# Efficient Functional Unification and Substitution

Atze Dijkstra and Arie Middelkoop and S. Doaitse Swierstra

Department of Information and Computing Sciences,
Universiteit Utrecht,
P.O.Box 80.089,
Padualaan 14, Utrecht, Netherlands,
{atze,ariem,doaitse}@cs.uu.nl,
WWW home page: http://www.cs.uu.nl

**Abstract.** Implementations of language processing systems often use unification and substitution to compute desired properties of input language fragments; for example when inferring a type for an expression. Purely functional implementations of unification and substitution usually directly correspond to the formal specification of language properties. Unfortunately the concise and understandable formulation comes with gross inefficiencies. A seond appoach is to focus on efficiency of implementation. However, efficient implementations of unification and substitution forgo pure functionality and rely on side effects. We present a third, 'best of both worlds', solution, which is both purely functional and efficient by simulating side effects functionally. We compare the three approaches side by side on implementation and performance.

## 1 Introduction

Although unification arises in many problem areas, for example in theorem proving systems and in Prolog implementations, our inspiration for this paper comes from its application in type checking and inferencing in a Haskell compiler (4; 6; 5). In Haskell we may write, for example:

$$first\ (a, b) = a$$
$$x_1 = first\ 3$$
$$x_2 = first\ (3, 4)$$
$$x_3 = first\ ((3, 4), 5)$$

For *first* we need to infer (or reconstruct) its type $\forall a\ b.(a, b) \rightarrow a$, whereas for $x_1, x_2$ and $x_3$ we need to check whether it is permitted to pass the given argument to *first*. Obviously this is not the case for $x_1$.

In implementations of type systems the reconstruction of yet unknown type information and the check whether known types match is usually done with the help of *unification* of types, the unification paradigm being one of many strategies to solve equations on types imposed by the formal specification of a type system. Types may contain type variables representing yet unknown type information; unification then either matches two types, possibly returning new

bindings for such type variables, referred to as *substitution*, or it fails with a type mismatch. For example, for the application of *first* to $(3, 4)$ in the definition of $x_2$ types $(Int, Int)$ and $(v_1, v_2)$ match with bindings for type variables $v_1$ and $v_2$; in the right hand side of the definition for $x_1$ the given argument type $Int$ and expected argument type $(v_1, v_2)$ do not match.

Formally, the unification problem is described as follows (see Knight (11)). We define a *term*, denoted by $\{s, t\}$, to be constructed from *function symbol*s $\{f, g\}$ and *variable symbol*s $\{v, w\}$:

$$t = f\ (t_1, \ldots, t_n)$$
$$\mid\ v$$

Function symbols take a possibly empty sequence of arguments; functions without arguments act as *constant symbol*s.

A *substitution* is a mapping from variables to terms: $\{v_1 \mapsto t_1, \ldots, v_n \mapsto t_n\}$. We wll use $\{\theta, \sigma, \vartheta\}$ to refer to substitutions. A substitution can be extended to a function from terms to terms via its application to terms, denoted by $\theta(t)$ or juxtaposition $\theta t$ when it is clear that substitution application is meant. The term $\theta t$ denotes the term in which each variable $v_i$ in $t$ in the domain of $\theta$ is replaced by $t_i = \theta(v_i)$:

$$\theta(f\ (t_1, \ldots, t_n)) = f\ (\theta t_1, \ldots, \theta t_n)$$
$$\theta(v \qquad\quad) = t, \{v \mapsto t\} \subset \theta$$
$$= v, otherwise$$

Substitutions can be composed: $\sigma\theta t$ denotes $t$ after the application of $\theta$ followed by $\sigma$. The application to a substitution $\theta = \{v_i \mapsto t_i\}$ is defined as $\sigma\theta = \sigma \cup \{v_i \mapsto \sigma t_i\}$. Composition of substitutions is associative, but in general not commutative.

Two terms $s$ and $t$ are *unifiable* if there exists a substitution $\theta$ such that $\theta s = \theta t$. The substitution $\theta$ is then called the *unifier*, $\theta t$ the *unification*. A unifier $\theta$ is called the *most general unifier* (MGU) if for any other unifier $\sigma$, there exists a substitution $\vartheta$ such that $\vartheta\theta = \sigma$. Two terms $s$ and $t$ may be *infinitely unifiable* if their unifier binds variables to infinitely long terms. In this paper we prevent this from happening.

