

First-order unification by structural recursion

CONOR MCBRIDE

Department of Computer Science, University of Durham, South Road Durham, DH1 3LE, UK
(e-mail: c.t.mcbride@durham.ac.uk)

Abstract

First-order unification algorithms (Robinson, 1965) are traditionally implemented via general recursion, with separate proofs for partial correctness and termination. The latter tends to involve counting the number of unsolved variables and showing that this total decreases each time a substitution enlarges the terms. There are many such proofs in the literature (Manna & Waldinger, 1981; Paulson, 1985; Coen, 1992; Rouyer, 1992; Jaume, 1997; Bove, 1999). This paper shows how a dependent type can relate terms to the set of variables over which they are constructed. As a consequence, first-order unification becomes a structurally recursive program, and a termination proof is no longer required. Both the program and its correctness proof have been checked using the proof assistant LEGO (Luo & Pollack, 1992; McBride, 1999).

Capsule Review

Using first-order unification as an example, McBride proselytizes for a style of programming using structural recursion over indexed families of datatypes. Although this style is not Turing complete, McBride shows that guaranteed termination and the expressiveness of dependent types, is a useful and beautiful paradigm.

1 Introduction

Inductive datatypes give more than just a machine representation for data. They equip it with a natural mode of inspection – case analysis – and a natural mode of terminating computation – structural recursion. Case analysis allows a function over an inductive datatype to decompose its input by stripping off the datatype’s constructor symbols, and structural recursion gives access to the outputs for the components so exposed. The structure of the computation is explained in terms of the structure of each input, and the well-foundedness of the latter ensures the termination of the former.

By contrast, functions defined by general recursion can make whatever recursive calls they like, regardless of the structure of their input. There is no *a priori* guarantee that such functions will ever deliver an output. We may need this extra freedom in order to write the programs we want, but we pay with the extra work required to convince ourselves that the structures we exploit, no longer those of the data, really

are well-founded. Such programs are still “effective procedures” in the sense that we can order a computer to execute them. They do not, however, fully express the understanding we have in our minds. What makes general recursion necessary is a “failure of language”: if we could find the structures we need amongst those we can express, we could give the machine not just orders, but an “effective explanation”.

Gödel’s incompleteness theorems tell us that this failure is sometimes inevitable for languages in which all programs terminate: a language restricted to structural recursion cannot encode its own interpreter. However, despite the fact that we cannot always write programs via structural recursion, there is no need to conceal the structural explanations of many common algorithms which do not involve such spectacular bootstrapping. Indeed, the languages we have today are already enough to make structural many algorithms where general recursion is popularly considered essential. The classic example is Quicksort, a non-structural recursion on lists. As David Turner points out (Turner, 1995), the real structure of Quicksort is the tree-like call structure between whose branches the partitioning operation distributes the list elements. If we turn this into a data structure, we acquire exactly the tree-sort algorithm with which Rod Burstall introduced us to case analysis and structural recursion (Burstall, 1969).

Nonetheless, there are many algorithms where the structure in play is not a simple inductive one. For example, first-order unification (Robinson, 1965) is not structurally recursive on terms, because the terms grow larger each time a variable is substituted. Typical implementations use general recursion. The termination proof is usually quite subtle, relying on the “occur-check” to ensure that each substitution reduces the number of distinct variables remaining in the problem. The usual representation of terms does not take account of the number of variables involved, hence this key piece of structure is hidden.

Dependent type systems give us a much richer language in which to describe data. In particular, we can stratify the type of terms into a family of types, $\text{Term } n$, representing terms over a set of n variables. This paper presents a first-order unification algorithm which is structurally recursive on this n : its correctness proof is available online (McBride, 2003a). Both program and proof have been coded and checked in the Extended Calculus of Constructions, as implemented in the proof assistant LEGO (Luo & Pollack, 1992; McBride, 1999) and are available online (McBride, 2003b). The point is this: dependent types make more recursion structural because dependent types can express more structure.

2 The programming language in this paper

The programs in this paper are all real terms in a real dependent type theory. The full LEGO development is available online (McBride, 2003b). However, LEGO was not designed for presenting programs, so for reasons of legibility, I have disguised this type theory as an imaginary functional language. In doing so, I have borrowed from Haskell and from Standard ML, as well as exploiting the cosmetic potential of L^AT_EX. Some less superficial aspects deserve more thorough discussion.

2.1 Datatype declarations and function definitions

Where, in Haskell, we might write

`data N = 0 | s N`

this paper presents datatypes by a formation rule and constructors, in natural deduction style:

$$\text{data} \quad \frac{}{\mathbf{N} : \text{Type}} \quad \text{where} \quad \frac{}{0 : \mathbf{N}} \quad \frac{n : \mathbf{N}}{sn : \mathbf{N}}$$

The more familiar notation suggests the form of a typical constructor application, using a type as a placeholder for its inhabitants – `s N`, for example. In a simply typed setting, there is no risk of mistaking `s N` for an actual application, because terms and types inhabit separate worlds. Dependent types break this separation, so there is perhaps less potential for misunderstanding if we avoid writing constructor declarations which resemble ill-typed expressions. The natural deduction style follows the same “typical application” idea, but with variables as the placeholders. The `s` constructor may still be used in partially applied form: its actual type is $\mathbf{N} \rightarrow \mathbf{N}$, as one might expect.

