

Bottom-up β -reduction: uplinks and λ -DAGs

Olin Shivers¹ and Mitchell Wand²

¹ Georgia Institute of Technology

² Northeastern University

Abstract. Representing a λ -calculus term as a DAG rather than a tree allows us to represent the sharing that arises from β -reduction, thus avoiding combinatorial explosion in space. By adding uplinks from a child to its parents, we can efficiently implement β -reduction in a bottom-up manner, thus avoiding combinatorial explosion in time required to search the term in a top-down fashion. We present an algorithm for performing β -reduction on λ -terms represented as uplinked DAGs; discuss its relation to alternate techniques such as Lamping graphs, explicit-substitution calculi and director strings; and present some timings of an implementation. Besides being both fast and parsimonious of space, the algorithm is particularly suited to applications such as compilers, theorem provers, and type-manipulation systems that may need to examine terms in-between reductions—*i.e.*, the “readback” problem for our representation is trivial. Like Lamping graphs, and unlike director strings or the suspension λ -calculus, the algorithm functions by side-effecting the term containing the redex; the representation is *not* a “persistent” one. The algorithm additionally has the charm of being quite simple: a complete implementation of the core data structures and algorithms is 180 lines of SML.

1 Introduction

The λ -calculus [2, 5] is a simple language with far-reaching use in the programming-languages and formal-methods communities, where it is frequently employed to represent, among other objects, functional programs, formal proofs, and types drawn from sophisticated type systems. Here, our particular interest is in the needs of client applications such as compilers, which may use λ -terms to represent both program terms as well as complex types. We are somewhat less focussed on the needs of graph-reduction engines, where there is greater representational license—a graph reducer can represent a particular λ -term as a chunk of machine code (*e.g.*, by means of supercombinator extraction), because its sole focus is on *executing* the term. A compiler, in contrast, needs to examine, analyse and transform the term in-between reductions, which requires the actual syntactic form of the term be available at the intermediate steps.

Of the three basic operations on terms in the λ -calculus— α -conversion, β -reduction, and η -reduction—it is β -reduction that accomplishes the “heavy lifting” of term manipulation. (The other two operations are simple to implement.) Unfortunately, naïve implementations of β -reduction can lead to exponential time and space blowup.

To appear in *Proceedings of the European Symposium on Programming (ESOP)*, Edinburgh, Scotland, April 2005.

There are only three forms in the basic language: λ expressions, variable references, and applications of a function to an argument:

$$t \in \text{Term} ::= \lambda x.t \mid x \mid t_f t_a$$

where “ x ” stands for a member of some infinite set of variables.

β -reduction is the operation of taking an application term whose function subterm is a λ -expression, and substituting the argument term for occurrences of the λ 's bound variable in the function body. The result, called the *contractum*, can be used in place of the original application, called the *redex*. We write

$$(\lambda x.b) a \Rightarrow [x \mapsto a]b$$

to express the idea that the redex applying function $\lambda x.b$ to argument a reduces to the contractum $[x \mapsto a]b$, by which we mean term b , with free occurrences of x replaced with term a .

We can define the core substitution function with a simple recursion:

$$\begin{aligned} [y \mapsto t][x] &= t & x = y \\ [y \mapsto t][x] &= x & x \neq y \\ [x \mapsto t][t_f t_a] &= ([x \mapsto t]t_f)([x \mapsto t]t_a) \\ [x \mapsto t][\lambda y.b] &= \lambda y'.([x \mapsto t][y \mapsto y']b) & y' \text{ fresh in } b \text{ and } t. \end{aligned}$$

Note that, in the final case above, when we substitute a term t under a λ -expression $\lambda y.b$, we must first replace the λ -expression's variable y with a fresh, unused variable y' to ensure that any occurrence of y in t isn't “captured” by the $[x \mapsto t]$ substitution. If we know that there are no free occurrences of y in t , this step is unnecessary—which is the case if we adopt the convention that every λ -expression binds a unique variable.

It is a straightforward matter to translate the recursive substitution function defined above into a recursive procedure. Consider the case of performing a substitution $[y \mapsto t]$ on an application $t_f t_a$. Our procedure will recurse on both subterms of the application... but we could also use a less positive term in place of “recurse” to indicate the trouble with the algorithmic handling of this case: search. In the case of an application, the procedure will blindly search *both* subterms, even though one or both may have no occurrences of the variable for which we search. Suppose, for example, that the function subterm t_f , is very large—perhaps millions of nodes—but contains no occurrences of the substituted variable y . The recursive substitution will needlessly search out the entire subterm, constructing an identical copy of t_f . What we want is some way to direct our recursion so that we don't waste time searching into subterms that do not contain occurrences of the variable being replaced.

2 Guided tree substitution

Let's turn to a simpler task to develop some intuition. Consider inserting an integer into a set kept as an ordered binary tree (Fig. 1). There are three things about this simple algorithm worth noting:

```

Procedure addItem(node, i)
  if node = nil then
    new := NewNode()
    new.val := i
    new.left := nil
    new.right := nil
  else if node.val < i then
    new := NewNode()
    new.val := node.val
    new.left := node.left
    new.right := addItem(node.right, i)
  else if node.val > i then
    new := NewNode()
    new.val := node.val
    new.right := node.right
    new.left := addItem(node.left, i)
  else new := node
  return new

```

Fig. 1. Make a copy of ordered binary tree *node*, with added entry *i*. The original tree is not altered.

