

A Type System with Usage Aspects

DAVID ASPINALL

LFCS Edinburgh, Mayfield Rd, Edinburgh EH9 3JZ, UK
(e-mail: David.Aspinall@ed.ac.uk)
<http://homepages.inf.ed.ac.uk/da/>

MARTIN HOFMANN

Institut für Informatik, Oettingenstraße 67, 80538 München, Germany
(e-mail: mhofmann@informatik.uni-muenchen.de)
<http://www.tcs.informatik.uni-muenchen.de/~mhofmann/>

MICHAL KONEČNÝ

Aston University, Aston Triangle, Birmingham, B4 7ET, UK
(e-mail: m.konecny@aston.ac.uk)
<http://www.aston.ac.uk/~konecny/>

Abstract

Linear typing schemes can be used to guarantee non-interference and so the soundness of in-place update with respect to a functional semantics. But linear schemes are restrictive in practice, and more restrictive than necessary to guarantee soundness of in-place update. This limitation has prompted research into static analysis and more sophisticated typing disciplines to determine when in-place update may be safely used, or to combine linear and non-linear schemes. Here we contribute to this direction by defining a new typing scheme that better approximates the semantic property of soundness of in-place update for a functional semantics. We begin from the observation that some data is used only in a “read-only” context, after which it may be safely re-used before being destroyed. Formalising the in-place update interpretation in a machine model semantics allows us to refine this observation, motivating three *usage aspects* apparent from the semantics that are used to annotate function argument types. The aspects are (1) used destructively, (2) used read-only but shared with result, and (3) used read-only and not shared with the result. The main novelty is aspect (2) that allows a linear value to be safely read and even aliased with a result of a function without being consumed. This novelty makes our type system more expressive than previous systems for functional languages in the literature. The system remains simple and intuitive, but it enjoys a strong soundness property whose proof is non-trivial. Moreover, our analysis features principal types and feasible type reconstruction, as shown in (Konecny, 2003b).

1 Introduction

The distinctive advantage of pure functional programming is that program functions may be viewed as ordinary mathematical functions. Powerful proof principles, such as equational reasoning with program terms and mathematical induction, are

```

reverse_aux : L(A), L(A) → L(A)
def reverse_aux(l, acc) =
  match l with
  | nil ⇒ acc
  | cons(h, t) ⇒ reverse_aux(t, cons(h, acc))

reverse : L(A) → L(A)
def reverse(l) = reverse_aux(l, nil)

```

Fig. 1. Functional list reverse

available. These principles are sound and do not need to use stores or other auxiliary entities, as is invariably the case when reasoning about imperative programs, as e.g. in Reynolds' separation logic (Reynolds, 2002).

Consider the functional implementation of linked list reversal, as shown in Fig. 1. This definition of reversal is readily verified by induction and equational reasoning over the set of finite lists. On the other hand, implementing reversal imperatively using pointers is (arguably) more cumbersome and error prone and, more seriously, would be harder to verify using complicated reasoning principles for imperative programs (see e.g. (Dor *et al.*, 2000)).

The advantage of the usual imperative implementation, of course, is that it modifies its argument in-place, whereas with the usual functional implementation the result must be created from scratch and garbage collection is necessary to salvage heap space. However, if by static analysis a clever compiler could determine that the original list `l` is not used after the call `reverse(l)`, it would be safe to instead use an optimised in-place update implementation, which avoids creating garbage. To be safe, one must trust both the compiler optimisation and the static analysis.

If we want to simplify the static analysis and ensure that an in-place update implementation is always possible, we can restrict programs using a linear typing scheme. With linear typing, every variable is used exactly once (or, in an affine linear system, at most once). However, it is well known that pure linear schemes are overly restrictive for real programming and need to be relaxed so that variables can be used more than once as far as possible. For example, consider the function `sumlist` that operates on lists of integers:

```

sumlist : L(N) → N
def sumlist(l) =
  match l with
  | nil ⇒ 0
  | cons(h, t) ⇒ h + sumlist(t)

```

In a purely linear type system, we would not be able to use the list `l` after a call to `sumlist(l)`, but this function merely examines its argument; the list would actually remain intact under any reasonable implementation so the linearity restriction that supposes it may be destroyed is overly restrictive.

Not only should the `sumlist` function inspect its argument list without modifying

it, but the result it returns no longer refers to the argument list. This means that we should be able to use `l` after the call to `sumlist`, so an expression like

```
cons(sumlist(l), reverse(l))
```

could be safely compiled even with the in-place implementation of `reverse`, if we assume that evaluation occurs from left to right and the list `l` is not used again afterwards.

There is less freedom with functions that return a result that contains some part of the argument. An example is the function `nth_tail` that returns the n th tail of a list:

```
nth_tail : N, L(A) → L(A)
def nth_tail(n, l) =
  if n <= 0 then l else match l with
    nil ⇒ nil
  | cons(h, t) ⇒ nth_tail(n-1, t)
```

Unlike `sumlist`, the result of `nth_tail` may be *shared* (aliased) with the argument. This aliasing means an expression

```
cons(nth_tail(2, l), nil)
```

will be sound, but

```
cons(nth_tail(2, l), cons(reverse(l), nil))
```

will *not* be soundly implemented using the in-place update version of `reverse`. If `l = [1, 2, 3]`, the expression should evaluate to the list `[[3], [3, 2, 1]]` but the in-place version would yield `[[1], [3, 2, 1]]` because of the aliasing. So, if we want to guarantee to use the in-place implementation of `reverse`, the second expression should not be allowed in the language. Simpler example functions in the same category as `nth_tail`, whose results share with the argument, include projection functions and the identity function.

As another example, consider the `append` function:

```
append : L(A), L(A) → L(A)
def append(l, m) =
  match l with
    nil ⇒ m
  | cons(h, t) ⇒ cons(h, append(t, m))
```

The desirable imperative implementation of `append` joins the second list to the first one in-place, by modifying the final pointer; it returns its first argument that now points to a new list. Thus, in `append(l, m)` the first list `l` has been *destroyed* so we should treat that in the same way as arguments to `reverse`. But the second list `m` is *shared* with the result, and so should be treated in the same way as arguments to `nth_tail`.

Usage aspects. Together with consideration of a machine model semantics for in-place update, these observations lead us to introduce three *usage aspects* for variables, which are the central innovation in our type system. The usage aspects are:

- Aspect 1: modifying use, e.g., `l` in `reverse(l)` and `append(l,m)`
- Aspect 2: non-modifying use, but shared with result, e.g., `m` in `append(l,m)`
- Aspect 3: non-modifying use, not shared with result, e.g., `l` in `sumlist(l)`.

As a first intuition, the numbers are in increasing order of “permissiveness to reuse” in the terms. Variables may only be accessed once with aspect 1. Variables can be used many times with aspect 2, but it prevents an aspect 1 usage later if intermediate results are retained. Finally, variables can be freely used with aspect 3, the pure “read-only” usage. It will become apparent later (e.g. Lemma 5.6) how the ordering can be specified precisely and that it demonstrates itself also in other forms than the one suggested by the intuition of “permissiveness to reuse”. The fact that the usage aspects are numbered, rather than named mnemonically, will allow us to make use of the ordering in the presentation of our type system.

Our type system decorates function arguments with usage aspects, and then tracks the way that variables are used. For example, we can give the following typings:

```
reverse : L(A)1 → L(A)
sumlist : L(A)3 → N
nth_tail : L(A)2, N3 → L(A)
append  : L(A)1, L(A)2 → L(A)
```

Types, such as `N`, that do not involve any heap storage in our machine model can always have the read-only aspect 3. We call these types *heap-free*. Functions that have a heap-free result (like `sumlist`) may have aspect 3 for their non-heap-free arguments, provided they are not modified when computing the result.

As will be seen later, these usage aspects have an intuitive interpretation in our semantics but, perhaps surprisingly, the exact distinctions appear to be novel. Although ours is not the first type system that relaxes strict linearity, we believe that it is more expressive than other systems for functional languages, while remaining simple and intuitive. Although even more expressive typing systems and specialised program logics for aliasing exist, especially for low-level imperative languages, we believe that our type structure has the advantage of being considerably less complex than others, so that types may be more easily understood by the programmer. The apparent simplicity does not mean that our system is easy to prove sound; our soundness proof (that goes beyond traditional type safety) is quite involved and reveals subtleties in the precise interpretation of the usage aspects. A further advantage of our system is that it enjoys principal types and a feasible type reconstruction algorithm; this property is demonstrated elsewhere (Konečný, 2003b). Specific comparison with other work appears at the end of the paper.

Outline. In Section 2 we introduce our setting for studying usage aspects, which is within a first-order linear functional programming language called LFPL, first introduced in (Hofmann, 2000). The particular appeal of LFPL for us here is that it has a canonical in-place update interpretation by using a special *resource type* that mediates construction and destruction of inductive datatypes. As well as adding usage aspects, we augment LFPL with a richer type structure, including two kinds of products (\otimes and \times), that allows data structures with or without sharing to be defined. We then give some larger examples that illustrate our usage aspects.

In Section 3 we present the formal syntax and typing rules. Although the types themselves are simple, some of our rules have slightly intricate side-conditions; but they can be explained by an intuitive reading of the meaning of the aspects. In Section 4 we give an imperative operational semantics for LFPL that formalises how its in-place interpretation works; practical implementations are described elsewhere.

In Section 5 we prove our central soundness result, which makes precise the meaning of usage aspects and establishes that the in-place operational semantics agrees with a safe, classical call-by-value operational semantics that is known to correspond to the functional set-theoretic interpretation. This result justifies reasoning about the functional semantics and also proves that the various possibilities for aspect annotations in our typing rules are sound. Significantly, this result goes beyond most other published type safety results for type systems that do not connect operational semantics with denotations. Finally, Section 6 concludes with a comparison to some of the related work.

2 LFPL with usage aspects

Whereas functional languages enjoy powerful proof principles, imperative implementations that directly manipulate pointers enjoy efficiency. We are interested in having the best of both worlds by using a semantics-preserving translation of functional programs into imperative ones that use in-place update and need no garbage collection, as far as possible. In previous work by the second author (Hofmann, 2000), a first-order functional language called LFPL was defined, together with such a translation into C.

LFPL relies on some programmer assistance to manage memory but without compromising the functional semantics in any way. This assistance is enabled by augmenting (non-nullary) constructors of inductive datatypes such as `cons` with an additional argument of an abstract “diamond” resource type \diamond . The elements of \diamond can be thought of as heap space areas; the diamond type corresponds exactly to the heap space for storing a single cell of some inductive datatype.

To construct an element of an inductive type, we must supply a value of this type. Values of type \diamond can be obtained from formal parameters of functions or by deconstructing elements of recursive types in a pattern match. Consider the functions `reverse` and `append` discussed earlier, and implemented in LFPL in Fig. 2. The functions have almost identical definitions as before, except that there is now an extra argument to `cons`, both in pattern matches and in constructor positions. The pattern match releases a \diamond cell used for a list node, whereas the

```

reverse_aux : L(A),L(A) → L(A)
def reverse_aux(l,acc) =
  match l with
  | nil ⇒ acc
  | cons(d,h,t) ⇒ reverse_aux(t,cons(d,h,acc))

reverse : L(A) → L(A)
def reverse(l) = reverse_aux(l, nil)

append : L(A),L(A) → L(A)
def append(l,m) =
  match l with
  | nil ⇒ m
  | cons(d,h,t) ⇒ cons(d,h,append(t,m))

```

Fig. 2. LFPL examples

constructor consumes one. The first argument to each use of `cons` in Fig. 2 is a value of type \diamond ; the `cons` on the right hand side of the match is “justified” by the preceding `cons` pattern. The correspondence need not always be as local; in particular, values of type \diamond may be passed as arguments to and returned by functions, as well as appear in data structures.

Although it may be realised as a pointer, there is no way to directly examine or manipulate an element of \diamond , and in fact values of \diamond are interchangeable in a program. In the functional semantics $\text{cons}(d, h, (\text{cons}(d', h', t)))$ is equivalent to $\text{cons}(d', h, (\text{cons}(d, h', t)))$.

Our LFPL implementations to date use a single \diamond type that is large enough to store a cell of any datatype in the program, so a `d` for a list `cons` constructor is also interchangeable with a `d` for a tree node constructor. This largest-sized \diamond approach is potentially wasteful but gives flexibility and helps avoid garbage collection. It could be improved by using several differently sized \diamond types. The sizes in particular positions could be inferred using a data-flow analysis of the program.

Linear typing in LFPL. In (Hofmann, 2000) it was shown that the semantics of a translation of LFPL into C preserves the functional semantics of the source program *provided* the latter admits an affine linear typing for inductive types and \diamond types, i.e., bound variables of inductive type are used at most once. The heap-free types were *not* subject to any linearity restriction. Linearity guarantees that the memory space pointed to by a \diamond value is not needed anywhere else, which controls the space usage of programs. It prevents function definitions such as:

```

twice : L(A) → L(N)
def twice(l) =
  match l with
  | nil ⇒ nil
  | cons(d,h,t) ⇒ cons(d,0,cons(d,0,twice(t)))

```

The functional semantics of `twice` maps a list l to a list twice as long as l with zero entries; on the other hand, the LFPL translation to C of the above code computes a circular list!¹

As one would hope, the translation of `append` in Fig. 2 appends one linked list to another in place; again, the translation of a non-linear phrase like `append(1, 1)` results in a circular list, disagreeing with the functional semantics.

Linear typing together with the resource type \diamond seems restrictive at first sight. In particular, without dynamic creation of memory in the translation, we are heap-bounded: no function can be written that increases the size of its input. Yet surprisingly, a great many standard examples from functional programming fit very naturally into the LFPL typing discipline, among them, insertion sort, quick sort, tree sort, breadth-first traversal of a tree and Huffman’s algorithm. Moreover, in (Hofmann, 2000) it was shown that every non-size-increasing function on lists over booleans in the complexity class ETIME can be represented. The language is therefore useful where careful control over space usage is needed. If the typing is made pure linear instead of affine, we could also prevent space leaks.

Nevertheless, if heap-boundedness is too restrictive for a particular domain, one can easily add constants to the language that generate and destroy elements of type \diamond , implemented using an external memory allocator; typings are shown later. If we moreover assumed existence of a garbage collector, we could alleviate the programmer from needing to always free memory. In this case we are getting back to the ordinary functional programming setting of the introduction, except that we have the new forms of pattern matching and data construction available, which can be used to give the compiler directions about ways to reuse memory.

2.1 Adding usage aspects

Whether or not we are heap-bounded, the linear restrictions remain for the in-place memory reuse interpretation. As with any other typing scheme, the linear restriction rejects many semantically valid programs. In our context a program is semantically valid if its translation to imperative code computes its functional semantics.² We cannot catch all semantically valid programs by a decidable typing discipline, of course, but we can try to refine the type system to reduce the “slack”, i.e., the discrepancy between the semantically valid programs and those that pass the typing discipline.