*The problem.* From the above definitions we already can see why a straightforward functional implementation will be inefficient. When we directly translate the definition of the substitution application $\theta t$ to a corresponding function application in a purely functional language like Haskell, each such application will construct a copy $s$ of $t$, differing only in the free $v_i$ for which $\theta$ has a binding. Furthermore, whenever $v_i$ occurs more than once in $t$, several copies of $\theta(v_i)$ will be present in $s$. This leads to duplication of work for a subsequent substitution, a situation which occurs when substitutions are composed. Substitution composition is done frequently; this then makes variable replacement in substitutions the culprit, and thus has to be avoided in more efficient implementations of the subsitution process.

*A solution with side effects and its derived problems.* The growth of terms via duplicate copies of substituted variables can be avoided by never replacing variables. Instead we let variables act as pointers to a possible replacement term. This is easily accomplished in imperative languages, but is more difficult in purely functional ones because of the side effects involved: initially a variable will have no replacement bound to it, and when later a replacement is found for the variable the pointer is made to point to the term replacing the variable.

In a functional language like Haskell we achieve this by leaving the side effect free functional world: the *IO* monad (Haskells imperative environment) and *IORef* s (Haskells pointer mechanism) are then used. This is the approach taken in the GHC (14; 20) by the type inferencer, with the following consequences:

– Side effects infect: term reconstruction (type inferencing) and related functionality all have to be aware of side effects and loose the benefits of pure functions.
– Once updated, a variable is changed forever after. This, for example, complicates the use of backtracking mechanisms that may need to undo substitutions.

How much we suffer from these consequences depends on the necessities of the program using unification. We found ourselves in a situation where we were hindered by the lack of efficiency of the basic functional solution, and did not want to corrupt the cleanliness of our compiler implementation (4; 6; 5). Furthermore we wanted the freedom to experiment with temporary assumptions about type variables, instead of fixing knowledge about such variables in one pass directly. So we designed a third solution which is both functionally transparant and efficient.

*Our contribution: a solution without side effects.* A solution infecting an otherwise functional program with side effects can be avoided by simulating side effects purely functionally. The essence of an efficient substitution mechanism is to share the binding of a variable instead of copying it. This can be implemented without relying on imperative constructs such as *IO* in Haskell. Our contribution thus is:

– Present our side effect free efficient functional unification and substitution.
– Compare our solution with the naive purely functional as well as the side effect solution. We look at both the implementation and performance.

Because of space limitations we focus on the core parts our solution and comparison, leaving a more detailed elaboration to the accompanying technical report (7).

*Related work.* Our work is closely related to *explicit substitution* s (1; 19) in which substitutions are modelled explicitly in $\lambda$-calculus for the same reason as we do, to avoid inefficient duplication of work. Explicit substitution also deals with garbage collection (of term variables), which we do not. On the other hand,

we are not aware of other published work describing a solution for unification and substitution in a practical and functional setting as ours; neither are we aware of side by side presentations with other solutions.

The purely functional solution is frequently used in textbook examples (10; 15), whereas the one with side effects is used when efficiency is important, such as in production quality compilers (14; 20).

Much work has been done on unification, in fact so much that we only mention some entry points into existing literature, amongst which some surveys (11; 2; 9) and seminal work by Robinson (16; 17; 18), Paterson and Wegman (13), and Martelli and Montanari (12).

Observable sharing (3) provides identity of values, allowing equality checking based on this identity. The low level implementation requires side effects, similar to the solution in this paper based on side effects.

The problem we encounter is a consequence of being purely functional. Hiding the problem and its solution can be done by offering unification as a language feature and building the implementation of unification into the language implementation, as done in Prolog and its implementations.

*Outline of the remainder of this paper.* In Section 2 we proceed with the preliminaries for our work, in particular a mini system, formally described, and implemented using the three variants of unification and substitution. In Section 3 we present the purely functional implementation, in Section 4 the one with side effects, and in Section 5 our solution, which we call FUNCTIONAL SHARING in this paper to emphasize the purely functional nature as well as sharing for efficiency. We look at performance results in Section 6and conclude in Section 7.

## 2   Preliminaries

*The essence of the problem: purely functional versus side effects.* A purely functional implementation of substitutions over a term is inefficient because of its repeated and expensive application to terms and other substitutions. Efficient solutions update the term itself, using some update mechanism with a side effect, thereby loosing purely functional behavior. Our solution is to model those side effects explicitly and parameterize functions with such reified side effects.