– **No search**

The pleasant property of ordered binary trees is that we have enough information as we recurse down into the tree to proceed only into subtrees that require copying.

– **Steer down; build up**

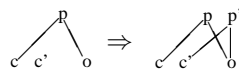
The algorithm’s recursive control structure splits decision-making and the actual work of tree construction: the downward recursion makes the decisions about which nodes need to be copied, and the upward return path assembles the new tree.

– **Shared structure**

We copy only nodes along the spine leading from targeted node to the root; the result tree shares as much structure as possible with the original tree.

3 Guiding tree search with uplinks

Unfortunately, in the case of β -reduction, there’s no simple, compact way of determining, as we recurse downwards into a tree, which way to go at application nodes—an application has two children, and we might need to recurse into one, the other, both, or neither. Suppose, however, that we represent our tree using not only down-links that allow us to go from a parent to its children, but also with redundant up-links that allow us to go from a child to its parent. If we can (easily) find the leaf node in the original tree we wish to replace, we can chase uplinks along the spine from the old leaf to the tree root, copying as we go (Fig. 2). The core iteration of this algorithm is the $c \mapsto c'$ upcopy:



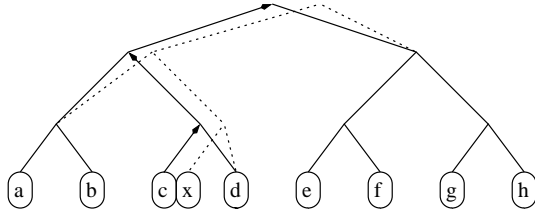


Fig. 2. Replacing a single leaf in a binary tree by following uplinks. Here, we make a copy of the original tree, replacing leaf c with x .

We take a child c and its intended replacement c' , and replicate the parent p of c , making the $c \mapsto c'$ substitution. This produces freshly-created node p' ; we can now iterate, doing a $p \mapsto p'$ upcopy into the parent of p at the next step, and so on, moving up through the original tree until we reach the root.

Note the similar properties this upcopy algorithm has with the previous algorithm: no search required; we build as we move upwards; we share as much structure as possible with the old tree, copying only the nodes along the “spine” leading from the leaf back to the root. For a balanced tree, the amount of copying is logarithmic in the total number of nodes. If we can somehow get our hands on the leaf node to be replaced in the old tree, the construction phase just follows uplinks to the root, instead of using a path saved in the recursion stack by the downwards search.

4 Upcopy with DAGs

We can avoid space blowup when performing β -reduction on λ -calculus terms if we can represent them as directed acyclic graphs (DAGs), not trees. Allowing sharing means that when we substitute a large term for a variable that has five or six references inside its binding λ -expression, we don’t have to create five or six distinct copies of the term (that is, one for each place it occurs in the result). We can just have five or six references to the same term. This has the potential to provide logarithmic compression on the simple representation of λ -calculus terms as trees. These term DAGs can be thought of as essentially a space-saving way to represent term trees, so we can require them, like trees, to have a single top or root node, from which all other nodes can be reached.

When we shift from trees to DAGs, however, our simple functional upcopy algorithm no longer suffices: we have to deal with the fact that there may be multiple paths from a leaf node (a variable reference) of our DAG up to the root of the DAG. That is, any term can have multiple parents. However, we can modify our upwards-copying algorithm in the standard way one operates on DAGs: we search upwards along all possible paths, marking nodes as we encounter them. The first time we copy up into a node n , we replicate it, as in the previous tree algorithm, and continue propagating the copy operation up the tree to the (possibly multiple) parents of n . However, before we move upwards from n , we first store the copy n' away in a “cache” field of n . If we later copy up into n via its other child, the presence of the copy n' in the cache slot of n will signal the algorithm that it should not make a second copy of n , and should not proceed

```

Procedure upcopy(childcopy, parent, relation)
  if parent.cache is empty then
    parcopy := NewNode()
    if relation is "left child" then
      parcopy.left := childcopy
      parcopy.right := parent.right
    else
      parcopy.right := childcopy
      parcopy.left := parent.left
    parent.cache := parcopy
    for-each <grandp,gprel> in parent.uplinks do
      upcopy(parcopy, grandp, gprel)
  else
    parcopy := parent.cache
    if relation is "left child"
    then parcopy.left := childcopy
    else parcopy.right := childcopy

```

Fig. 3. Procedure upcopy makes a copy of a binary DAG, replacing the *relation* child (left or right) of *parent* with *childcopy*.

upwards from n —that has already been handled. Instead, it mutates the existing copy, n' , and returns immediately.

The code to copy a binary DAG, replacing a single leaf, is shown in Fig. 3. Every node in the DAG maintains a set of its uplinks; each uplink is represented as a $\langle parent, relation \rangle$ pair. For example, if node c is the left child of node p , then the pair $\langle p, \text{left-child} \rangle$ will be one of the elements in c 's uplink set.