Another benefit of illustrating typical usage in a declaration is that we can indicate that an argument is to be left implicit – and inferred by the machine – by leaving it as a schematic variable in a rule. Consider, for example

$$\text{data} \quad \frac{T : \text{Type}}{\text{Maybe } T : \text{Type}} \quad \text{where} \quad \frac{}{\text{no} : \text{Maybe } T} \quad \frac{t : T}{\text{yes } t : \text{Maybe } T}$$

Here, the actual type of `no` is a dependent function space: $\forall T : \text{Type}. \text{Maybe } T$, where the return type depends on the type passed as an argument. However, this dependency enables the machine to infer the argument of `no`, given a desired return type. Similarly, the type of `yes` is $\forall T : \text{Type}. T \rightarrow \text{Maybe } T$. Here, $T \rightarrow \text{Maybe } T$ may be used as syntactic sugar for $\forall t : T. \text{Maybe } T$, as t does not occur in its scope. Correspondingly, we must supply t as an explicit argument when we apply `yes`, but we may still hide T . Schematic variables not explicitly declared in the premises of a rule yield \forall 's in the type of the symbol declared in the conclusion, but the corresponding arguments are left implicit in applications. This notation only gives explicit types to the explicit objects. Arguments which are usually implicit will be written as subscripts if they deserve attention.

Function definitions are introduced by the keyword “let”, with a type signature also in natural deduction style, then an equational presentation of the function. For example, here are two function-transforming operations capturing the monadic behaviour which `Maybe` supports:

$$\begin{array}{l} \text{let} \quad \frac{f : S \rightarrow T}{f \downarrow : S \rightarrow \text{Maybe } T} \quad \quad \quad f \downarrow \mapsto \text{yes} \cdot f \\[10pt] \frac{f : S \rightarrow \text{Maybe } T}{f \downarrow : \text{Maybe } S \rightarrow \text{Maybe } T} \quad \quad \quad \begin{array}{l} f \downarrow \text{ no} \mapsto \text{no} \\ f \downarrow (\text{yes } s) \mapsto f s \end{array} \end{array}$$

The \cdot symbol stands for function composition. Less conventionally, I write \mapsto to indicate a directed computation rule, keeping $=$ for equational propositions like

this one:

$$\text{fact} \quad (f \cdot g) \Downarrow s = ((f \Downarrow) \cdot \Downarrow g \Downarrow) s$$

This is a proposition – indeed, a type – which I assert with the keyword ‘fact’, and which we might prove for every f , g and s , but it is not a rule used for evaluating \Downarrow . Dependent types give us ready access to a logic for reasoning about programs. I use this programming language to define logical notions and its type system to state theorems. The typechecker also checks proofs!

Every monad gives rise to a functor composing two ‘half-lifting’ operations – those defined for *Maybe* above. By careful choice of operator precedence, we may write the full functorial lifting of some f by ‘bracketing’ it thus: $\Downarrow f \Downarrow$. It is also possible to generalise $\Downarrow \cdot \Downarrow$ to a $\Downarrow \cdot \Downarrow^n$ operation which lifts an n -ary function to its n -ary *Maybe* counterpart. The entire family can be defined once for all by recursion on n , but it is not worth the trouble here, as the largest n used in this paper is 2.

$$\text{let} \quad \frac{f : R \rightarrow S \rightarrow T}{\Downarrow f \Downarrow^2 : \text{Maybe } R \rightarrow \text{Maybe } S \rightarrow \text{Maybe } T} \quad \dots$$

2.2 Inductive families, pattern matching

Dependent types systems allow us to talk about collections of types in a systematic way by creating type *families*, typically with a signature resembling this:

$$\frac{t : T}{F t : \text{Type}}$$

We say that F is a family *indexed* by T . We may refer to a particular instance of the family by applying F to a particular t , yielding $F t$. We may abstract over the whole family by abstracting over all t ’s, as in $\forall t : T. F t \rightarrow \dots$. We may abstract over a *section* of the family by applying F to an expression containing abstracted variables, thus $\forall s : S. F (g s) \rightarrow \dots$.

Type families support fine analysis of data structures. Many common datatypes can be redesigned as carefully stratified families, with their operations acquiring more precise types. For example, indexing matrices by their dimensions allows a multiplier whose type requires that the column-count on the left equals the row-count on the right. The termination of unification relies on an analysis of the number of variables over which terms are constructed: later, we shall make that analysis explicit by indexing the representation of terms.

Let us first look more closely at how to declare type families and define operations over them. Most dependent type systems allow the introduction of families by inductive definition, although there is considerable variation in the styles supported. Systems like LEGO, COQ and ALF (Luo, 1994; Pfenning & Paulin-Mohring, 1990; Magnusson, 1994) adopt minor variations on the scheme set out in (Dybjer, 1994). They view an inductive family as a collection of mutually defined datatypes, with each choice of index yielding a branch of the definition. We declare a formation rule for the new family, together with constructors which are free to target and to source their recursive arguments from *any* section of it. For example, we may declare a

family of finite types, $\text{Fin } n$ – the index n gives the size of each instance:

$$\text{data } \frac{n : \mathbb{N}}{\text{Fin } n : \text{Type}} \quad \text{where } \frac{}{\text{f0}_n : \text{Fin } sn} \quad \frac{x : \text{Fin } n}{\text{fs}_n x : \text{Fin } sn}$$

The fs constructor embeds $\text{Fin } n$ as the “old” elements of $\text{Fin } sn$, whilst f0 makes a “new” element: the indices n are details which can be inferred mechanically, so I shall usually omit the arguments which I have subscripted above. Observe that both constructors target a restricted section of the family – the types with at least one element. $\text{Fin } 0$ is quite rightly uninhabited.