Here we address one particular source for slack, namely the implicit assumption

¹ Incidentally, we can implement `twice` in LFPL if we use another typing:

$$L(N \otimes \diamond) \rightarrow L(N)$$

where the argument list provides the right amount of extra space. In recent work with Steffen Jost (Hofmann & Jost, 2003) and Dilsun Kırli, we have shown that such additions of \diamond types can be automatically inferred, using integer linear programming.

² We could show a slightly stronger soundness condition: a well-typed program in fact evaluates with *benign sharing* as defined in (Hofmann & Jost, 2003). It means that the program will pass certain conservative run-time checks of non-interference and thus the agreement with functional semantics is not merely accidental.

```

sort : L(L(N))1 → L(L(N))
def sort(l) =
  match l with
    nil ⇒ nil
  | cons(d,h,t) ⇒ insert(d,h,sort(t))

insert : ◇1,L(N)2,L(L(N))1 → L(L(N))
def insert(d,l,ll) =
  match ll with
    nil ⇒ cons(d,l,nil)
  | cons(d',h,t) ⇒
    if sumlist(l) > sumlist(h)
    then cons(d',h,insert(d,l,t))
    else cons(d,l,cons(d',h,t))

```

Fig. 3. Example: insertion sort

that every access to a variable is potentially destructive, i.e., changes the memory reachable from this variable. This assumption is overly conservative: multiple uses of a variable need not compromise semantic validity, as long as only the last one is potentially destructive, and provided the results of the earlier accesses do not interfere with an ultimate destructive access.

We model this observation by relaxing the linearity by adding the usage aspects motivated in the introduction and using them to control the typing. The starting point is with the typing for `cons`:

$$\text{cons}_A \quad : \quad \diamond^1, A^2, L(A)^2 \rightarrow L(A)$$

The aspect 1 annotation indicates that the diamond argument `d` is destroyed in `cons(d, h, t)` and must not be used again. Any function that consumes a diamond element will give it aspect 1. The aspect 2 annotations on `cons` indicate that the head and tail arguments of the constructor are not destroyed, and moreover, they share with the result of the construction.

The typing rule for pattern matches allows some flexibility in usage aspect labels:

$$\frac{\Gamma \vdash e_{\text{nil}} : B \quad \Gamma, x_d \overset{i_d}{:} \diamond, x_h \overset{i_h}{:} A, x_t \overset{i_t}{:} L(A) \vdash e_{\text{cons}} : B \quad i = \min(i_d, i_h, i_t)}{\Gamma, x \overset{i}{:} L(A) \vdash \text{match } x \text{ with nil} \Rightarrow e_{\text{nil}} \mid \text{cons}(x_d, x_h, x_t) \Rightarrow e_{\text{cons}} : B}$$

This rule says that the usage aspect for a list `x` in an expression matching against this list will be the most destructive aspect among the aspects exhibited for the pieces of a `cons` cell in the `econs` expression. This rule now allows both destructive matches (e.g., where the diamond `xd` is used to construct a new cell in `econs`, as in `reverse`) and non-destructive matches (e.g., where the diamond `xd` is not used and the list is only examined, as in `sumlist`).

The typing of `cons` and the rule for list elimination allows us to derive the usage aspect typings given in the introduction for the corresponding LFPL programs.

A longer example is the LFPL implementation of insertion sort in Fig. 3. The

`sort` function here sorts a list of integer lists in ascending order of their sums. In the comparison `test sumlist(l) > sumlist(h)` both variables are used with aspect 3 and thus these variables can be safely re-used in the branches of the `if` statement. The type of the `insert` function indicates that `insert(d,a,l)` consumes the diamond `d`, inserts the data item `a` without modifying it, but sharing with the result; moreover, it destroys the input list `l`.

2.2 Sharing data

The strict linear type system of LFPL as presented in (Hofmann, 2000) prevents sharing in data structures, which can lead to bad space behaviour in some programs. The `append` function shows how we might be able to allow some limited sharing within data structures but still use an in-place update implementation, provided we take care over when modification is allowed. For example, we would like to allow the expression

```
let x=append(u,w) and y=append(v,w) in e
```

provided that we don't refer to `x` or `y` in `e` after the other has been modified. Similarly, we would like to allow a tree that has sharing amongst sub-trees, in the simplest case a node constructed as follows:

```
let u=sharednode(d, (a,t,t)) in e
```

(where `d:◇` and `a` is a label). This data structure should be safe so long as we do not modify both branches of `u`. The kinds of data structure we are considering here are DAGs.

The “not modifying both parts” flavour of these examples leads us to include two kinds of products in our augmented version of LFPL. Consider binary trees. In a linear setting we have two kinds of trees; one corresponding to trees laid out in full in memory (\otimes -trees), the other corresponding to trees with sharing of sub-trees (\times -trees). In ordinary functional programming these two are extensionally equivalent; in the presence of linearity constraints they differ considerably. The \otimes -trees allow simultaneous access to all their components thus encompassing, e.g. computing the list of leaf labellings, whereas access to the \times -trees is restricted to essentially search operations. Conversely, \otimes -trees are more difficult to construct because of their storage requirements, whereas the typing rules allow easy construction of a full binary \times -tree that is represented as a rather small DAG. Thus the type system can reflect the kind of choices that a programmer would normally make in selecting the best data representation for a purpose.

The tensor product (denoted by \otimes) is accessed using the pattern matching construct

```
match p with (x $\otimes$ y)  $\Rightarrow$  e
```

that allows both `x` and `y` to be accessed simultaneously in `e`. Given a \otimes -product of two lists in a pure linear type system like that of LFPL, we can access (maybe modify) both components; to be sound, the two lists must have no sharing.

With usage aspects, this constraint is relaxed somewhat: we can construct \otimes -products *within* a program whose components share in certain cases when it is safe to do so, so long as the sharing does not escape the scope of an expression. A contrived but short example is the expression:

```
let l2 = l
    p = l $\otimes$ l2
in match p with (x $\otimes$ y)  $\Rightarrow$  sumlist(x) + sumlist(y)
```

Here, the \otimes -product p is allowed to exhibit sharing because both of its components are used read-only. (The typing derivation for this term can be understood in conjunction with the typing rules that will be explained in Section 3. In particular, the rules do not permit to form a tensor product $l\otimes l$ directly.)

The cartesian product (denoted \times), which corresponds to the $\&$ connective of linear logic, has a different behaviour: we may access one component or the other, but not both; this means that the two components may have sharing. Again, with usage aspects, we can be more permissive than allowing just access to one component of the product. We can safely allow access to both components, so long as at most one component is modified, and if it is, the other one is not referenced thereafter. (The typing rules for `let` and cartesian product have special side conditions that allow this behaviour.) Cartesian products are accessed via projection functions:

$$\begin{aligned} \text{fst} & : (A \times B)^2 \rightarrow A \\ \text{snd} & : (A \times B)^2 \rightarrow B \end{aligned}$$

The usage aspect 2 here indicates that the result shares with the argument and allows us to use *both* components later in the program. Such usage would not be possible in a pure linear type system with cartesian products.

To allow sharing to persist within data structures used in a program, we may give constructors arguments of cartesian product types. Ideally, we would allow the user to choose exactly where cartesian products are used and where \otimes -products are used, to allow the user to define datatypes appropriate for their application. For the purpose of exposition in this paper, however, we will include only tensor lists and tensor tree types as primitives (as in LFPL); however, the language and results are easily extended to their variants using cartesian products, or indeed, a general datatype mechanism that would allow the user to choose.

2.3 Further examples

We will now describe two prototypical examples of aliasing that can be proved safe by our usage aspects and are not too far away from real-life programming scenarios. In Figure 4 a number of functions are defined that process lists of numbers. There are several supporting functions as well as the two example functions, called S1 and S2. The flow of data in the terms defining S1 and S2 is illustrated in Fig. 5, including aspect annotations. The detail of the labelling for the annotations will be better understood by the reader after considering the typing rules given in the next section; however, the motivation of the examples can be explained first.

```

highest_sum : L(L(N))2 → L(N)      hs_aux : L(N)2, L(L(N))2 → L(N)
def highest_sum(l) =                def hs_aux(prev, l) =
  match l with                      match l with
  nil ⇒ nil                          nil ⇒ prev
  | cons(d, h, t) ⇒                 | cons(d, h, t) ⇒
    hs_aux(h, t)                     if sumlist(h) > sumlist(prev)
                                      then hs(h, t)
                                      else hs(prev, t)

first_similar : L(N)3, L(L(N))2 → L(N)  is_similar : L(N)3, L(N)3 → N
def first_similar(m, l) =           def is_similar(l1, l2) =
  match l with                      match l1 with
  nil ⇒ nil                          nil ⇒
  | cons(d, h, t) ⇒                 match l2 with
    if is_similar(m, h)              nil ⇒ 1
    then h                            | cons(d, h, t) ⇒ 0
    else first_similar(t)            | cons(d1, h1, t1) ⇒
                                      match l2 with
                                      nil ⇒ 0
                                      | cons(d2, h2, t2) ⇒
                                        similar_number(h1, h2)
                                        * is_similar(t1, t2)

filter_similar : L(N)3, L(L(N))1 → L(L(N))
def filter_similar(m, l) =
  match l with
  nil ⇒ nil
  | cons(d, h, t) ⇒
    let r = filter_similar(m, t)
    in if is_similar(m, h)
       then cons(d, h, r)
       else r

similar_number : N3, N3 → N
def similar_number(n1, n2) =
  (n1 - n2 < 2)
  * (n2 - n1 < 2)

S1 : L(L(N))2 → L(N)                S2 : L(N)1, L(L(N))1, ◇1 → L(L(N))
def S1(l) =                          def S2(m, l, d) =
  let b = highest_sum(l)              let l' = filter_similar(m, l)
  in first_similar(b, l)              in map_reverse(cons(d, m, l'))

```

Fig. 4. Functions demonstrating the usefulness of usage aspects 2 and 3

Function S1 takes a list of lists of numbers and returns one element of this list. The returned element is the first element of the input list that is similar to another element that is, in some sense, best. For brevity, we use simple notions of measure and similarity: we compare lists of numbers by their sums and consider two lists similar if they have the same length and all their elements are similar, i.e., differ by at most one. The important notion in this example is the following data flow: the input is first processed read-only, creating an aliased piece of data, namely the “best” element. Then both the original list and its “best” element are used as arguments to `first_similar`. This sharing of arguments is safe because neither of them is destroyed in the final step and only one of them is aliased with the final result.

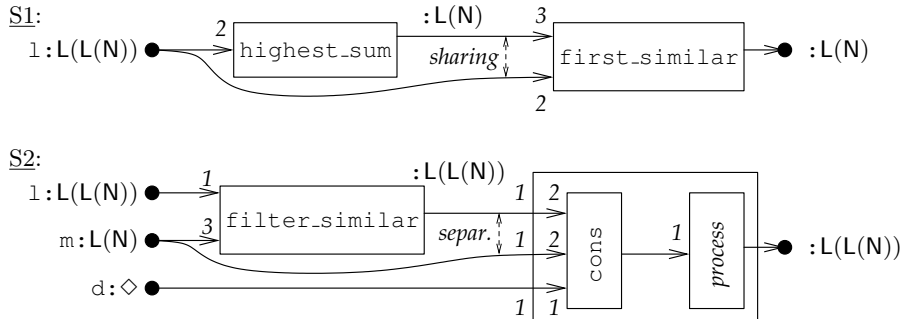


Fig. 5. The data-flow and usage aspects within the terms S1 and S2

Function S2 is more involved but has less intricate effects. It, again, takes a list of lists of numbers l as its primary argument, but it takes also a list of numbers m and a single diamond d . The function S2 filters from l all elements that are similar to m and then adds m back into this list. Finally it processes the resulting list in some destructive way, e.g. reverses all its elements in place. The important notion for S2 is the guarantee that after filtering, m is not only kept intact but is also separate from the result of the filtering, shown by its aspect 3 on its input to `filter_similar`.

These examples justify the need for all three aspects. If one would have only two aspects, leaving out, say, aspect 2, then S1 would not type-check. If one would leave out aspect 3, then S2 would not type-check. The examples give some hint about why the soundness proof for our system becomes involved: one has to reason about the annotations in different ways when they appear in different positions, in particular, considering both aliasing and separation of portions of heap data as shown in Fig. 5. We will return to these examples again when we compare our system to the related work in Section 6.

3 Syntax and typing

The grammar for the types and terms of our language is given in Fig. 6.

Types. Types consist of integers N , the resource type \diamond , lists $L(-)$ and trees $T(-)$, and the binary \otimes -product and \times -product. For brevity, we will use N also for the type of booleans. Types not containing diamonds \diamond or $L(-)$ or $T(-)$ are called *heap-free*, e.g. N and $N \otimes N$ are heap-free. (A pair is stored in a single logical memory cell as explained in Section 4.)

Terms and programs. A program consists of a series of (possibly mutually recursive) function definitions of the form $f(x_1, \dots, x_n) = e_f$. We use x and variants to range over variables and f to range over function symbols. To simplify the presentation, we restrict the syntax so that most term formers can be applied only to variables, like in the standard A-normal form (wik, 2007; Sabry & Felleisen, 1993). In practice, we can define the more general forms such as $f(e_1, \dots, e_n)$ as syntactic sugar for

$$\begin{aligned}
A & ::= \mathbf{N} \mid \diamond \mid \mathbf{L}(A) \mid \mathbf{T}(A) \mid A_1 \otimes A_2 \mid A_1 \times A_2 \\
e & ::= c \mid f(x_1, \dots, x_n) \mid x \mid \mathbf{let} \ x = e_x \ \mathbf{in} \ e \mid \mathbf{if} \ x \ \mathbf{then} \ e_1 \ \mathbf{else} \ e_2 \\
& \mid x_1 \otimes x_2 \mid \mathbf{match} \ x \ \mathbf{with} \ (x_1 \otimes x_2) \Rightarrow e \mid (e_1, e_2) \mid \mathbf{fst}(x) \mid \mathbf{snd}(x) \\
& \mid \mathbf{nil} \mid \mathbf{cons}(x_d, x_h, x_t) \mid \mathbf{match} \ x \ \mathbf{with} \ \mathbf{nil} \Rightarrow e_n \mid \mathbf{cons}(x_d, x_h, x_t) \Rightarrow e_c \\
& \mid \mathbf{leaf}(x_d, x_a) \mid \mathbf{node}(x_d, x_a, x_l, x_r) \\
& \mid \mathbf{match} \ x \ \mathbf{with} \ \mathbf{leaf}(x_d, x_a) \Rightarrow e_l \mid \mathbf{node}(x_d, x_a, x_l, x_r) \Rightarrow e_n
\end{aligned}$$

Fig. 6. LFPL abstract syntax grammar

nested let expressions. Using the same trick, we also impose a convention that in each application $f(x_1, \dots, x_n)$ the variables x_1, \dots, x_n are equal to the formal parameters used in e_f . This convention greatly simplifies definitions and various rules related to function applications.

