

Reference Counting with Frame Limited Reuse

ANTON LORENZEN, University of Bonn, Germany

DAAN LEIJEN, Microsoft Research, USA

The recently introduced *Perceus* algorithm can automatically insert reference count instructions such that the resulting (cycle-free) program is *garbage free*: objects are freed at the very moment they can no longer be referenced. An important extension is reuse analysis. This optimization pairs objects of known size with fresh allocations of the same size and tries to reuse the object in-place at runtime if it happens to be unique. Unfortunately, current implementations of reuse analysis are fragile with respect to small program transformations, or can cause an arbitrary increase in the peak heap usage. We present a novel *drop-guided* reuse algorithm that is simpler and more robust than previous approaches. Moreover, we generalize the linear resource calculus to precisely characterize garbage-free and frame-limited evaluations. On each function call, a frame-limited evaluation may hold on to memory longer if the size is bounded by a constant factor. Using this framework we show that our drop-guided reuse is frame-limited and find that an implementation of our new reuse approach in Koka can provide significant speedups.

CCS Concepts: • **Software and its engineering** → **Control structures; Polymorphism; • Theory of computation** → **Type theory.**

Additional Key Words and Phrases: Reference Counting, Reuse, Frame Limited, Koka

ACM Reference Format:

Anton Lorenzen and Daan Leijen. 2022. Reference Counting with Frame Limited Reuse. *Proc. ACM Program. Lang.* 6, ICFP, Article 103 (August 2022), 24 pages. <https://doi.org/10.1145/3547634>

1 INTRODUCTION

Reference counting [Collins 1960] is a technique for automatic memory management where each allocated object stores the number of references that point to it. Reinking, Xie, de Moura, and Leijen [2021] describe the *Perceus* algorithm for automatically inserting reference count instructions such that the resulting (cycle-free) program is *garbage free*: objects are freed at the very moment they can no longer be referenced. Even though the *Perceus* algorithm itself needs an internal calculus with explicit control flow, the authors apply this work in the context of the Koka language [2021] which supports full algebraic effect handlers. These can be used to define various kinds of implicit control flow, including features like exceptions, *async/await*, and backtracking [Leijen 2017; Plotkin and Power 2003; Plotkin and Pretnar 2009; Xie and Leijen 2021].

An important extension is *reuse analysis*. This optimization pairs objects of known size with fresh allocations of the same size and tries to reuse the object in-place at runtime if it happens to be unique. This was first described in earlier work by Ullrich and de Moura [2019] in the context of the Lean theorem prover [Moura and Ullrich 2021]. Unfortunately, the published algorithms for reuse analysis all have various weaknesses; for example, they are fragile with respect to small program transformations, where inlining or rearranging expressions can cause reuse analysis to fail unexpectedly.

Authors' addresses: Anton Lorenzen, University of Bonn, Bonn, Germany, anton.lorenzen@uni-bonn.de; Daan Leijen, Microsoft Research, Redmond, WA, USA, daan@microsoft.com.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2022 Copyright held by the owner/author(s).

2475-1421/2022/8-ART103

<https://doi.org/10.1145/3547634>

Moreover, while Perceus itself is garbage-free, reuse analysis does not have this property. By construction, it holds on to memory that is to be reused later, which can lead to an increased peak memory usage. We find the maximum increase is not just a constant factor but can be much larger, and is generally not *safe for space* [Appel 1991; Paraskevopoulou and Appel 2019].

In this work we improve upon this with a novel reuse algorithm and a formal framework for reasoning about heap bounds that can be applied to general transformations (including reuse analysis). In particular:

- We define a new approach to reuse called *drop guided reuse* (Section 3.2). In contrast to earlier techniques we perform the reuse *after* Perceus has inserted reference count instructions. This both simplifies the analysis and makes it more robust with respect to small program transformations. We formalize drop-guided reuse generally in the form of declarative derivation rules where we can discuss clearly the various choices an algorithm can make (Section 5.4).
- We illustrate this with a practical example in Section 2.2 where we see significantly better reuse for red-black tree balanced insertion in comparison with the previous algorithms. In combination with the tail-recursion-modulo-cons (TRMC) optimization, we show how for red-black tree insertion we get a highly optimized code path reusing a (unique) tree in-place while using minimal stack space. Our straightforward and purely functional implementation is about 19% faster than the manually optimized in-place mutating red-black tree implementation in the C++ STL library (`std::map`) (Section 6).
- We reformulate the original linear resource calculus λ^1 [Reinking, Xie et al. 2021] in a normalized form (λ^{1n}) where we can reason precisely about reuse. As for λ^1 , the declarative derivation rules for the λ^{1n} calculus are non-deterministic and can derive many programs with reference counting instructions that are all correct but may differ in their memory consumption. However, in the normalized reformulation, it now is possible to add a single logical side condition (\star) to the LET rule that captures various important variants concisely (Section 5.3):
 - If the (\star) condition is unrestricted, the resulting programs are *sound* – that is, the reference counting is correct and the final heap contains no garbage.
 - By restricting (\star) in a particular way, we can show that all derived programs are *garbage-free* where at every allocating evaluation step there is no garbage (and all programs resulting from the Perceus algorithm fall in this class).
 - Finally, we can weaken the (\star) condition of the garbage-free system to also allow derivations that we call *frame-limited*, where every function call uses at most a constant factor c more memory.
- Transformations like reuse and borrowing no longer have the garbage-free property as they hold on to some memory for a bit longer (the cell to reuse, or the data that is borrowed). With the new formalization, we can now show that some of these transformations are still *frame-limited*, and we prove that our new *drop-guided* reuse analysis is frame-limited (Section 5.4). In contrast, we show that some previous reuse algorithms and unrestricted borrow inference [Ullrich and de Moura 2019] are *not* frame-limited transformations (and we argue such transformations should therefore be avoided in practice).
- Building on robust reuse analysis, we can often express imperative style algorithms in a functional way. Such algorithms are called *Functional But In-Place* (FBIP). In Section 7 we use this technique to create a faster version of the parallel `binarytrees` benchmark. We then generalize this approach to visitor data types (as a *derivative* of the original data type) to derive a novel FBIP algorithm for red-black tree balanced insertion that further improves upon the standard Okasaki style implementation [Okasaki 1999a].

The work described in this paper applies in particular to systems that use Perceus style reference counting: both drop-guided reuse and the (\star) condition to reason about space usage rely on the linear resource calculus. As noted by Reinking, Xie et al. [2021], Perceus cannot collect cycles, and works best in a language with limited use of (concurrent) mutable references and where potential sharing of structures across threads is detectable. Currently, Perceus is implemented in the Lean and Koka compilers where these assumptions hold, but we believe there are many more (ML-style applicative) languages that could take advantage of Perceus and simplify their runtime systems.

Proofs and extended sources can be found in the associated technical report [Lorenzen and Leijen 2021].

2 OVERVIEW AND BACKGROUND

We start with a short overview of background material and related work; in particular the Koka language, Perceus-style reference counting, and reuse analysis. All our examples use the Koka language [Leijen 2014 2017 2021; Xie and Leijen 2021] – a strongly typed functional language with effect handlers which tracks (side) effects in the type of every function. Koka uses algebraic data types extensively. For example, we can define a polymorphic list of elements of type a as:

```
type list(a)
  Cons( head : a, tail : list(a) )
  Nil
```

We can match on a list to define a polymorphic `map` function that applies a function f to each element of a list xs :

```
fun map( xs : list(a), f : a -> e b ) : e list(b)
  match xs
  Cons(x,xx) -> Cons(f(x), map(xx,f))
  Nil       -> Nil
```

Here we transform the list of generic elements of type a to a list of generic elements of type b . Since `map` itself has no intrinsic effect, the overall effect of `map` is polymorphic, and equals the effect e of the function f as it is applied to every element.

2.1 Perceus

By starting from a language with strong static guarantees (like Koka), the Perceus algorithm [Reinking, Xie et al. 2021] can insert optimized reference count instructions during compilation. Note though it still needs separate mechanisms to address cyclic data and mitigate the impact of thread shared reference counts – we refer to the Perceus paper for a in-depth discussion of these.

The main attribute that sets Perceus apart from most automatic memory management systems is that it is *garbage-free*: for a cycle-free program, an object is freed as soon as no more references remain. Consider for example the following function:

```
fun main()
  val xs = list(1,1000000) // allocate a large list
  val ys = map(xs,inc)    // increment each element
  println(ys)
```

Many reference count systems would drop the references to xs and ys based on the lexical scope; for example:

```
fun main()
  val xs = list(1,1000000)
  val ys = map(xs,inc)
  println(ys)
  drop(xs)
  drop(ys)
```

where we use a gray background for generated operations. The `drop(xs)` operation decrements the reference count of an object and, if it drops to zero, recursively drops all children of the object and frees its memory. These “scoped lifetime” reference counts are for example used by a C++ `shared_ptr(T)` (calling the destructor at the end of the scope), Rust’s `Rc(T)` (using the `Drop` trait), and Nim (using a `finally` block to call `destroy`) [Yarantsev 2020]. It is not required by the semantics, but Swift typically emits code like this as well [Gallagher 2016].

Implementing reference counting this way is straightforward and integrates well with exception handling where the drop operations are performed as part of stack unwinding. But from a performance perspective, the technique is not always optimal: in the previous example, the large list `xs` is retained in memory while a new list `ys` is built. Moreover, at the end of the scope, a long cascading chain of drop operations happens for each element in both lists.

Ownership. Perceus takes a more aggressive approach where *ownership* of references is passed down into each function: now `map` is in charge of freeing `xs`, and `ys` is freed by `print`: no `drop` operations are emitted inside `main` as all local variables are *consumed* by other functions, while the `map` and `print` functions drop the list elements as they go. Let’s take a look at what reference count instructions Perceus generates for the `map` function:

```
fun map(xs : list<a>, f : a -> e b) : e list<b>
  match xs
  Cons(x, xx) ->
    dup(x); dup(xx); drop(xs)
    Cons( dup(f)(x), map(xx, f))
  Nil ->
    drop(xs); drop(f)
  Nil
```

In the `Cons` branch, first the head and tail of the list are *dupped*, where a `dup(x)` operation increments the reference count of an object and returns itself. The `drop(xs)` then frees the initial list node. We need to `dup f` as well as it is used twice, while `x` and `xx` are consumed by `f` and `map` respectively.

Transferring ownership, rather than retaining it, means we can free an object immediately when no more references remain. This both increases cache locality and decreases memory usage. For `map`, the memory usage is halved: the list `xs` is deallocated while the new list `ys` is being constructed.