The ability to restrict the target sections of constructors makes pattern matching over inductive families both more complex and more powerful than the standard notion found in Hindley–Milner languages. If we are defining a function whose domain is a particular section of a family, then, in principle, we need only write patterns for those constructors which target that section. The machine determines which these are by attempting to unify the target indices of each constructor with the indices of the section being matched. Unification is undecidable in general, but we still get a very powerful language when we restrict to decidable first-order problems. Whenever a constructor’s target section does not overlap the function’s domain, we do not need to write a pattern for it. For example, given that no constructor targets $\text{Fin } 0$, we can define the following function with a type signature but no equations!

$$\text{let } \frac{}{\text{empty} : \text{Fin } 0 \rightarrow T}$$

The designers of dependently typed languages continue to debate which notions of inductive family it is sensible to support. This paper exploits one particular choice, but it is reasonable to consider others. Both Cayenne and Agda (Augustsson, 1998; Coquand & Coquand, 1999) are at present less liberal. Whilst they still allow a family’s constructors to source recursive arguments from any section, they demand that every constructor targets the whole family. This removes the complexity from pattern matching: no matter which section of the family we match over, every constructor always applies. However, it substantially restricts the families we can readily define. The above definition of Fin , for example, is no longer available. Fortunately, there is an alternative in this case – we may compute a type of size n , using Maybe and the empty type, \perp :

$$\text{let } \frac{n : \mathbb{N}}{\text{Fin}' n : \text{Type}} \quad \begin{array}{l} \text{Fin}' 0 \mapsto \perp \\ \text{Fin}' (sn) \mapsto \text{Maybe} (\text{Fin}' n) \end{array}$$

One can also imagine systems which allow an even more liberal notion of inductive family than that used in this paper. *Inductive-recursive* definitions (Dybjer & Setzer, 1999) allow datatype families and operations over them to be mutually defined. We might even allow inductive families to be defined mutually with the types which index them. Making structure explicit and precise is habit forming – failure to express oneself brings a ready hunger for more language!

2.3 Structural recursion

Pattern matching gives us access to the components from which a function’s inputs are constructed. Structural recursion gives us access to the function’s outputs for the

components so exposed. It is, perhaps, helpful to view the constructors of a datatype as building larger “new” elements from smaller “old” ones. Structural recursion can be rationalised correspondingly as computing “new” outputs for larger inputs, given that we can already see the “old” outputs for smaller inputs. The fact that recursive calls are initiated after the call which requires them is a consequence of the “just in time” mechanisms which sensible implementations employ. It is this shift, interpreting recursion not as “seeing prior results” but as “doing subsequent computations”, which can also give a meaning to general recursion.

Rod Burstall (1987) describes a variation of ML which restores the “seeing” interpretation. Each pattern variable x standing for a subterm of the input is automatically accompanied by its “inductive hypothesis”, a variable $\$x$ bound to the output for x . There are no recursive calls – return values are built from objects already known and named. The resulting programs may look a little odd, but they make their structural behaviour, and hence guaranteed termination, clear to behold.

We can present structural recursion with a more conventional notation if we restrict the arguments on which recursive calls can be made. In particular, we must be able to identify an argument position for which the recursive calls on the right use only strict subterms exposed by pattern matching on the left. For example, the first argument in the following definition of “reversing append”:

$$\text{let } \frac{xs, ys : \text{List } T}{\text{revapp } xs \ ys : \text{List } T} \quad \begin{array}{l} \text{revapp } [] \quad ys \mapsto ys \\ \text{revapp } (x :: xs) \ ys \mapsto \text{revapp } xs \ (x :: ys) \end{array}$$

We may extend this notion to cover datatypes with higher-order constructors (e.g. Brouwer ordinals, infinitely branching trees, etc.) by considering a functional pattern variable to expose everything in its image. We may liberalise further, allowing *nested* structural recursion. Here we identify a sequence of argument positions, where recursive calls may preserve the first n positions as long as the argument in position $n + 1$ decreases – Ackermann’s function is a typical example:

$$\text{let } \frac{m, n : \mathbb{N}}{\text{ack } m \ n : \mathbb{N}} \quad \begin{array}{l} \text{ack } 0 \quad n \mapsto sn \\ \text{ack } (sm) \ 0 \mapsto \text{ack } m \ (s0) \\ \text{ack } (sm) \ (sn) \mapsto \text{ack } m \ (\text{ack } (sm) \ n) \end{array}$$

Structural recursion thus delivers more functions than primitive recursion (in its original first-order sense). Further, allowing recursive calls on any exposed subterm loosens primitive recursion’s requirement to peel off exactly one constructor at each step. The direct (but inefficient) definition of the Fibonacci function is thus permitted:

$$\text{let } \frac{n : \mathbb{N}}{\text{fib } n : \mathbb{N}} \quad \begin{array}{l} \text{fib } 0 \mapsto 0 \\ \text{fib } (s0) \mapsto s0 \\ \text{fib } (s(sn)) \mapsto \text{fib } n + \text{fib } (sn) \end{array}$$

A dependent type system such as LEGO’s, which only provides higher-order primitive recursion, can nonetheless support all the programs definable by pattern

matching and structural recursion in the above sense. The full translation procedure is given in my thesis (McBride, 1999), but the basic idea is to use primitive recursion to build an auxiliary data structure memoizing all “prior” outputs for each input – the original function’s recursive calls are replaced by look-up operations on this structure. “Doing” becomes “seeing”, just as in Burstall’s language!

Dependent types add more than just detail to the class of structurally recursive function definitions. An inductive family indexed by a datatype has two notions of structural recursion: the constructors of the index type give one, and the constructors of the family itself give another. The two may have quite different behaviours, and they may be combined by nesting. Inductive families thus allow us to avoid general recursion in situations where we can see how to turn hidden recursive structure into the explicit inductive structure of an index type. That is how we shall write first-order unification by structural recursion.