The function definitions e_f must be well-typed. To make this typing possible in the presence of recursion, a program is given together with a signature Σ , which is a finite function from *function symbols* to first-order function types with usage aspects, i.e. of the form $A_1^{i_1}, \dots, A_n^{i_n} \rightarrow A$. In the typing rules we will assume a fixed program with signature Σ .

We also treat constructors as function symbols declared in the signature and include primitive arithmetic and comparison operations in the signature. Specifically, we assume that Σ contains a number of declarations:

$$\begin{aligned}
+, -, *, <, > & : \mathbf{N}^3, \mathbf{N}^3 \rightarrow \mathbf{N} \\
\mathbf{nil}_A & : \mathbf{L}(A) \\
\mathbf{cons}_A & : \diamond^1, A^2, \mathbf{L}(A)^2 \rightarrow \mathbf{L}(A) \\
\mathbf{leaf}_A & : \diamond^1, A^2 \rightarrow \mathbf{T}(A) \\
\mathbf{node}_A & : \diamond^1, A^2, \mathbf{T}(A)^2, \mathbf{T}(A)^2 \rightarrow \mathbf{T}(A) \\
\mathbf{fst}_{A \times B} & : (A \times B)^2 \rightarrow A \\
\mathbf{snd}_{A \times B} & : (A \times B)^2 \rightarrow B
\end{aligned}$$

for suitable types A as used in the program³. In examples where the types can be easily deduced from the context, we omit the type subscripts in \mathbf{nil} , \mathbf{cons} , etc.

The comma between argument types is treated as a \otimes -product, which means that these typings, and the corresponding elimination rules in Fig. 7, describe lists and trees with simultaneous access to sub-components. Hence they must be implemented without sharing, unless the access is guaranteed to be read-only.

Typing contexts. We keep track of *usage aspects* for variables as introduced above. We write $x:^i A$ to mean that $x:A$ will be used with aspect $i \in \{1, 2, 3\}$ in the subject of the typing judgement. A *typing context* Γ is a finite function from identifiers to types A with usage aspects; $|\Gamma|$ denotes the domain of Γ . If $x:^i A \in \Gamma$ we write $\Gamma(x) = A$ and $\Gamma[x] = i$.

We use familiar notation for extending contexts. If $x \notin |\Gamma|$, then we write $\Gamma, x:^i A$ for the extension of Γ with $x:^i A$. More generally, if $|\Gamma| \cap |\Delta| = \emptyset$ then we write

³ Thus the term constructor formers in Fig. 6 are actually superfluous but make the grammar more comprehensible.

Γ, Δ for the disjoint union of Γ and Δ . If such notation appears in the premise or conclusion of a rule below it is implicitly understood that these disjointness conditions are met.

In a couple of the typing rules we need some additional notation for manipulating usage aspects on variables. The “committed to i ” context Δ^i is the same as Δ , but each declaration $x:^2A$ of an aspect 2 (aliased) variable is replaced with $x:^iA$. If we have two contexts Δ_1, Δ_2 that differ only in usage aspects, so $|\Delta_1| = |\Delta_2|$ and $\Delta_1(x) = \Delta_2(x)$ for all x , then we define the merged context $\Gamma = \Delta_1 \wedge \Delta_2$ by $|\Gamma| = |\Delta_1|$, $\Gamma(x) = \Delta_1(x)$, $\Gamma[x] = \min(\Delta_1(x), \Delta_2(x))$. The merged context takes the “worst” usage aspect of each variable.

Richer types and signatures. Although we fix particular list and tree types and their interpretation for this paper, it should be clear that the language and our results can easily be extended to a general inductive datatype mechanism. Moreover, we can introduce datatypes that admit sharing by giving suitable typings for constructors in the signature. For example, to add a type of “sharing trees” $\text{ST}(A)$ with unrestricted sharing between components, we could use the constructor typing:

$$\text{sharednode}_A : \diamond^1, (A \times \text{ST}(A) \times \text{ST}(A))^2 \rightarrow \text{ST}(A).$$

In this typing, there can be sharing amongst the label and sub-trees, but still no sharing with the \diamond argument, of course, since the region pointed to by the \diamond argument is overwritten to store the constructed cell.

LFPL and our examples do not use it, but it is also straightforward to add explicit allocation and deallocation to the language. One simply needs to include in the program signature two further constants:

$$\begin{aligned} \text{new} & : \rightarrow \diamond \\ \text{dispose} & : \diamond^1 \rightarrow \mathbf{N} \end{aligned}$$

These built-ins can be implemented by interfacing to an external memory manager; of course, without further restriction, this extension breaks the heap-bounded nature of the system.⁴

Typing rules. Now we explain the typing rules, shown in Fig. 7, which define a judgement of the form $\Gamma \vdash e : A$. Most rules are straightforward. Our system is affine linear, so it includes a rule for weakening, **WEAK**. In the rules **VAR** and **⊗-PAIR**, variables are given the default aspect 2, to indicate their obvious sharing with the result. If the result is a value of a heap-free type, then with **RAISE** we can promote variables of aspect 2 to aspect 3 to reflect that they manifestly do not share with the result. The rule **DROP** goes the other way and allows us to assume that a variable is used in a more destructive fashion than it actually is.

The two non-standard structural rules give insight into the semantics and allow

⁴ In other work (Hofmann & Jost, 2003; MacKenzie & Wolverson, 2004) we have implemented a similar language which does have mechanisms for predicting heap usage in the presence of allocation and deallocation.

$$\begin{array}{c}
 \frac{}{\vdash c : \mathbb{N}} \text{ (CONST)} \qquad \frac{}{x \overset{?}{:} A \vdash x : A} \text{ (VAR)} \qquad \frac{\Gamma \vdash e : A}{\Gamma, \Delta \vdash e : A} \text{ (WEAK)} \\
 \\
 \frac{\Gamma \vdash e : A \quad A \text{ heap-free}}{\Gamma^3 \vdash e : A} \text{ (RAISE)} \qquad \frac{\Gamma, x \overset{i}{:} A \vdash e : B \quad j \leq i}{\Gamma, x \overset{j}{:} A \vdash e : B} \text{ (DROP)} \\
 \\
 \frac{f : A^{i_1}, \dots, A^{i_n} \rightarrow B \text{ in } \Sigma}{x_1 \overset{i_1}{:} A_1, \dots, x_n \overset{i_n}{:} A_n \vdash f(x_1, \dots, x_n) : B} \text{ (APP)} \\
 \\
 \frac{\Gamma \vdash e_1 : A \quad \Gamma \vdash e_2 : A}{\Gamma, x \overset{?}{:} \mathbb{N} \vdash \text{if } x \text{ then } e_1 \text{ else } e_2 : A} \text{ (IF)} \\
 \\
 \frac{\Gamma_1, \Delta_1 \vdash e_1 : A_x \quad \Gamma_2, \Delta_2, x \overset{i}{:} A_x \vdash e_2 : B \quad \text{condition } \star \text{ (see Table 8)}}{\Gamma_1^i, \Gamma_2, \Delta_1^i \wedge \Delta_2 \vdash \text{let } x = e_1 \text{ in } e_2 : B} \text{ (LET)} \\
 \\
 \frac{}{x_1 \overset{?}{:} A_1, x_2 \overset{?}{:} A_2 \vdash x_1 \otimes x_2 : A_1 \otimes A_2} \text{ (\otimes-PAIR)} \\
 \\
 \frac{\Gamma_1, \Delta_1 \vdash e_1 : A_1 \quad \Gamma_2, \Delta_2 \vdash e_2 : A_2 \quad \text{condition } \star\star \text{ (see page 17)}}{\Gamma_1, \Gamma_2, \Delta_1 \wedge \Delta_2 \vdash (e_1, e_2) : A_1 \times A_2} \text{ (\times-PAIR)} \\
 \\
 \frac{\Gamma, x_1 \overset{i_1}{:} A_1, x_2 \overset{i_2}{:} A_2 \vdash e : B \quad i = \min(i_1, i_2)}{\Gamma, x \overset{i}{:} A_1 \otimes A_2 \vdash \text{match } x \text{ with } (x_1 \otimes x_2) \Rightarrow e : B} \text{ (PAIR-ELIM)} \\
 \\
 \frac{\Gamma \vdash e_{\text{nil}} : B \quad \Gamma, x_d \overset{i_d}{:} \diamond, x_h \overset{i_h}{:} A, x_t \overset{i_t}{:} \mathbb{L}(A) \vdash e_{\text{cons}} : B \quad i = \min(i_d, i_h, i_t)}{\Gamma, x \overset{i}{:} \mathbb{L}(A) \vdash \text{match } x \text{ with nil} \Rightarrow e_{\text{nil}} | \text{cons}(x_d, x_h, x_t) \Rightarrow e_{\text{cons}} : B} \text{ (LIST-ELIM)} \\
 \\
 \frac{\Gamma, x_d \overset{i_d}{:} \diamond, x_a \overset{i_a}{:} A \vdash e_{\text{leaf}} : B \quad \Gamma, x_d \overset{i_d}{:} \diamond, x_a \overset{i_a}{:} A, x_l \overset{i_l}{:} \mathbb{T}(A), x_r \overset{i_r}{:} \mathbb{T}(A) \vdash e_{\text{node}} : B \quad i = \min(i_a, i_d, i_l, i_r)}{\Gamma, x \overset{i}{:} \mathbb{T}(A) \vdash \text{match } x \text{ with leaf}(x_d, x_a) \Rightarrow e_{\text{leaf}} | \text{node}(x_d, x_a, x_l, x_r) \Rightarrow e_{\text{node}} : B} \text{ (TREE-ELIM)}
 \end{array}$$

Note: $\Delta^i[x] = \begin{cases} i & \text{if } \Delta[x] = 2 \\ \Delta[x] & \text{otherwise} \end{cases}$ and $(\Delta_1 \wedge \Delta_2)[x] = \min(\Delta_1[x], \Delta_2[x])$.

Fig. 7. Typing rules

$\Delta_1[z] \rightarrow$	1			2			3		
$\downarrow i/\Delta_2[z] \rightarrow$	1	2	3	1	2	3	1	2	3
1	\times	\times	\times	\times	\times	\times	\checkmark	\checkmark	\checkmark
2	\times	\times	\times	\times	\times	\checkmark	\checkmark	\checkmark	\checkmark
3	\times	\times	\times	\times	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark

Fig. 8. The side condition \star of LET in Fig. 7: $\forall z \in \Delta_1$ the table has to show a tick.

us to give simplified versions of the other rules that would otherwise require greater flexibility. For example, the rule IF assumes the same context Γ when typing both branches but in fact the usage aspects of variables may differ between e_1 and e_2 — a variable may be used destructively in one branch but not the other; the rule DROP is used in this case to unite the contexts. The rule RAISE allows us to simplify VAR and \otimes -PAIR that would otherwise need to be generalised to also allow the aspect 3 in the case that the associated type is heap-free.

The rule LET is somewhat intricate. To type $\text{let } x = e_1 \text{ in } e_2$, the context is split into three pieces: variables specific to the definition e_1 , in Γ_1 ; variables specific to the body e_2 , in Γ_2 ; and common variables, in Δ_1 and Δ_2 , which may be used with different aspects in e_1 and e_2 .

First, we type-check the definition to find its type A_x . Then we type-check the body using some usage aspect i for the bound variable x . The usage aspect the bound variable x has in the body is used for any aliased variables belonging to e_1 . For example, if x is used destructively in e_2 , then all aliased variables in Γ_1 and Δ_1 are used destructively in the overall expression; this observation accounts for the use of Γ_1^i and Δ_1^i in the conclusion. The aspects in Δ_1 and Δ_2 are merged in the overall expression, taking into account the way that x is used in e_2 .

The side condition “ \star ” is defined by Fig. 8.⁵ It prevents any common variable z being modified in e_1 or e_2 before it is referenced in e_2 . More exactly, $\Delta_1[z] = 1$ is not allowed (the value of the variable would be destroyed in the binding); $\Delta_1[z] = 3$ is always allowed (the value of the variable has no heap overlap with the binding value), and $\Delta_1[z] = 2$ is allowed provided neither $i = 1$ nor $\Delta_2[z] = 1$ (the value of the common variable may have aliasing with x , provided it is not partly or completely destroyed in e_2 : the modification may happen before the reference). Also, it is not safe to have $\Delta_1[z] = \Delta_2[z] = 2$ when $i = 2$ as such combination would allow one to create an illegally sharing tensor pair by $\text{let } x = z \text{ in } x \otimes z$.

As an instance of LET, we get a derived rule of contraction for aspect 3 variables:

$$\frac{\Gamma, x^i : A, y^3 : A \vdash e : B \quad i \geq 2}{\Gamma, x^i : A \vdash e[x/y] : B} \quad (\text{CONTR})$$

where $e[x/y]$ stands for $\text{let } y = x \text{ in } e$. This rule is comparable to the rule of

⁵ The table can be derived from the proof of soundness in Subsection 5.2, see Fig. 14 on page 27.

contraction given in (Hofmann, 2000) that allows contraction only for heap-free types A .

The only constructor rules we need are for the two kinds of pairs because the other constructors are function symbols in the signature. The rule for constructing a \times -pair ensures that all variables which are shared between the components have aspect at least 2. The “condition $\star\star$ ” in rule \times -PAIR is:

- $\Delta_1[z] \geq 2$ and $\Delta_2[z] \geq 2$ for all $z \in |\Delta_1| = |\Delta_2|$. ($\star\star$)

which ensures that no part of memory shared between the components is destroyed when the pair is constructed. If we made the evaluation or compilation dependent on the typing derivation, choosing which component to evaluate first, we could generalise condition $\star\star$ slightly, to say that either

- $\Delta_1[z] \geq 2$ and $\Delta_1[z] = 3 \implies \Delta_2[z] \geq 2$ for all $z \in |\Delta_1| = |\Delta_2|$, or
- $\Delta_2[z] \geq 2$ and $\Delta_2[z] = 3 \implies \Delta_1[z] \geq 2$ for all $z \in |\Delta_1| = |\Delta_2|$.

Intuitively, the reason for this condition is that we must be able to evaluate e_1 after e_2 or vice versa without corrupting the previous result.

In the destructor rules we type-check the branches in possibly extended contexts, and then pass the worst-case usage aspect as the usage for the term being destructed. For example, if we destroy one half of a pair in PAIR-ELIM, so x_1 has usage aspect 1, then the whole pair is considered destroyed in the conclusion. These rules could be simplified by unifying all usage aspects and relying on the RAISE rule. Nevertheless, we feel that the explicit calculation of the result aspect better expresses the intuition.

