow D]$  iff  $e = f(d)$  implies  $\hat{e} = \hat{f}(\hat{d})$  for any  $d$  and  $e$  in  $D$ . Not every function in  $[\hat{D} \rightarrow \hat{D}]$  is an extension of one in  $[D \rightarrow D]$ ; those that are, we call conservative.

A function  $g$  in  $[\hat{D} \rightarrow \hat{D}]$  is said to be tight iff

$$g(x) = \bigcup_{d \in x} g(\hat{d})$$

for any  $x$  in  $\hat{D}$ . In other words, a tight function is determined by its values over  $D$ . For example, if  $g$  is tight (and  $\hat{D}$  is as in the introduction) then  $g(\underline{\text{integer}})$  must be the lub of  $g(0), g(1), g(-1), \dots$ . Given a function  $h$  in  $[\hat{D} \rightarrow \hat{D}]$  we define the tightening  $\bar{h}$  of  $h$  as follows:

$$\bar{h}(x) = \bigcup_{d \in x} h(\hat{d})$$

for any  $x$  in  $\hat{D}$ .

THEOREM II. For any  $h$  in  $[\hat{D} \rightarrow \hat{D}]$ ,

- (i)  $\bar{h}$  is the least function in  $[\hat{D} \rightarrow \hat{D}]$  which agrees with  $h$  on  $D$ ;
- (ii)  $\bar{h}$  is a tight function;
- (iii)  $h$  is tight iff  $h = \bar{h}$ .

Given any  $f$  in  $[D \rightarrow D]$ , we define the tight extension  $\bar{f}$  of  $f$

to  $[\hat{D} \rightarrow \hat{D}]$  as follows:

$$\bar{f}(x) = \bigcup_{d \in x} f(\bar{d})$$

for any  $x$  in  $\hat{D}$ . The tight extension of a function is that extension which, in a sense, no new possibilities not already inherent in the function itself. If, for example,  $f(i) = 0$  for every integer  $i$ , then  $\bar{f}(\text{integer})$  must also be 0.

THEOREM III. For any  $f$  in  $[D \rightarrow D]$ :

- (i)  $\bar{f}$  is a tight function;
- (ii)  $\bar{f}$  is the least extension of  $f$ ;
- (iii)  $\bar{f} = \bar{\bar{f}}$  for any extension  $\bar{f}$  of  $f$ .

It might seem that nontight functions serve no purpose, and could be eliminated. There are three reasons why this is not the case:

- (i) the composition of tight functions may not be tight;
- (ii) the tight extension of a given function may be very complex, while at the same time simple and adequate nontight extensions exist;
- (iii) the least fixedpoint of  $t_{\hat{b}}$  may not be tight, even when  $\hat{b}$  consists of the tight extensions of the functions of  $b$ .

Finally, we should mention why we define  $[D \rightarrow E]$  to consist of all monotonic functions from  $D$  to  $E$ , and not just the continuous ones. The reason is that it is possible to find examples of a domain  $D$ , a type extension  $\hat{D}$  of  $D$  and a continuous function  $f$  from  $D$  to  $D$  which cannot be extended to a continuous function from  $\hat{D}$  to  $\hat{D}$ . It is possible that adding more restrictions on  $D$  and  $\hat{D}$  would eliminate this difficulty.

### 3. Arrow Functions

We now investigate more closely the properties of the arrow operator defined in the introduction. Our first result justifies the claim that  $x \rightarrow y$  represents those functions which, given an argument of type  $x$ , return a result of type  $y$ .

THEOREM IV. For any  $x$  and  $y$  in  $\hat{D}$  and any  $h$  in  $[\hat{D} \rightarrow \hat{D}]$ , the following are equivalent:

- (i)  $h \sqsubseteq x \rightarrow y$ ;
- (ii)  $h(x) \sqsubseteq y$ ;
- (iii)  $h(z) \sqsubseteq y$  whenever  $z \sqsubseteq x$ .



The arrow operation is increasing in  $y$  but decreasing in  $x$ .

THEOREM V. For any  $x, x', y,$  and  $y'$  in  $D$ :

if  $x' \subseteq x$  and  $y \subseteq y'$  then  $x \rightarrow y \subseteq x' \rightarrow y'$ .