2.2 Reuse

Reuse analysis [Reinking, Xie et al. 2021; Ullrich and de Moura 2019] is an optimization that takes advantage of precise reference counts to try to reuse objects in-place. We can pair objects of known size with same sized allocated constructors and try to reuse these in-place at runtime. Reuse analysis rewrites `map` into:

```
fun map(xs : list<a>, f : a -> e b) : e list<b>
  match xs
  Cons(x, xx) ->
    dup(x); dup(xx); val r = dropru(xs)
    Cons@r( dup(f)(x), map(xx, f))
```

The *reuse token* `r` becomes the address of the `Cons` cell `xs` if `xs` happens to be unique, and `NULL` otherwise. The `Cons@r` allocation reuses `xs` in-place if `r` is non `NULL`, and otherwise allocates a fresh `Cons` cell. In case we `map` over a unique list, the list elements are updated in-place. Given the cost of allocation versus a single uniqueness check, reuse optimization (almost) certainly improves performance when it can apply, and often with a significant factor [Reinking, Xie et al. 2021].

A further rewriting technique called *drop specialization* can further optimize this by inlining the `dropru` operation and simplifying such that no reference count operations are necessary in the case that the list is unique:

```

type color { R; B }

type tree(a)
  Node( clr:color, l:tree(a), key:int, value:a, r:tree(a) )
  Leaf

fun is-red(t : tree(a)) : bool
  match t
  Node(R,_,_,_,_) -> True
  _ -> False

fun lbal(l :tree(a), k : int, v : a, r : tree(a)) : tree(a)
  match l
  Node(_, Node(R, lx, kx, vx, rx), ky, vy, ry) ->
    Node(R, Node(B,lx,kx,vx,rx), ky, vy, Node(B,ry,k,v,r))
  Node(_, ly, ky, vy, Node(R, lx, kx, vx, rx)) ->
    Node(R, Node(B,ly,ky,vy,lx), kx, vx, Node(B,rx,k,v,r))
  Node(_, lx, kx, vx, rx) ->
    Node(B, Node(R,lx,kx,vx,rx), k, v, r)
  Leaf -> Leaf

fun rbal(l : tree(a), k : int, v : a, r : tree(a)) : tree(a)
  ...

fun ins(t :tree(a), k : int, v : a) : tree(a)
  match t
  Leaf -> Node(R, Leaf, k, v, Leaf)
  Node(B, l, kx, vx, r) ->
    if k < kx then
      if is-red(l) then lbal(ins(l,k,v), kx, vx, r)
      else Node(B, ins(l,k,v), kx, vx, r)
    elif k > kx then
      if is-red(r) then rbal(l, kx, vx, ins(r,k,v))
      else Node(B, l, kx, vx, ins(r,k,v))
    else Node(B, l, k, v, r)
  Node(R, l, kx, vx, r) ->
    if k < kx then Node(R, ins(l,k,v), kx, vx, r)
    elif k > kx then Node(R, l, kx, vx, ins(r,k,v))
    else Node(R, l, k, v, r)

```

Fig. 1. Balanced red-black tree insertion

```

fun map(xs : list<a>, f : a -> e b) : e list<b>
  match xs
  Cons(x,xx) ->
    val r = if unique(xs) then &xs // reuse the address of xs directly
            else dup(x); dup(xx); decref(xs); NULL
    Cons@r( dup(f)(x), map(xx,f) )
  ...

```

This is an important optimization in practice and our new reuse algorithm is compatible with it, but for clarity we generally leave it out in the following examples.

Balanced Trees. Reinking, Xie et al. [2021] provide an appealing example of the effectiveness of reuse analysis using balanced insertion in red-black trees [Guibas and Sedgewick 1978]. Figure 1 shows the implementation in Koka based on Okasaki’s algorithm [Okasaki 1999a]. A red-black tree has the invariant that the number of black nodes from the root to any of the leaves are the same, and that a red node is never a parent of red node. Together this ensures that the trees are always balanced. When inserting nodes, the invariants are maintained by rebalancing the nodes when needed.

If we look closely at all the `match` branches in Figure 1, (and assume that the `lbal` and `rbal` functions get inlined), then we can see that we always either match one `Node` and allocate one, or we match three nodes deep and allocate three (when rebalancing). In the case that the tree is unique, the reuse analysis reuses every `Node` along the spine without doing any further allocations! Moreover, if we use the tree *persistently* [Okasaki 1999b] where the tree is shared (or has shared subtrees), it adapts to copying exactly the shared spine of the tree (and no more).

3 DROP-GUIDED REUSE

Previous work has shown that reuse analysis can be very effective, and has been implemented in both the Koka and Lean languages. Unfortunately, it turns out that both previously published algorithms for reuse analysis are flawed: the Koka algorithm, which we call algorithm K [Reinking, Xie et al. 2021], is fragile with respect to small program transformations, where rearranging expressions can cause reuse analysis to fail unexpectedly. The Lean algorithm, called algorithm D [Ullrich and de Moura 2019 (Fig 3)], is more robust but can lead to an arbitrary increase in peak memory usage. In this section we look at each algorithm, and propose a new approach, called *drop-guided reuse*, that improves upon both the previous techniques.

3.1 Problems with Reuse

Both Reinking, Xie et al. [2021] and Ullrich and de Moura [2019] describe reuse algorithms as a pass *before* the main Perceus algorithm. The reason they chose this approach is two-fold: the `drop` function can be seen as consuming its argument and thus needs no special treatment from Perceus, and similarly, a reuse token can be deallocated by Perceus if it is not used (for example if the constructor is only allocated in one branch of a nested `match`-statement but not another). However, in practice we observed that this can lead to situations where the reuse is not optimal.

Algorithm K as implemented in the Koka compiler tries to reuse at the start of a branch whenever a matching constructor size can be found in the branch body. That seems reasonable at first, but consider the following example (A):

```
match x
  Just(_) -> // x is still live here
    match y
      0 -> x
      _ -> Just(y)
```

In this case, the `Just(y)` matches with the deconstructed `x`, where reuse analysis will try to reuse `x`, and the generated code becomes:

```
match x
  Just(_) ->
    val r = drop(ru(dup(x)))
    match y
      0 -> drop(y); drop(r); x
      _ -> drop(x); Just@r(y)
```

Unfortunately, since `x` is still live in the scope it prevents any reuse at runtime due to the inserted `dup` operation, and `r` will always be `NULL` at runtime.

An obvious way to improve on this, is by pushing down reuse operations into branches behind the last use of an object. This is the approach used by algorithm D, where example (A) can reuse `x` now effectively:

```
match x
  Just(_) ->
    match y
      0 -> drop(y); x
      _ -> val r = drop(ru(dup(x))); Just@r(y)
```

This approach has a different weakness though, and can lead to an arbitrary increase in peak memory usage. Consider the following example (B):

```
match xs
  Cons(_,_) ->
    val y = f(xs)
    Cons(y,Nil)
```

Using algorithm D, a reuse is inserted right after the call to `f`, which leads to Perceus (running afterwards) inserting a dup on the `xs` parameter:

```
match xs
  Cons(_,_) ->
    val y = f(dup(xs))
    val r = dropru(xs)
    Cons@r(y,Nil)
```

Even though the `Cons` cell of `xs` is actually available for reuse, it now also holds on to the full `xs` list during evaluation of `f`. This not only means that we may use `sizeof(xs)` more memory than necessary, but also that any reuse by `f` of `xs` is prevented as well (as `xs` is now certainly not unique).

Reuse in Balanced Trees. The previous patterns actually occur regularly in practice. In fact, both algorithm D and K also fail to effectively reuse, even under small rewrites, the balanced tree insertion example. Consider again the code in Figure 1, and let's focus on the second branch in the `ins` function:

```
match t
  Node(B, l, kx, vx, r) ->
    if k < kx then
      if is-red(l) then lbal(ins(l,k,v), kx, vx, r)
      else Node(B, ins(l,k,v), kx, vx, r)
```

With both reuse algorithms, the `Node` allocation in the `else` branch will reuse the outer matched node `t`. It gets more interesting though, if both `lbal` and `is-red` get inlined (and simplified); we focus just on the `then` case and the first branch of `lbal` here:

```
match t
  Node(B, l, kx, vx, r) ->
    if k < kx then
      match l
        Node(R,_,_,_,_) ->
          val t2 = ins(l,k,v)
          match t2
            Node(_, Node(R, ...), ...) ->
              Node(R, Node(B, ...), ..., Node(B, ...))
```

Due to the inlined `is-red` function, we now get a situation where both algorithm D and K fall short. This is exactly like example (B): algorithm K will try to reuse `l` even though it is still used later on, leading to:

```
match l
  Node(R,_,_,_,_) ->
    val ru2 = dropru(dup(l))
    val u = ins(l,k,v)
```

where `ru2` is never available for reuse at runtime. Algorithm D does not fare much better either as it holds on to `l` preventing the recursive `ins` from reusing the tree:

```
match l
  Node(R,_,_,_,_) ->
    val t2 = ins(dup(l),k,v)
    val ru2 = dropru(l)
```

As shown, both published algorithms are quite fragile with respect to small program transformations.

3.2 Drop Guided Reuse

We believe the essential weakness of both previous algorithms is that they fail to take liveness into account. As such, we propose to perform reuse analysis *after* doing Perceus dup-drop insertion.

Indeed, Perceus already performs precise liveness analysis and it has been shown it inserts optimal dup/drop operations in the sense that the resulting program is garbage-free. This means that the drop operations signify precisely when a cell may become available for reuse – and this is exactly the point where we should rewrite the drop to a `dropru` if there is any potential for reuse. The call to `dropru` cannot hold on to objects any longer than necessary since by the garbage-free property an object must be live until directly before it is dropped.

With *drop-guided* reuse analysis, we keep track of the currently known sizes of each variable (updated at each branch pattern) and if we encounter a drop we can statically determine if it can pair with a later allocation of the same size. As type information is erased in Koka's runtime, it does not matter if the values are of different types. Looking at our earlier example (A) from the previous section, Perceus generates first:

```
match x
  Just(_) ->
    match y
      0 -> drop(y); x
      _ -> drop(x); Just(y)
```

The *drop-guided* reuse analysis can now rewrite the `drop(x)` into a `dropru(x)` as we know that the size of `x` matches the size of the following `Just` allocation:

```
match x
  Just(_) ->
    match y
      0 -> drop(y); x
      _ -> val r = dropru(x); Just@r(y)
```

For the other example (B) in the previous section, no drop operations are generated in the first place, and no `dropru` is inserted either – perfect!