4 Imperative operational semantics

To establish the soundness of our typing rules, we need to formalise the intended in-place update interpretation of the language. In (Hofmann, 2000), a translation to C and a semantics for the target sub-language of C was used. Here we use an abstract machine model instead; this approach allows us to more easily consider alternative translations to other languages, such as the typed assembly language interpretation, as given in (Aspinall & Compagnoni, 2003).

Let Loc be a set of *locations* that model memory addresses on a heap. We use l to range over elements of Loc . Next we define two sets of values, *stack values* SVal , ranged over by v , and *heap values* HVal , ranged over by h , thus:

$$\begin{aligned} v & ::= c \mid l \mid \text{NULL} \mid (v, v) \\ h & ::= \{f_1 = v_1 \dots f_n = v_n\} \end{aligned}$$

A stack value is either a constant c (in our case an integer), a location l , a null value NULL , or a pair of stack values (v, v) . A heap value is an n -ary record consisting of named fields with stack values. The operational semantics is based around an abstract notion of stack and heap. A stack $S: \text{Var} \rightarrow \text{SVal}$ is a partial mapping from variables to stack values, and a heap $\sigma: \text{Loc} \rightarrow \text{HVal}$ is a partial mapping from locations to heap values. Evaluation of an expression e takes place with a given

stack and heap, and yields a stack value and a possibly updated heap. Thus we have a relation of the form

$$S, \sigma \vdash e \rightsquigarrow v, \sigma'$$

expressing that the evaluation of e under stack S and heap σ terminates and results in stack value v . As a side effect the heap is modified to σ' .

The only way the heap is modified is as a side effect of evaluating constructors that take \diamond arguments, following our in-place update interpretation of the language.

The stack is extended with additional variable bindings whenever we enter a new scope, inside sub-terms in the premises of the evaluation rules. When we evaluate a function body, we use a stack that mentions only the actual parameters, intuitively preventing access beyond the stack frame. Notice that the stack may contain pointers into the heap (i.e. locations) but there are no pointers going from the heap into the stack.

Most of the rules defining the evaluation relation are in Fig. 9. The rule APP does not apply to the built-in functions, such as `nil`, `cons`, because they have no bodies. Thus we also need the obvious rules for evaluating arithmetic and comparison operators, conditional expressions, and rules for trees, which are similar to those for lists. To represent trees on the heap, we store leaves as records $\{\text{label} = a\}$ and nodes as records $\{\text{label} = a, \text{left} = t_l, \text{right} = t_r\}$. The only interesting cases in the operational semantics are the ones for the heap datatypes, which make use of \diamond values as heap locations. In the CONS case, the first argument x_d of `cons` is a variable of \diamond type. The result is the location $S(x_d)$ where we make the cons cell by updating the heap, using a record with `hd` and `tl` fields. The match rule LIST-ELIM-CONS performs the opposite operation, exposing the contents and location of a cons cell.

5 Soundness

In this section we will prove that, for a well-typed program, the imperative operational semantics is sound with respect to a “safe” operational semantics that does not update in-place and is therefore equivalent to the usual functional denotational semantics. The safe operational semantics is expressed by the relation:

$$\eta \vdash e \rightsquigarrow_{\text{SF}} a$$

It is defined by very similar rules as the semantics \rightsquigarrow described in the previous section except that all references to a heap are dropped, stack values v are replaced with semantic values a and stacks S are replaced with valuations η , which map variables to semantic values.

We do not explicitly define semantic values because they are simply symbolic expressions built from the operators `null`, `cons(h, t)` for lists, `leaf(a)`, `node(a, l, r)` for trees and constants c, \diamond . A formal definition of semantic values is implicit in Definition 5.2 where it is shown how these values can be represented on the heap.

The only rules that differ significantly are the ones that overwrite a heap location, i.e. CONS and similar rules for trees:

$$\eta \vdash \text{cons}(x_d, x_h, x_t) \rightsquigarrow_{\text{SF}} \text{cons}(\eta(x_h), \eta(x_t)) \quad (\text{CONS})$$

$$\begin{array}{c}
 S, \sigma \vdash c \rightsquigarrow c, \sigma \quad (\text{CONST}) \\
 \\
 S, \sigma \vdash x \rightsquigarrow S(x), \sigma \quad (\text{VAR}) \\
 \\
 \frac{S(x_1) = v_1 \cdots S(x_n) = v_n \quad [x_1 \mapsto v_1, \dots, x_n \mapsto v_n], \sigma \vdash e_f \rightsquigarrow v, \sigma'}{S, \sigma \vdash f(x_1, \dots, x_n) \rightsquigarrow v, \sigma'} \quad (\text{APP}) \\
 \\
 \frac{S, \sigma \vdash e_1 \rightsquigarrow v_x, \sigma' \quad S[x \mapsto v_x], \sigma' \vdash e_2 \rightsquigarrow v, \sigma''}{S, \sigma \vdash \text{let } x = e_1 \text{ in } e_2 \rightsquigarrow v, \sigma''} \quad (\text{LET}) \\
 \\
 \frac{}{S, \sigma \vdash x_1 \otimes x_2 \rightsquigarrow (S(x_1), S(x_2)), \sigma} \quad (\otimes\text{-PAIR}) \\
 \\
 \frac{S(x) = (v_1, v_2) \quad S[x_1 \mapsto v_1][x_2 \mapsto v_2], \sigma \vdash e \rightsquigarrow v, \sigma'}{S, \sigma \vdash \text{match } x_p \text{ with } (x_1 \otimes x_2) \Rightarrow e \rightsquigarrow v, \sigma'} \quad (\text{PAIR-ELIM}) \\
 \\
 \frac{S, \sigma \vdash e_1 \rightsquigarrow v_1, \sigma' \quad S, \sigma' \vdash e_2 \rightsquigarrow v_2, \sigma''}{S, \sigma \vdash (e_1, e_2) \rightsquigarrow (v_1, v_2), \sigma''} \quad (\times\text{-PAIR}) \\
 \\
 \frac{S(x) = (v_1, v_2)}{S, \sigma \vdash \text{fst}(x) \rightsquigarrow v_1, \sigma} \quad (\text{FST}) \qquad \frac{S(x) = (v_1, v_2)}{S, \sigma \vdash \text{snd}(x) \rightsquigarrow v_2, \sigma} \quad (\text{SND}) \\
 \\
 S, \sigma \vdash \text{nil} \rightsquigarrow \text{NULL}, \sigma \quad (\text{NIL}) \\
 \\
 \frac{}{S, \sigma \vdash \text{cons}(x_d, x_h, x_t) \rightsquigarrow S(x_d), \sigma[S(x_d) \mapsto \{\text{hd} = S(x_h), \text{tl} = S(x_t)\}]} \quad (\text{CONS}) \\
 \\
 \frac{S(x) = \text{NULL} \quad S, \sigma \vdash e_{\text{nil}} \rightsquigarrow v, \sigma'}{S, \sigma \vdash \text{match } x \text{ with nil} \Rightarrow e_{\text{nil}} \mid \text{cons}(x_d, x_h, x_t) \Rightarrow e_{\text{cons}} \rightsquigarrow v, \sigma'} \quad (\text{LIST-ELIM-NIL}) \\
 \\
 \frac{S(x) = l \quad \sigma(l) = \{\text{hd} = v_h, \text{tl} = v_t\} \quad S[x_d \mapsto l, x_h \mapsto v_h, x_t \mapsto v_t], \sigma \vdash e_{\text{cons}} \rightsquigarrow v, \sigma'}{S, \sigma \vdash \text{match } x \text{ with nil} \Rightarrow e_{\text{nil}} \mid \text{cons}(x_d, x_h, x_t) \Rightarrow e_{\text{cons}} \rightsquigarrow v, \sigma'} \quad (\text{LIST-ELIM-CONS})
 \end{array}$$

Fig. 9. Evaluation relation with in-place update

Notice that x_d is completely ignored by this evaluation because no in-place update could possibly take place at this level of abstraction.

We now attempt to formulate the soundness theorem:

Preliminary Theorem 5.1 (Soundness)

Assume the following data and conditions:

- H1. a well-typed program P over some signature Σ ;
- H2. a well-typed term $\Gamma \vdash e : A$ over Σ for some Γ, e, A ;
- H3. a heap σ , a stack S and a valuation η such that each value $\eta(x)$ is *appropriately represented* by $S(x)$ and σ .

Then we have also:

- C1. $S, \sigma \vdash e \rightsquigarrow v, \sigma'$ implies there is a such that $\eta \vdash e \rightsquigarrow_{\text{sf}} a$ and v, σ' *appropriately represent* value a (i.e. in-place update evaluation is correct);
- C2. $\eta \vdash e \rightsquigarrow_{\text{sf}} a$ implies there are v and σ' such that $S, \sigma \vdash e \rightsquigarrow v, \sigma'$ and v, σ' *appropriately represent* value a (i.e. in-place update evaluation is complete).

The statement of the theorem is not yet in a form ready to be proved because assumption H3 as well as the conclusions are rather vague. We need to clarify what kinds of heap representations of values are appropriate for this theorem. Contrary to our initial expectation the soundness proof is highly non-trivial. In particular, the just mentioned “appropriateness” assumption needs to be elaborated in a non-obvious way with added separation conditions, as expressed by our separation theorem (Theorem 5.4).

5.1 Meaningful stack values in a heap

A stack value is *meaningful* for a particular type and heap if it has a sensible interpretation in the heap for that type, i.e. it *appropriately represents* a semantic value of this type. For instance, if $\sigma(v) = \{\text{hd} = 1, \text{tl} = \text{NULL}\}$ then v would be a meaningful stack value of type $L(\mathbf{N})$ with respect to σ and it would represent the singleton semantic list $[1]$ (using the usual notation for list expressions). In that same heap (v, v) would be a meaningful stack value of type $L(A) \times L(A)$ representing the semantic pair $([1], [1])$. Perhaps surprisingly, the value (v, v) will also be a meaningful stack value of type $L(A) \otimes L(A)$ in case it is used in a read-only fashion. This occurs for example in the term $f(x \otimes x)$ when $f : (A \otimes A)^3 \rightarrow B$. This means that “meaningfulness” of a stack value depends on the aspect with which the value is going to be used. We will parametrise our interpretation on a product separation flag that reflects this intuition. The flag value ∞ indicates that heap separation of the components of all tensor products is required while flag value ∞ indicates that sharing is allowed even inside tensor products.

To express heap separation, we need first to define the *region* $R_A(v, \sigma)$ of a stack value v of type A in heap σ . It is defined as the least set of locations satisfying the

$$\begin{aligned}
R_{\mathbf{N}}(c, \sigma) &= \emptyset & R_{\mathbf{L}(A)}(\mathbf{NULL}, \sigma) &= \emptyset & R_{\diamond}(l, \sigma) &= \{l\} \\
R_{A \times B}((v_1, v_2), \sigma) &= R_{A \otimes B}((v_1, v_2), \sigma) = R_A(v_1, \sigma) \cup R_B(v_2, \sigma) \\
\frac{\sigma(l) = \{\mathbf{hd} = h, \mathbf{tl} = t\}}{R_{\mathbf{L}(A)}(l, \sigma) = \{l\} \cup R_A(h, \sigma) \cup R_{\mathbf{L}(A)}(t, \sigma)} \\
\frac{\sigma(l) = \{\mathbf{label} = v\}}{R_{\mathbf{T}(A)}(l, \sigma) = \{l\} \cup R_A(v, \sigma)} \\
\frac{\sigma(l) = \{\mathbf{label} = v, \mathbf{left} = t_l, \mathbf{right} = t_r\}}{R_{\mathbf{T}(A)}(l, \sigma) = \{l\} \cup R_A(v, \sigma) \cup R_{\mathbf{T}(A)}(t_l, \sigma) \cup R_{\mathbf{T}(A)}(t_r, \sigma)}
\end{aligned}$$

Fig. 10. Heap region of a stack value

$$\begin{aligned}
& c, \sigma \Vdash_{\mathbf{N}}^p c & l, \sigma \Vdash_{\diamond}^p \diamond & \mathbf{NULL}, \sigma \Vdash_{\mathbf{L}(A)}^p \text{null} \\
\frac{v_k, \sigma \Vdash_{A_k}^p a_k \text{ for } k = 1, 2}{(v_1, v_2), \sigma \Vdash_{A_1 \times A_2}^p (a_1, a_2)} & \frac{v_k, \sigma \Vdash_{A_k}^p a_k \text{ for } k = 1, 2, \\ (p = \infty \Rightarrow R_{A_1}(v_1, \sigma), R_{A_2}(v_2, \sigma) \text{ are disjoint})}{(v_1, v_2), \sigma \Vdash_{A_1 \otimes A_2}^p (a_1, a_2)} \\
\frac{\sigma(l) = \{\mathbf{hd} = v_h, \mathbf{tl} = v_t\}, \quad v_h, \sigma \Vdash_A^p h, \quad v_t, \sigma \Vdash_{\mathbf{L}(A)}^p t, \\ (p = \infty \Rightarrow R_{\diamond}(l, \sigma), R_A(v_h, \sigma), R_{\mathbf{L}(A)}(v_t, \sigma) \text{ are pairwise disjoint})}{l, \sigma \Vdash_{\mathbf{L}(A)}^p \text{cons}(h, t)} \\
\frac{\sigma(l) = \{\mathbf{label} = v_a\}, \quad v_a, \sigma \Vdash_A^p a, \\ (p = \infty \Rightarrow R_{\diamond}(l, \sigma), R_A(v_a) \text{ are disjoint})}{l, \sigma \Vdash_{\mathbf{T}(A)}^p \text{leaf}(a)} \\
\frac{\sigma(l) = \{\mathbf{label} = v_a, \mathbf{left} = v_l, \mathbf{right} = v_r\}, \quad v_a, \sigma \Vdash_A^p a, \quad v_l, \sigma \Vdash_{\mathbf{T}(A)}^p l, \quad v_r, \sigma \Vdash_{\mathbf{T}(A)}^p r \\ (p = \infty \Rightarrow R_{\diamond}(l, \sigma), R_A(v_a, \sigma), R_{\mathbf{T}(A)}(v_l, \sigma), R_{\mathbf{T}(A)}(v_r, \sigma) \text{ are pairwise disjoint})}{l, \sigma \Vdash_{\mathbf{T}(A)}^p \text{node}(a, l, r)}
\end{aligned}$$

Fig. 11. Heap representation relation, with separation flag $p = \infty$ or ∞

rules in Figure 10 (if such a set exists). This is a type sensitive notion of reachability. It should be clear that $R_A(v, \sigma)$ is the part of the domain of σ that is relevant for v . According to this intuition, we have $R_A(v, \sigma) = R_A(v, \sigma')$ whenever $\sigma(l) = \sigma'(l)$ for all $l \in R_A(v, \sigma)$. Also, if A is a heap-free type, then $R_A(v, \sigma) = \emptyset$.