This brings out the crucial difference between our arrow operation and the function domain constructor. The domain  $[D \rightarrow E]$  is the collection of all monotonic functions from  $D$  to  $E$ , and the larger are  $D$  and  $E$ , the larger is  $[D \rightarrow E]$ . But functions below  $x \rightarrow y$  are functions which, given something in  $x$ , return something in  $y$ . Increasing  $x$  makes the condition  $h \subseteq x \rightarrow y$  more restrictive, and decreasing  $x$  makes it less restrictive, because  $h \subseteq x \rightarrow y$  says nothing about what  $h$  does to arguments not of type  $x$ . Our arrow function is similar to Yeung's construct (see [3]), but Scott's arrow operation on retracts [4] acts as a domain construction operation. Of course, it is only the fact that one type can be a subtype of another that brings about the distinction; in most systems, being of type  $x \rightarrow y$  is the same as being an element of the domain  $[x \rightarrow y]$ .

Compound function types can be formed by  $\cap$ -ing together (taking glb's of) arrow functions, and these combinations obey certain intuitively plausible rules.

THEOREM VI. For any  $x, x', y, y'$  and  $z$  in  $D$ :

- (i)  $x \rightarrow y \cap y \rightarrow z \subseteq x \rightarrow z$ ;
- (ii)  $x \rightarrow y \cap x' \rightarrow y' \subseteq (x \cap x') \rightarrow (y \cap y')$ .

Thus  $\cap$  acts very much as intersection does in a set-based system, such as Yeung's (indeed his has much in common with ours). On the other hand,  $\sqcup$ , since it works pointwise, is very different from set union; in fact it is really unnecessary because

$$(x \rightarrow y) \sqcup (x' \rightarrow y') = (x \cap x') \rightarrow (y \sqcup y')$$

for any  $x, x', y$  and  $y'$ .

Of special interest are compound types which result from splitting an arrow type into cases. More precisely, a case analysis of the arrow function  $x \rightarrow y$  is a function of the form

$$(x_0 \rightarrow y) \cap (x_1 \rightarrow y) \cap \dots \cap (x_{n-1} \rightarrow y)$$

such that the type  $x$  is the set union of the types  $x_0, x_1, \dots, x_{n-1}$ . For example, eveninteger  $\rightarrow$  real  $\cap$  oddinteger  $\rightarrow$  real is a case analysis of integer  $\rightarrow$  real. Note that these two functions are not equal; the first when applied to integer gives U, whereas the second gives real.

If  $g$  is a case analysis of  $x \rightarrow y$ , it can be shown that  $\bar{g} = x \rightarrow y$  and  $x \rightarrow y \subseteq g$ .

Our careful distinction between a function and its case analysis may seem pointless to those who think about types in terms of sets. But as it turns out, it is exactly the lack of such distinctions which inhibits proofs by case analysis in simple type checking systems. A special rule which handles the problem will be developed in section 6.

#### 4. Recursion Over the Extended Domain

We now give more details about the relationship between the meanings of recursive programs interpreted over  $D$  and over  $\hat{D}$ . As in the introduction, we assume that  $t$  is a term in some pure recursive programming language (we will not go into details on this point) and that we are interested in the recursive program

$$F = t[F].$$

For example, if  $t$  is the term

$$\lambda n \text{ if } n < 1 \text{ then } 1 \text{ else } n * F(n-1)$$

we are dealing with a program for the factorial function. In this example the base functions symbols are "if-then-else", "<", "\*" and "-".

Our most important result is that the least fixedpoint of  $t_{\hat{b}}$  (as defined in the introduction) is an extension of the least fixedpoint of  $t_b$ . We prove first the following result to the effect that  $t_{\hat{b}}$  is in a sense an extension of  $t_b$ .

LEMMA I. For any domain  $D$ , any type extension  $\hat{D}$  of  $D$ , any term  $t$ , any sequence  $b$  of operations over  $D$  with extensions  $\hat{b}$  over  $\hat{D}$ , any element  $g$  of  $[D \rightarrow D]$  with extension  $\hat{g}$  in  $[\hat{D} \rightarrow \hat{D}]$ ,

$$t_{\hat{b}}(\hat{g}) \text{ is an extension of } t_b(g).$$

PROOF (sketch). Let  $d$  be in  $D$ . Then in the process of evaluating  $t_{\hat{b}}(g)(\hat{d})$ , new elements of  $D$  will never arise; thus the evaluation of  $t_{\hat{b}}(g)(\hat{d})$  will be completely analogous to that of  $t_b(g)(d)$ , so that if the result of the latter is  $e$ , the result of the former will be  $\hat{e}$ .