A drawback of the new approach is that we need to explicitly free newly created reuse tokens if these happen to be unused in some branch, and we can no longer rely on Perceus doing this for us. It may seem we need to perform another mini Perceus pass to address this, but it can be done in a simpler way: since reuse tokens are only generated in a specific way, we can show they are only used either never or once, and never captured under a lambda or passed as an argument. As such, it suffices to locally check at each branch of a match expression if a given reuse token `r` occurs in this branch and insert a `free(r)` instruction if this is not the case. In the new Koka implementation we combine this in one pass with the reuse analysis.

Finally, with balanced tree insertion, none of the earlier problems with algorithm D and K occur. The drop-guided approach is robust and leads to optimal reuse:

```
match t
  Node(B, l, kb, vb, r) ->
    dup(l); dup(kb); dup(vb); dup(r)
    val ru = dropru(t)
    if k < kb then
      match l
        Node(R, _, _, _, _) ->
          val t2 = ins(l, k, v)
          match t2
            Node(_, l2 as Node(R, ...), ..., r2) ->
              val ru2 = dropru(t2)
              val ru3 = dropru(l2)
              Node@ru(R, Node@ru2(B, ...), ..., Node@ru3(B, ...))
```


Again, the reuse analysis will reuse every `Node` along the spine without doing any further allocations if the tree is unique. As we show in the benchmarking section, with the new improved reuse analysis the performance on unique trees rivals that of the manually optimized in-place mutating red-black tree implementation in the C++ STL library (`std::map`).

TRMC: Tail-Recursion-Modulo-Cons. We can make the red-black tree insertion a bit faster still with *tail-recursion-modulo-cons* (TRMC) optimization. Usually, with *tail-recursion* any functions that call themselves recursively in a tail position are transformed into a loop instead (using no extra stack space). With TRMC, such function can make its tail call inside any tail-position *expression* consisting of just constructors and non-allocating total expressions. For example, `map` is such function where the recursion in `Cons(f(x), map(xx, f))` is transformed into a loop. This is done by pre-allocating the `Cons` node ahead of the recursive call with a *hole* in the tail field, which is later assigned by the recursive call.

TRMC interacts well with reuse analysis as often the recursive call is inside a constructor that is reused. In the `map` function for example, this will result in traversing the list in a tight loop while updating each element in-place when the list `xs` happens to be unique.

For red-black tree rebalancing, we can see in Figure 1 that there are four TRMC recursive `ins` calls and only in the `rbal/lbal` cases this does not hold. When we study the generated C code the final result is quite sophisticated: there is an outer TRMC loop for the four TRMC `ins` calls, but it is interspersed at runtime with the two actual recursive `ins` calls. Moreover, due to reuse, the code updates unique nodes in place when rebalancing the spine, but adapts to copying for shared subtrees. Overall this is very efficient code that would be difficult to write directly by hand.

4 REASONING ABOUT SPACE

As we have seen, liveness analysis helps us avoid the problems of algorithm K, but have we also avoided the problem of algorithm D of using too much space? Indeed, since reuse analysis keeps heap cells alive until they can be reused, it means that no reuse analysis can preserve the garbage-free property! Is there a way to still characterize the space usage of such transformations that is more restrictive than allowing arbitrary increases, but also more permissive than garbage-free?

4.1 Frame-Limited Transformations

Drop guided reuse analysis can, at any evaluation step, hold on to a single cell per reuse token `r` that was created by `dropru`, but not used at a constructor yet. Since any function can only contain a constant number of `dropru` calls, one might expect the total overhead to be constant as well, which would make drop guided reuse *safe for space* [Appel 1991; Paraskevopoulou and Appel 2019], in the sense that the maximum peak memory increase is bounded by a constant. However, before a reuse token is used there might be a recursive call. Consider the `map` function we first viewed:

```
val r = dropru(xs)
Cons@r( dup(f)(x), map(xx, f))
```

Here, `r` is live during the recursive call and so reuse analysis can hold on to as many `Cons` cells as either the list is long or the stack allows. In other words, we can only hope to bound the extra memory needed for reuse analysis by a constant factor times the current number of stack frames – we call this *frame-limited*. We formalize this notion in the next section, and formally prove in Section 5.4 that our new drop guided reuse *is* a frame-limited transformation. In contrast, as we showed in Section 3.1, the reuse algorithm D [Ullrich and de Moura 2019] is not frame limited and can lead to an arbitrary increase in memory usage.

Even though weaker than being garbage-free or safe for space, we argue that frame-limited transformations are still good in practice. First of all, programmers are already aware of recursion

and take steps to avoid unbounded recursion. Secondly, in practice the stack size is usually already bounded – in such case, that makes the frame-limited bound constant (and thus safe for space).

Note that in practice, backend optimizations like tail-recursion may optimize stack frames away. However, we formalize frame-limited in terms of the size of the evaluation context (i.e. the recursion depth) instead of actual stack frames and thus the bound still applies.

4.2 Borrowing

Another example of a transformation that does not preserve the garbage-free property is *borrowing* [Ullrich and de Moura 2019]: even though the first example in Section 2.1 argues that arguments should be passed owned to the callee, this is not always optimal – sometimes it is better to pass arguments as *borrowed* instead. Consider converting a list into an unbalanced binary tree:

```
fun make-tree( xs : list(a) ) : tree(a)
  match xs
    Cons(x, xx) -> Node( R, Leaf, 0, x, make-tree(xx) )
    Nil -> Leaf
```

Perceus inserts `dup(x)`; `dup(xx)`; `drop(xs)`; in the `Cons` branch, which causes some overhead, especially since there are no reuse opportunities. When we annotate a parameter like `xs` to be borrowed (as `^xs`), the caller keeps ownership of the parameter. As a result, no reference count operations need to be performed at all for the `xs` parameter in our example. Note that borrow annotations are strictly a performance hint, and do not change semantics or whether a program is well-typed (in contrast to the notion of borrowing in a language like Rust for example).

Borrow annotations are not always beneficial though: if `xs` happens to be unique, we allocate the tree while the full `xs` list is still live – exactly the situation we wanted to avoid in Section 2.1. In general, borrowing can increase the memory usage of a program by an arbitrary amount and it is generally *not* frame-limited. Ullrich and de Moura [2019] describe a borrow inference algorithm that marks `xs` automatically as borrowed. However, given that this is not safe for space, we argue that automatic borrow inference should be further restricted to guarantee it is at least frame-limited. Therefore, Koka currently has no automatic borrow inference and generally only uses borrowing for built-in primitives (like big integer operations).

5 FORMALIZATION

We formalize our results using the linear resource calculus λ^1 as given by Reinking, Xie et al. [2021] (see figure 2). This is essentially just lambda calculus extended with let bindings and pattern matching. We assume that the patterns p_i in a match are all distinct. The semantics for λ^1 is standard using strict evaluation contexts E [Wright and Felleisen 1994]. The evaluation contexts uniquely determine where to apply an evaluation step using the `eval` rule. As such, evaluation contexts neatly abstract from the usual implementation context of a stack and program counter. The small step evaluation rules perform function application (*app*), let-binding (*let*), and pattern matching (*match*).

5.1 Heap Semantics

To reason precisely about reference counting, we need a semantics with an explicit heap. Reinking, Xie et al. [2021] define a heap semantics directly over λ^1 . However, since we aim to reason precisely about the space behavior, this is not quite sufficient for our case as we need to be more explicit about sharing and evaluation order. We therefore translate any expression e into *normalized form* $[e]$, as defined in Figure 3, where all arguments become variables instead of values [Flanagan et al. 1993].

$$\begin{array}{l}
e ::= v \\
\quad | \quad e e \\
\quad | \quad \text{let } x = e \text{ in } e \\
\quad | \quad \text{match } e \{ \overline{p_i \rightarrow e_i} \} \\
v ::= x \\
\quad | \quad \lambda x. e \\
\quad | \quad C \bar{v} \\
p ::= C \bar{x}
\end{array}$$

Semantics:

$$\begin{array}{l}
E ::= \square \mid E e \mid v E \\
\quad | \quad \text{let } x = E \text{ in } e \\
\quad | \quad \text{match } E \{ \overline{p_i \rightarrow e_i} \} \\
\frac{e \longrightarrow e'}{E[e] \mapsto E[e']} \quad [\text{EVAL}]
\end{array}$$

Small step transitions:

$$\begin{array}{l}
(\text{app}) \quad (\lambda x. e) v \quad \longrightarrow \quad e[x:=v] \\
(\text{let}) \quad \text{let } x = v \text{ in } e \quad \longrightarrow \quad e[x:=v] \\
(\text{match}) \quad \text{match } (C \bar{v}) \{ \overline{p_i \rightarrow e_i} \} \quad \longrightarrow \quad e_i[\bar{x}:=\bar{v}] \\
\quad \text{where } p_i = C \bar{x}
\end{array}$$

Fig. 2. Syntax and semantics of λ^1

Figure 4 defines the syntax of the *normalized linear resource calculus* λ^{ln} where all arguments are now variables. Moreover the syntactic constructs in gray are only generated in derivations of the calculus and are not exposed to users. Among those constructs, dup and drop form the basic instructions of reference counting, while $r \leftarrow \text{dropru}$ is used for reuse. Also, every lambda $\lambda_{\bar{z}}x. e$ is annotated with its free variables \bar{z} which becomes important during evaluation.

Contexts Δ, Γ are *multisets* containing variable names. We use the compact comma notation for summing (or splitting) multisets. For example, (Γ, x) adds x to Γ , and (Γ_1, Γ_2) appends two multisets Γ_1 and Γ_2 . The set of free variables of an expression e is denoted by $\text{fv}(e)$, and the set of bound variables of a pattern p by $\text{bv}(p)$.

Using our normalized calculus, Figure 5 defines the semantics in terms of a reference counted heap, where sharing of values is explicit, and substitution only substitutes variables. Here, each heap entry $x \mapsto^n v$ points to a value v with a reference count of n (with $n \geq 1$). In these semantics, values other than variables are allocated in the heap with rule (lam_r) and rule (con_r) . The evaluation rules discard entries from the heap when the reference count drops to zero. Any allocated lambda is annotated as $\lambda_{\bar{z}}x. e$ to clarify that these are essentially *closures* holding an environment \bar{z} and a code pointer $\lambda x. e$. Note that it is important that the environment \bar{z} is a multi-set. After the initial translation, \bar{z} will be equivalent to the free variables in the body (see rule LAM), but during evaluation substitution may substitute several variables with the same reference. To keep reference counts correct, we need to keep considering each one as a separate entry in the closure environment.