Definition 5.2

Given a stack value v , a heap σ , a semantic value a of type type A and a product separation flag $p \in \{\infty, \infty\}$, we define a five-place relation $v, \sigma \Vdash_A^p a$ which expresses that the semantic value a is *appropriately represented* by a *meaningful* stack value

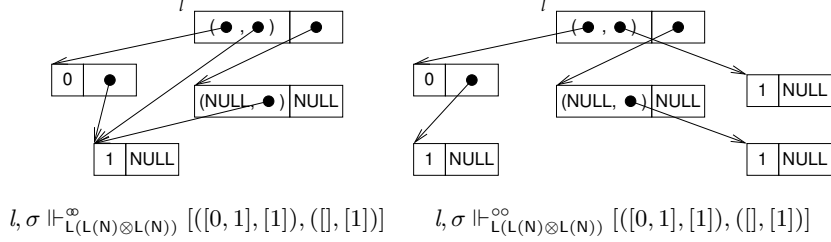


Fig. 12. Heap representation of $[([0, 1], [1]), ([], [1])]$ with and without sharing

v and heap σ with or without the condition of separation of tensor products. It is defined inductively by the rules shown in Figure 11. \square

When we want to ignore the semantic value, we can leave it out and write $v, \sigma \Vdash_A^p$.

The difference between \Vdash^{∞} and \Vdash° is that the latter prevents any “internal sharing” within \otimes -product types in the heap representation (see Fig. 12 for an illustration). We extend this relation to stacks, valuations and typing contexts:

Definition 5.3

For a stack S , heap σ , typing context Γ , a separation flag p and a valuation η , it holds that $S, \sigma \Vdash_{\Gamma}^p \eta$ if

- $S(x), \sigma \Vdash_{\Gamma(x)}^{p_{\Gamma(x)}} \eta(x)$ for each $x \in |\Gamma|$ where $p_1 = \circ, p_2 = p, p_3 = \infty$
- If for any two variables $x \neq y$ the regions $R_{\Gamma, x}(S, \sigma), R_{\Gamma, y}(S, \sigma)$ are *not* disjoint, then $\Gamma[x] \geq 2, \Gamma[y] \geq 2$ and if $p = \circ$, then moreover either $\Gamma[x] = 3$ or $\Gamma[y] = 3$.

\square

This definition amounts to saying that $S, \sigma \Vdash_{\Gamma}^{\circ} \eta$ holds if stack S and heap σ are meaningful at appropriate types and aspects, and moreover, the region for each aspect 1 variable does not overlap with the region for any other variable. (Informally: the aspect 1 variables are safe to update.) The case that $S, \sigma \Vdash_{\Gamma}^{\infty} \eta$ is stronger: variables with aspect 2 have internally separated tensor products and cannot share with each other. These extra separation conditions will guarantee that the result of a computation typed using this judgement is represented with its tensor products separated on the heap.

5.2 Separation theorem

With Definitions 5.2 and 5.3 in place, we can prove that the evaluation of a well-typed term under a meaningful stack and heap gives a meaningful result. This is the hardest part of the soundness theorem. Instead of finishing the formulation and proof of the soundness theorem, we will therefore first prove a theorem that states the separation properties of the in-place update evaluation and ignores the values themselves. It is easier to focus on separation before considering other aspects of soundness. Our soundness proof in Subsection 5.3 is an extension of the separation proof shown here.

Theorem 5.4 (Separation)

Assume the following data and conditions:

- H1. a well-typed program P over some signature Σ ;
- H2. a well-typed term $\Gamma \vdash e : A$ over Σ for some Γ, e, A ;
- H3. a heap σ and a stack S with $S, \sigma \Vdash_{\Gamma}^{\infty}$
(i.e. arguments are meaningful, separated according to aspects);
- H4. $S, \sigma \vdash e \rightsquigarrow v, \sigma'$.

Then we have also:

- C1. $R_A(v, \sigma') \subseteq \bigcup_{x \in \Gamma, \Gamma[x] < 3} R_{\Gamma, x}(S, \sigma)$
(i.e. the result is contained entirely within the heap space of arguments with aspect 1 or 2);
- C2. if $l \in R_{\Gamma, x}(S, \sigma)$ and $\Gamma[x] \geq 2$ then $\sigma(l) = \sigma'(l)$
(i.e. heap occupied by arguments with aspect 2 or 3 is not modified);
- C3. $v, \sigma' \Vdash_A^{\infty}$ (i.e. the result is meaningful);
- C4. $S, \sigma \Vdash_{\Gamma}^{\circ\circ}$ implies $v, \sigma' \Vdash_A^{\circ\circ}$
(i.e. if tensors are respected in the arguments, so they are in the result).

Specialising this, perhaps daunting, theorem to the particular case of a unary function on lists yields the following representative corollary:

Corollary 5.5

Let P be a well-typed program having a function symbol $f : \mathbf{L}(\mathbf{N})^i \rightarrow \mathbf{L}(\mathbf{N})$.

If σ is a store and l is a location such that l points in σ to a linked list with integer entries $w = [x_1, \dots, x_n]$ in σ and $[x \mapsto l], \sigma \vdash f(x) \rightsquigarrow v, \sigma'$ for some v, σ' , then v points in σ' to a linked list with integer entries.

Moreover, if $i \neq 1$, then the argument list w remains intact in σ' ; if $i = 3$, then the heap region $R_{\mathbf{L}(\mathbf{N})}(v, \sigma')$ of the result list is disjoint from the heap region $R_{\mathbf{L}(\mathbf{N})}(l, \sigma)$ of the argument.

Since our system has no memory allocation, we can conclude that if $i = 2$, the result must be some tail of the argument list; and if $i = 3$, the result must be the empty list.

To prove the theorem we need the following two elementary observations:

Lemma 5.6

Whenever we have two contexts Γ, Γ' that differ only in usage aspects and $\Gamma[x] \leq \Gamma'[x]$ for all x in Γ then

- 1. $S, \sigma \Vdash_{\Gamma}^{\infty}$ implies $S, \sigma \Vdash_{\Gamma'}^{\infty}$
- 2. $S, \sigma \Vdash_{\Gamma}^{\circ\circ}$ implies $S, \sigma \Vdash_{\Gamma'}^{\circ\circ}$

In plain words, raising usage aspects weakens the separation conditions.

Lemma 5.7

For any evaluation $S, \sigma \vdash e \rightsquigarrow v, \sigma'$ such that $S, \sigma \Vdash_{\Gamma}^{\infty}$ and $v, \sigma' \Vdash_A^{\infty}$, it holds

- 1. $R_A(v, \sigma') \subseteq R_{\Gamma}(S, \sigma)$
- 2. $\forall \ell \in \text{Dom}(\sigma) \setminus R_{\Gamma}(S, \sigma), \sigma(\ell) = \sigma'(\ell)$

In plain words, operational semantics neither overwrites memory that is not referenced via the context nor uses it in the result.

Further, we define a shortcut notation for several assertions about which we reason in the situation of the above lemma. The name *rg* stands for *heap region*, *sep* stands for *heap region separation*, *tens* for *tensor product separation* and *pres* for *heap preservation*.

$$\begin{aligned} |\Gamma|_i &= \{x \in \Gamma \mid \Gamma[x] = i\}, & |\Gamma|_{i,j} &= |\Gamma|_i \cup |\Gamma|_j \\ \text{rg}(x) &= R_{\Gamma,x}(S, \sigma) \\ \text{sep}(x, y) &\equiv \text{rg}(x) \cap \text{rg}(y) = \emptyset, & \text{sep}(T, T') &\equiv \bigwedge_{x \in T, x' \in T', x \neq x'} \text{sep}(x, x') \\ \text{tens}(x) &\equiv S(x), \sigma \Vdash_{\Gamma(x)}^{\circ}, & \text{tens}(T) &= \bigwedge_{x \in T} \text{tens}(x) \\ \text{pres}(x) &\equiv (\forall \ell \in \text{rg}(x)) (\sigma(\ell) = \sigma'(\ell)), & \text{pres}(T) &= \bigwedge_{x \in T} \text{pres}(x) \end{aligned}$$

Moreover, the symbol *res* is treated as a special variable representing the result on the heap σ' , i.e. $S(\text{res}) = v$, $\Gamma(\text{res}) = A$ and in relation to *res* the heap σ' is used instead of σ . When we want to relate these assertions to a particular evaluation $S, \sigma \vdash e \rightsquigarrow v, \sigma'$, we can write $v, \sigma', S, \sigma \Vdash$ in front of it. If there is no occurrence of *res* in it, we can leave v and σ' out. The relevant typing context Γ and result type A will always be clear from the context.

Using this notation we can express more succinctly the defining properties of heap representation for contexts:

$$\begin{aligned} S, \sigma \Vdash_{\Gamma}^{\circ} &\implies \text{sep}(|\Gamma|_1, |\Gamma|) \wedge \text{tens}(|\Gamma|_1) \\ S, \sigma \Vdash_{\Gamma}^{\circ} &\implies \text{sep}(|\Gamma|_1, |\Gamma|) \wedge \text{sep}(|\Gamma|_2, |\Gamma|_2) \wedge \text{tens}(|\Gamma|_{1,2}) \end{aligned}$$

Proof of the separation theorem (Theorem 5.4)

We proceed by induction on the lexicographic product of the depth of the derivations of $S, \sigma \vdash e \rightsquigarrow v, \sigma'$, i.e., computation time (1st priority) and of $\Gamma \vdash e : A$ (2nd priority). In most cases, we consider a derivation step of the operational semantics together with one corresponding typing rule. To be able to derive valid typing judgements for the premises of operational semantics rules, sometimes it is necessary to additionally use the WEAK, DROP and RAISE typing rules on their own.

CONST: Conclusions hold trivially because $\sigma = \sigma'$ and v is heap-free.

VAR: C1 holds because both sides are clearly equal—the region of the result is equal to the region of the only argument and this argument has aspect 2. C2 holds since $\sigma = \sigma'$. C3 follows directly from H3. C4 is similarly trivial.

APP: We assume $e = f(x_1, \dots, x_n)$ and $\Sigma(f) = A_1^{i_1}, \dots, A_n^{i_n} \rightarrow A$. Thus we have $x_1 : A_1^{i_1}, \dots, x_n : A_n^{i_n} \vdash e_f : A$ and $S, \sigma \vdash e_f \rightsquigarrow v, \sigma'$, using the operational and typing rules. Now the induction hypothesis applies and gives all the desired conclusions verbatim.

IF-TRUE: We assume $e = \text{if } x \text{ then } e_1 \text{ else } e_2$, $\Gamma = \Gamma_1, x : \mathbb{N}$, $\Gamma \vdash e_i : A$ for $i = 1, 2$ and $S(x) \neq 0$. Thus we must have $S, \sigma \vdash e_1 \rightsquigarrow v, \sigma'$ and the induction hypothesis can be applied to this judgement, and yields the desired conclusions verbatim.

IF-FALSE: Symmetric to IF-TRUE.

WEAK: Assume $\Gamma = \Gamma_1, \Gamma_2$ and $\Gamma \vdash e : A$ can be derived from a valid judgement $\Gamma_1 \vdash e : A$. Then conditions H2–H4 still hold after replacing Γ with Γ_1 and S with $S_1 = S|_{\Gamma_1}$. Thus we can apply the induction hypothesis and it yields the desired conclusions except that instead of Γ, S , they refer to Γ_1, S_1 respectively. The desired conclusions follow: C1 because the right hand side set of locations increases with Γ instead of Γ_1 ; C2 because heap space occupied by Γ_2 cannot share with Γ_1 except on the region of $|\Gamma_1|_{2,3}$, which is not modified, and the rest is not modified thanks to Lemma 5.7; C3 is identical to C3 in the induction hypothesis; C4 is weakened by a strengthened premise with S instead of S_1 .

DROP: The induction hypothesis can be applied thanks to Lemma 5.6. The resulting conclusions imply the desired conclusions because C1, C2, C4 become weaker and C3 remains the same when the aspect is dropped back to its original value.

RAISE: Assume $\Gamma_* \vdash e : A$ for some context Γ_* with $\Gamma = \Gamma_*^3$ and let A be heap-free. Thus H2 holds when replacing Γ by Γ_* . Also H3 holds after the change because the weaker form of heap representation does not make any distinction between aspects 2 and 3. We can therefore apply the induction hypothesis and obtain the desired conclusions with Γ replaced by Γ_* . The desired conclusion C3 is identical in this transformation and is thus proved. C4 follows trivially from C3 because A is heap-free. C2 does not make distinction between aspects 2 and 3 and is thus proved directly from the induction hypothesis. Finally, C1 is trivial for a heap-free type A .

⊗-PAIR, NIL, CONS, LEAF and NODE: Conclusions follow directly from the assumptions, like in the case of VAR.

PAIR-ELIM: We have $S, \sigma \vdash e \rightsquigarrow v, \sigma'$ where $e = \text{match } x \text{ with } x_1 \otimes x_2 \Rightarrow e'$, which implies $S(x) = (v_1, v_2)$ and $S'', \sigma \vdash e' \rightsquigarrow v, \sigma'$ where $S'' = S[x_1 \mapsto v_1][x_2 \mapsto v_2]$. Further we have $S, \sigma \Vdash_{\Gamma}^{\infty}$ where $\Gamma = \Gamma', x :^i A_1 \otimes A_2$. Put $\Gamma'' = \Gamma', x_1 :^{i_1} A_1, x_2 :^{i_2} A_2$.

We aim at proving $S'', \sigma \Vdash_{\Gamma''}^{\infty}$ in order to be able to apply the induction hypothesis. Basic separation conditions confined within Γ' translate directly from Γ to Γ'' . Since $i \leq i_1, i_2$, any separation required in Γ'' between x_1 or x_2 and some variable $y \in \Gamma'$ follows from the separation between x and y in Γ . Also, if $\text{tens}(x_1)$ or $\text{tens}(x_2)$ is required in Γ'' , then it can be obtained from $\text{tens}(x)$, which must be required in Γ in this case. Finally, if $\text{sep}(x_1, x_2)$ is required in Γ'' then $i = 1$, which implies that $\text{tens}(x)$ is required in Γ , which is sufficient for $\text{sep}(x_1, x_2)$ to hold. This concludes the proof of $S'', \sigma \Vdash_{\Gamma''}^{\infty}$.

From the induction hypothesis we get $v, \sigma' \Vdash_A^{\infty}$. It remains to prove the preservation and separation conclusions. Again, conditions confined to Γ' and the result translate from Γ'' trivially. If $i > 1$, we need to show $\text{pres}(x)$ and we can derive it from $\text{pres}(x_1)$ and $\text{pres}(x_2)$, which hold thanks to $i_1, i_2 \geq i > 1$. Similarly, we get $\text{rg}(\text{res}) \subseteq \text{rg}(|\Gamma''|_{1,2}) \subseteq \text{rg}(|\Gamma|_{1,2})$.