COROLLARY I. Let  $t$ ,  $b$ ,  $\hat{b}$ ,  $g$  and  $\hat{g}$  be as above, and let  $\bar{b}$  be the sequence of tight extensions of the operations in  $b$ . Then

$$\overline{t_{\hat{b}}(\hat{g})} = \overline{t_{\bar{b}}(g)} \subseteq t_{\bar{b}}(\bar{g}) \subseteq t_{\bar{b}}(\hat{g}) \subseteq t_{\hat{b}}(\hat{g}).$$

Note that if  $f$  and  $h$  are in  $[D \rightarrow D]$  and  $[\hat{D} \rightarrow \hat{D}]$  respectively,

then  $h$  is an extension of  $f$  iff  $\bar{f} \subseteq h$ .

THEOREM VII. Let  $D, \hat{D}, t, b$  and  $\hat{b}$  be as in the previous lemma. Then the least fixedpoint of  $t_{\hat{b}}$  is an extension of the least fixedpoint of  $t_b$ .

PROOF (sketch). Let  $f_0 = \lambda x. 1$  and for any positive ordinal  $k$  let  $f_k = \bigcup_{j < k} t_b(f_j)$ , define the sequence  $\hat{f}_k$  analogously. Then a simple induction on  $k$  shows that  $\hat{f}_k$  is, for each  $k$ , an extension of  $f_k$ . Since the two least fixedpoints are the limits of the respective sequences, our result follows. The only complication is that since the functions involved may not be continuous, we must consider infinite as well as finite ordinals.

The following corollary justifies our approach to type checking, as described in the next section.

COROLLARY I. Let  $D, \hat{D}, b$  and  $\hat{b}$  be as before, let  $x$  and  $y$  be types in  $\hat{D}$  and let  $f$  and  $\hat{f}$  be the least fixed points of  $t_b$  and  $t_{\hat{b}}$  respectively. If either  $\hat{f} \subseteq x \rightarrow y$  or  $\bar{f} \subseteq x \rightarrow y$  then  $f(d) \in y$  for any  $d$  in  $x$ .

There is also a refinement of theorem VII analogous to that of lemma I.

COROLLARY II. Let  $t, b, \hat{b}$  and  $\bar{b}$  be as before. Then

$$\overline{Yt_{\hat{b}}} = \overline{Yt_b} \subseteq Yt_{\bar{b}} \subseteq Yt_{\hat{b}}.$$

## 5. Type Checking

In our system, type checking the recursive program  $F = t[F]$  is reduced to proving an inclusion of the form

$$\overline{Yt_b} \subseteq (x_0 \rightarrow y_0) \cap (x_1 \rightarrow y_1) \cap \dots \cap (x_n \rightarrow y_n)$$

and, as we saw in the last section, it is sufficient to prove the weaker version in which  $Yt_{\hat{b}}$  replaces  $\overline{Yt_b}$ . For example, suppose that our program is

$$F(n) = \text{if } n=0 \text{ then } 0 \text{ else } 3 * F(n-1) + 1,$$

and that we wish to show that its least fixed point maps even integers to even integers and odd to odd. Then we must show

$$Yt_{\hat{b}} \subseteq (\underline{\text{eveninteger} \rightarrow \text{eveninteger}}) \cap (\underline{\text{oddinteger} \rightarrow \text{oddinteger}}).$$

The important point is that type checking is now just a case of the general and well studied problem of proving assertions about least fixed-points, and so we have at our disposal several useful methods.

One of the simplest of these is direct evaluation: we evaluate the term  $F(x_i)$  (using the standard substitution/simplification method) and try to obtain something below the corresponding  $y_i$ . Sometimes this results in a set of mutually dependant equations for  $F(x_0), F(x_1), \dots$  which can be solved by computing successive approximations. For example, with the program given above, direct evaluation gives the equations

$$\begin{aligned} F(\text{eveninteger}) &= 0 \sqcup 3 * F(\text{oddinteger}) + 1 \\ F(\text{oddinteger}) &= 3 * F(\text{eveninteger}) + 1 \end{aligned}$$

and the approximations settle down after five steps to the pair of values  $[\text{eveninteger}, \text{oddinteger}]$ .