When applying an abstraction, rule (app_r) needs to satisfy the assumptions made when deriving the abstraction in rule LAM (shown in Figure 6). First, the (app_r) rule inserts dup to duplicate the free closure variables \bar{z} , as these are owned in rule LAM . It then drops the reference to the closure itself.

A difference between (app_r) and (match_r) is that for applications the free variables \bar{z} are dynamic and thus the duplication must be done at runtime. In contrast, a match knows the bound variables in a pattern statically and we therefore generate the required dup and drop operations statically during elaboration for each branch (as shown in Figure 6) – this is essential as that enables the further static optimizations of dup/drop pairs and reuse analysis.

We discuss the reuse evaluation rules later in Section 5.4.

$$\begin{aligned}
[x] &= x \\
[\lambda x. e] &= \lambda x. [e] \\
[e e'] &= \text{let } f = [e] \text{ in let } x = [e'] \text{ in } f x \\
[C v_1 \dots v_n] &= \text{let } x_1 = [v_1] \text{ in } \dots \text{let } x_n = [v_n] \text{ in } C x_1 \dots x_n \\
[\text{match } e \{ \overline{p_i \rightarrow e_i} \}] &= \text{let } x = [e] \text{ in match } x \{ \overline{p_i \rightarrow [e_i]} \} \\
[\text{let } x = e_1 \text{ in } e_2] &= \text{let } x = [e_1] \text{ in } [e_2]
\end{aligned}$$
Fig. 3. Normalization. All f and x are fresh
$$\begin{array}{ll}
e ::= v & v ::= x \\
| e x & | \lambda_{\bar{z}} x. e \\
| \text{let } x = e \text{ in } e & | C \bar{x} \\
| \text{match } x \{ \overline{p_i \rightarrow e_i} \} & \\
| \text{dup } x; e & p ::= C \bar{x} \\
| \text{drop } x; e & \\
| r \leftarrow \text{dropru } x; e & \Delta, \Gamma ::= \emptyset \mid \Delta \cup x \\
\end{array}$$

$$\lambda x. e \doteq \lambda_{\bar{z}} x. e \quad (\bar{z} = \text{fv}(e))$$
Fig. 4. The normalized linear resource calculus λ^{ln} .
$$\begin{array}{l}
H : x \mapsto (\mathbb{N}^+, v) \\
E ::= \square \mid E x \mid \text{let } x = E \text{ in } e
\end{array}
\quad
\frac{H \mid e \longrightarrow_r H' \mid e'}{H \mid E[e] \mapsto_r H' \mid E[e']} \text{ [EVAL]}$$

$$\begin{array}{ll}
(\text{lam}_r) & H \mid \lambda_{\bar{z}} x. e \longrightarrow_r H, f \mapsto^1 \lambda_{\bar{z}} x. e \quad | f \quad \text{fresh } f \\
(\text{con}_r) & H \mid C x_1 \dots x_n \longrightarrow_r H, z \mapsto^1 C x_1 \dots x_n \quad | z \quad \text{fresh } z \\
(\text{app}_r) & H \mid f y \longrightarrow_r H \mid \text{dup } \bar{z}; \text{drop } f; e[x:=y] \quad (f \mapsto^n \lambda_{\bar{z}} x. e) \in H \\
(\text{match}_r) & H \mid \text{match } y \{ \overline{p_i \rightarrow e_i} \} \longrightarrow_r H \mid e_i[\bar{x}:=\bar{z}] \quad \text{with } p_i = C \bar{x} \text{ and } (y \mapsto^n C \bar{z}) \in H \\
(\text{let}_r) & H \mid \text{let } x = z \text{ in } e \longrightarrow_r H \mid e[x:=z] \\
(\text{dup}_r) & H, x \mapsto^n v \quad | \text{dup } x; e \longrightarrow_r H, x \mapsto^{n+1} v \quad | e \\
(\text{drop}_r) & H, x \mapsto^{n+1} v \quad | \text{drop } x; e \longrightarrow_r H, x \mapsto^n v \quad | e \quad \text{if } n \geq 1 \\
(\text{dlam}_r) & H, x \mapsto^1 \lambda_{\bar{x}} y. e' \quad | \text{drop } x; e \longrightarrow_r H \mid \text{drop } \bar{x}; e \\
(\text{dcon}_r) & H, x \mapsto^1 C \bar{x} \quad | \text{drop } x; e \longrightarrow_r H \mid \text{drop } \bar{x}; e
\end{array}$$

Extension with reuse (with fresh z, \bar{z}):

$$\begin{array}{ll}
(\text{drop}_{ru}) & H, x \mapsto^{n+1} v \quad | r \leftarrow \text{dropru } x; e \longrightarrow_r H, x \mapsto^n v, z \mapsto^1 () \mid e[r:=z] \quad \text{if } n \geq 1 \\
(\text{dlam}_{ru}) & H, x \mapsto^1 \lambda_{\bar{x}} y. e' \quad | r \leftarrow \text{dropru } x; e \longrightarrow_r H, z \mapsto^1 () \mid \text{drop } \bar{x}; e[r:=z] \\
(\text{dcon}_{ru}) & H, x \mapsto^1 C \bar{x} \quad | r \leftarrow \text{dropru } x; e \longrightarrow_r H, \bar{z} \mapsto^1 (), z \mapsto^1 C \bar{z} \mid \text{drop } \bar{x}; e[r:=z]
\end{array}$$
Fig. 5. Reference-counted heap semantics for λ^{ln} .

5.2 Dup-Drop Insertion in λ^{ln}

Figure 6 defines the logical derivation rules over λ^{ln} for inserting reference count instructions such that the resulting expression can be soundly evaluated by the heap semantics of Figure 5. The judgement $\Delta \mid \Gamma \vdash e \rightsquigarrow e'$ in Figure 6 reads as follows: given a *borrowed environment* Δ , a *linear environment* Γ , an expression e is translated into an expression e' with explicit reference counting

$\frac{\Delta \mid \Gamma \vdash e \rightsquigarrow e'}{\uparrow \quad \uparrow \quad \uparrow \quad \downarrow} \quad (\uparrow \text{ is input, } \downarrow \text{ is output})$	
$\Delta \text{ and } \Gamma \text{ are multisets of the borrowed and owned variables in scope}$	
$\frac{}{\Delta \mid x \vdash x \rightsquigarrow x} \text{ [VAR]}$	$\frac{}{\Delta \mid \bar{x} \vdash C \bar{x} \rightsquigarrow C \bar{x}} \text{ [CON]}$
$\frac{\Delta, x \mid \Gamma \vdash e \rightsquigarrow e'}{\Delta \mid \Gamma, x \vdash e x \rightsquigarrow e' x} \text{ [APP]}$	$\frac{\emptyset \mid \Gamma, x \vdash e \rightsquigarrow e' \quad \Gamma = \text{fv}(\lambda x. e)}{\Delta \mid \Gamma \vdash \lambda x. e \rightsquigarrow \lambda_{\Gamma} x. e'} \text{ [LAM]}$
$\frac{\Delta \mid \Gamma, x \vdash e \rightsquigarrow e' \quad x \in \Delta, \Gamma}{\Delta \mid \Gamma \vdash e \rightsquigarrow \text{dup } x; e'} \text{ [DUP]}$	$\frac{\Delta \mid \Gamma \vdash e \rightsquigarrow e'}{\Delta \mid \Gamma, x \vdash e \rightsquigarrow \text{drop } x; e'} \text{ [DROP]}$
$\frac{\Delta \mid \Gamma, r \vdash e \rightsquigarrow e' \quad \text{fresh } r}{\Delta \mid \Gamma, x \vdash e \rightsquigarrow r \leftarrow \text{dropru } x; e'} \text{ [DROPRU]}$	
$\frac{\Delta, \Gamma_2 \mid \Gamma_1 \vdash e_1 \rightsquigarrow e'_1 \quad \Delta \mid \Gamma_2, x \vdash e_2 \rightsquigarrow e'_2 \quad (\star)}{\Delta \mid \Gamma_1, \Gamma_2 \vdash \text{let } x = e_1 \text{ in } e_2 \rightsquigarrow \text{let } x = e'_1 \text{ in } e'_2} \text{ [LET]}$	
$\frac{\Delta \mid \Gamma, \bar{z}_i \vdash e_i \rightsquigarrow e'_i \quad p_i = C_i \bar{z}_i \quad x \in \Delta, \Gamma}{\Delta \mid \Gamma \vdash \text{match } x \{ \bar{p}_i \mapsto \bar{e}_i \} \rightsquigarrow \text{match } x \{ p_i \mapsto \text{dup}(\bar{z}_i); e'_i \}} \text{ [MATCH]}$	

Fig. 6. Logical derivation rules of λ^{ln} . The rules are parameterized by the (\star) condition on LET. We write \vdash_{GF} for garbage-free derivations, and \vdash_{FL} for derivations that are frame-limited.

instructions. We call variables in the linear environment *owned*.

The key idea is that each resource (i.e., owned variable) is consumed *exactly* once. That is, a resource needs to be explicitly duplicated (in rule DUP) if it is needed more than once; or be explicitly dropped (in rule DROP) if it is not needed. The rules are closely related to linear typing.

The rules are close to the derivation rules by Reinking, Xie et al. [2021] but differ in important details. In particular, by using a normalized form we only split the owned environment Γ in the LET rule, and no longer in the APP and CON rules which are now much simpler. This in turn allows us to parameterize the system by a single side condition (\star) that allows us to concisely capture *garbage-free*, *frame-limited*, and *sound* transformations as shown in Section 5.3.

The LET rule splits the owned environment Γ into two separate contexts Γ_1 and Γ_2 for expression e_1 and e_2 respectively. Each expression then consumes its corresponding owned environment. Since Γ_2 is consumed in the e_2 derivation, we know that resources in Γ_2 are surely alive when deriving e_1 , and thus we can *borrow* Γ_2 in the e_1 derivation. The rule is quite similar to the [LET!] rule of Wadler’s linear type rules [Wadler 1990, pg.14] where a linear type can be “borrowed” as a regular type during evaluation of a binding.