From the separation precondition for $\text{tens}(\text{res})$ present in $S, \sigma \Vdash_{\Gamma}^{\infty}$ we can deduce $S'', \sigma \Vdash_{\Gamma''}^{\infty}$ analogously to deducing the ∞ version above. Thus the desired guarantee $\text{tens}(\text{res})$ follows from $\text{tens}(\text{res})$ for the induction hypothesis, proving C4.

×-PAIR: From $S, \sigma \vdash (e_1, e_2) \rightsquigarrow v, \sigma''$ we get $S, \sigma \vdash e_1 \rightsquigarrow v_1, \sigma', S, \sigma' \vdash e_2 \rightsquigarrow v_2, \sigma''$ and $v = (v_1, v_2)$. It is straightforward to derive the assumptions of the first

induction hypothesis, in particular $S_1, \sigma \Vdash_{\Gamma_1, \Delta_1}^\infty$ where $S_1 = S|_{\Gamma_1, \Delta_1}$, from the present assumption $S, \sigma \Vdash_\Gamma^\infty$ where $\Gamma = \Gamma_1, \Gamma_2, \Delta_1 \wedge \Delta_2$. Thus the conclusions of the first induction hypothesis hold, in particular $v_1, \sigma' \Vdash_{A_1}^\infty$ where $A = A_1 \times A_2$ and $\text{pres}(|\Gamma_1, \Delta_1|_{2,3})$.

Analogously, for the second induction hypothesis we can define S_2 and obtain $S_2, \sigma \Vdash_{\Gamma_2, \Delta_2}^\infty$. We need to show that this representation still holds after the first evaluation, i.e. $S_2, \sigma' \Vdash_{\Gamma_2, \Delta_2}^\infty$. To make this conclusion, we realise that the region of Γ_2, Δ_2 on σ is preserved in the first evaluation, i.e. this region is separated from $|\Gamma_1, \Delta_1|_1$, which follows from $\text{sep}(|\Gamma|_1, (|\Gamma_2, \Delta_2|_1))$ and the side condition that Δ_2 does not contain any variable that has aspect 1 in Δ_1 . The other assumptions of the second induction hypothesis are straightforward. Thus we can use the conclusions of the second induction hypothesis, including $v_2, \sigma'' \Vdash_{A_2}^\infty$.

Since $R_{A_1}(v_1, \sigma') \subseteq R_{|\Gamma_1, \Delta_1|_{1,2}}(S|_{\Gamma_1, \Delta_1|_{1,2}}, \sigma)$, the region $R_{A_1}(v_1, \sigma')$ has not been modified by the evaluation of e_2 thanks to $\text{sep}(|\Gamma_2, \Delta_2|_1, \Gamma)$ and the side condition. Thus we also have $v_1, \sigma'' \Vdash_{A_1}^\infty$. Now we can deduce $(v_1, v_2), \sigma'' \Vdash_{A_1 \times A_2}^\infty$. From the conclusions of both hypotheses it is easy to deduce also $\text{pres}(|\Gamma|_{2,3})$ and $\text{rg}(\text{res}) \subseteq \text{rg}(|\Gamma|_{1,2})$ using the straightforward properties

$$|\Gamma|_{2,3} \subseteq |\Gamma_1, \Delta_1|_{2,3} \cup |\Gamma_2, \Delta_2|_{2,3} \quad \text{and} \quad |\Gamma_1, \Delta_1|_{1,2} \cup |\Gamma_2, \Delta_2|_{1,2} \subseteq |\Gamma|_{1,2}.$$

When $\text{sep}(|\Gamma|_2, |\Gamma|_2)$ and $\text{tens}(|\Gamma|_2)$ hold on σ then the same holds for Γ_1, Δ_1 and Γ_2, Δ_2 on their respective heaps σ and σ' . Thus v_1 and v_2 are represented on σ'' with separated tensor products. As there is no tensor separation needed between v_1 and v_2 , we get $(v_1, v_2) = v, \sigma'' \Vdash_{A_1 \times A_2}^\infty$, i.e. the final conclusion C4.

LIST-ELIM and TREE-ELIM Proofs for these cases can be obtained by a straightforward combination and adaptation of the proofs for IF and PAIR-ELIM.

LET: The skeleton of this part of the proof is in Fig. 13. A typical row of the table contains a statement preceded by its code and followed by a list of references to statements and definitions from which it has been derived. (We will elaborate on these, often non-trivial, deductions shortly.) The statements whose names end with *i1* are valid under the condition that $i = 1$, etc. A proof of conclusion C4 includes the dotted parts while the proof of C1–C3 is obtained by ignoring the dotted parts. To improve clarity, the table also includes the definitions of certain important symbols, each definition immediately preceding the first use of the symbol.

The horizontal lines divide the proof into three parts. The first and second parts prove that the induction hypothesis can be applied to the first and second sub-term (i.e. e_1 and e_2), respectively. The final part shows the validity of conclusions C1–C4 for the whole let expression.

The proof is much more involved than originally anticipated. One unexpected aspect of this proof is that the stronger conclusion C4 of the induction hypothesis for e_1 is needed even when proving conclusions C1–C3 of the theorem when $i = 1$ (see condition *Hr1*). Moreover, even when proving C4, we need C3 and not C4 for the induction hypothesis for e_1 when $i = 3$. Intuitively, we need the stronger ($p = \infty$) interpretation of the theorem even when proving the weaker version ($p = \infty$) and vice versa. This means that we cannot decouple the two

Implicit conditions: $(\dots i1) \sim (i = 1)$, $(\dots i12) \sim (i \in \{1, 2\})$, etc.

Proof of conclusion C4 additionally includes all the dotted material.

(Op)	$S, \sigma \vdash \text{let } x = e_1 \text{ in } e_2 \rightsquigarrow v, \sigma''$	assumption H4
$(Op1,2)$	$S, \sigma \vdash e_1 \rightsquigarrow v_x, \sigma' \quad S[x \mapsto v_x], \sigma' \vdash e_2 \rightsquigarrow v, \sigma''$	Op, \rightsquigarrow
(Hs)	$S, \sigma \Vdash_{\Gamma}^{\circ} \quad S, \sigma \Vdash_{\Gamma}^{\circ\circ}$	assumption H3 $\overset{\circ}{C4a}$
(Sep)	$S, \sigma \Vdash \text{sep}(\Gamma _1, \Gamma) \wedge \text{sep}(\Gamma _2, \Gamma _2) \wedge \text{tens}(\Gamma _{1,2})$	Hs
	$(\Gamma = \Gamma_1^i, \Gamma_2, \Delta_1^i \wedge \Delta_2 \text{ as used in the derivation of } \Gamma \vdash e : A)$	
$(1-23)$	$ \Gamma_1, \Delta_1 _1 \cup \Gamma_2, \Delta_2 _1 \subseteq \Gamma _1 \quad \Gamma _{2,3} \subseteq \Gamma_1 _{2,3} \cup \Gamma_2 _{2,3} \cup (\Delta_1 _{2,3} \cap \Delta_2 _{2,3})$	
	$(S_1 = S _{\Gamma_1, \Delta_1})$	
$(Hs1)$	$S_1, \sigma \Vdash_{\Gamma_1, \Delta_1}^{\circ}$	$Hs, Sep, \Gamma_1^i, \Delta_1^i$
$(Hs1i2)$	$S_1, \sigma \Vdash_{\Gamma_1, \Delta_1}^{\circ\circ}$	$Hs, Sep, \Gamma_1^i, \Delta_1^i$
<hr/>		
$(Hr1i2)$	$v_x, \sigma' \Vdash_{A_x}^{\circ\circ}$	induction 1, C4
$(Hr1i23)$	$v_x, \sigma' \Vdash_{A_x}^{\circ}$	induction 1, C3
$(G1)$	$v_x, \sigma', S_1, \sigma \Vdash \text{pres}(\Gamma_1, \Delta_1 _{2,3}) \wedge \text{rg}(\text{res}) \subseteq \text{rg}(\Gamma_1, \Delta_1 _{1,2})$	ind. 1, C1, C2
	$(S_2 = S _{\Gamma_2, \Delta_2})$	
$(Hs2p)$	$S_2, \sigma \Vdash_{\Gamma_2, \Delta_2}^{\circ}$	$S_2, \sigma \Vdash_{\Gamma_2, \Delta_2}^{\circ\circ} \quad Hs$
$(Hs2g)$	$S_2, \sigma' \Vdash_{\Gamma_2, \Delta_2}^{\circ}$	$S_2, \sigma' \Vdash_{\Gamma_2, \Delta_2}^{\circ\circ} \quad Hs2p, Sep, 1-23, G1, \text{ side cond.}$
	$(S_x = S_2[x \mapsto v_x])$	
$(S2xi1)$	$S_x, \sigma' \Vdash \text{sep}(x, (\Gamma_2, \Delta_2)) \wedge \text{tens}(x)$	$Sep, G1, Hr1i2, \text{ side cond.}$
$(S2xi2)$	$S_x, \sigma' \Vdash \text{sep}(x, \Gamma_2, \Delta_2 _{1,2}) \wedge \text{tens}(x)$	$Sep, Hr1i2, \Delta_1^i \wedge \Delta_2, \text{ side cond.}$
$(S2xi3)$	$S_x, \sigma' \Vdash \text{sep}(x, \Gamma_2, \Delta_2 _1)$	$Sep, \text{ side cond.}$
	$(\Gamma_x = \Gamma_2, \Delta_2, x : A_x)$	
$(Sep2)$	$S_x, \sigma' \Vdash \text{sep}(\Gamma_x _1, \Gamma_x) \wedge \text{sep}(\Gamma_x _2, \Gamma_x _2) \wedge \text{tens}(\Gamma_x _{1,2})$	$Hs2g, S2x$
$(Hs2)$	$S_x, \sigma' \Vdash_{\Gamma_x}^{\circ}$	$Hs2g, Hr1, Sep2$
<hr/>		
$(Hr2)$	$v, \sigma'' \Vdash_{A'}^{\circ}$	induction 2, C3 $\overset{\circ}{C4}$
$(G2)$	$v, \sigma'', S_x, \sigma' \Vdash \text{pres}(\Gamma_x _{2,3}) \wedge \text{rg}(\text{res}) \subseteq \text{rg}(\Gamma_x _{1,2})$	induction 2, C1, C2
$(Ri12)$	$\text{rg}(\text{res}) \subseteq \text{rg}(\Gamma_2, \Delta_2 _{1,2}) \cup \text{rg}(\Gamma_1, \Delta_1 _{1,2}) = \text{rg}(\Gamma _{1,2})$	$G2, G1$
$(Ri3)$	$\text{rg}(\text{res}) \subseteq \text{rg}(\Gamma_2, \Delta_2 _{1,2}) \subseteq \text{rg}(\Gamma _{1,2})$	$G2$
$(a23i1)$	$ \Gamma _{2,3} = \Gamma_1 _3 \cup \Gamma_2 _{2,3} \cup (\Delta_1 _3 \cap \Delta_2 _{2,3})$	$\Delta_1^i \wedge \Delta_2$
$(S1-23)$	$(\Delta_1 _{2,3} \cap \Delta_2 _{2,3}) \text{ disjoint from } \Delta_1 _1 \cup \Delta_2 _1$	$\Delta_1^i \wedge \Delta_2$
(P)	$v, \sigma'', S, \sigma \Vdash \text{pres}(\Gamma _{2,3})$	$G1, G2, 1-23, a23i1, S1-23, Sep$
(G)	$v, \sigma'', S, \sigma \Vdash \text{pres}(\Gamma _{2,3}) \wedge \text{rg}(\text{res}) \subseteq \text{rg}(\Gamma _{1,2})$	P, R

Fig. 13. Evaluation of let expressions preserves separation properties

		$\Delta_1[z] \rightarrow$	1			2			3		
		$\downarrow i/\Delta_2[z] \rightarrow$	1	2	3	1	2	3	1	2	3
$Hs2g$	dashed region	1	✗	✗	✗	✗	✗	✗	✓	✓	✓
$S2xi1$	dotted region $i = 1$	2	✗	✗	✗	✗	✗	✓	✓	✓	✓
$S2xi2$	dotted regions $i = 2$	3	✗	✗	✗	✗	✓	✓	✓	✓	✓
$S2xi3$	dotted regions $i = 3$	3	✗	✗	✗	✗	✓	✓	✓	✓	✓

Fig. 14. Condition \star and proof of the LET case.

interpretations from each other and prove them separately; we have no choice but to prove them both in parallel.

The core of the proof lies in the statements $Hs2g$, $S2xi1$, $S2xi2$ and $S2xi3$, which claim that the evaluation of e_1 has not interfered with the context for e_2 and prepared the interim result in such a way that it fits the representation criteria of the typing for e_2 . These statements directly rely on the side condition. Fig. 14 is a copy of Fig. 8 in which it is shown on which forbidden cases each of these four statements relies.

Individual statements up to induction 1. Op is an elaborated version of assumption H4. Statements $Op1,2$ follow from the operational LET rule that we assume was used to derive Op . Statement Hs (i.e. heap representation of the stack) combines H3 with the assumption of C4. Hs is further elaborated as Sep . Observing the particular form Γ takes thanks to the typing rule LET, we make the elementary observation that any variable with aspect 1 in either of the two sub-expressions has aspect 1 in Γ (1-23). The second inclusion in 1-23 is a rephrasing of the first inclusion, focusing on aspects 2,3 instead of aspect 1.

At this point all assumptions of the induction hypothesis for e_1 are established by the typing and operational rules except for H3 (i.e. statement $Hs1$). When $i = 1$, we will also need conclusion C4 and will therefore need to strengthen $Hs1$ as shown in the table. The stronger stack representation clause is also needed when $i = 2$ while proving C4 for the let expression. This is because separation of tensors will be required within x in these cases.

Why does $Hs1$ follow from Hs ? All separation requirements of Hs (i.e. Sep) restricted to the variables in Γ_1, Δ_1 are stronger than the separation requirements in $Hs1$ by Lemma 5.6 because the aspects of these variables are lower in Hs than they are in $Hs1$. If $i = 1$, then all variables with aspect 2 in Γ_1, Δ_1 gain aspect 1 in Γ and thus, by Hs , have their tensors separated and are separated from each other. Thus we have the stronger version of $Hs1$ when $i = 1$. When observing the dotted reading, the stronger assumption holds in Hs and thus also $Hs1$ has separated tensors and variables with aspect 2, unless $i = 3$ because then variables with aspect 2 get aspect 3 in Γ and we do not know anything about them from Hs . This concludes the proof of $Hs1$ and the induction hypothesis can be used as required, yielding, in particular, statements $Hr1$ and $G1$.

Individual statements up to induction 2. The next major goal is to derive the assumptions of the induction hypothesis for e_2 and its corresponding context $\Gamma_x = \Gamma_2, \Delta_2, x.^iA_x$ and stack S_x . Again, the only hard part of this part of the proof is to derive that the stack values are represented properly with all the separation conditions as dictated by the aspects in the context (i.e. assumptions H3 and C4a, formulated in statement $H2s$, elaborated as $Sep2$).