*Experimental environment.* Our experimental environment consists of an implementation resembling structures found in many compilers. We thus mimic the actual runtime environment we are interested in, while keeping things as simple as possible. Fig. 1 shows the algorithmic rules for our system (we omitted the declarative variant because of space limitations); it should be familiar to those acquainted with type systems. Since we want to focus on unification mechanisms without wandering off to type systems, our example system neutrally specifies which values *Val* are to be associated with a tree *Tree*.

A *Tree* offers constructs for binding and using program identifiers, as well as constructing and deconstructing pairs of (ultimately) some constant. The

$$\boxed{\theta_{\mathbf{in}}; \Gamma \vdash Tree : Val \rightsquigarrow \theta_{out}}$$

$$\frac{}{\theta; \Gamma \vdash C : c \rightsquigarrow \theta} \text{ T.CON}_A \qquad \frac{(n \mapsto v) \in \Gamma}{\theta; \Gamma \vdash n : \theta \; v \rightsquigarrow \theta} \text{ T.USEB}_A$$

$$\frac{\begin{array}{c} \theta; \Gamma \vdash x : v \rightsquigarrow \theta_x \\ \theta_x; n \mapsto v, \Gamma \vdash y : w \rightsquigarrow \theta_y \end{array}}{\theta; \Gamma \vdash \mathbf{bind} \; n = x \; \mathbf{in} \; y : w \rightsquigarrow \theta_y} \text{ T.DEFB}_A \qquad \frac{\begin{array}{c} \theta; \Gamma \vdash x : v \rightsquigarrow \theta_x \\ \theta_x; \Gamma \vdash y : w \rightsquigarrow \theta_y \end{array}}{\theta; \Gamma \vdash (x, y) : (\theta_y v, w) \rightsquigarrow \theta_y} \text{ T.TUP}_A$$

$$\frac{\begin{array}{c} \theta; \Gamma \vdash x : vw \rightsquigarrow \theta_x \\ v, w \; \text{fresh} \\ (v, w) \equiv vw \rightsquigarrow \theta_m \end{array}}{\theta; \Gamma \vdash \mathbf{fst} \; x : \theta_m \theta_x v \rightsquigarrow \theta_m \theta_x} \text{ T.FST}_A \qquad \frac{\begin{array}{c} \theta; \Gamma \vdash x : vw \rightsquigarrow \theta_x \\ v, w \; \text{fresh} \\ (v, w) \equiv vw \rightsquigarrow \theta_m \end{array}}{\theta; \Gamma \vdash \mathbf{snd} \; x : \theta_m \theta_x w \rightsquigarrow \theta_m \theta_x} \text{ T.SND}_A$$

**Fig. 1.** Rules for Val of Tree (A)

concrete syntax is included in comment, the exclamation mark enforces strictness and can be ignored for the purpose of understanding:

```
data Tree                      -- concrete syntax:
  = Constant                   -- C
  | UseBind  String            -- n
  | DefBind  String Tree Tree  -- bind n = x in y
  | Tuple    Tree    Tree      -- (x,y)
  | First    Tree              -- fst x
  | Second   Tree              -- snd x
```

The rules associate a *Val* with a *Tree*. Again, a *Val* is inspired by type systems, but for the purposes of this paper it is just some structure, complex enough to discuss unification and substitution. Therefore, in the remainder of this paper a *Val* is a term participating in unification and substitution.

```
data Val              -- concrete syntax:
  = Pair  Val  Val    -- (v,w)
  | Const             -- c
  | Var   VarId
  | Err   String
type VarId = Int
```

A *Val* has two alternatives in its structure which do not have a *Tree* constructor as counterpart: a construct *Var* for encoding variables as used in unification and substitution, and a construct *Err* for signalling errors.

*Test examples.* For example, with the following tree:

```
bind v₁ = C                in
bind v₂ = (v₁, v₁)          in
bind v₃ = (snd v₂, fst v₂) in v₃
```

the rules associate the value $(c, c)$. This example is one of the test cases we use, where we also vary in the number of bindings similar to $v_3$. The value of the tree is always $(c, c)$.

The second example we use for testing infers a *Val* of exponential size in terms of the number of bindings similar to $v_4$, yielding values $((c, c), ((c, c), (c, c)))$ and so forth for increasing numbers of similar bindings:

```
bind v₁ = C                    in
bind v₂ = (v₁, v₁)              in
bind v₃ = (fst v₂, v₂)          in
bind v₄ = (snd v₃, (v₂, v₂)) in v₄
```

The first example provides typical programming language input, with many small definitions, whereas the second example provides a worst case scenario. We label the tests respectively LINEAR and EXPONENTIAL.