3 First-order terms, renaming and substitution

Let us begin our development of unification with a type of terms over n variables:

$$\begin{array}{l} \text{data} \quad \frac{n : \mathbb{N}}{\text{Term } n : \text{Type}} \\ \\ \text{where} \quad \frac{x : \text{Fin } n}{\iota x : \text{Term } n} \quad \frac{}{\text{leaf} : \text{Term } n} \quad \frac{s, t : \text{Term } n}{s \text{ fork } t : \text{Term } n} \end{array}$$

Here $\text{Fin } n$ represents the variable names, embedded into $\text{Term } n$ with the ι constructor. We may compute on $\text{Term } n$ by nesting recursion, first on n (allowing us to grow the terms if we eliminate a variable), then on the terms themselves. Building the number of variables into the definition of terms explains the structure by which unification operates.

We can represent a *renaming* by a function between variable sets, and a *substitution* by a function from one variable set to terms over another. Substitution has a monadic behaviour and renaming is the associated functor. Just as we did with *Maybe*, we may define a pair of half-lifting operators, with \flat turning a renaming into a substitution, and \flat applying a substitution to a term. The functorial map is again given by “bracketing”, $\flat _ \flat$.

$$\begin{array}{l} \text{let} \quad \frac{r : \text{Fin } m \rightarrow \text{Fin } n}{\flat r : \text{Fin } m \rightarrow \text{Term } n} \quad \flat r \mapsto \iota \cdot r \\ \\ \frac{f : \text{Fin } m \rightarrow \text{Term } n}{f \flat : \text{Term } m \rightarrow \text{Term } n} \quad \begin{array}{l} f \flat (\iota x) \mapsto f x \\ f \flat \text{ leaf} \mapsto \text{leaf} \\ f \flat (s \text{ fork } t) \mapsto (f \flat s) \text{ fork } (f \flat t) \end{array} \end{array}$$

The same substitution, in terms of effect, can have different functional encodings. We shall need to work with a *pointwise* (or *extensional*) equality for substitutions:

$$\text{let} \quad \frac{f, g : \text{Fin } m \rightarrow \text{Term } n}{f \doteq g : \text{Type}} \quad f \doteq g \mapsto \forall x : \text{Fin } m. f x = g x$$

All of the reasoning in this paper relies only on the pointwise behaviour of substitutions. I shall omit the details and freely rewrite with “equations” of the form $f \doteq g$.

Inspecting the definition of \Downarrow , we can see that ι plays the rôle of the identity substitution. We may define the usual composition and check useful monadic properties:

$$\text{let } \frac{f : \text{Fin } m \rightarrow \text{Term } n \quad g : \text{Fin } l \rightarrow \text{Term } m}{f \diamond g : \text{Fin } l \rightarrow \text{Term } n} \quad f \diamond g \mapsto (f \Downarrow) \cdot g$$

$$\text{fact } \quad \iota \Downarrow t = t \quad (f \diamond g) \Downarrow t = f \Downarrow (g \Downarrow t) \quad f \diamond (\Downarrow r) \doteq f \cdot r$$

4 The occur-check, through thick and thin

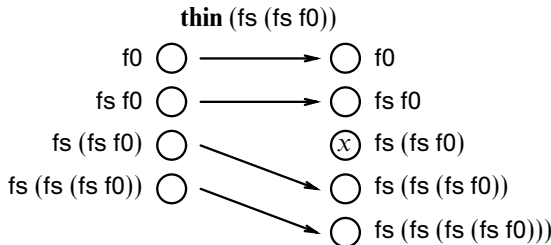
The notion of *occurrence* plays a key rôle in both the termination and correctness of first-order unification. If we are attempting to unify a variable x with a non-variable term t , then we must check whether x is used in t . If so, there is no unifier; if not, then the substitution “ t for x ” is not only a most general unifier, but it is also a means to eliminate x from the remaining parts of the problem: the corresponding reduction in the number of variables is sufficient to guarantee termination, even if the substitution enlarges the terms. The usual boolean occur-check does not expose this aspect of the algorithm: that is the crux of the problem this paper solves.

Our more stratified representation gives us the extra structure we need. If $x : \text{Fin } sn$ and $t : \text{Term } sn$, then any unifying substitution $f : \text{Fin } sn \rightarrow \text{Term } n$ would yield terms with structurally smaller indices. This section, presents an occur-check re-designed more positively and informatively as the search for such a substitution – the effective witness that x does not occur in t .

The boolean occur-check is built from the boolean equality test on variable names. A more informative occur-check needs a more informative equality test, returning a witness when its arguments are different. Hence we need a more concrete way to describe difference in variable sets. If $x : \text{Fin } sn$, we may collect the other n elements of $\text{Fin } sn$ as the image of an embedding, **thin** x from $\text{Fin } n$, “diluting” $\text{Fin } n$ with x :

$$\text{let } \frac{x : \text{Fin } sn \quad y : \text{Fin } n}{\text{thin}_n x y : \text{Fin } sn} \quad \begin{array}{ll} \text{thin}_n \text{ f0} & y \mapsto \text{fs } y \\ \text{thin}_{sn} (\text{fs } x) & \text{f0} \mapsto \text{f0} \\ \text{thin}_{sn} (\text{fs } x) (\text{fs } y) & \mapsto \text{fs } (\text{thin}_n x y) \end{array}$$

The indices are highly significant in the pattern-matching behaviour of **thin**: in order to type the recursive call, we must ensure that there really is an s to strip from the index. However, there is no need for a **thin**₀ (fs x)... case, as there is no such x . The following diagram illustrates **thin**’s behaviour:



The key properties of thinning are as follows:

$$\text{fact} \quad \frac{\mathbf{thin} \ x \ y = \mathbf{thin} \ x \ z}{y = z} \quad \mathbf{thin} \ x \ y \neq x \quad \frac{x \neq y}{\exists y'. \mathbf{thin} \ x \ y' = y}$$