The upcopy algorithm explores each edge on all the paths from the root of the DAG to the copied leaf exactly once; marking parent nodes by depositing copies in their cache slots prevents the algorithm from redundant exploration. Hence this graph-marking algorithm runs in time proportional to the number of edges, *not* the number of paths (which can be exponential in the number of edges). Were we to “unfold” the DAG into its equivalent tree, we would realise this exponential blowup in the size of the tree, and, consequently, also in the time to operate upon it. Note that, analogously to the tree-copying algorithm, the new DAG shares as much structure as possible with the old DAG, only copying nodes along the spine (in the DAG case, spines) from the copied leaf to the root.

After an upcopy has been performed, we can fetch the result DAG from the cache slot of the original DAG's root. We must then do another upwards search along the same paths to clear out the cache fields of the original nodes that were copied, thus resetting the DAG for future upcopy operations. This cache-clearing pass, again, takes time linear in the number of edges occurring on the paths from the copied leaf to the root. (Alternatively, we can keep counter fields on the nodes to discriminate distinct upcopy operations, and perform a global reset on the term when the current-counter value overflows.)

5 Operating on λ -DAGs

We now have the core idea of our DAG-based β -reduction algorithm in place, and can fill in the details specific to our λ -expression domain.

Basic representation We will represent a λ -calculus term as a rooted DAG.

Sharing Sharing will be generally allowed, and sharing will be *required* of variable-reference terms. That is, any given variable will have no more than one node in the DAG representing it. If one variable is referenced by (is the child of) multiple parent nodes in the graph, these nodes will simply all contain pointers to the same data structure.

Bound-variable short-cuts Every λ -expression node will, in addition to having a reference to its body node, also have a reference to the variable node that it binds. This, of course, is how we navigate directly to the leaf node to replace when we begin the upcopy for a β -reduction operation. Note that this amounts to an α -uniqueness condition—we require that every λ -expression bind a unique variable.

Cache fields Every application node has a cache field that may either be empty or contain another application node. λ -expression nodes do not need cache fields—they only have one child (the body of the λ -expression), so the upcopy algorithm can only copy up through a λ -expression once during a β -reduction.

Uplinks Uplinks are represented by $\langle \text{parent}, \text{relation} \rangle$ pairs, where the three possible relations are “ λ body,” “application function,” and “application argument.” For example, if a node n has an uplink $\langle l, \lambda\text{-body} \rangle$, then l is a λ -expression, and n is its body.

Copying λ -expressions With all the above structure in place, the algorithm takes shape. To perform a β -reduction of redex $(\lambda x.b) a$, where b and a are arbitrary subterms, we simply initiate an $x \mapsto a$ upcopy. This will copy up through all the paths connecting top node b and leaf node x , building a copy of the DAG with a in place of x , just as we desire.

Application nodes, having two children, are handled just as binary-tree nodes in the general DAG-copy algorithm discussed earlier: copy, cache & continue on the first visit; mutate the cached copy on a second visit. λ -expression nodes, however, require different treatment. Suppose, while we are in the midst of performing the reduction above, we find ourselves performing a $c \mapsto c'$ upcopy, for some internal node c , into a λ parent of c : $\lambda y.c$. The general structure of the algorithm calls for us to make a copy of the λ -expression, with body c' . But we must also allocate a fresh variable, y' , for our new λ -expression, since we require all λ -expressions to bind distinct variables. This gives us $\lambda y'.c'$. Unfortunately, if old body c contains references to y , these will also occur in c' —not y' . We can be sure c' contains no references to y' , since y' was created after c' ! We need to fix up body c' by replacing all its references to y with references to y' .

Luckily, we already have the mechanism to do this: before progressing upwards to the parents of $\lambda y.c$, we simply initiate a $y \mapsto y'$ upcopy through the existing DAG. This

upcopy will proceed along the paths leading from the y reference, up through the DAG, to the $\lambda y.c$ node. If there are such paths, they *must* terminate on a previously-copied application node, at which point the upcopy algorithm will mutate the cached copy and return.

Why must these paths all terminate on some previously copied application node? Because we have already traversed a path from x up to $\lambda y.c$, copying and caching as we went. Any path upwards from the y reference must eventually encounter $\lambda y.c$, as well—this is guaranteed by lexical scope. The two paths must, then, converge on a common application node—the only nodes that have two children. That node was copied and cached by the original x -to- $\lambda y.c$ traversal.

When the $y \mapsto y'$ upcopy finishes updating the new DAG structure and returns, the algorithm resumes processing the original $c \mapsto c'$ upcopy, whose next step is to proceed upwards with a $(\lambda y.c) \mapsto (\lambda y'.c')$ upcopy to all of the parents of $\lambda y.c$, secure that the c' sub-DAG is now correct.