Park's method is, as was mentioned in the introduction, also quite useful. For the program given above, we must show  $t_b^{\wedge}(g) \sqsubseteq g$  where  $g$  is  $\text{eveninteger} \rightarrow \text{eveninteger} \sqcap \text{oddinteger} \rightarrow \text{oddinteger}$ . This is equivalent to showing the two inclusions

$$\begin{aligned} t_b^{\wedge}(g)(\text{eveninteger}) &\sqsubseteq \text{eveninteger} \\ t_b^{\wedge}(g)(\text{oddinteger}) &\sqsubseteq \text{oddinteger}. \end{aligned}$$

These calculations are straightforward, e.g.

$$\begin{aligned} t_b^{\wedge}(g)(\text{eveninteger}) &= \text{if } \text{eveninteger}=0 \text{ then } 0 \text{ else } 3 * g(\text{eveninteger}-1) + 1 \\ &= \text{if } \text{boolean} \text{ then } 0 \text{ else } 3 * g(\text{eveninteger}-1) + 1 \\ &= 0 \sqcup 3 * g(\text{eveninteger}-1) + 1 \\ &= 0 \sqcup 3 * g(\text{oddinteger}) + 1 \\ &= 0 \sqcup 3 * \text{oddinteger} + 1 \\ &= 0 \sqcup \text{eveninteger} \\ &= \text{eveninteger}. \end{aligned}$$

We might at this point comment briefly on a peculiar property of the ordering on  $\hat{D}$ : the (vaguely defined) notion of "amount of information" changes in two opposing directions. If  $x \sqsubseteq y$  in  $\hat{D}$ , then  $y$  is a more defined data object than  $x$ , but a less precise (larger) type. In particular, the most defined object,  $U$ , gives the least type information - none at all.

It might seem, then, that it would be a good idea to change the definition of  $\hat{D}$  by placing the types below the appropriate data objects (let us call the resulting domain  $\check{D}$ ). This is possible, and the analogs of the theorems we have so far proved are also valid; in particular, the least fixedpoint  $\check{f}$  of  $t_b^{\vee}$  a monotonic extension of  $f$ . But the problem is that in most cases  $\check{f}$  gives us no type information at all, because  $\check{f}$  applied to any type is simply  $\perp$  (this is the case with the example above). The reason we use  $\hat{D}$  instead of  $\check{D}$  is that for the purposes of type checking, we want the least fixedpoint to draw

out the maximum type information from the program - and this is possible because  $\sqsubseteq_{\hat{D}}$  orders types in the reverse order from set inclusion.

## 6. Case Analysis

The type checking methods just discussed fall down when some sort of analysis by cases is required. As a very simple example, consider the following program

$$F(n) = \text{if } n=0 \text{ then } n \text{ else } 0$$

with least fixedpoint  $f$ , and suppose that we are trying to prove that  $f$  applied to any integer is 0. It is futile to try to show that  $\hat{f}$  is of type integer $\rightarrow$ 0 because it is not true:  $\hat{f}(\text{integer})$  is integer (with  $D$  and  $\hat{D}$  as in the introduction).

Now suppose that  $\hat{D}$  has an extra type nonzerointeger (abbreviated nzi) and that nzi $\rightarrow$ 0 is ff. Then simple calculation will show that

$$f \sqsubseteq (0 \rightarrow 0) \sqcap (\text{nzi} \rightarrow 0).$$

The right hand side is a case analysis of integer $\rightarrow$ 0, and so its tightening is integer $\rightarrow$ 0; thus

$$\bar{f} \sqsubseteq \text{integer} \rightarrow 0$$

and from this we can conclude that  $f$  applied to any integer is 0.

This type of case analysis extends the power of the system, but it usually breaks down for more realistic, genuinely recursive programs. Consider the following program defining a function which flattens binary trees into strings:

$$F(u) = \text{if } \text{isatom}(u) \text{ then } u \text{ else } F(\text{left}(u)) * F(\text{right}(u)).$$

Assume that the domain  $D$  has binary trees (some of which are atoms) and strings of atoms, that  $b$  interprets "\*" as string concatenation, and "left", "right" and the predicate "isatom" in the usual way.

We are trying to prove that the least fixedpoint  $f$  of this program takes trees into strings. If we construct  $\hat{D}$  simply by adding the types tree and string (and define  $\hat{b}$  appropriately) we will fail for reasons similar to those in the first example. But even if we add two more types, atom and nonatom, both below tree, and let

$$g = (\text{atom} \rightarrow \text{string}) \sqcap (\text{nonatom} \rightarrow \text{string})$$

we still will not succeed. If we use Park's method, the evaluation of  $t_{\hat{D}}(g)(\text{nonatom})$  gives us  $g(\text{tree}) * g(\text{tree})$  and since  $g(\text{tree}) = U$ , the result is  $U$ .