In effect, every y apart from x is the image under $\mathbf{thin} \ x$ of a unique y' . Hence the terms t without x are exactly those expressible as $\mathbf{thin} \ x \ t'$. The search for a unifying, variable-removing substitution thus reduces to the search for such a t' . This reduces in turn to finding the appropriate y' for each variable y we encounter in t . The refinement of boolean equality we need is thus the partial inverse to $\mathbf{thin} \ x$, here named $\mathbf{thick} \ x$.

$$\begin{array}{l} \text{let} \quad \frac{x, y : \text{Fin } sn}{\mathbf{thick}_n \ x \ y : \text{Maybe } (\text{Fin } n)} \quad \begin{array}{l} \mathbf{thick}_n \quad f0 \quad f0 \mapsto \text{no} \\ \mathbf{thick}_n \quad f0 \quad (fs \ y) \mapsto \text{yes } y \\ \mathbf{thick}_{sn} \ (fs \ x) \quad f0 \mapsto f0 \\ \mathbf{thick}_{sn} \ (fs \ x) \ (fs \ y) \mapsto \mathbf{fs}_n \ (\mathbf{thick}_n \ x \ y) \end{array} \end{array}$$

An induction on x shows that \mathbf{thick} really is the partial inverse of \mathbf{thin} :

$$\text{fact} \quad \frac{\mathbf{thick} \ x \ y = r}{\bigvee \left\{ \begin{array}{l} y = x \wedge r = \text{no} \\ \exists y'. y = \mathbf{thin} \ x \ y' \wedge r = \text{yes } y' \end{array} \right.}$$

We can see $\mathbf{thick} \ x$ as a “failure-prone renaming”. The occur-check just propagates $\mathbf{thick} \ x$ monadically through the structure of terms:

$$\begin{array}{l} \text{let} \quad \frac{x : \text{Fin } sn \quad t : \text{Term } sn}{\mathbf{check} \ x \ t : \text{Maybe } (\text{Term } n)} \quad \begin{array}{l} \mathbf{check} \ x \quad (\iota \ y) \mapsto \mathbf{thick} \ x \ y \\ \mathbf{check} \ x \quad \text{leaf} \mapsto \text{yes leaf} \\ \mathbf{check} \ x \ (s \ \text{fork } t) \mapsto \\ \quad (\mathbf{check} \ x \ s) \mathbf{fork} \ \ell^2 \ (\mathbf{check} \ x \ t) \end{array} \end{array}$$

If $\mathbf{check} \ x \ t$ returns no, then x occurs in t ; if we get $\text{yes } t'$, then $\mathbf{thin} \ x \ t' = t$. We may now define the corresponding unifier, $t' \text{ for } x$.

$$\text{let} \quad \frac{x : \text{Fin } sn \quad t' : \text{Term } n}{t' \text{ for } x : \text{Fin } sn \rightarrow \text{Term } n} \quad (t' \text{ for } x) \ y \mapsto \begin{array}{l} \text{case } \mathbf{thick} \ x \ y \\ \text{of } \text{yes } y' \Rightarrow \iota \ y' \\ \quad \text{no} \Rightarrow t' \end{array}$$

As $t' \text{ for } x$ “thickens” everything apart from x itself, we may show that

$$\text{fact} \quad t' \text{ for } x \cdot \mathbf{thin} \ x \doteq \iota$$

Hence, $t' \text{ for } x$ really does unify x with t :

$$\begin{array}{l} \text{fact} \quad \frac{\mathbf{check} \ x \ t = \text{yes } t'}{t' \text{ for } x \ t = t' \text{ for } x \ (\mathbf{thin} \ x \ t) = (t' \text{ for } x \cdot \mathbf{thin} \ x) \ t' = t' \\ \quad = (t' \text{ for } x) \ x} \end{array}$$

5 A refinement of substitution

In order to exploit recursion on n when unifying over $\text{Term } n$, we need to know that each step of the substitution we accumulate really does get rid of a variable. It is not obvious how a functional substitution breaks down into individual steps, but we can define an inductive representation of substitutions built by solving for a

sequence of variables in the manner described above.

$$\begin{array}{l} \text{data} \quad \frac{m, n : \mathbb{N}}{\text{AList } m \ n : \text{Type}} \\ \text{where} \quad \frac{}{\text{anil} : \text{AList } n \ n} \quad \frac{\sigma : \text{AList } m \ n \quad t' : \text{Term } m \quad x : \text{Fin } sm}{\sigma \text{ asnoc } t'/x : \text{AList } sm \ n} \end{array}$$

Looking just at the indices, $\text{AList } m \ n$ behaves very like $m \geq n$. Looking at the data, we see that $\text{AList } m \ n$ stores association lists, linking each x with a t' over the set with x removed, until the number of variables has been whittled down from m to n . The family $\text{AList } m \ n$ captures exactly the subset of $\text{Fin } m \rightarrow \text{Term } n$ we shall actually use, and in a form which admits structural decomposition, not just application. We can easily recover the functional form:

$$\begin{array}{l} \text{let} \quad \frac{\sigma : \text{AList } m \ n}{\text{sub } \sigma : \text{Fin } m \rightarrow \text{Term } n} \quad \text{sub} \quad \text{anil} \quad \mapsto \iota \\ \quad \text{sub } (\sigma \text{ asnoc } t'/x) \mapsto \text{sub } \sigma \diamond t' \text{ for } x \end{array}$$

Observe that anil gives a concrete name to the identity substitution (for each n). We can also define composition:

$$\begin{array}{l} \text{let} \quad \frac{\rho : \text{AList } m \ n \quad \sigma : \text{AList } l \ m}{\rho \mathbin{++} \sigma : \text{AList } l \ n} \\ \quad \rho \mathbin{++} \text{anil} \quad \mapsto \rho \\ \quad \rho \mathbin{++} (\sigma \text{ asnoc } t/x) \mapsto (\rho \mathbin{++} \sigma) \text{ asnoc } t/x \end{array}$$

It is easy to check that **sub** respects $++$ in the appropriate way:

$$\text{fact} \quad \text{sub } (\rho \mathbin{++} \sigma) \doteq \text{sub } \rho \diamond \text{sub } \sigma$$

That is, the $\text{AList } m \ n$ form a category whose objects are just those types in the image of Term and whose collections of arrows are a subset of the usual function spaces over those objects – just those functions in the image of $(\cdot) \cdot \text{sub}$.