*The Val computation algorithm.* Typical of algorithmic rules restrictions on *Val*'s are enforced by unification. For example, the argument of **fst** is constrained to have a *Val* of the form $(v, w)$. To achieve this, the argument and a *Val* of the form $(v, w)$ are unified. The constraining *Val* is built from variables guaranteed to be unique (called *fresh*), whereas the extraction is done by simply using the unique variables together with a substitution $\theta$ holding possible additional information about the variables.

The algorithm threads a substitution $\theta$ through its computation, while gathering information about the *Var*s participating in the construction of the *Val* associated with the root of the tree. The rules maintain the invariant that $\theta$ is already taken into account in resulting $t$'s, that is $\theta t = t$, where $t$ refers to the *Val* component of the conclusion.

A substitution $\theta$ is represented by a variable mapping *VMp*, mapping identifiers *VarId* of variables to terms *Val*:

**newtype** *VMp* = *VMp* (*Map VarId Val*)

We need the usual functions for constructing and querying which we assume to be self-explanatory.

Contextual information $\Gamma$ holding assumptions for program identifiers is encoded by an environment *Env*:

**newtype** *Env* = *Env* (*Map String Val*)

We omit definitions for functions on *Env* and assume their names are understandable enough to indicate their meaning.

Finally, in the following we restrict ourselves to first order unification, and do not allow infinite values.

## 3   Substitution by copying

We first discuss the purely functional reference implementation to which we compare the others. We present the overall computational structure on which we vary in the subsequent alternate implementations. We label this solution by FUNCTIONAL.

The rules in Fig. 1 strongly suggest a direction in which information flows over a tree, upward or synthesized for e.g. *Val*, downward or inherited for e.g. $\Gamma$, and chained for $\theta$. We use a state monad to encode this flow:

```
data St = St{ stUniq :: ! VarId
            , stEnv   :: ! Env
            , stVMp :: ! VMp
            }
type Compute v = State St v
```

The *Compute* state monad threads the following three values through the computation:

- a counter used for creating fresh variables,
- an environment *Env* holding $\Gamma$,
- and a variable mapping *VMp* corresponding to both the inherited and synthesized substitution $\theta$.

*Substituting.* In a *Val* substitutable variables may occur, and thus also in *Env*. Substitutabilty is expressed by the class *Substitutable*:

```
class Substitutable x where
  (|@) :: VMp → x → x
  ftv  :: x → Set VarId
```

The application $\theta x$ of a substitution $\theta$ to some $x$ is expressed by the function |@. The function *ftv* computes the free variables of a $x$, however we omit its implementations.

Substitution over a *Val* is straightforwardly encoded as a recursive replacement:

```
instance Substitutable Val where
  s |@ v
     = sbs s v
       where sbs s (Pair v w) = Pair (sbs s v) (sbs s w)
             sbs s v@(Var i)  = case i |? s of
                                    Just v' → v'
                                    _        → v
             sbs s v          = v
```

The composition of two substitutions, that is, substituting over a substitution itself means taking the union of two *VMp*s and ensuring that all *Val*s in the

```
valUnify :: Val → Val → Compute Val
valUnify v w
  = uni v w where
  uni    v@(Const) (Const)        = return v
  uni    v@(Var i) (Var j) | i == j = return v
  uni    (Var i)     w            = bindv i w
  uni    v           w@(Var _)    = uni w v
  uni    (Pair p q) (Pair r s)    =
    do pr   ← uni p r
       st1  ← get
       qs   ← uni (stVMp st1 |@ q) (stVMp st1 |@ s)
       st2  ← get
       return (Pair (stVMp st2 |@ pr) qs)
  uni    _           _            = err "fail"
  bindv i v
    | Set.member i (ftv v)        = err "inf"
    | otherwise                   =
        do st ← get
           put (st{stVMp = vmUnit i v |@ stVMp st})
           return v
  err x = return (Err x)
```

**Fig. 2.** Val unification in the FUNCTIONAL solution

previous substitution are substituted over as well, the previous substitution being the second operand to |@:

**instance** *Substitutable VMp* **where**
  $s |@ (VMp\ m) = s\ `vmUnion`\ VMp\ (Map.map\ (s |@)\ m)$

Applying the |@ from this instance over and over again makes the update of a substitution with new bindings for variables a costly operation, and alone is responsible for a major part of the efficiency loss of this solution.