It is quite easy to see that, apart from x , all the values used in e_2 are properly separated on the original heap σ (statement $Hs2p$) because their aspects are not raised in Γ , although they may be lowered. To lift this condition to heap σ' (statement $Hs2g$), we need to establish that the entire region of these values was

left unmodified by the evaluation of e_1 , i.e. there is no aliasing between aspect 1 variables in Γ_1, Δ_1 and the region of S_2 on heap σ . This claim is trivial for Γ_2 because of its inclusion in Γ (statement *Sep*) and the observation 1-23. For Δ_2 , however, we have to additionally use the fact that the side condition forbids a shared variable to appear in e_1 with aspect 1. This concludes the proof of *Hs2g*. Now we need to prove that x on σ' is separated from the remainder of the context Γ_x as required by *Sep2*. To do so, we will consider the three possible values of i separately.

$i = 1$ (*S2xi1*): The required internal separation of x follows straight from *Hr1*.

The separation of x from the rest of the context follows from the fact that the value of x is contained within the region of variables with aspect 1 or 2 in Γ_1, Δ_1 . All these variables are separated from everything else thanks to *Sep* and the definition of Γ_1^1, Δ_1^1 and are forbidden to appear in Δ_2 by the side condition.

$i = 2$ (*S2xi2*): The minimum requirement here is that x is separated from all variables with aspect 1 in Γ_2, Δ_2 . This claim can be argued in the same way as in the previous case from *Sep* and the side condition.

In the dotted reading, we also need separation from variables with aspect 2 and the separation of tensors in x , which follows straight from *Hr1*. The separation of x (i.e. variables with aspect 1 or 2 in Γ_1, Δ_1) from variables with aspect 2 in Γ_2, Δ_2 can be derived from *Sep* and the side condition as follows. No variable can appear in both sets under consideration because the side condition does not allow any variable to have aspect 1 or 2 in Δ_1 and aspect 2 in Δ_2 . The rest follows directly from the dotted version of *Sep*, considering the fact that the aspects of these variables cannot change to 3 in Γ .

$i = 3$ (*S2xi3*): The side condition implies that variables with aspect 1 in Δ_2 must have aspect 3 in Δ_1 . The variables with aspect 1 in Γ_2, Δ_2 still have aspect 1 in Γ . Thus *Sep* guarantees that variables with aspect 1 or 2 in Γ_1, Δ_1 (which contain the region of x) are separated from the variables with aspect 1 in Γ_2, Δ_2 because the two sets of variables are disjoint.

Putting together all statements from *S2x* and *Hs2g*, we get *Sep2*, which is a rephrasing of *Hs2*. This statement allows us to use the induction hypothesis for expression e_2 , yielding the conclusions *Hr2* and *G2*.

Individual statements after induction 2. *Hr2* is equal to the desired conclusions C3 and C4. The remainder of the proof is concerned with proving C1 and C2, which have been rephrased as statement *G*.

To prove that the result is contained in $|\Gamma|_{1,2}$, consider that we know from *G2* that it is contained in $rg(|\Gamma_x|_{1,2})$, which itself is contained in $rg(|\Gamma_2, \Delta_2|_{1,2}) \cup rg(|\Gamma_1, \Delta_1|_{1,2})$ if $i < 3$ (thanks to *G1*) and is equal to $rg(|\Gamma_2, \Delta_2|_{1,2})$ if $i = 3$. The proof of *Ri3* is concluded when we take into account $|\Gamma_2, \Delta_2|_{1,2} \subseteq |\Gamma|_{1,2}$. The other case (i.e. *Ri12*) needs additionally $rg(|\Gamma_1, \Delta_1|_{1,2}) \subseteq |\Gamma|_{1,2}$, which holds whenever $i < 3$ because in this case the aspect of a variable cannot be higher in Γ than in Γ_1, Δ_1 .

It remains to prove the preservation of the heap contents for variables $|\Gamma|_{2,3}$ (statement P). First consider which of these variables play which role in the two sub-evaluations, using elementary statements $1-23$ and $a23i1$. Next we will show that these preserved variables are not aliasing with what may be modified by the sub-evaluations. This claim follows almost straightforwardly from Sep because the modified variables have aspect 1, but we also need to confirm that the two sets of variables do not overlap. The only danger of overlap is for variables within Δ_1 but the elementary statement $S1-23$ shows that no overlap happens for these variables either. Now we can see that the variables $|\Gamma|_{2,3}$ are not modified in either of the two evaluations.

The last danger comes from the fact that variables with aspect 2 in e_1 may share with x , which may be destroyed if $i = 1$. Nevertheless, in this case the variables that have aspect 2 in e_1 have the aspect raised to 1 in Γ and are therefore not obliged to be preserved (which is implicit in $a23i1$). This concludes the proof of P , which, as stated earlier, concludes the proof of the theorem for let expressions.

□

5.3 Soundness theorem

Now we are ready to complete Preliminary Theorem 5.1 by interpreting its assumption H3 as $S, \sigma \Vdash_{\Gamma}^{\circ} \eta$ and its conclusions as $v, \sigma' \Vdash_A^{\circ} a$:

Theorem 5.8 (Soundness)

Assume the following data and conditions:

- H1. a well-typed program P over some signature Σ ;
- H2. a well-typed term $\Gamma \vdash e : A$ over Σ for some Γ, e, A ;
- H3. a heap σ , a stack S and a valuation η such that $S, \sigma \Vdash_{\Gamma}^{\circ} \eta$;

Then we have also:

- C1. $S, \sigma \vdash e \rightsquigarrow v, \sigma'$ implies there is a such that $\eta \vdash e \rightsquigarrow_{\text{sf}} a$ and $v, \sigma' \Vdash_A^{\circ} a$ (i.e. in-place update evaluation is correct);
- C2. $\eta \vdash e \rightsquigarrow_{\text{sf}} a$ implies there are v and σ' such that $S, \sigma \vdash e \rightsquigarrow v, \sigma'$ and $v, \sigma' \Vdash_A^{\circ} a$ (i.e. in-place update evaluation is complete).

Proof

Conclusion C1 is obtained straightforwardly by adapting the separation theorem (Theorem 5.4) and its proof to include semantic values alongside the separating heap representation statements, e.g. $v, \sigma' \Vdash_A^{\circ} a$ instead of $v, \sigma' \Vdash_A$. We thus consider conclusion C2 proved.

The assumptions of this theorem for proving conclusion C2 are almost identical to the assumptions of the separation theorem. The only difference is in switching from ordinary operational semantics to safe operational semantics. We will therefore be able to reuse the structure and certain parts of the proof of the separation theorem.

Like in the separation theorem, the proof of C2 is by induction simultaneously on the derivation of the operational semantics and on the derivation of the typing

judgement. The rules for safe operational semantics are in an obvious one-to-one correspondence with the rules for in-place updating operational semantics and the corresponding rules are quite similar.

Whenever we need to use the induction hypothesis, one of the most laborious steps is to prove assumption H3 for the sub-expression. Fortunately, this step can always be reused from the proof of the separation theorem without any modifications. We will therefore omit this step in our proof description. The reason why we can reuse the proof of H3 for sub-expressions from the separation theorem is that the proofs made no use of the operational semantics. In cases where there are two sub-expressions, when proving H3 for the second one, the conclusions of the separation theorem are needed for the first sub-expression. This can be obtained here by the separation theorem applied on the first sub-expression with the missing operational semantics condition (H4 in Theorem 5.4) supplied by the first induction hypothesis.

CONST: In the ordinary operational semantics, we get $S, \sigma \vdash c \rightsquigarrow c, \sigma$, trivially proving the conclusion.

VAR: This case is very similar to **CONST** except that we need to realise that $S(x), \sigma \Vdash_{\Gamma}^{\infty} \eta(x)$ by H4 to prove that the value is equivalent.

APP, IF-TRUE, IF-FALSE, DROP, RAISE: We can apply the induction hypothesis with an updated typing judgement and operational rule but the same S , η and a to give us the desired conclusions verbatim.

WEAK: Assume $\Gamma = \Gamma_1, \Gamma_2$ and $\Gamma \vdash e : A$ can be derived from a valid judgement $\Gamma_1 \vdash e : A$. We can apply the induction hypothesis on the same assumptions except for replacing Γ with Γ_1 , S with $S_1 = S|_{\Gamma_1}$ and η with $\eta_1 = \eta|_{\Gamma_1}$. In the conclusions that we get we can now substitute back S for S_1 because extending the stack in operational semantics does not invalidate it.

⊗-PAIR, NIL, LEAF and NODE: Conclusions follow directly from the assumptions, like in the case of **VAR** and **CONST**.

CONS: Aligning the **CONS** rules for safe and in-place updating operational semantics allows us to draw conclusion C1, provided we can show that $S(x_d)$ is not in $R_{\Gamma(x_h)}(S(x_h), \sigma)$ and $R_{\Gamma(x_t)}(S(x_t), \sigma)$. This is guaranteed by condition H3 and the fact that x_d has aspect 1.

PAIR-ELIM: We apply the induction hypothesis on the only operational premise, substituting $\Gamma, x{:}^i A_1 \otimes A_2$ with $\Gamma, x_1{:}^{i_1} A_1, x_2{:}^{i_2} A_2$ and η with $\eta[x_1 \mapsto a_1, x_2 \mapsto a_2]$. We get the desired conclusions except for having to apply the imperative operational rule to turn $S[x_1 \mapsto v_1, x_2 \mapsto v_2]$ back to S . The remainder of the induction hypothesis is identical to the desired statement.

LET: We can trivially apply the induction hypothesis on the first operational premise and also on the second premise thanks to the ability to reuse the proof from Theorem 5.4 to get H3 for it. These hypotheses yield the two premises of the imperative operational rule **⊗-PAIR**, which then derives the operational statement in conclusion C1. The representation of the final value follows directly from the second induction hypothesis.

⊗-PAIR: This is very similar to **LET**, except that we need to work harder to prove that the final value is equivalent to the semantic value.

To obtain $(v_1, v_2), \sigma'' \Vdash_{A_1 \times A_2}^\infty (a_1, a_2)$, we need $v_1, \sigma'' \Vdash_{A_1}^\infty a_1$, which follows from $v_1, \sigma' \Vdash_{A_1}^\infty a_1$ and the fact that evaluation of e_2 does not modify the region of v_1 , which follows from:

- conclusion C2 of the separation theorem applied on the second premise;
- conclusion C1 of the separation theorem applied on the first premise;
- the side condition of \times -PAIR.

□

6 Conclusions and related work

In this paper we introduced a typing scheme that relaxes linear or affine typing schemes, allowing some cases of contraction. Variables can be used more than once in certain points in program expressions. The points of re-use are determined by three *usage aspects* that characterise variable usages as *destructive* (aspect 1), *read-only shared* (aspect 2) and *read-only unshared* (aspect 3). Aspect 2 is the most interesting and novel: it captures when the machine representations of function arguments and result may share memory (i.e., exhibit aliasing). It allows the aliasing to be tracked through expressions.

Our typing scheme was demonstrated in a version of the resource-controlled linear functional programming language LFPL (Hofmann, 2000), which has an affine typing scheme that additionally distinguishes between heap-allocated datatypes and register or stack-allocated ones. These “heap-free” types are always given aspect 3 here. LFPL has an intended implementation using an in-place update interpretation for datatype construction; the resource type \diamond is used to represent a unit of heap space that can be destructively updated to construct a cell of a data structure. With usage aspects, a datatype constructor typing always attaches aspect 1 to its resource argument.

We gave an operational semantics to formalise a machine model of the in-place update interpretation of LFPL. We proved that usage aspect typing ensures that evaluation in the language is both type sound for a memory model and equivalent to the standard call-by-value semantics, which corresponds to the usual denotational functional semantics. The denotational correctness establishes that the program is always properly evaluated; it justifies reasoning about programs using their functional semantics, without caring about the underlying imperative behaviour. The proof is highly non-trivial because of the semantics of aspects in contexts, especially the need to establish separation and preservation properties between memory areas occupied by variables in the context and the results of expression evaluation. This suggests that a simple usage aspect typing can express a rather complex aliasing and separation relation. We believe that the result here goes beyond most existing proofs of type safety in the literature for related type systems.

6.1 Related work

The closest strand of work begins with Wadler’s introduction of the idea of a sequential let (Wadler, 1990). If we assume that e_1 is evaluated before e_2 in the

expression

let $x=e_1$ in e_2 ,

then we can allow sharing of a variable z between e_1 and e_2 , as long as z is not modified in e_1 and the type of e_1 is such that its value x cannot possibly be aliased with z . This condition is strictly stronger than z having aspect 3 in our system: e.g. our example term $S2$ in Subsection 2.3 would not type-check translated into Wadler's language in a straightforward way:

let! (m) $l' = \text{filter_similar}(m, l)$ **in** $\text{map_reverse}(\text{cons}(m, l'))$

because the type of l' , i.e. $L(L(N))$ is potentially sharing with the type of m , i.e. $L(N)$.

Odersky (1992) improved Wadler's system by introducing *observer* type annotations, which semantically correspond to our aspect 3: example $S2$ works in his system. In certain ways, Odersky's system is more powerful than ours: it supports polymorphism, higher order functions and observer status can be independently granted to individual components of an argument's type rather than the argument as a whole. Nevertheless, the observer annotations cannot model aspect 2, which can be demonstrated with our example $S1$. In Odersky's system, when typing the term

let! $b = \text{find_best}(l)$ **in** $\text{first_similar}(b, l)$

the inner list within the type of l , i.e. $L(L(N))$, will have to be non-linear because it cannot be made an observer in the first term where it shares with the result b .

Kobayashi (1999) introduces *quasi-linear* types. This typing scheme also allows sharing in let expressions. It has a δ use, which corresponds roughly to our aspect 3 usage. Kobayashi's motivation was to detect statically points where deallocation occurs; this requires stack-managed extra heap, augmenting region analysis (Tofte & Talpin, 1997). To illustrate Kobayashi's system, we quote the typing of the append function from (Kobayashi, 1998), page 30:

$$\forall \alpha, \beta, \gamma, i, j, k :: \{\alpha \geq \gamma, \beta \geq \gamma, i \geq \delta, j \geq k, j \geq \delta, k \geq [k], l \geq \delta\}.$$

$$((\alpha \text{ list}^i \times^l \beta \text{ list}^j) \rightarrow^{\omega, \omega} \gamma \text{ list}^k)$$

The judgement entails that the cons cells of both input lists have been used but are not bound to the result ($i \geq \delta, j \geq \delta$) and thus can be safely deallocated and re-used. In particular, it is safe to implement the append in-place. Also, the judgement says that the elements of the lists are not touched ($\alpha \geq \gamma, \beta \geq \gamma$) and there are no restrictions on how many times the function may be safely used ($\rightarrow^{\omega, \omega}$).