Nevertheless, it is still true that  $\bar{f} \sqsubseteq g$  and, as we have seen, that is all we need. Fortunately, there is a variant of the Park fixed-point method with which we can prove many inclusions of the form  $\overline{Yt_b} \sqsubseteq g$  without actually determining  $\overline{Yt_b}$ .

THEOREM VIII. Let  $D$ ,  $\hat{D}$ ,  $t$ ,  $b$  and  $\hat{b}$  be as in the previous theorems, and let  $g$  be any element of  $[\hat{D} \rightarrow \hat{D}]$ . Then

$$t_{\hat{b}}(\bar{g}) \sqsubseteq g \text{ implies } \overline{Yt_b} \sqsubseteq g.$$

PROOF (sketch). For any ordinal  $k$  define  $f_k$  as in theorem vii. We prove by induction on  $k$  that  $\bar{f}_k \sqsubseteq g$ .

The base step is straight forward. Now let  $k$  be a positive ordinal and let  $m < k$ . By induction we have  $\bar{f}_m \sqsubseteq g$ . Thus  $\bar{g} \sqsupseteq \bar{f}_m = \bar{f}_m$  and so  $t_{\hat{b}}(\bar{f}_m) \sqsubseteq t_{\hat{b}}(\bar{g}) \sqsubseteq g$ . This implies  $t_{\hat{b}}(\bar{f}_m) \sqsubseteq g$  and  $t_{\hat{b}}(\bar{f}_m) = \overline{t_{\hat{b}}(\bar{f}_m)}$  by corollary ii of theorem vii. Thus  $\bar{f}_{m+1} \sqsubseteq g$  and since this is true of every  $m$  less than  $k$ , we can conclude that  $\bar{f}_k \sqsubseteq g$  and our result follows easily.

We illustrate the method on our tree flattening example. Since the function  $g$  as defined above is a case analysis of tree $\rightarrow$ string, we have  $\bar{g} = \text{tree} \rightarrow \text{string}$ . Thus we need only show that

$$t_{\hat{b}}(\text{tree} \rightarrow \text{string}) \sqsubseteq (\text{atom} \rightarrow \text{string}) \sqcap (\text{nonatom} \rightarrow \text{string}).$$

This is equivalent to showing that

$$\begin{aligned} t_{\hat{b}}(\text{tree} \rightarrow \text{string})(\text{atom}) &\sqsubseteq \text{string} \\ t_{\hat{b}}(\text{tree} \rightarrow \text{string})(\text{nonatom}) &\sqsubseteq \text{string} \end{aligned}$$

and the calculations are straight forward. With this method, as well as with the ordinary Park method, its actual use involves no induction.

As another example, consider the program

$$F(n) = \text{if } n < 100 \text{ then } n-10 \text{ else } F(F(n+11))$$

for the well known 91-function. If  $\hat{D}$  contains the type integer and also the types ge91 (of integers greater than or equal to 91), gr100 (of integers greater than 100) and le100 (of integers less than or equal to 100) then our method shows that

$$\bar{f} \sqsubseteq \text{integer} \rightarrow \text{ge91}$$

using the case analysis

$$(\text{gr100} \rightarrow \text{ge91}) \sqcap (\text{le100} \rightarrow \text{ge91}).$$

A natural question concerning this method is how to find the appropriate partition of given types into complementary subtypes. Our experience indicates that the tests of the if-then-else's are a good guide.

### Acknowledgements

We would like to thank Robin Milner and Lockwood Morris for a very useful discussion which provided the original inspiration for this work.

The research was financed in part by a grant from the Science Research Council of the United Kingdom.

### References

1. R. Milner, L. Morris and M. Newey, "A logic for computable functions with reflexive and polymorphic types", proceedings of the Symposium on Proving and Improving Programs, Arc et Senans 1975, pp371-394.
2. A. Church, "A formulation of the simple theory of types", Journal of Symbolic Logic, v5 (1940), pp56-68.
3. Yeung, PhD thesis, Queen Mary College, London, 1976.
4. Dana Scott, "Data types as lattices", SIAM Journal on Computing, v5 (1976), pp522-587.