*Value unification.* Unification tells us whether two values can be made syntactically equal, and a substitution tells us which variables in these values have to be bound to another value to make this happen. Fig. 2 shows the code for *valUnify*, which unifies two *Val*s. Function *valUnify* applied to $t$ and $s$ yields the unification $\theta t$ directly and the substitution $\theta$ via the state of *Compute*. A unification may also fail, which we simply signal by the *Err* alternative of *Val*.

We note that always returning the unification $\theta t$ is convenient but strictly not necessary, as $\theta$ and $t$ can also be combined outside *valUnify*. Now additional *Val*s are constructed, however, we could not observe an effect on performance (see Section 6 for further discussion). Encoding an error as part of *Val* is also a matter of convenience, and merely to show where errors arise; we do not report those errors and in our test cases no errors arise.

The function *valUnify* assumes that its *Val* parameters do not contain free variables bound by the substitution *stVMp* passed via the *Compute* state. When-

ever a variable is encountered during the comparison of the two types being unified, it is bound to the other comparand. We prevent recursive bindings causing infinite values, like $v \mapsto (v, v)$, from occurring by performing the so called *occurs check* done in *bindv*, and by checking on the trivial unification of $v$ with $v$.

Unification proceeds recursively over *Pair*s. We ensure the invariant that *Val*s passed for further comparison always have the most recent substitution already applied to them.

*Fresh variables.* Besides the environment and the current substitution, the state *St* contains a counter for the generation of fresh variables. Function *newVar* increments the counter *stUniq* in the *Compute* state and returns *Var*s with unique *VarId*s, *newVars* conveniently returns a group of such variables:

> *newVar* :: *Compute Val*
> *newVars* :: *Int* → *Compute* [ *Val* ]

*Computing a Val over a Tree.* Of the computation of a *Val* over a *Tree* we show the implementation for rule T.FST:

> *First* $x$ →
>   **do** *vw*      ← *treeCompute x*
>        [ *v, w* ] ← *newVars* 2
>        *valUnify* (*Pair v w*) *vw*
>        *st*      ← *get*
>        *return* (*stVMp st* |@ *v*)

We closely follow the rule by recursing over the $x$ component of **fst** $x$, allocating fresh variables, using these to match the value of $x$ and returning its first component with the most recent substitution applied. We also slightly deviate from the rule by threading the full *Compute* state through *valUnify* instead of computing additional bindings only.

This completes our basic reference implementation, often used for its simplicity in explanations, but avoided in real world systems because of the time and memory spent in copying and substituting over the content pointed to by variables.

## 4   Substitution by sharing

We can avoid the copying of *Val*s during substitution in the previous solution by sharing the content bound to variables. Variables become pointers[1] in a directed acyclic graph (DAG) representation of *Val* instead of a tree representation as used by the FUNCTIONAL solution (13). We use an *IORef* to encode such a pointer (14), with utility functions like *newRef* for hiding its use. Note that

---

[1] We still need the *VarId* fields because of the computation of *ftv* returning a *Set*; *IORef* is not an instance of *Ord* required for *Set*.

*refRead* is not returning a *Compute* monad; a tricky point we come back to at the end of this section. We label this solution SHARING.

```
data Val                 -- concrete syntax:
  = Pair  Val    Val   -- (v,w)
  | Const              -- c
  | Var   VarId Ref
  | Err   String
type     RefContent = Maybe Val
newtype Ref          = Ref (IORef RefContent)
data St = St{ stUniq ::  VarId
            , stEnv  ::  Env
            }
type Compute v = StateT St IO v
newRef  :: Compute Ref
refRead :: Ref → RefContent
refWrite :: Ref → RefContent → Compute ()
newRef = do r ← lift $ newIORef Nothing
               return (Ref r)
```

In essence, we now store the substitution which maps variables to values directly in a *Var*. Hence we do not need the *VMp* in the *Compute* state anymore. On the other hand, we need to combine the *State* monad with the *IO* monad because of the use of *IORef*. A fresh variable now also gets a fresh shared memory location *Ref*, initialized to hold nothing.