The most we can say in advance about the unifying substitution we hope to find for two terms in $\text{Term } m$ is that it inhabits $\text{AList } m \ n$ for *some* n . We can express this existential notion via *dependent pairing*, which generalizes the cartesian product, \times , in the same way that \forall generalizes \rightarrow .

$$\text{data} \quad \frac{T : S \rightarrow \text{Type}}{\exists T : \text{Type}} \quad \text{where} \quad \frac{s : S \quad t : T s}{\langle s, t \rangle : \exists T}$$

If we wish to express existential propositions, we may regard the $\exists x : S. T$ as syntactic sugar for $\exists(\lambda x : S. T)$. The non-sugared syntax is also useful: $\exists(\text{AList } m)$ is exactly the type of “substitutions for *some* target”.

Let us overload $_ \text{ asnoc } _ / _$ with its “uncurried” counterpart, extending substitutions paired with their targets:

$$\begin{array}{l} \text{let} \quad \frac{a : \exists(\text{AList } m) \quad t' : \text{Term } m \quad x : \text{Fin } sm}{a \text{ asnoc } t'/x : \exists(\text{AList } sm)} \\ \quad \langle n, \sigma \rangle \text{ asnoc } t'/x \mapsto \langle n, \sigma \text{ asnoc } t'/x \rangle \end{array}$$

The effect of \exists here is to hide AList ’s “target” index. Would it not be simpler, one might reasonably ask, to define AList with just a “source” index in the first place?

Given such a definition, we could still compute the target, should we need it:

$$\begin{array}{l}
 \text{data} \quad \frac{m : \mathbb{N}}{\text{AList}' m : \text{Type}} \\
 \text{where} \quad \frac{}{\text{anil}' : \text{AList}' m} \quad \frac{\sigma : \text{AList}' m \quad t' : \text{Term } m \quad x : \text{Fin } sm}{\sigma \text{ acons}' t' / x : \text{AList}' sm} \\
 \text{let} \quad \frac{\sigma : \text{AList}' m}{\text{targ } \sigma : \mathbb{N}} \quad \text{targ} \quad \text{anil}'_m \mapsto m \\
 \quad \text{targ } (\sigma \text{ acons}' t' / x) \mapsto \text{targ } \sigma
 \end{array}$$

However, this would make types which do not involve targets marginally less complex at the cost of complicating those which do:

$$\frac{}{\text{sub}' \sigma : \text{Fin } m \rightarrow \text{Term } (\text{targ } \sigma)} \quad \frac{\rho : \text{AList}' (\text{targ } \sigma) \quad \sigma : \text{AList}' m}{\rho ++' \sigma : \text{AList}' m}$$

Moreover, commuting interpretation with composition now changes the type of the resultant substitution:

$$\begin{array}{l}
 \text{sub}' (\rho ++' \sigma) : \text{Fin } m \rightarrow \text{Term } (\text{targ } (\rho ++' \sigma)) \\
 \text{sub}' \rho \diamond \text{sub}' \sigma : \text{Fin } m \rightarrow \text{Term } (\text{targ } \rho)
 \end{array}$$

It requires a proof by induction to show that the targets will be equal for all values of ρ and σ . Not all programs exploit the same structural aspects of a datatype. It costs us little to discard explicit structure with \exists ; it costs rather more to recover structure which we have failed to express.

6 First-order unification

We are now ready to write the first-order unification algorithm. It takes two terms s and t and attempts to find a most general unifying substitution f . That is, f must satisfy:

1. $f \Downarrow s = f \Downarrow t$
2. If $g \Downarrow s = g \Downarrow t$ then, for some h , $g \doteq h \diamond f$

In particular, if $s, t : \text{Term } m$, then f should be $\text{sub } \sigma$ where $\sigma : \text{AList } m n$ for some n . Hence we may give unification the type signature

$$\text{let} \quad \frac{s, t : \text{Term } m}{\text{mgu}_m s t : \text{Maybe } (\exists (\text{AList } m))}$$

One might wonder whether this type for mgu is a little tight, in that it apparently forces s and t to “have the same number of variables”, but that is not what the type means. The point is that s and t are expressed over a common set of unknowns, not all (or indeed any) of which they are obliged to mention. It only makes sense to unify expressions from a common language.

Let us employ the usual technique of introducing an accumulator:

$$\begin{array}{l}
 \text{let} \quad \frac{s, t : \text{Term } m \quad \text{acc} : \exists (\text{AList } m)}{\text{amgu}_m s t \text{ acc} : \text{Maybe } (\exists (\text{AList } m))} \\
 \quad \text{mgu}_m s t \mapsto \text{amgu}_m s t \langle m, \text{anil} \rangle
 \end{array}$$

The idea behind the accumulator technique is to proceed *optimistically*, guessing that the unifier is the most general solution to the subproblems examined thus far.