The example typing shows that our usage aspect annotation is considerably more concise than Kobayashi's typings, although we assume a more restricted scenario where an in-place update interpretation is intended rather than inferred. Kobayashi proves a traditional type soundness (subject reduction) property, which shows an internal consistency of his system, whereas we have characterised and proved equivalence with an independently meaningful semantic property. It might well be possible to prove similar results to ours for Kobayashi's system, but we believe that

by considering the semantical property at the outset, we have introduced a rather more natural syntactic system, with simpler types and typing rules.

To summarise so far, in the realm of functional programming, the work that is related to ours mostly encompasses our usage aspect 3, but does not have anything analogous to our aspect 2. There is also a long history of work on aliasing in the realm of imperative programming, most lately introducing expressive type systems and specialised logics, which we outline next. Although type systems in this setting may appear to be closely comparable to ours, there are important differences. For example, it seems hard to concisely represent aspect 2 in an imperative setting: an imperative procedure may have many non-local effects during its execution that could introduce aliasing at almost arbitrary positions in the store, leading to the need to introduce names for locations (and then polymorphism over those names). Things are more controlled in a functional setting, since a function returns a single result, so we can express the aspect 2 concisely as aliasing with this single result and monitor the effect of the introduced aliasing more easily.

Some formalisms for reasoning about aliasing in imperative programs are flexible enough to express and prove safe an aliasing corresponding to our aspect 2. Probably the closest of these to our system are *fractional permissions* (Boyland, 2003), which track all aliases of a linear value across procedures (and even parallel threads) and statically detect when the value becomes unique and can be safely overwritten. The tool for analysing C programs described in (Evans, 1996) has several annotations whose meaning is related to our aspects. Although this system has been popular and useful for finding errors, its theoretical underpinning is limited and its guarantees are therefore much weaker than those of the other systems mentioned. Also, separation logic (Reynolds, 2002) and alias types (Smith *et al.*, 2000) can be used to express and verify safety of aliasing across procedures.

In separation logic and the logic of bunched implications (Ishtiaq & O’Hearn, 2001), it is natural to express aliasing properties but not a read-only behaviour. This means that without amendments neither aspect 2 or 3 can be expressed. One way to express read-only behaviour that has been suggested is to combine these logics with passive types from syntactic control of interference (Reynolds, 1978; O’Hearn *et al.*, 1995). Another way to express read-only properties in separation logic is by the use of auxiliary variables that provide a link between the heap before and after the execution of a command. But again, such systems provide much more complicated assertions than ours, and need sophisticated machinery to reason with them.

Separation logic and alias types and are complex general mechanisms in which automatic checking of higher-level type safety (such as aliasing over values of recursive datatypes) is difficult to obtain. Some works close to alias types, e.g. (Fahndrich & DeLine, 2002; Aiken *et al.*, 2003) and the work by Boyland already mentioned, do make steps in this direction.

All the mentioned methods for analysing memory aliasing in imperative programs make assertions about individual memory cells or unstructured heap blocks. Defining recursive datatypes and expressing aliasing properties for them as a whole is not straightforward. It is apparent that, e.g. for alias types (Walker & Morrisett, 2001),

the resulting assertions are complex and it will be more difficult to infer them than to infer our usage aspects.

We would like to stress again that imperative languages do not admit an intuitive denotational semantics comparable to that of our functional language. Consequently, it does not make as much sense to prove the correctness of the execution of well-typed programs in the quoted systems as much as it does in ours. Indeed, the mentioned works prove only type safety of the execution, not relating the behaviour of the programs to their specifications. Our approach gives the programmer a transparent denotational semantics that can be used to check that programs meet their specifications. Reasoning at this level, usage aspects and diamond values can be completely ignored; our proof guarantees that the imperative interpretation with aliasing is sound. Diamond values or other means of explicit memory deallocation or reuse allow the programmer to safely optimise the memory usage of their programs without changing its functional behaviour.

There is further less-closely related work on formal systems for reasoning or type-checking in the presence of aliasing, including for example work on the imperative λ -calculus (Yang & Reddy, 1997); uniqueness types (Barendsen & Smetsers, 1996), usage types for optimised compilation of lazy functional programs (Peyton Jones & Wansbrough, 2000) and program analyses for destructive array updates (Draghicescu & Purushothaman, 1993; Wand & Clinger, 1998) as automated in PVS (Shankar, 1999). There is also related work in the area of compiler construction and typed assembly languages, where researchers have investigated static analysis techniques for determining when optimisations such as in-place update or compile-time garbage collection are admissible; recent examples include shape analysis (Wilhelm *et al.*, 2000), already mentioned alias types (Smith *et al.*, 2000), and static capabilities (Crary *et al.*, 1999), which are an alternative and more permissive form of region-based memory management.

One of our future goals is to relate our work back to research on compiler optimisations and typed low-level languages in the hope that we can *guarantee* that certain optimisations will always be possible in LFPL by virtue of its type system. This is in contrast to the behaviour of many present optimising compilers where it is often difficult for the programmer to be sure if a certain desirable optimisation will be performed by the compiler or not. Work in this directions has begun in (Aspinall & Compagnoni, 2003) where a typed assembly language is developed that has high-level types designed to support compilation from LFPL to obviate the need for garbage collection.

6.2 Outlook

We see the work reported here as a step along the way towards a powerful high-level language equipped with notions of resource control. There are more steps to take. We want to consider richer type systems closer to those used in present functional programming languages, in particular, including polymorphic and higher-order types. We do not expect significant problems integrating Milner-style polymorphism into our language and higher-order functions should be possible using the

standard types and effects technique (Gifford & Lucassen, 1986; Talpin & Jouvelot, 1994).

Another step is to consider inference mechanisms for adding resource annotations, including the \diamond arguments (we mentioned some progress on such inference in Section 1) and usage aspects, as well as the possibility of automatically choosing between \otimes -types and \times -types. Inference of our usage aspects has been addressed in (Konečný, 2003b) and a generalisation of the inference of the product types is included in (Konečný, 2003a) where usage aspects are assigned to certain type sub-terms instead of variables.

In (Konečný, 2003b), the third author showed that for every program typable in our system, there is a typing with best (i.e. largest) usage aspects. These usage aspects can be automatically reconstructed using an iterative search for a fixed point, starting with the most optimistic typing of all functions in the program, i.e. assuming that every argument in every function is used with aspect 3.

We have supported our theoretical work with the development of experimental prototype compilers for LFPL. Type-checking and usage aspect inference for the present system have been implemented by the third author as a front-end to a compiler by Robert Atkey. More information about the compiler can be found via (Aspinall & Konečný, 2003).

Acknowledgements. This work was supported by the EPSRC grant GR/N28436/01. The authors are grateful to Robert Atkey for discussion and comments and indebted to the referees who provided useful suggestions for improvement. The authors owe a number of good ideas to them. Especially they would like to thank the referee who suggested the use of safe operational semantics to simplify the presentation and proof of soundness. An earlier version of this work was reported as (Aspinall & Hofmann, 2002). The details of the soundness proof reported there turned out to be considerably more intricate than first anticipated; they were provided by the third author in this paper.

References

- (2007). *Wikipedia article on the administrative normal form*. http://en.wikipedia.org/wiki/Administrative_Normal_Form. Accessed on 27 March 2007.
- Aiken, Alex, Foster, Jeffrey S., Kodumal, John, & Terauchi, Tachio. (2003). Checking and inferring local non-aliasing. *Pages 129–140 of: PLDI '03: Proceedings of the ACM SIGPLAN 2003 conference on programming language design and implementation*. New York, NY, USA: ACM Press.
- Aspinall, David, & Compagnoni, Adriana. (2003). Heap-bounded assembly language. *Journal of automated reasoning*, **31**(3–4), 261–302.
- Aspinall, David, & Hofmann, Martin. (2002). Another type system for in-place update. *Pages 36–52 of: Métayer, D. Le (ed), Programming languages and systems, proceedings of 11th european symposium on programming*. Springer-Verlag. Lecture Notes in Computer Science 2305.
- Aspinall, David, & Konečný, Michal. 2003 (February). *Type Systems for Resource Bounded Programming and Compilation project homepage*. <http://homepages.inf.ed.ac.uk/da/resbnd>.

- Barendsen, E., & Smetsers, S. (1996). Uniqueness typing for functional languages with graph rewriting semantics. *Mathematical structures in computer science*, **6**, 579–612.
- Boyland, John. (2003). Checking interference with fractional permissions. *Pages 55–72 of: Cousot, R. (ed), Static analysis: 10th international symposium*. Lecture Notes in Computer Science, vol. 2694. Berlin, Heidelberg, New York: Springer.
- Crary, Karl, Walker, David, & Morrisett, Greg. (1999). Typed memory management in a calculus of capabilities. *Pages 262–275 of: Proceedings ACM principles of programming languages*.
- Dor, Nurit, Rodeh, Michael, & Sagiv, Mooly. (2000). Checking cleanness in linked lists. *Pages 115–134 of: Proceedings of the seventh international static analysis symposium*. Springer Berlin / Heidelberg. Lecture Notes in Computer Science 1824.
- Draghicescu, M., & Purushothaman, S. (1993). A uniform treatment of order of evaluation and aggregate update. *Theoretical computer science*, **118**(2), 231–262.
- Evans, David. (1996). Static detection of dynamic memory errors. *Pages 44–53 of: PLDI '96: Proceedings of the ACM SIGPLAN 1996 conference on Programming language design and implementation*. New York, NY, USA: ACM Press.
- Fahndrich, Manuel, & DeLine, Robert. (2002). Adoption and focus: practical linear types for imperative programming. *Pages 13–24 of: Pldi '02: Proceedings of the acm sigplan 2002 conference on programming language design and implementation*. New York, NY, USA: ACM Press.
- Gifford, David K., & Lucassen, John M. (1986). Integrating functional and imperative programming. *Pages 28–38 of: LFP '86: Proceedings of the 1986 ACM conference on LISP and functional programming*. New York, NY, USA: ACM Press.
- Hofmann, Martin. (2000). A type system for bounded space and functional in-place update. *Nordic journal of computing*, **7**(4), 258–289. An extended abstract has appeared in *Programming Languages and Systems*, G. Smolka, ed., Springer LNCS, 2000.
- Hofmann, Martin, & Jost, Steffen. (2003). Static prediction of heap space usage for first-order functional programs. *Pages 185–197 of: 30th ACM SIGPLAN-SIGACT symposium on principles of programming languages (POPL '03)*.
- Ishtiaq, Samin, & O'Hearn, Peter W. (2001). BI as an assertion language for mutable data structures. *Pages 14–26 of: The 28th ACM SIGPLAN-SIGACT symposium on Principles of programming languages (POPL'01)*.
- Kobayashi, Naoki. (1998). *Quasi-linear types*. Tech. rept. 98–02. Department of Information Science, University of Tokyo.
- Kobayashi, Naoki. (1999). Quasi-linear types. *Pages 29–42 of: 26th ACM SIGPLAN-SIGACT symposium on Principles of programming languages (POPL'99)*.
- Konečný, Michal. (2003a). Functional in-place update with layered datatype sharing. *Pages 195–210 of: TLCA 2003, Valencia, Spain, proceedings*. Springer-Verlag. Lecture Notes in Computer Science 2701.
- Konečný, Michal. (2003b). Typing with conditions and guarantees for functional in-place update. *Pages 182–199 of: TYPES 2002 workshop, Nijmegen, proceedings*. Springer-Verlag. Lecture Notes in Computer Science 2646.
- MacKenzie, Kenneth, & Wolverson, Nicholas. (2004). Camelot and grail: resource-aware functional programming on the JVM. *Pages 29–46 of: Trends in functional programming*, vol. 4. Intellect.
- Odersky, Martin. (1992). Observers for linear types. *Pages 390–407 of: Krieg-Brückner, B. (ed), 4th European symposium on programming (ESOP'92)*. Rennes, France: Springer-Verlag. Lecture Notes in Computer Science 582.
- O'Hearn, P. W., Takeyama, M., Power, A. J., & Tennent, R. D. (1995). Syntactic control

- of interference revisited. *MFPS XI, Conference on mathematical foundations of program semantics*. Electronic Notes in Theoretical Computer Science, vol. 1. Elsevier.
- Peyton Jones, Simon, & Wansbrough, Keith. 2000 (September). Simple usage polymorphism. *Proc. 3rd ACM SIGPLAN Workshop on types in compilation (TIC 2000)*. Technical Report CMU-CS-00-161.
- Reynolds, J. C. (1978). Syntactic control of interference. *Pages 39–46 of: Proceedings of the Fifth ACM symposium on Principles of programming languages (POPL)*. Tucson, Arizona: ACM Press.
- Reynolds, John C. (2002). Separation logic: A logic for shared mutable data structures. *Pages 55–74 of: Proceedings of 17th annual IEEE symposium on Logic in computer science (LICS'02)*.
- Sabry, Amr, & Felleisen, Matthias. (1993). Reasoning about programs in continuation-passing style. *LISP and symbolic computation*, **6**(3–4), 289–360.
- Shankar, Natarajan. 1999 (November). *Efficiently executing PVS*. Project report, Computer Science Laboratory, SRI International, Menlo Park, CA.
- Smith, Frederick, Walker, David, & Morrisett, Greg. (2000). Alias types. *Pages 366–381 of: Smolka, G. (ed), 9th European symposium on programming (ESOP'00)*. Springer-Verlag. Lecture Notes in Computer Science 1782.
- Talpin, Jean-Pierre, & Jouvelot, Pierre. (1994). The type and effect discipline. *Inf. comput.*, **111**(2), 245–296.
- Tofte, M., & Talpin, J.-P. (1997). Region-based memory management. *Information and computation*, **132**(2), 109–176.
- Wadler, Philip. (1990). Linear types can change the world. *Pages 561–581 of: Broy, M., & Jones, C. B. (eds), IFIP TC 2 Working conference on Programming concepts and methods*. Sea of Galilee, Israel: North-Holland.
- Walker, David, & Morrisett, J. Gregory. (2001). Alias types for recursive data structures. *Pages 177–206 of: TIC '00: Selected papers from the third international workshop on types in compilation*. London, UK: Springer-Verlag.
- Wand, Mitchell, & Clinger, William D. (1998). Set constraints for destructive array update optimization. *Pages 184–193 of: Proc. IEEE international conference on computer languages (ICCL'98)*.
- Wilhelm, Reinhard, Sagiv, Mooly, & Reps, Thomas. (2000). Shape analysis. *Pages 1–17 of: Proc. 9th International conference on Compiler construction (CC 2000)*. Springer-Verlag. Lecture Notes in Computer Science 1781.
- Yang, H., & Reddy, U. (1997). *Imperative lambda calculus revisited*. Electronic manuscript.