The LAM rule is interesting as it essentially derives the body of the lambda independently. The premise $\Gamma = \text{fv}(\lambda x. e)$ requires that exactly the free variables in the lambda are owned – this corresponds to the notion that a lambda is allocated as a closure at runtime that holds all free variables of the lambda (and thus the lambda expression *consumes* the free variables). The body of a lambda is evaluated only when applied, so it is derived under an empty borrowed environment only owning the argument and the free variables (in the closure). The translated lambda is also annotated with Γ , as $\lambda_{\Gamma} x. e$, so we know precisely the resources the lambda should own when

evaluated in a heap semantics. We often omit the annotation when it is irrelevant.

Another important difference from earlier work is that the `MATCH` rule now statically generates `dup` instructions for the pattern bindings (since the new ($match_r$) rule no longer dups the fields at runtime). Inserting `dup` instructions statically is essential to actually perform further `dup/drop` optimizations (like reuse!) on the final derived expression. Finally, we also added the `DROPRU` rule to reason precisely about reuse analysis described later.

Properties. Many of the properties proven for λ^1 [Reinking, Xie et al. 2021] carry over to λ^{1n} . In particular, the logical derivation rules precisely elaborate expressions with reference count operations such that they can be correctly evaluated by the target heap semantics, as stated in the following theorem:

Theorem 1. (*Reference-counted heap semantics is sound*)

If we have $\emptyset \mid \emptyset \vdash e \rightsquigarrow e'$ and $e \mapsto^* v$, then we also have $\emptyset \mid e' \mapsto^*_r H \mid x$ with $[H]x = v$.

We prove this theorem in separate steps: first we show soundness in a heap semantics that ignores reference count instructions (App. C.2 in the tech report), then use a separate resource calculus to show reference counts are correct (App. C.3 in the tech report), and finally combine these results for the final proof (App. C.4 in the tech report).

Second, we prove that the reference counting semantics never *hold on* to unused variables. We first define the notion of *reachability*.

Definition 1. (*Reachability*)

We say a variable x is reachable in terms of a heap H and an expression e , denoted as $\text{reach}(x, H \mid e)$, if (1) $x \in \text{fv}(e)$; or (2) for some y , we have $\text{reach}(y, H \mid e) \wedge y \mapsto^n v \in H \wedge \text{reach}(x, H \mid v)$.

With reachability, we can show (see App. C.5 in the tech report):

Theorem 2. (*Reference counting leaves no garbage*)

Given $\emptyset \mid \emptyset \vdash e \rightsquigarrow e'$, and $\emptyset \mid e' \mapsto^*_r H \mid x$, then for every intermediate state $H_i \mid e_i$, we have for all $y \in \text{dom}(H_i)$, $\text{reach}(y, H_i \mid e_i)$.

Note that similar to λ^1 , λ^{1n} does not model *mutable references*. A natural extension of the system is to include mutable references and thus cycles. In that case, we could generalize Theorem 2, where the conclusion would be that for all resources in the heap, it is either reachable from the expression, or it is part of a cycle.

These theorems establish the correctness of the reference-counted heap semantics. However, correctness does not imply that evaluation is *garbage-free*. Eventually all live data is discarded but an evaluation may well hold on to live data too long by delaying drop operations. As an example, consider $y \mapsto^1 () \mid (\lambda x. x) (\text{drop } y; ())$, where y is reachable but dropped too late: it is only dropped after the lambda gets allocated. In contrast, a garbage-free algorithm would produce $y \mapsto^1 () \mid \text{drop } y; (\lambda x. x) ()$.

5.3 Reasoning about Space with the “Star” Condition

Reinking, Xie et al. [2021] give declarative derivation rules for reference counting but then provide a separate *algorithm* that is then proven to be garbage-free. This is not ideal as it does not provide any particular insight why the algorithm is garbage-free, and if other approaches may exist as well.

By making evaluation order explicit in the normalized λ^{1n} calculus, we found a way to capture the garbage-free property *declaratively* as a single side condition on the `LET` rule in our new derivation rules. Moreover, by weakening the condition, this can also be used to characterize other interesting points in the design space, and provide a general framework to reason about memory consumption. In particular, we can concisely characterize frame-limited derivations. We can thus instantiate the

rules by giving a specific (\star) condition:

- *General derivations* \vdash . When we define (\star) to be true, the evaluation of any derived expression is *sound* (Theorem 1).
- *Garbage-free derivations* \vdash_{GF} . When we define the (\star) condition as $\Gamma_2 \subseteq \text{fv}(e_2)$, then the evaluation of any derived expression is *garbage-free* (Definition 2 and Theorem 3). At the LET rule we have the freedom to split Γ into Γ_1 and Γ_2 in any way. For garbage-free derivations though we try to minimize borrowing in the e_1 derivation (of Γ_2) and thus the condition captures intuitively that we should use the smallest Γ_2 possible, and not include variables that are not needed for the e_2 derivation.
- *Frame-limited derivations* \vdash_{FL} . When we define (\star) as $\Gamma_2 = \Gamma', \Gamma''$ where $\Gamma' \subseteq \text{fv}(e_2)$ and $\text{sizeof}(\Gamma'') \leq c$ for some constant c , then the evaluation for any derived expression is *frame-limited* (Definition 3 and Theorem 4). We define $\text{sizeof}(\Gamma'')$ as the sum of the sizes of each element: $\sum_{y \in \Gamma''} \text{sizeof}(y)$. This weakens the garbage-free condition to allow borrowing of any y where the runtime size of y is known to be limited by a constant. This is just enough for transformations like reuse where we borrow a reuse token until it can be reused.

Garbage-Free Evaluation. We define *garbage-free* evaluation formally as:

Definition 2. (*Garbage free evaluation*)

An evaluation $\emptyset \mid e' \mapsto^*_r H \mid x$ is called *garbage-free* iff for every intermediate state $H_i \mid E[v]$ in the evaluation, we have that for all $y \in \text{dom}(H_i)$, $\text{reach}(y, H_i \mid \lceil E[v] \rceil)$.

where we use the notation $\lceil e \rceil$ to erase all drop and dup in the expression e . This is a refinement of the definition given by Reinking, Xie et al. [2021], which considered any non dup/drop steps, while we weaken this and consider only value steps v . By using $E[v]$ we consider exactly those points where we are at an allocation step ($C \bar{x}$ or $\lambda_{\bar{x}}x. e$), and this is exactly the point where we want to ensure that there is no garbage. In particular, match and let are heap invariant, and applications just expand to dups, drops, and substitution. This also gives more freedom to garbage-free algorithms as it becomes possible for example to push a drop down into the branches of a match which was not possible before. Using our new definition, we can then prove (App. C.6 in the tech report):

Theorem 3. (*\vdash_{GF} derivations are garbage-free*)

If $\emptyset \mid \emptyset \vdash_{\text{GF}} e \rightsquigarrow e'$ and $\emptyset \mid e' \mapsto^*_r H \mid x$, then the evaluation is garbage-free.

Even with the garbage-free side-condition, there are still many choice points in the derivations which gives us freedom to consider various algorithms. Generally though, when implementing Perceus one wants to dup as late as possible (push up dup into the leaves of the derivation), and do drops as early as possible. The original Perceus algorithm does this, and also trivially satisfies our garbage-free condition as it has the invariant $\Gamma \subseteq \text{fv}(e)$ for any derivation step.

Frame-Limited Evaluation. Similar to garbage-free evaluations, we can define frame-limited evaluations:

Definition 3. (*Frame-limited evaluation*)

An evaluation $\emptyset \mid e' \mapsto^*_r H \mid x$ is called *frame-limited* iff for every intermediate state $H_i \mid E[v]$ in the evaluation, we have that H_i equals (H_1, H_2) such that for all $y \in \text{dom}(H_1)$, $\text{reach}(y, H_1 \mid \lceil E[v] \rceil)$ and $|H_2| \leq c \cdot |E|$ for some constant c .

This expresses that at every allocation step, we may now hold on to extra heap space H_2 , but the size of H_2 is limited by a constant amount times the size of the evaluation context (i.e. the stack). We can now prove (App. C.7 in the tech report):

$\begin{array}{c} S \mid R \Vdash e \rightsquigarrow e' \quad (\uparrow \text{ is input, } \downarrow \text{ is output}) \\ \uparrow \quad \uparrow \quad \uparrow \quad \downarrow \\ S \text{ maps variables to their heap size (if known)} \\ R \text{ maps reuse variables } r \text{ to their available heap size, } \text{dom}(R) \cap \text{fv}(e) \end{array}$	
$\frac{}{S \mid \emptyset \Vdash x \rightsquigarrow x} \text{ [RVAR]}$	$\frac{S \mid R \Vdash e \rightsquigarrow e'}{S \mid R, r : n \Vdash e \rightsquigarrow \text{drop } r; e'} \text{ [RDROPR]}$
$\frac{}{S \mid \emptyset \Vdash C \bar{x} \rightsquigarrow C \bar{x}} \text{ [RCON]}$	$\frac{x : n \in S \quad S \mid R, r : n \Vdash e \rightsquigarrow e' \quad \text{fresh } r}{S \mid R \Vdash \text{drop } x; e \rightsquigarrow r \leftarrow \text{dropru } x; e'} \text{ [RDROP-REUSE]}$
$\frac{S \mid R \Vdash e \rightsquigarrow e'}{S \mid R \Vdash e x \rightsquigarrow e' x} \text{ [RAPP]}$	$\frac{S \mid R_1 \Vdash e_1 \rightsquigarrow e'_1 \quad S \mid R_2 \Vdash e_2 \rightsquigarrow e'_2}{S \mid R_1, R_2 \Vdash \text{let } x = e_1 \text{ in } e_2 \rightsquigarrow \text{let } x = e'_1 \text{ in } e'_2} \text{ [RLET]}$
$\frac{\emptyset \mid \emptyset \Vdash e \rightsquigarrow e'}{S \mid \emptyset \Vdash \lambda_{\Gamma} x. e \rightsquigarrow \lambda_{\Gamma} x. e'} \text{ [RLAM]}$	$\frac{S \mid R, r \Vdash e \rightsquigarrow e'}{S \mid R \Vdash r \leftarrow \text{dropru } x; e \rightsquigarrow r \leftarrow \text{dropru } x; e} \text{ [RDROPRU]}$
$\frac{S \mid R \Vdash e \rightsquigarrow e'}{S \mid R \Vdash \text{dup } x; e \rightsquigarrow \text{dup } x; e'} \text{ [RDUP]}$	$\frac{S \mid R \Vdash e \rightsquigarrow e'}{S \mid R \Vdash \text{drop } x; e \rightsquigarrow \text{drop } x; e'} \text{ [RDROP]}$
$\frac{S, x : n \mid R \Vdash e_i \rightsquigarrow e'_i \quad p_i = C x_1 \dots x_n}{S \mid R \Vdash \text{match } x \{ \overline{p_i \mapsto e_i} \} \rightsquigarrow \text{match } x \{ \overline{p_i \mapsto e'_i} \}} \text{ [RMATCH]}$	

Fig. 7. Declarative rules for drop-guided reuse analysis

Theorem 4. (\vdash_{FL} derivations are frame-limited)

If $\emptyset \mid \emptyset \vdash_{\text{FL}} e \rightsquigarrow e'$ and we have $\emptyset \mid e' \mapsto_r^* H \mid x$, then the evaluation is frame-limited with respect to the constant c chosen in the (\star) condition.