We thus initialize the accumulator with the identity substitution and extend it only when disillusioned. This optimism is justified: once we have a most general unifier for part of the problem, we can show that extending it just enough to solve the rest yields a most general unifier for the whole thing. An abstract justification of *optimistic optimization* can be found in my thesis. However, the correctness proof, available online, illustrates the technique.

As suggested above, **amgu** works by nested structural recursion, first on m , then on s . There is no reason why **amgu** should perform its *case analyses* in that order. Traditional primitive recursors and fold operators make a tight connection between recursion and case analysis, in that they force recursion on some x to be conducted via case analysis on x exactly once and first of all. In fact, **amgu** does case analysis on s and t , and then only if a variable is encountered will it examine acc . If acc is $\langle n, \sigma \text{ asnoc } r/z \rangle$, then we know we can get rid of z . That is, although our outermost recursion is on m , we never examine m directly – the constructor types of AList tell us that nontrivial substitutions have nonempty domains.

I have written **amgu** with *prioritized* patterns, saving space by combining clauses of the case analysis which operate uniformly. The algorithm is fairly straightforward, dividing unification problems into three classes:

- *Rigid-rigid* problems: each term begins with leaf or fork. Like constructors yield subproblems for corresponding subterms; otherwise, unification fails.
- *Flexible* problems: at least one term is a variable and the accumulator holds `anil`. The appropriate helper function solves these immediately.
- *Postponed* problems: at least one term is a variable, but the accumulated σ might instantiate that variable. We peel off one step from σ and expand the terms – removing one variable justifies a recursive call. If successful, we glue the step back on to the resulting substitution.

let $\frac{s, t : \text{Term } m \quad acc : \exists(\text{AList } m)}{\text{amgu}_m s t acc : \text{Maybe } (\exists(\text{AList } m))}$

| | | | | |
|--|---|---------------------------|--|---|
| amgu _{m} | leaf | leaf | acc | $\mapsto \text{yes } acc$ |
| amgu _{m} | leaf | $(t_1 \text{ fork } t_2)$ | acc | $\mapsto \text{no}$ |
| amgu _{m} | $(s_1 \text{ fork } s_2)$ | leaf | acc | $\mapsto \text{no}$ |
| amgu _{m} | $(s_1 \text{ fork } s_2)$ | $(t_1 \text{ fork } t_2)$ | acc | \mapsto |
| | amgu _{m} $s_2 t_2 \nmid \text{amgu}_m s_1 t_1 acc$ | | | |
| amgu _{m} | (ιx) | (ιy) | $\langle m, \text{anil} \rangle$ | $\mapsto \text{yes } (\text{flexFlex}_m x y)$ |
| amgu _{m} | (ιx) | t | $\langle m, \text{anil} \rangle$ | $\mapsto \text{flexRigid}_m x t$ |
| amgu _{m} | t | (ιx) | $\langle m, \text{anil} \rangle$ | $\mapsto \text{flexRigid}_m x t$ |
| amgu _{sm} | s | t | $\langle n, \sigma \text{ asnoc } r/z \rangle$ | \mapsto |
| | $\nmid \lambda \sigma. \sigma \text{ asnoc } r/z \nmid \text{amgu}_m (r \text{ for } z \nmid s) (r \text{ for } z \nmid t) \langle n, \sigma \rangle$ | | | |

$\frac{x, y : \text{Fin } m}{\text{flexFlex}_m x y : \exists(\text{AList } m)}$

flexFlex _{sm} $x y \mapsto$ case **thick** $x y$
 of **yes** $y' \Rightarrow \langle m, \text{anil asnoc } \iota y'/x \rangle$
 no $\Rightarrow \langle sm, \text{anil} \rangle$

$$\frac{x : \text{Fin } m \quad t : \text{Term } m}{\text{flexRigid}_m x t : \text{Maybe } (\exists (\text{AList } m))}$$

$$\begin{aligned} \text{flexRigid}_{sm} x t \mapsto & \text{case } \text{check } x t \\ & \text{of } \text{yes } t' \Rightarrow \text{yes } \langle m, \text{anil asnoc } t' / x \rangle \\ & \text{no} \Rightarrow \text{no} \end{aligned}$$

Prioritizing the patterns has saved us six lines in all: the two lines which call **flexRigid** apply whether the non-variable argument is a leaf or a fork, whilst the last line applies in each of the five *postponed* cases.

flexFlex and **flexRigid** each take a variable x (perforce from a nonempty set) and attempt to express their other argument as an image of $\triangleright \text{thin } x \downarrow$. If successful, we have found a unifier. Otherwise, in the **flexFlex** case, the two variables were already the same, so the identity substitution suffices, and in the **flexRigid** case, the occur-check has failed, so there is no unifier.

The online supplement to this paper (McBride, 2003a) contains a proof that the answer returned by this program, either a substitution or an indication of failure, is correct. It follows the same lines as the traditional proofs of partial correctness found in many papers about first-order unification, but with two modest novelties:

- The step cases of the algorithm are justified by a proof technique which I call *optimistic* optimization. The algorithm works by recursively decomposing a problem and refining an over-general approximation to its solution only when absolutely necessary. The key lemma in the proof shows that this strategy produces most general solutions because unifiers are closed under post-composition.
- The notion of *occurrence* has a *data* representation, inspired by Huet’s notion of “zipper” (Huet, 1997). This allows a more direct characterization of **check**, showing exactly its relationship to unification, and comparing favourably with the collection of variable-counting lemmas previously required.

The major novelty is that there is no need for a termination proof: by expressing the structure which unification exploits, **mgu** becomes a recognizable program in a language of total functions.