Unification now has to be aware that variables are pointers: the SHARING solution is presented in Fig. 3. Relative to the FUNCTIONAL solution we need to modify the following:

- When comparing a variable *Var* we no longer can assume that the variable is still unbound. Hence we need to inspect its *Ref* and use it for further comparison.
- Binding a variable in *bindv* now also involves updating the reference with the bound value.
- There is no *VMp* threaded through the *Compute* state, hence we need not maintain the invariant that it is always applied, for example when comparing *Pair*s. This is now guaranteed via the *Ref* mechanism.

The implementation of *treeCompute* becomes simpler, because we need not apply the *VMp* here either. As before, we highlight the *First* case branch for rule T.FST; also for the other alternatives the only difference with the FUNC-TIONAL solution is the removal of the application of *VMp*.

```
First x →
  do vw     ← treeCompute x
     [v, w] ← newVars 2
```

```
valUnify :: Val → Val → Compute Val
valUnify v w
  = uni v w where
  uni   v@(Const)   (Const)                = return v
  uni   v@(Var i _) (Var j _) | i == j     = return v
  uni   (Var _ r)    w        | isJust mbv = uni v' w
    where mbv = refRead r
          v'  = fromJust mbv
  uni   v@(Var _ _) w                      = bindv v w
  uni   v            w@(Var _ _)           = uni w v
  uni   (Pair p q)   (Pair r s)            =
    do pr   ← uni p r
       qs   ← uni q s
       return (Pair pr qs)
  uni   _            _                     = err "fail"
  bindv (Var i r) v
    | Set.member i (ftv v)                 = err "inf"
    | otherwise                            =
        do refWrite r (Just v)
           return v
  err x = return (Err x)
```

**Fig. 3.** Val unification in the SHARING solution

```
valUnify (Pair v w) vw
return v
```

The substitution mechanism is completely hidden as a side effect throughout the *Compute* state.

Finally, when computing free variables one also has to be aware of *Ref*s. Function *ftv* has to recurse over the value bound to such a *Ref* to compute its free variables.

The price we have to pay for this solution is that we only may have at most one binding for a *Var*, the one stored in the *Var* itself. This is problematic if we want to have more than one binding during the computation, for example when we want to compute a tentative value and later backtrack on it (5; 8). We have lost the parameterizability of the binding by introducing side effects and giving up purely functional behavior of substitutions.

The use of *IORef* has other, more subtle, consequences typical of the use of monads, which we only mention briefly here because of space limitations. We are forced to either use and hide the side effect of *unsafePerformIO* in *refRead*, as done in Fig. 3, or make this side effect visible and force the *Var* case alternative to be aware of it, thereby requiring a significant rewrite because the use of *IO* necessitates a too early commit to the *Var* case alternative.

```
valUnify :: Val → Val → Compute Val
valUnify v w
   = do { st ← get; uni st v w } where
   uni st v@(Const) (Const)            = return v
   uni st v@(Var i)  (Var j) | i == j  = return v
   uni st (Var i)     w      | isJust mbv = uni st v' w
        where mbv = i |? stVMp st
              v'  = fromJust mbv
   uni st (Var i)     w                = bindv st i w
   uni st v           w@(Var _)        = uni st w v
   uni st (Pair p q) (Pair r s)        =
      do pr  ← uni st p r
         st2 ← get
         qs  ← uni st2 q s
         return (Pair pr qs)
   uni _ _            _                = err "fail"
   bindv st i v                        =
      do put (st{ stVMp = vmUnit i v |@ stVMp st })
         return v
   err x = return (Err x)
```

<div align="center">

**Fig. 4.** Val unification in the FUNCTIONAL SHARING solution

</div>

## 5   Substitution by functional shared memory

We regain purely functional behavior of the unification and substitution machinery by letting a *Var* itself –once again– be unaware of its content, and thus decouple it from the particular baked-in way *IORef*s implement the notion of pointers to memory content. Instead we implement our own dereferencing mechanism by combining *VMp*s from the FUNCTIONAL solution with the pointer based approach of the SHARING solution. We use the *Val* definition of the FUNCTIONAL solution, and adapt the *valUnify* function of the SHARING solution: instead of *IORef*s we create 'do it yourself' memory in the *VMp* as shown in Fig. 4. The key difference with SHARING is that the dereferencing required for a variable now is implemented via a lookup in the threaded *stVMp*. The key commonality with SHARING is that we do not replace a variable; we do not apply the substution to a variable but only use the variable itself.