Borrowing a variable for a function call can also be frame-limited: Whenever an owned variable is passed as borrowed, we need to move it from the linear environment to the borrowed environment. But this can only happen in the `APP` and `LET` rules. With the `APP` rule, we can safely borrow since we use the variable later and with the `LET` rule this is frame-limited whenever the size of the borrowed variables is bounded by a constant.

5.4 Drop-Guided Reuse

Figure 7 shows the declarative derivation rules for drop-guided reuse analysis. A rule $S \mid R \Vdash e \rightsquigarrow e'$ states we can derive e' from e given a mapping S from variables to their heap cell size (if known), and a mapping R from reuse tokens r to their available heap size. Again, we use a multi-set for S and R ; we need to ensure we only use a reuse token r once, and similar to the owned environment Γ the leaf derivations require R to be empty (as in `RVAR`, `RCON`, and `RLAM`).

We state drop-guided reuse in terms of derivation rules instead of a specific algorithm in order to clearly expose the choice points. At any `drop` $x; e$ expression we can choose to either leave it as is, or use `RDROP-REUSE` to try to reuse it at runtime. We can only use `RDROP-REUSE` though if the size of x is known statically – in our rules (and in our implementation) this only happens by matching on a particular constructor (in `RMATCH`) but in principle we could also use other sources, like type information for example where all constructors happen to have the same size. Furthermore, in an implementation we also would only apply `RDROP-REUSE` if there is an actual opportunity in e for

reuse to occur – that is, we look first if e contains an occurrence of $C\ x_1 \dots x_n$. This is what we do in the Koka implementation where we actually do this in a single pass by passing all available reuse tokens R statefully and after transforming the sub expression we only apply `DROP-REUSE` if the reuse token was actually used (and otherwise we leave it unchanged as `RDROP`).

Another choice point is in the `RLET` rule where we can freely split R into R_1 and R_2 , where we may either reuse early or late if reuse is possible in both e_1 and e_2 . Also, rule `RDROPR` is not syntax directed and thus we have a choice in what reuse token to use if there are multiple reuse tokens of the right size. In the Koka implementation we always reuse early and in the `RLET` rule we first process e_1 and then e_2 with an R_2 consisting of any unused reuse tokens. If there is a choice of reuse tokens of the same size, we prefer a token of the same constructor and with the most fields that are unchanged, as that is beneficial for reuse specialization [Reinking, Xie et al. 2021, section 2.5].

To simplify the formalization and proofs, we do not introduce a new form of constructor reuse that we used before (as `C@r`) but instead assume the runtime evaluation recognizes a pair `drop r; C x` as a reuse opportunity, that is:

$$H, r \mapsto^1 C' \perp_1 \dots \perp_n \mid \text{drop } r; C\ x_1 \dots x_n \longrightarrow_r H, x \mapsto^1 C\ x_1 \dots x_n \mid x$$

Also, to simplify the formalization, we do not add \perp either and use allocated unit constructors `()` instead. The rules for `drop` evaluation are given in Figure 5. We can now show that `drop`-guided reuse is sound and frame-limited (App. C.8 in the tech report):

Theorem 5. (*Reuse is frame-limited*)

If $\Delta \mid \Gamma \vdash_{\text{GF}} e \rightsquigarrow e'$, and reuse derives $S \mid R \Vdash e' \rightsquigarrow e''$, then $\Delta \mid \Gamma, R \vdash_{\text{FL}} e \rightsquigarrow e''$.

That is, if we have a garbage free derivation followed by `drop`-guided reuse, we could have derived that same expression directly as well using a frame-limited derivation. This is how Koka implements this as well: first a Perceus algorithm (that satisfies garbage-free derivations) followed by a `drop`-guided reuse algorithm (that satisfies reuse derivations).

6 BENCHMARKS

We measured the performance of the new `drop`-guided reuse in Koka, versus the previous algorithm K. To get a sense of the absolute performance of Koka with the new reuse and as evidence that our compilation techniques (Perceus, `drop`-guided reuse, TRMC) are viable and can be competitive, we also included comparisons against other mature systems that use a range of memory reclamation techniques and are considered best-in-class. However, when comparing across systems we should interpret those results with care; for example, a particular benchmark may not be idiomatic for that particular language (like a pure tree versus using in place updates) etc. Similarly, we cannot draw conclusions from these small benchmarks on the relation of Perceus reference counting versus GC in general. However, the difference between the Koka variants is of course carefully controlled and supports the claims in this paper.

We use the following functional programming languages:

- “Koka”: Koka v2.3.3 with `drop`-guided reuse, compiling the generated C code with gcc 9.4.0 using a customized version of the `mimalloc` allocator [Leijen et al. 2019].
- “Koka, no trmc”: As “Koka”, but we disable TRMC (with “`-fno-opttrmc`”) to measure the impact of this optimization.
- “Koka, old”: As “Koka”, but we use Algorithm K instead of `drop`-guided reuse and no borrowing (on a branch of Koka: v2.3.3-old).
- “Koka, fbip”: As “Koka”, but with the implementations discussed in section 7.
- Multi-core OCaml 4.12. This has a concurrent generational collector with a minor and major heap. The minor heap uses a copying collector, while non-copying mark-sweep collector is

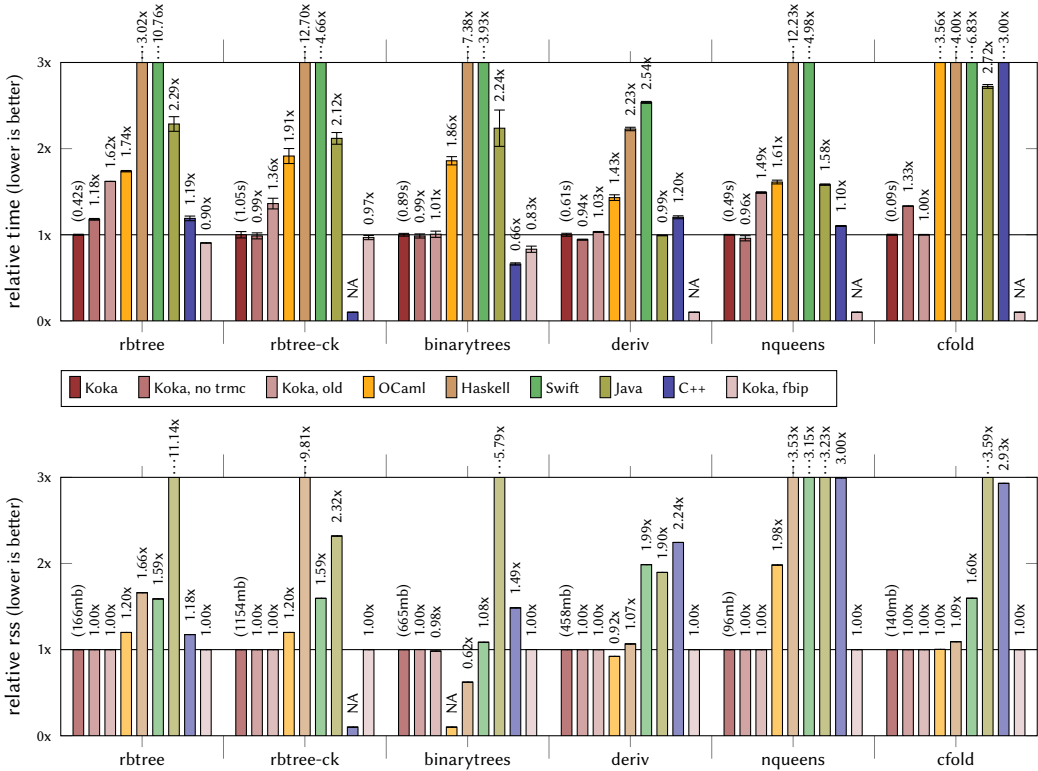


Fig. 8. All benchmarks with relative execution time and peak working set with respect to Koka. Using a 32-core x64 AMD5950X 3.4Ghz with 32GiB 3600Mhz memory (32KiB L1, 512KiB L2), Ubuntu 20.04.

used for the major heap [Doligez and Leroy 1993; Minsky et al. 2012, Chap.22; Sivaramakrishnan et al. 2020].

- Haskell, GHC 8.6.5. A highly optimizing compiler with a multi generational garbage collector. We used strictness annotations in the data structures to speed up the benchmarks, as well as to ensure that the same amount of work is done.

We also compare against Swift 5.6.1 and Java 17.0.1 LTS. Here, we keep the benchmarks in a functional style (without direct mutation) and only replace tail-calls with explicit loops. While Swift also uses reference counting [Choi et al. 2018; Ungar et al. 2017], Java uses the HotSpot JVM and the G1 concurrent, low-latency, generational garbage collector.

Finally, we use C++ as our performance baseline (with gcc 9.4.0 and using the standard libc allocator): For the `rbtree` benchmark we used the standard STL `std::map` implementation that uses a highly optimized in-place updating version of red-black trees [Free Software Foundation et al. 1994]. The `binarytrees` benchmark use a monotonic buffer resource for memory management, while the other benchmarks (`nqueens`, `deriv`, and `cfold`) do not reclaim memory at all (for C++).

The results of all benchmarks on an AMD5950x x64 system with Ubuntu 20.04 are shown in Figure 8. The most important conclusion is that over these benchmarks the new drop-guided reuse is always as fast (and sometimes up to 1.6× faster) than the old reuse algorithm.

6.1 Benchmarking Balanced Trees

We use the same red-black tree benchmarks as used by Reinking, Xie et al. [2021] (which makes use of the insertion algorithm in our running example from section 2.2). There are two versions:

- `rbtree`: inserts 4200000 elements in tree and afterwards folds over the tree counting the `True` elements.
- `rbtree-ck`: like `rbtree` but also keeps a list of every 5th tree that is generated, effectively sharing many subtrees. This implies many subtrees have a non-unique reference count which causes many more slow paths to be taken in the Koka code. This is not done for C++ as `std::map` does not support persistence.