The single-DAG requirement We've glossed over a limitation of the uplink representation, which is that a certain kind of sharing is not allowed: after a β -reduction, the original redex must die. That is, the model we have is that we start with a λ -calculus term, represented as a DAG. We choose a redex node somewhere within this DAG, reduce it, and *alter the original DAG to replace the redex with the contractum*. When done, the original term has been changed—where the redex used to be, we now find the contractum. What we *can't* do is to choose a redex, reduce it, and then continue to refer to the redex or maintain an original, unreduced copy of the DAG. Contracting a redex kills the redex; the term data structure is not “pure functional” or “persistent” in the sense of the old values being unchanged. (We can, however, first “clone” a multiply-referenced redex, splitting the parents between the original and the clone, and then contract only one of the redex nodes.)

This limitation is due to the presence of the uplinks. They mean that a subterm can belong to only one rooted DAG, in much the same way that the backpointers in a doubly-linked list mean that a list element can belong to only one list (unlike a singly-linked list, where multiple lists can share a common tail). The upcopy algorithm assumes that the uplinks exactly mirror the parent→child downlinks, and traces up through all of them. This rules out the possibility of having a node belong to multiple distinct rooted DAGs, such as a “before” and “after” pair related by the β -reduction of some redex occurring within the “before” term.

Hence the algorithm, once it has finished the copying phase, takes the final step of disconnecting the redex from its parents, and replacing it with the contractum. The redex application node is now considered dead, since it has no parents, and can be removed from the parent/uplink sets of its children and deallocated. Should one of its two children thus have its parent set become empty, it, too, can be removed from the parent sets of its children and deallocated, and so forth. Thus we follow our upwards-recursive construction phase with a downwards-recursive deallocation phase.

It's important to stress, again, that this deallocation phase is not optional. A dead node must be removed from the parent sets of its children, lest we subsequently waste time doing an upcopy from a child up into a dead parent during a later reduction.

Termination and the top application Another detail we’ve not yet treated is termination of the upcopy phase. One way to handle this is simply to check as we move up through the DAG to see if we’ve arrived at the λ -expression being reduced, at which point we could save away the new term in some location and return without further upward copying. But there is an alternate way to handle this. Suppose we are contracting $\text{redex}(\lambda x.b)n$, for arbitrary sub-terms b and n . At the beginning of the reduction operation, we first check to see if x has no references (an easy check: is its uplink set empty?). If so, the answer is b ; we are done.

Otherwise, we begin at the λ -expression being reduced and scan downwards from λ -expression to body, until we encounter a non- λ -expression node—either a variable or an application. If we halt at a variable, it *must* be x —otherwise x would have no references, and we’ve already ruled that out. This case can also be handled easily: we simply scan back through this chain of nested λ -expressions, wrapping fresh λ -expressions around n as we go.

Finally, we arrive at the general case: the downward scan halts at the topmost application node a of sub-term b . We make an identical copy a' of a , *i.e.* one that shares both the function and argument children, and install a' in the cache slot of a .

Now we can initiate an $x \mapsto n$ upcopy, knowing that all upwards copying must terminate on a previously-copied application node. This is guaranteed by the critical, key invariant of the DAG: all paths from a variable reference upward to the root *must* encounter the λ -node binding that variable—this is simply lexical-scoping in the DAG context. The presence of a' in the cache slot of a will prevent upward copying from proceeding above a . Node a acts as a sentinel for the search; we can eliminate the root check from the upcopy code, for time savings.

When the upcopy phase finishes, we pass a' back up through the nested chain of λ -expressions leading from a back to the top $\lambda x.b$ term. As we pass back up through each λ -expression $\lambda y.t$, we allocate a fresh λ -expression term and a fresh variable y' to wrap around the value t' passed up, then perform a $y \mapsto y'$ upcopy to fix up any variable references in the new body, and then pass the freshly-created $\lambda y'.t'$ term on up the chain. (Note that the extended example shown in Sec. 7 omits this technique to simplify the presentation.)

6 Fine points

These fine points of the algorithm can be skipped on a first reading.

Representing uplinks A node keeps its uplinks chained together in a doubly-linked list, which allows us to remove an uplink from a node’s uplink set in constant time. We will need to do this, for example, when we mutate a previously copied node n to change one of its children—the old child’s uplink to n must be removed from its uplink set.

We simplify the allocation of uplinks by observing that each parent node has a fixed number of uplinks pointing to it: two in the case of an application and one in the case of a λ -expression. Therefore, we allocate the uplink nodes along with the parent, and thread the doubly-linked uplink lists through these pre-allocated nodes.

An uplink doubly-linked list element *appears* in the uplink list of the child, but the element *belongs* to the parent. For example, when we allocate a new application

node, we simultaneously allocate two uplink items: one for the function-child uplink to the application, and one for the argument-child uplink to the application. These three data structures have identical lifetimes; the uplinks live as long as the parent node they reference. We stash them in fields of the application node for convenient retrieval as needed. When we mutate the application node to change one of its children, we also shift the corresponding uplink structure from the old child’s uplink list to the new child’s uplink list, thus keeping the uplink pointer information consistent with the downlink pointer information.

The single-reference fast path Consider a redex $(\lambda x.b) n$, where the λ -expression being reduced has exactly one parent. We know what that parent must be: the redex application itself. This application node is about to die, when all references to it in the term DAG are replaced by references to the contractum. So the λ -expression itself is about to become completely parentless—*i.e.*, it, too, is about to die. This means that any node on a path from x up to the λ -expression will also die. Again, this is the key invariant provided by lexical scope: all paths from a variable reference upward to the root *must* encounter the λ -expression binding that variable. So if the λ -expression has no parents, then all paths upwards from its variable must terminate at the λ -expression itself.