   We now also can avoid the expensive copying because we follow pointers instead of accessing a copied value directly. The implementation of the *Substitutable VMp* instance no longer needs to update the 'previous' VMp, a subtle but most effective memory saving change:

```
instance Substitutable VMp where
   s |@ s₂ = s 'vmUnion' s₂
```

*Dereferencing and infinite values.* The consequence of derefencing via a table lookup is a performance loss because such a lookup is expensive compared to

| test | depth | FUNCTIONAL | | SHARING | | FUNCTIONAL SHARING | | FUNCTIONAL SHARING NO TOP SUBST | | FUNCTIONAL SHARING OCCUR CHECK | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | sec | Mb | sec | Mb | sec | Mb | sec | Mb | sec | Mb |
| LINEAR | 500 | 0.67 | 61.7 | 0.07 | 1.8 | 0.03 | 2.8 | 0.04 | 2.8 | 0.52 | 2.9 |
| | 1100 | 4.10 | 391.3 | 0.30 | 3.2 | 0.08 | 5.5 | 0.10 | 5.5 | 3.14 | 5.8 |
| | 1600 | 8.60 | 687.5 | 0.63 | 4.9 | 0.13 | 7.4 | 0.14 | 7.4 | 7.67 | 7.5 |
| EXPONENTIAL | 20 | 0.04 | 4.4 | 0.00 | 1.3 | 0.01 | 2.1 | 0.00 | 1.3 | 0.01 | 1.3 |
| | 25 | 0.89 | 60.7 | 0.11 | 1.3 | 0.21 | 13.7 | 0.09 | 1.3 | 0.08 | 1.3 |
| | 28 | 5.63 | 438.7 | 0.58 | 1.3 | 1.38 | 107.7 | 0.42 | 1.3 | 0.44 | 1.3 |

**Fig. 5.** Performance results

a plain memory dereference. Both *valUnify* and its use of *ftv* now require such table lookups. Our design choice is to avoid excessive dereferencing by not using *ftv* at all during unification, and consequently omitting the occurs check from unification. In turn this means that unification may return a substitution with cycles, and we have to deal with infinite values and the occurs check elsewhere, that is, all functions traversing a *Val* need to be aware that an infinite value may indirectly occur via a substitution.

For example, we need to check during application of a substitution to a *Val*. We adapt the application of a substitution to a *Val* to implement the occurs check: we return an error whenever a substitution for a variable occurs twice.

Actually, the necessity for such a check depends on the context in which unification and substitution are used. In this case we could have done without the check because a binding for a variable leading to an infinite value, like $v \mapsto (v, v)$, only arises when we would have had recursive references to bindings in the *Tree* language. We come back to its effect on performance (by putting the occurs check back into *valUnify*) when discussing performance (Section 6).

Finally, for the result to be usable without being aware of a *VMp*, we apply the substitution outside *treeCompute*, in the toplevel test function. For example, our pretty printing is unaware of a *VMp*. Again, we come back to this because of its degrading effect on performance.

## 6 Performance results

We compared the three solutions, FUNCTIONAL, SHARING and FUNCTIONAL SHARING, by running two test trees, LINEAR and EXPONENTIAL, with various depths. Both tests are described in Section 2 and are characterized by manipulation of *Val*s, linear and exponential in the number of bindings introduced (which equals the depth of the tree) by the test *Tree*s. The results are shown in Fig. 5. The FUNCTIONAL, SHARING and FUNCTIONAL SHARING variants are already described; the remaining variants are introduced and discussed hereafter as part of

the performance analysis. The memory sizes in Fig. 5 correspond to the maximum resident set size as reported by the Unix time command, and is because of the GHC garbage collection an overestimate of the actual memory requirements. However, it still gives an indication of the proportial memory use. Tests were run on a MacBook Pro 2.2Ghz Intel Core 2 Duo with 2GB memory, MacOS X 10.5.4, the programs compiled with GHC 6.8.2 without optimization flags. Each test was run twice, the results taken from the second run. Further runs did not give significant variation in the results.