If we look at the results in Figure 8, we see that Koka performs surprisingly well in comparison with other mature systems, and outperforms the C++ implementation by almost 20% – how is that even possible? We conjecture this is mainly due to two factors: 1) the TRMC optimization leads to using less stack space than the C++ implementation (and perhaps less register-spills), and 2) the allocator that Koka uses can use 8-byte alignment while C++ requires 16-byte minimal alignment which leads to less allocated memory. If we compare the non-TRMC optimized version it performs only slightly better than C++ which seems to support our thesis. Even though this is a typical functional style algorithm both OCaml and Haskell are quite a bit slower (1.7× and 3×).

The `rbtree-ck` version with many shared subtrees increases the relative speed of Koka even further. This is also somewhat surprising as it is clearly much worse for reuse (with many shared subtrees), but of course it similarly causes more pressure on garbage collectors as well as more objects get promoted to older generations. Still, with only every 5th tree shared there is some opportunity for reuse left. The Multi-core OCaml GC performs especially well here staying within 20% of Koka’s *garbage-free* memory usage.

If we compare Koka with drop-guided reuse against Koka “old” which uses algorithm K, we can see the new algorithm does about 1.6× better on this benchmark due to the improved reuse for the `is-red` test as explained in Section 2.2. Finally, what is the fastest “Koka, fbip” version? We will discuss this in Section 7.

6.2 Binary Trees

The `binarytrees` benchmark comes from the Computer Language Benchmark Game [Game 2021] which is an adaptation of Hans Boehm’s `GCbench` benchmark [Boehm 2000] (which in turn was adapted from a benchmark by John Ellis and Pete Kovac). This is an interesting benchmark as it uses *concurrent* allocation of many binary trees and calculates their checksums. Moreover, we can compare against the best performing implementations that were created by experts in each of our comparison languages.

The top 17 implementations are all languages with manual allocation (C++, C, Rust, and Free Pascal) with the top entry being C++ (#7) followed very closely by Rust and two other C++ entries (#5 and #4). For our purposes, we use C++ benchmark #5 since the performance of #7 was uneven across systems. The C++ entries all use very efficient allocation by using a `monotonic_buffer_resource` for bump-pointer allocation where all nodes are freed at once at every iteration.

Figure 8 shows the benchmark results of `binarytrees` (the Koka implementation can be found in App. A.1 in the tech report). Besides C++, Koka outperforms all other languages here even though it uses a simple thread pool implementation without work-stealing. The C++ implementation is still quite a bit faster though (0.66×). Since the benchmark does concurrent allocations reference counting generally needs to be atomic. However, due to careful language design, Koka can avoid most of the overhead by checking upfront if a reference count needs to be atomic or not [Reinking, Xie et al. 2021].

6.3 Other Benchmarks

The `nqueens`, `cfold`, and `deriv` benchmarks are the same as used in the Perceus paper. We included them here for completeness but for these benchmarks the new drop guided reuse gives very similar results to algorithm K.

The `nqueens` benchmark computes a list of all solutions to the N-queens problem of size 21. The large speedup here compared to Koka “old” is actually due to borrowing in the `safe` function. The `cfold` benchmark performs constant folding in a program, and the `deriv` benchmark computes a symbolic derivative of a large expression. Again, in these benchmarks the difference with the previous reuse algorithm is minimal.

7 FBIP: FUNCTIONAL BUT IN-PLACE ALGORITHMS

Just like tail-recursion let us express loops as recursive functions, reuse analysis can be used to express imperative algorithms in a functional style. We call such algorithms *Functional But In-Place* (FBIP) [Reinking, Xie et al. 2021]. In particular, it is often possible to reformulate an algorithm with non-tail-recursive calls into one that is tail recursive using an explicit visitor data type. By relying on drop-guided reuse analysis, we can now make the allocation of the visitor data type “free” by ensuring we can reuse existing objects. In this section we illustrate this technique to the `rbtree` and `binarytrees` benchmarks and show we can improve their performance even further.

7.1 Binary Trees

Most of the work in the `binarytrees` benchmark is in creating the trees and calculating their size using the `check` function:

```
type tree
  Node( left: tree, right: tree )
  Tip

fun check( t : tree ) : int
  match t
  Node(l,r) -> check(l) + check(r) + 1
  Tip      -> 0
```

This consists of two non-tail-recursive calls to `check`. We can improve upon this using FBIP where we use a *visitor* data type that tracks where we are in the tree. For our purposes we define:

```
type visit
  NodeR( right: tree, rest: visit )
  Done
```

We can use the new visitor datatype to write a `check` function that uses no extra stack space as all calls are tail-recursive:

```
fun check-fbip( t : tree, v : visit, acc : int ) : int
  match t
  Node(l,r) -> check-fbip( l, NodeR(r,v), acc + 1) // (A)
  Tip      -> match v
  NodeR(r,v') -> check-fbip( r, v', acc) // (B)
  Done       -> acc // (C)
```

At every `Node` we directly go down the left branch but remember that we still need to visit the right node by extending our `visit` datatype (A). When we reach a `Tip` we go through our visitor to now visit the saved right nodes (B) until we are done (C).

This may look more expensive, but when `t` happens to be unique at runtime, the `NodeR` allocations in (1) will reuse the `Node` that is matched – effectively updating these nodes in place to create a list of nodes that still need to be visited. At runtime this becomes tight loop that directly reuses the

memory of the tree it is checking. Effectively, the generated code uses in-place *pointer reversal* to visit the tree without using further stack space much like stackless marking in garbage collectors with the Schorr-Waite algorithm [Schorr and Waite 1967]. In Figure 8 this is the *Koka fbip* variant and with this implementation of check Koka becomes 17% faster and within 25% of the performance of the C++ implementation (0.8× versus Koka fbip).

7.2 Balanced Trees

We can apply the same FBIP technique on our previous red-black tree balanced insertion: as we saw in Section 3.2, due to TRMC most calls to `ins` are tail recursive but not the ones that needed rebalancing. We can define a visitor for red-black trees as the *derivative* of the `tree` data type [Huet 1997; McBride 2001]:

```
type zipper(a)
  NodeR(clr:color, l:tree(a), key:int, value:a, zip:zipper(a))
  NodeL(clr:color, zip:zipper(a), key:int, value:a, r:tree(a))
  Done
```

Using this data type we can now traverse down a tree in tail-recursive way to the insertion point, while building up the zipper that tracks where we are in the tree:

```
fun ins(t : tree(a), k : int, v : a, z : zipper(a)) : tree(a)
  match t
  Node(c, l, kx, vx, r)
    -> if k < kx then ins(l, k, v, NodeL(c, z, kx, vx, r))
        elif k > kx then ins(r, k, v, NodeR(c, l, kx, vx, z))
        else rebuild(z, Node(c, l, kx, vx, r)) // A
  Leaf -> balance(z, Leaf, k, v, Leaf) // B
```

If the element is already present (A), we can use the tail-recursive `rebuild` function to reconstruct the tree using our just constructed zipper:

```
fun rebuild( z : zipper(a), t : tree(a) ) : tree(a)
  match z
  NodeR(c, l, k, v, z1) -> rebuild(z1, Node(c, l, k, v, t))
  NodeL(c, z1, k, v, r) -> rebuild(z1, Node(c, t, k, v, r))
  Done -> t
```

If we reach a leaf node though (B), we use `balance` to rebalance the tree going upward. Rebalancing is also tail-recursive now:

```
fun balance( z : zipper(a), l : tree(a), k : int, v : a, r : tree(a) ) : tree(a)
  match z
  Done -> Node(Black, l, k, v, r)
  NodeR(Black, l1, k1, v1, z1) -> rebuild( z1, Node(Black, l1, k1, v1, Node(Red, l, k, v, r)) )
  NodeR(Red, l1, k1, v1, z1) -> match z1
    Done -> Node(Black, l1, k1, v1, Node(Red, l, k, v, r))
    NodeR(_, l2, k2, v2, z2) -> balance(z2, Node(Black, l2, k2, v2, l1), k1, v1, Node(Black, l, k, v, r))
    NodeL(...) -> ...
  NodeL(Black, ...) -> ...
  NodeL(Red, ...) -> ...
```

As before, besides the inserted node, every `Node` allocation in `rebuild` and `balance` can be paired with a matched `NodeR` or `NodeL` (and the other way around in `ins`), and all can be reused in-place at runtime if the tree happens to be unique. Furthermore, for a shared tree we only allocate the zipper upfront, but the zipper itself is always unique and reused in-place by `rebuild` and `balance`.

Our novel FBIP algorithm for balanced insertion improves further upon the standard Okasaki style algorithm [Okasaki 1999a] since we stop rebalancing as soon as we reach a `Black` node, and switch to `rebuild` instead (which is also done in the usual imperative algorithms [Guibas and Sedgewick 1978]). The full implementation can be found in App. A.2 in the tech report. This is the *Koka fbip* variant in Figure 8 which is about 10% faster than the regular version and now around 30% faster than the C++ STL version.

8 RELATED WORK

Our work is closely based earlier work by Reinking, Xie et al. [2021] and Ullrich and de Moura [2019] (in the context of the Lean theorem prover [Moura and Ullrich 2021]). In this work we improve upon the both of the two earlier reuse algorithms, and show that drop-guided reuse is strictly better as it is frame-limited and can find more opportunities for reuse. Our λ^{ln} calculus, heap semantics, and derivation rules differ in important details from [Reinking, Xie et al. 2021]: the normalization simplifies the derivation rules and allows us to concisely express the garbage-free and frame-limited side conditions (and no longer as a particular property of a specific algorithm). It is also now immediate that the previous published Perceus algorithm is garbage-free. The new frame-limited condition allows us to reason about various transformations that are no longer garbage-free but can still be bounded.

The notion of safe-for-space was introduced by Appel [1991] and Paraskevopoulou and Appel [2019] studied this further. Similar to our reuse transformation and frame-limited notion, the latter work introduces a general framework to show that the closure conversion transformation with flat environments is safe-for-space, while linked closure conversion is not. Other examples where a program transformation was proven to respect a resource bound include Crary and Weirich [2000] and Minamide [1999].

Using explicit reference count instructions in order to optimize them via static analysis is described as early as Barth [1975]. Mutating unique references in place has traditionally focused on array updates [Hudak and Bloss 1985], as in functional array languages like Sisal [McGraw et al. 1983] and SaC [Grelck and Trojahnner 2004; Scholz 2003]. Férey and Shankar [2016] provide functional array primitives that use in-place mutation if the array has a unique reference which is also present in the Koka implementation. We believe this would work especially well in combination with reuse-analysis for BTree-like structures using trees of small functional arrays.