This opens up the possibility of an alternate, fast way to produce the contractum: when the λ -expression being reduced has only one parent, mutate the λ -expression’s body, altering all of x ’s parents to refer instead to n . We do no copying at all, and may immediately take the λ -expression’s body as our answer, discarding the λ -expression and its variable x (in general, a λ -expression and its variable are always allocated and deallocated together).

Opportunistic iteration The algorithm can be implemented so that when a node is sequencing through its list of uplinks, performing a recursive upcopy on each one, the final upcopy can be done with a tail recursion (or, if coded in a language like C, as a straight iteration). This means that when there is no sharing of nodes by parents, the algorithm tends to iteratively zip up chains of single-parent links without pushing stack frames.

7 Extended example

We can see the sequences of steps taken by the algorithm on a complete example in Fig. 4. Part 4(a) shows the initial redex, which is $(\lambda x.(x(\lambda y.x(uy)))(\lambda y.x(uy))) t$, where the $(\lambda y.x(uy))$ subterm is shared, and t and u are arbitrary, unspecified subterms with no free occurrences of x or y . To help motivate the point of the algorithm, imagine that the sub-terms t and u are enormous—things we’d like to avoid copying or searching—and that the λx node has other parents besides application 1—so we cannot blindly mutate it at will, without corrupting what the other parents see. (If the λx node *doesn’t* have other parents, then the single-reference fast-path described in the previous section applies, and we *are* allowed to mutate the term, for a very fast reduction.)

In the following subfigure, 4(b), we focus in on the body of the λ -expression being reduced. We iterate over the parents of its variable-reference x , doing an $x \mapsto t$ upcopy;

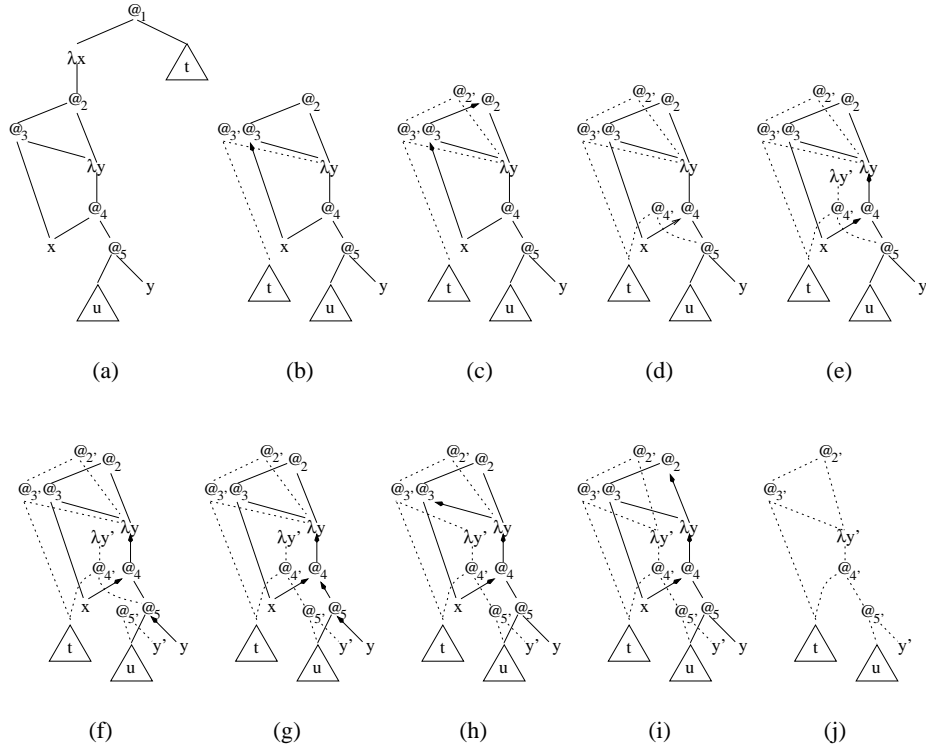


Fig. 4. A trace of a bottom-up reduction of the term $(\lambda x.(x(\lambda y.x(uy)))(\lambda y.x(uy)))t$, where the $(\lambda y.x(uy))$ term is shared, and sub-terms t and u are not specified.

this is the redex-mandated substitution that kicks off the entire reduction. The first parent of x is application 3, which is copied, producing application 3', which has function child t instead of the variable reference x , but has the same argument child as the original application 3, namely the λy term. The copy 3' is saved away in 3's cache slot, in case we copy into 3 from its argument child in the future.

Once we've made a copy of a parent node, we must recursively perform an upcopy for it. That is, we propagate a $3 \mapsto 3'$ upcopy to the parents of application 3. There is only one such parent, application 2. In subfigure 4(c), we see the result of this upcopy: the application 2' is created, with function child 3' instead of 3; the argument child, λy , is carried over from the original application 2. Again, application 2' is saved away in the cache slot of application 2.