We observe the following:

– On the linear test cases all but the FUNCTIONAL variant perform equally well, using small amounts of memory.
– On the exponential test case the SHARING variant runs best, the FUNCTIONAL variant worst, especially in terms of memory. The FUNCTIONAL SHARING variant sits in between. It turned out this was caused by the substitution still applied in the toplevel test function. Variant FUNCTIONAL SHARING NO TOP SUBST has this substitution removed and replaced by code forcing a deep evaluation over the *Val* and substitution jointly. The results are now similar to those of SHARING, even a bit faster.
– Omitting the occurs check in FUNCTIONAL SHARING is worthwhile. Variant FUNCTIONAL SHARING OCCUR CHECK includes the occurs check relative to the fastest variant FUNCTIONAL SHARING NO TOP SUBST: it is significantly slower for the linear test. This is an apparent trade-off between efficiency and responsibility of doing the occurs check: encapsulated in unification or outside of unification. Carrying the 'occurs check' responsibility implies additional program complexity, but, in the light of variant FUNCTIONAL SHARING NO TOP SUBST, no loss of efficiency. We did not further experiment and measure this. In our real world use (6) of our solution only a limited number of functions is aware of substitutions, yielding a sufficient gain in efficiency.

## 7  Conclusion

Avoiding copying and the resulting memory allocation, and using sharing mechanisms instead, pays off. This is the overall conclusion which can be drawn. Furthermore, using a solution with *IORef* based side effect can be avoided without performance penalties; there is no need to fall back on the *IO* monad to achieve acceptable levels of performance.

Our 'best of both worlds' solution has been implemented as part of EHC, the essential Haskell Compiler (4; 6); The programs discussed in this paper can also be found there as part of its experiments subdirectory. Because we have based our EHC implementation on attribute grammars, avoiding the dependency on *IO* and side effects, the efficient functional solution was critical to the success of the implementation of the type system in EHC.

# Bibliography

[1] M. Abadi, L. Cardelli, and P-L. Curien. Explicit Substitutions. *Journal of Functional Programming*, 1(4):375–416, Oct 1991.

[2] Franz Baader and Wayne Snyder. *Unification Theory*, chapter 8, pages 447–531. Elsevier Science Publishers, 2001.

[3] Koen Claessen and David Sands. Observable Sharing for Functional Circuit Description. In *Asian Computing Science Conference*, number 1742 in LNCS, pages 62–73, 1999.

[4] Atze Dijkstra. EHC Web. `http://www.cs.uu.nl/groups/ST/Ehc/WebHome`, 2004.

[5] Atze Dijkstra. *Stepping through Haskell*. PhD thesis, Utrecht University, Department of Information and Computing Sciences, 2005.

[6] Atze Dijkstra, Jeroen Fokker, and S. Doaitse Swierstra. The Structure of the Essential Haskell Compiler, or Coping with Compiler Complexity. In *Implementation of Functional Languages*, 2007.

[7] Atze Dijkstra, Arie Middelkoop, and S. Doaitse Swierstra. Efficient Functional Unification and Substitution. `http://www.cs.uu.nl/wiki/Atze/WebHome`, 2008.

[8] Atze Dijkstra and S. Doaitse Swierstra. Exploiting Type Annotations. Technical Report UU-CS-2006-051, Department of Computer Science, Utrecht University, 2006.

[9] Jean H. Gallier. Unification Procedures in Automated Deduction Methods Based on Matings: A Survey. Technical Report MS-CIS-91-76, University of Pennsylvania, Department of Computer and Information Science, 1991.

[10] Mark P. Jones. Typing Haskell in Haskell. In *Haskell Workshop*, 1999.

[11] Kevin Knight. Unification: a multidisciplinary survey. *ACM Computing Surveys*, 21(1):93–124, Mar 1989.

[12] Alberto Martelli and Ugo Montanari. An Efficient Unification Algorithm. *ACM TOPLAS*, 4(2):258 – 282, April 1982.

[13] M.S. Paterson and M.N. Wegman. Linear unification. *Journal of Computer and System Sciences*, 16(2):158–167, Apr 1978.

[14] Simon Peyton Jones and Mark Shields. Practical type inference for arbitrary-rank types, 2004.

[15] Benjamin C. Pierce. *Types and Programming Languages*. MIT Press, 2002.

[16] J.A. Robinson. A machine-oriented logic based on the resolution principle. *Journal of the ACM*, 12(1):23–41, Jan 1965.

[17] J.A. Robinson. Computational logic: The unification computation. In B. Meltzer and D. Michie, editors, *Machine Intelligence*, pages 63–72, 1971.

[18] J.A. Robinson. Fast unification. In *Theorem Proving Workshop*, 1976.

[19] Kristoffer H. Rose. Explicit Substitution - Tutorial and Survey. Technical Report BRICS-LS-96-3, Department of Computer Science, University of Aarhus, 1996.

[20] GHC Team. GHC Developer Wiki. `http://hackage.haskell.org/trac/ghc`, 2007.