7 Discussion

This development of unification is another in a long line. From Zohar Manna and Richard Waldinger’s pioneering hand-synthesis (Manna & Waldinger, 1981), through Larry Paulson’s machine verification in LCF (Paulson, 1985) to the more recent work in diverse proof systems (Coen, 1992; Rouyer, 1992; Jaume, 1997; Bove, 1999), all have faced the same inherent problem of explaining programs which simply could not and did not make all of the sense which their makers had in mind. Understanding the size of the variable signature over which terms are constructed is crucial to the termination of first-order unification. All the above treatments have been forced to fill in this missing structure after the fact, with complex lemmas relating substitution, the occur-check and the function which counts the number of distinct variables used in a term.

Ana Bove's treatment (Bove, 1999), like mine, yields a unification function which is an executable term in a strongly normalising type theory. She sets out to show that a standard Haskell function using general recursion can be imported systematically into type theory, acquiring an extra "termination" argument (a proof of termination for the actual input) and a proof obligation (that there is such a proof for each input). The extra argument belongs to an inductive family whose constructors and indices reflect exactly the case analyses and recursive calls in the algorithm: the function is thus a structural recursion on termination proofs. The original Haskell code survives almost intact, but the real work now lies in the proof obligation, and it is here that the usual variable-counting lemmas reappear. The "program" is still missing its explanation, but Bove's inductive presentation of termination proofs adds it very cleanly, with the bonus of providing exactly the right induction principle with which to prove partial correctness!

A dependent type system thus serves well as an "explanation language" for a standard, simply-typed, general recursive program. However, this paper gives strong evidence that dependent types can be used to improve the language of programs themselves, reducing the amount of subsequent explanation required by making more structure explicit in advance. This is not just the shifting of a comparable burden from proof to program: it is genuinely easier to exploit structure which is present in data than to recover structure which is absent.

What I hope this paper shows is that the apparent need for general recursion is evidence not that structural recursion is a flawed ideal, but rather that we need a better language with which to express structure. With dependent types, we can express our *understanding*, not just our *procedure*. That is the very purpose of declarative programming – to make it more likely that we mean what we say by improving our ability to say what we mean.

Acknowledgements

This work was supported in part by grants from the UK EPSRC (GR/N 24988/01) and from the Leverhulme Trust. The original development was supervised by Rod Burstall, Healfdene Goguen and James McKinna. Many thanks also to Ana Bove, Paul Callaghan and the anonymous referees for helpful discussions, advice and encouragement.

References

- Augustsson, L. (1998) Cayenne – a language with dependent types. *ACM International Conference on Functional Programming*. ACM.
- Bove, A. (1999) *Programming in Martin-löf Type Theory. Unification: A non-trivial Example*. Licentiate Thesis, Chalmers University of Technology, Sweden.
- Burstall, R. (1969) Proving Properties of Programs by Structural Induction. *Computer J.* **12**(1), 41–48.
- Burstall, R. (1987) Inductively Defined Functions in Functional Programming Languages. *J. Comput. Syst. Sci.* **34**, 409–421.
- Coen, M. (1992) *Interactive Program Derivation*. PhD thesis, University of Cambridge.

- Coquand, C. and Coquand, T. (1999) Structured Type Theory. *Workshop on Logical Frameworks and Metalanguages*.
- Dybjer, P. (1994) Inductive families. *Formal Aspects of Comput.* **6**, 440–465.
- Dybjer, P. and Setzer, A. (1999) A finite axiomatization of inductive-recursive definitions. *Proceedings of Typed Lambda Calculi and Applications: LNCS*, pp. 129–146. Springer-Verlag.
- Huet, G. (1997) The Zipper. *J. Functional Program.* **7**(5), 549–554.
- Jaume, M. (1997) Unification: a case study in transposition of formal properties. In: Gunter, E. L. and Felty, A., editors, *Supplementary Proceedings of the 10th International Conference on Theorem Proving in Higher Order Logics: Poster session TPHOLs'97*, pp. 79–93.
- Luo, Z. (1994) *Computation and Reasoning: A Type Theory for Computer Science*. Oxford University Press.
- Luo, Z. and Pollack, R. (1992) *LEGO Proof Development System: User's Manual*. Technical report ECS-LFCS-92-211, Laboratory for Foundations of Computer Science, University of Edinburgh.
- Magnusson, L. (1994) *The implementation of ALF – A Proof Editor based on Martin-Löf's Monomorphic Type Theory with Explicit Substitution*. PhD thesis, Chalmers University of Technology, Göteborg.
- Manna, Z. and Waldinger, R. (1981) Deductive Synthesis of the Unification Algorithm. *Sci. Comput. Program.* **1**, 5–48. North-Holland.
- McBride, C. (1999) *Dependently Typed Functional Programs and their Proofs*. PhD thesis, University of Edinburgh.
- McBride, C. (2003a) *Proving First-Order Unification Correct*. CUP website.
- McBride, C. (2003b) *LEGO development of First-Order Unification*. CUP website.
- Paulson, L. (1985) Verifying the Unification Algorithm in LCF. *Sci. f Comput. Program.* **5**, 143–169.
- Pfenning, F. and Paulin-Mohring, C. (1990) Inductively defined types in the calculus of constructions. In: Main, M., Melton, A., Mislove, M. and Schmidt, D., editors, *Mathematical Foundations of Programming Semantics 1989: LNCS 442*, pp. 209–228. Springer-Verlag.
- Robinson, A. (1965) A Machine-oriented Logic Based on the Resolution Principle. *J. ACM*, **12**, 23–41.
- Rouyer, J. (1992) *Développement de l'algorithme d'unification dans le Calcul des Constructions avec types inductifs*. Technical report 1795, INRIA-Lorraine.
- Turner, D. (1995) Elementary strong functional programming. *Functional Programming Languages in Education, First International Symposium: LNCS 1022*. Springer-Verlag.