Application 2 is the root of the upcopy DAG, so once it has been copied, control returns to application 3 and its $3 \mapsto 3'$ upcopy. Application 3 has only one parent, so it is done. Control returns to x and its $x \mapsto t$ upcopy, which proceeds to propagate upwards to the second parent of x , application 4.

We see the result of copying application 4 in subfigure 4(d). The new node is $4'$, which has function child t where 4 has x ; $4'$ shares its argument child, application 5, with application 4. Once again, the copy $4'$ is saved away in the cache slot of application 4.

Having copied application 4, we recursively trigger a $4 \mapsto 4'$ upcopy, which proceeds upwards to the sole parent of application 4. We make a copy of λy , allocating a fresh variable y' , with the new body $4'$. This is shown in subfigure 4(e).

Since the new $\lambda y'$ term binds a fresh variable, while processing the λy term we must recursively trigger a $y \mapsto y'$ upcopy, which begins in subfigure 4(f). We iterate through the parents of variable-reference y , of which there is only one: application 5. This is copied, mapping child y to replacement y' and sharing function child u . The result, $5'$, is saved away in the cache slot of application 5.

We then recursively trigger a $5 \mapsto 5'$ upcopy through the parents of application 5; there is only one, application 4. Upon examining this parent (subfigure 4(g)), we discover that 4 already has a copy, $4'$, occupying its cache slot. Rather than create a second, new copy of 4, we simply mutate the existing copy so that its argument child is the new term $5'$. Mutating rather than freshly allocating means the upcopy proceeds no further; responsibility for proceeding upwards from 4 was handled by the thread of computation that first encountered it and created $4'$. So control returns to application 5, which has no more parents, and then to y , who also has no more parents, so control finally returns to the λy term that kicked off the $y \mapsto y'$ copy back in subfigure 4(f).

In subfigure 4(h), the λy term, having produced its copy $\lambda y'$, continues the upcopy by iterating across its parents, recursively doing a $\lambda y \mapsto \lambda y'$ upcopy. The first such parent is application 3, which has already been copied, so it simply mutates its copy to have argument child $\lambda y'$ and returns immediately.

The second parent is application 2, which is handled in exactly the same way in subfigure 4(i). The λy term has no more parents, so it returns control to application 4, who has no more parents, and so returns control to variable reference x . Since x has no more parents, we are done. The answer is application $2'$, which is shown in subfigure 4(j). We can change all references to application 1 in the DAG to point, instead, to application $2'$, and then deallocate 1. Depending on whether or not the children of application 1 have other parents in the DAG, they may also be eligible for deallocation. This is easily performed with a downwards deallocation pass, removing dead nodes from the parent lists of their children, and then recursing if any child thus becomes completely parentless.

8 Experiments

To gain experience with the algorithm, a pair of Georgia Tech undergraduates implemented three β -reduction algorithms: the bottom-up algorithm (BUBS), a reducer based on the suspension λ -calculus (SLC, see Sec. 9.1), and a simple, base-line reducer, based on the simple top-down, blind-search recursive procedure described in Sec. 1. For a toy client application that would generate many requests for reduction, we then built a pair of simple normalisers (one total and one weak-head) on top of the reducers. We did two independent implementations, the first in SML, and a second, in C; the C implementation gave us tighter control over the algorithm and data structures for the purposes of

	CPU time (ms)			# reductions	
	BUBS	SLC	Simple	BUBS	Tree
(fact 2)	0	10	10	123	180
(fact 3)	0	20	20	188	388
(fact 4)	0	40	∞	286	827
(fact 5)	0	160	∞	509	2045
(fact 6)	10	860	∞	1439	7082
(fact 7)	20	5620	∞	7300	36180
(fact 8)	190	48600	∞	52772	245469
nasty-I	30	740	∞	7300	8664
pearl10	0	N/A	N/A	10	N/A
pearl18	0	N/A	N/A	18	N/A
tree10	0	0	0	1023	1023
tree18	740	2530	1980	262143	262143

Fig. 5. Timings for three different implementations of reduction. The system gave us a measurement precision of 10 ms; an entry of 0ms means below the resolution of the timer—*i.e.*, less than 10ms; a measurement of ∞ means the measurement was halted after several cpu-minutes.

measurement. The SLC and simple reducers managed storage in the C implementation with the Boehm-Demers-Weiser garbage collector, version 6.2; the BUBS algorithm requires no garbage collector.

Space limitations restrict us to presenting a single set of comparisons from these tests (Fig. 5). The “fact” entries are factorial computations, with Church-numeral encodings. “Nasty-I” is a 20,152-node tree of S and K combinators that reduces to I. A “tree *i*” entry is a full binary tree of applications, *i* deep, with I combinators at the leaves; a “pearl *i*” is this tree collapsed to a DAG—a linear sequence of *i* application nodes with a single I leaf. We compiled the code with gcc 2.95.4 -g -O2 -Wall and performed the test runs on an 800 MHz PIII (256 KB cache), 128 MB RAM, Debian GNU/Linux 3.0 system. These measurements are fairly minimal; we are currently porting Shao’s FLINT [14] system to BUBS to get a more realistic test of the algorithm in actual practice.

One of the striking characteristics of the bottom-up algorithm is not only how fast it is, but how well-behaved it seems to be. The other algorithms we’ve tried have fast cases, but also other cases that cause them to blow up fairly badly. The bottom-up algorithm reliably turns in good numbers. We conjecture this is the benefit of being able to exploit both sharing and non-sharing as they arise in the DAG. If there’s sharing, we benefit from re-using work. If there’s no sharing, we can exploit the single-parent fast path. These complementary techniques may combine to help protect the algorithm from being susceptible to particular inputs.

9 Related work

A tremendous amount of prior work has been carried out exploring different ways to implement β -reduction efficiently. In large part, this is due to β -reduction lying at the heart of the graph-reduction engines that are used to execute lazy functional languages. The text by Peyton Jones *et al.* [13] summarises this whole area very well.

However, the focus of the lazy-language community is on representations tuned for *execution*, and the technology they have developed is cleverly specialised to serve this need. This means, for example, that it’s fair game to fix on a particular reduction order. For example, graph reducers that overwrite nodes rely on their normalisation order to keep the necessary indirection nodes from stacking up pathologically. A compiler, in contrast, is a λ -calculus client that makes reductions in a less predictable order, as analyses reveal opportunities for transformation.

Also, an implementation tuned for execution has license to encode terms, or parts of terms, in a form not available for examination, but, rather, purely for execution. This is precisely what the technique of supercombinator compilation does. Our primary interest at the beginning of this whole effort was instead to work in a setting where the term being reduced is always directly available for examination—again, serving the needs of a compiler, which wants to manipulate and examine terms, not execute them.

9.1 Explicit-substitution calculi

One approach to constructing efficient λ -term manipulators is to shift to a language that syntactically encodes environments. The “suspension λ -calculus” developed by Nadathur *et al.* [12] is one such example that has been used with success in theorem provers and compilers. However, these implementations are quite complex, inflict de Bruijn encodings on the client, and their “constant-time” reductions simply shift the burden of the reduction to readback time. In the terms we’ve defined, these technologies use “blind search” to find the variables being substituted. Also, their use of de Bruijn encodings is a barrier to sharing internal structure: de Bruijn-index references are context dependent. *E.g.*, if a term $\lambda x.y$ appears multiple times underneath a λy parent, the index used for the y reference can vary.

One of the major algorithmic payoffs of these representations, lazy reduction, is not so useful for compilers, which typically must examine all the nodes of a term in the course of processing a program. SLC has been successfully employed inside a compiler to represent Shao’s FLINT typed intermediate language [14], but the report on this work makes clear the impressive, if not heroic, degree of engineering required to exploit this technology for compiler internals—the path to good performance couples the core SLC representation with hash consing as well as memoisation of term reductions.

The charm of the bottom-up technique presented here is its simplicity. The data structure is essentially just a simple description of the basic syntax as a datatype, with the single addition of child→parent backpointers. It generalises easily to the richer languages used by real compilers and other language-manipulation systems. It’s very simple to examine this data structure during processing; very easy to debug the reduction engine itself. In contrast to more sophisticated and complex representations such as SLC, there are really only two important invariants on the structure: (1) all variables are in scope (any path upwards from a variable reference to the root must go through the variable’s binding λ -expression), and (2) uplink backpointers mirror downlink references.

9.2 Director strings

Director strings [7] are a representation driven by the same core issue that motivates our uplink-DAG representation: they provide a way to guide search when performing β -reduction. In the case of director strings, however, one can do the search top-down. Unfortunately, director strings can impose a quadratic space penalty on our trees. Uplinked λ -DAGs are guaranteed to have linear space requirements. Whether or not the space requirements for a director strings representation will blow up in practice depends, of course, on the terms being manipulated. But the attraction of a linear-space representation is knowing that blow-up is completely impossible.

Like the suspension λ -calculus, director strings have the disadvantage of not being a direct representation of the original term; there is some translation involved in converting a λ -calculus term into a director strings.

Director strings can be an excellent representation choice for graph-reducing normalising engines. Again, we are instead primarily focussed on applications that require fine-grained inter-reduction access to the term structure, such as compilers.

9.3 Optimal λ reduction

The theory of “optimal λ reduction” [10, 9, 6] (or, OLR), originated by Lévy and Lamping, and developed by Abadi, Asperti, Gonthier, Guerrini, Lawall, Mairson *et al.*, is a body of work that shares much with bottom-up β -reduction. Both represent λ -terms using graph structure, and the key idea of connecting variable-binders directly to value-consumers of the bound variable is present in both frameworks—and for the same reason, namely, from a desire that substitution should be proportional to the number of references to the bound variable, removing the need to blindly search a term looking for these references.

However, the two systems are quite different in their details, in fairly deep ways. Lamping graphs allow *incremental reduction* by means of adding extra “croissant,” “bracket” and “fan” nodes to the graph. This exciting alternative model of computation, however, comes with a cost: the greatly increased complexity of the graph structure and its associated operations. In particular, in actual use, the croissant and bracket marks can frequently pile up uselessly along an edge, tying up storage and processing steps. It also makes it difficult to “read” information from the graph structure. As Gonthier, Abadi and Lévy state [6], “it seems fair to say that Lamping’s algorithm is rather complicated and obscure.” The details of this complexity have prevented OLR-based systems from widespread adoption.

9.4 Two key issues: persistence and readback

Our comparisons with other techniques have repeatedly invoked the key issues of persistence and readback. Our data structure is not a “persistent” one—performing a reduction inside a term changes the term. If an application needs to keep the old term around, then our algorithm is not a candidate (or, at least, not without some serious surgery). So perhaps it is unfair to compare our algorithm’s run times to those of persistent algorithms, such as SLC or director strings.

However, we can turn this around, and claim that the interesting feature of our algorithm is that it *exploits* lack of persistence. If an application doesn't need persistence, it shouldn't have to pay for it. The standard set of choices are invariably persistent; our algorithm provides an alternative design point. (Note that reduction on Lamping graphs is also not persistent, which is, again, either a limitation or a source of efficiency, depending on your point of view.)

The other key, cross-cutting issue is readback. An application that doesn't need to examine term structure in-between reductions has greater flexibility in its requirements. If readback is a requirement, however, then Lamping graphs and the SLC are much less attractive. Readback with our representation is free: one of the pleasant properties of a DAG is that it can be viewed just as easily as a tree; there is no need to convert it.

Thus, bottom-up β -reduction is a technology which is well suited to applications which (1) don't need persistence, but (2) do need fine-grained readback.

10 Conclusion

We certainly are not the first to consider using graph structure to represent terms of the λ -calculus; the ideas go back at least to 1954 [4, 15]. The key point we are making is that two of these ideas work together:

- representing λ -terms as DAGS to allow sharing induced by β -reduction, and
- introducing child \rightarrow parent backpointers and $\lambda\rightarrow$ variable links to efficiently direct search and construction.

The first idea allows sharing *within* a term, while the second allows sharing *across* a reduction, but they are, in fact, mutually enabling: in order to exploit the backpointers, we need the DAG representation to allow us to build terms without having to replicate the subterm being substituted for the variable. This is the source of speed and space efficiency.

The algorithm is simple and directly represents the term without any obscuring transform, such as combinators, de Bruijn indices or suspensions, a pleasant feature for λ -calculus clients who need to examine the terms. It is also, in the parlance of the graph-reduction community, fully lazy.

11 Acknowledgements

Bryan Kennedy and Stephen Strickland, undergraduates at Georgia Tech, did the entire implementation and evaluation reported in Sec. 8. We thank, of course, Olivier Danvy. Zhong Shao provided helpful discussions on the suspension λ -calculus. Chris Okasaki and Simon Peyton Jones tutored us on director strings. Harry Mairson and Alan Bawden provided lengthy and patient instruction on the subtleties of optimal lambda reduction and Lamping graphs.

References

1. Andrea Asperti and Stefano Guerrini. *The Optimal Implementation of Functional Programming Languages*. Cambridge University Press, 1999.
2. Henk Barendregt. *The Lambda Calculus*. North Holland, revised edition, 1984.
3. Alan Bawden. Personal communication, November 2002. Alan wrote the compiler, a toy exercise for Scheme, sometime in the late 1980's.
4. N. Bourbaki. *Théorie des ensembles*. Hermann & C. Editeurs, 1954.
5. Alonzo Church. *The Calculi of Lambda Conversion*. Princeton University Press, 1941.
6. Georges Gonthier, Martín Abadi and Jean-Jacques Lévy. The geometry of optimal lambda reduction. In *Conference Record of the Nineteenth Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 15–26, January 1992.
7. J. R. Kennaway and M. R. Sleep. Director strings as combinators. *ACM Transactions on Programming Languages and Systems*, 10, pages 602–626, (October 1988).
8. Richard A. Kelsey. A correspondence between continuation-passing style and static single assignment form. In *ACM SIGPLAN Workshop on Intermediate Representations, SIGPLAN Notices*, vol. 30, no. 3, pages 13–22, January 1995.
9. John Lamping. An algorithm for optimal lambda-calculus reduction. In *Proceedings of the Seventeenth Annual ACM Symposium on Principles of Programming Languages*, pages 16–30, January 1990.
10. Jean-Jacques Lévy. *Réductions Correctes et Optimales dans le Lambda-calcul*. Ph.D. thesis, Université Paris VII, 1978.
11. R. Milner, M. Tofte, R. Harper, D. MacQueen. *The Definition of Standard ML (Revised)*. MIT Press, 1997.
12. Gopalan Nadathur and Debra Sue Wilson. A notation for lambda terms: A generalization of environments. *Theoretical Computer Science* 198(1–2):49–98, May 1998.
13. Simon L. Peyton Jones. *The Implementation of Functional Programming Languages*. Prentice-Hall, 1987.
14. Zhong Shao, Christopher League, and Stefan Monnier. Implementing typed intermediate languages. In *Proceedings of the 1998 ACM SIGPLAN International Conference on Functional Programming Languages*, September 1998.
15. C. P. Wadsworth. *Semantics and pragmatics of the lambda-calculus*. PhD dissertation, Oxford University, 1971.