The λ^1 and λ^{ln} calculus are closely based on linear logic. Turner and Wadler [1999] give a heap-based operational interpretation which does not need reference counts as linearity is tracked by the type system. In contrast, Chirimar et al. [1996] give an interpretation of linear logic in terms of reference counting, but in their system, values with a linear type are not guaranteed to have a unique reference at runtime.

Generally, a system with linear types [Wadler 1990], like linear Haskell [Bernardy et al. 2017], the uniqueness typing of Clean [Barendsen and Smetsers 1996; Vries et al. 2008], or the quantitative type theory of Idris [Brady 2021], can offer *static* guarantees that the corresponding objects are unique at runtime, so that destructive updates can always be performed safely. However, this also requires writing multiple versions of a function for each case (unique- versus shared argument, or an in-place mutating data structure versus a persistent one). In contrast, reuse analysis relies on dynamic runtime information, and thus reuse can be performed generally. This is also what enables FBIP to use a single function that can be used for both unique or shared objects (since the uniqueness property is *not* part of the type). These two mechanisms can be combined: if our system is extended with unique types, then reuse analysis can eliminate corresponding uniqueness checks.

9 CONCLUSION AND FUTURE WORK

In this work we show the effectiveness of drop-guided reuse, and give a precise characterization of garbage-free and frame-limited evaluations. However, this works well partly because in a functional style language like Koka it is uncommon to create cycles (which can only be created through mutable references). We would like to see if we can combine some of the static analysis with cycle collection. Also, we are interested in language support to guarantee that reuse is happening at compile time.

REFERENCES

- Andrew W. Appel. 1991. *Compiling with Continuations*. Cambridge University Press. <https://doi.org/10.1017/CBO9780511609619>.
- Erik Barendsen, and Sjaak Smetsers. 1996. Uniqueness Typing for Functional Languages with Graph Rewriting Semantics. *Mathematical Structures in Computer Science* 6 (6). Cambridge University Press: 579–612. <https://doi.org/10.1017/S0960129500070109>.
- Jeffrey M. Barth. Jun. 1975. *Shifting Garbage Collection Overhead to Compile Time*. UCB/ERL M524. EECS Department, University of California, Berkeley. <https://doi.org/10.1145/359636.359713>.
- Jean-Philippe Bernardy, Mathieu Boespflug, Ryan R. Newton, Simon Peyton Jones, and Arnaud Spiwack. Dec. 2017. Linear Haskell: Practical Linearity in a Higher-Order Polymorphic Language. *Proc. ACM Program. Lang.* 2 (POPL). <https://doi.org/10.1145/3158093>.
- Hans Boehm. 2000. GC Bench. https://hboehm.info/gc/gc_bench.
- Edwin Brady. 2021. Idris 2: Quantitative Type Theory in Practice. In *35th European Conference on Object-Oriented Programming (ECOOP 2021)*, edited by Anders Møller and Manu Sridharan, 194:9:1–9:26. Leibniz International Proceedings in Informatics (LIPIcs). Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl, Germany. <https://doi.org/10.4230/LIPIcs.ECOOP.2021.9>.
- Jawahar Chirimar, Carl A. Gunter, and Jon G. Riecke. 1996. Reference Counting as a Computational Interpretation of Linear Logic. *Journal of Functional Programming* 6: 6–2. <https://doi.org/10.1017/S0956796800001660>.
- Jiho Choi, Thomas Shull, and Josep Torrellas. 2018. Biased Reference Counting: Minimizing Atomic Operations in Garbage Collection. In *Proceedings of the 27th International Conference on Parallel Architectures and Compilation Techniques*. PACT '18. Limassol, Cyprus. <https://doi.org/10.1145/3243176.3243195>.
- George E Collins. 1960. A Method for Overlapping and Erasure of Lists. *Communications of the ACM* 3 (12). ACM New York, NY, USA: 655–657. <https://doi.org/10.1145/367487.367501>.
- Karl Cray, and Stephanie Weirich. 2000. Resource Bound Certification. In *Proceedings of the 27th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, 184–198. POPL '00. Boston, MA, USA. <https://doi.org/10.1145/325694.325716>.
- Damien Doligez, and Xavier Leroy. Jan. 1993. A Concurrent, Generational Garbage Collector for a Multithreaded Implementation of ML. In *Proceedings of the 20th ACM Symposium on Principles of Programming Languages (POPL)*, 113–123. ACM press. <https://doi.org/10.1145/158511.158611>.
- Gaspard Férey, and Natarajan Shankar. 2016. Code Generation Using a Formal Model of Reference Counting. In *NASA Formal Methods*, edited by Sanjai Rayadurgam and Oksana Tkachuk, 150–165. Springer International Publishing. https://doi.org/10.1007/978-3-319-40648-0_12.
- Cormac Flanagan, Amr Sabry, Bruce F. Duba, and Matthias Felleisen. 1993. The Essence of Compiling with Continuations. In *Proceedings of the ACM SIGPLAN 1993 Conference on Programming Language Design and Implementation*, 237–247. PLDI '93. Albuquerque, New Mexico, USA. <https://doi.org/10.1145/155090.155113>.
- Free Software Foundation, Silicon Graphics, and Hewlett–Packard Company. 1994. Internal Red-Black Tree Implementation for "Stl:map." <https://github.com/gcc-mirror/gcc/tree/master/libstdc++v3/src/c++98/tree.cc>.
- Matt Gallagher. Dec. 2016. Reference Counted Releases in Swift. <https://www.cocoawithlove.com/blog/resources-releases-reentrancy.html>. Blog post.
- The Computer Language Benchmark Game. Nov. 2021. Binarytrees. <https://benchmarksgame-team.pages.debian.net/benchmarksgame/performance/binarytrees.html>.
- Clemens Grelck, and Kai Trojahnner. Sep. 2004. Implicit Memory Management for SAC. In *6th International Workshop on Implementation and Application of Functional Languages (IFL'04)*. Lübeck, Germany.
- Leo J Guibas, and Robert Sedgwick. 1978. A Dichromatic Framework for Balanced Trees. In *19th Annual Symposium on Foundations of Computer Science (sfcs 1978)*, 8–21. IEEE. <https://doi.org/10.1109/SFCS.1978.3>.
- Paul Hudak, and Adrienne Bloss. 1985. The Aggregate Update Problem in Functional Programming Systems. In *Proceedings of the 12th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages*, 300–314. POPL '85. ACM, New Orleans, Louisiana, USA. <https://doi.org/10.1145/318593.318660>.
- Gérard P. Huet. 1997. The Zipper. *Journal of Functional Programming* 7 (5): 549–554. <https://doi.org/10.1017/S0956796897002864>.
- Daan Leijen. 2014. Koka: Programming with Row Polymorphic Effect Types. In *MSFP'14, 5th Workshop on Mathematically Structured Functional Programming*. <https://doi.org/10.4204/EPTCS.153.8>.
- Daan Leijen. Jan. 2017. Type Directed Compilation of Row-Typed Algebraic Effects. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages (POPL'17)*, 486–499. Paris, France. <https://doi.org/10.1145/3009837.3009872>.
- Daan Leijen. 2021. The Koka Language. <https://koka-lang.github.io>.

- Daan Leijen, Zorn Ben, and Leo de Moura. 2019. Mimalloc: Free List Sharding in Action. *Programming Languages and Systems*, LNCS, 11893. Springer International Publishing. https://doi.org/10.1007/978-3-030-34175-6_13. APLAS'19.
- Anton Lorenzen, and Daan Leijen. Nov. 2021. *Reference Counting with Frame Limited Reuse (extended Version)*. MSR-TR-2021-30. Microsoft Research.
- Conor McBride. 2001. The Derivative of a Regular Type Is Its Type of One-Hole Contexts. <http://strictlypositive.org/diff.pdf>. (Extended Abstract).
- J. McGraw, S. Skedzielewski, S. Allan, D. Grit, R. Oldehoeft, J. Glauert, I. Dobes, and P. Hohensee. Jul. 1983. *SISAL: Streams and Iteration in a Single-Assignment Language*. *Language Reference Manual, Version 1.1*. LLL/M-146, ON: DE83016576. Lawrence Livermore National Lab., CA, USA.
- Yasuhiko Minamide. 1999. Space-Profiling Semantics of the Call-by-Value Lambda Calculus and the CPS Transformation. *Electronic Notes in Theoretical Computer Science* 26: 105–120. [https://doi.org/10.1016/S1571-0661\(05\)80286-5](https://doi.org/10.1016/S1571-0661(05)80286-5). HOOTS '99, Higher Order Operational Techniques in Semantics.
- Yaron Minsky, Anil Madhavapeddy, and Jason Hickey. 2012. *Real World OCaml: Functional Programming for the Masses*. <https://dev.realworldocaml.org>.
- Leonardo de Moura, and Sebastian Ullrich. 2021. The Lean 4 Theorem Prover and Programming Language. In *Automated Deduction – CADE 28*, edited by André Platzer and Geoff Sutcliffe, 625–635. https://doi.org/10.1007/978-3-030-79876-5_37.
- Chris Okasaki. 1999. Red-Black Trees in a Functional Setting. *Journal of Functional Programming* 9 (4). Cambridge University Press: 471–477. <https://doi.org/10.1017/S0956796899003494>.
- Chris Okasaki. Jun. 1999. *Purely Functional Data Structures*. Colombia University, New York.
- Zoe Paraskevopoulou, and Andrew W Appel. 2019. Closure Conversion Is Safe for Space. *Proceedings of the ACM on Programming Languages* 3 (ICFP). ACM New York, NY, USA: 1–29. <https://doi.org/10.1145/3341687>.
- Gordon D. Plotkin, and John Power. 2003. Algebraic Operations and Generic Effects. *Applied Categorical Structures* 11 (1): 69–94. <https://doi.org/10.1023/A:1023064908962>.
- Gordon D. Plotkin, and Matija Pretnar. Mar. 2009. Handlers of Algebraic Effects. In *18th European Symposium on Programming Languages and Systems*, 80–94. ESOP'09. York, UK. https://doi.org/10.1007/978-3-642-00590-9_7.
- Reinking, Xie, de Moura, and Leijen. 2021. Perceus: Garbage Free Reference Counting with Reuse. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, 96–111. PLDI 2021. ACK, New York, NY, USA. <https://doi.org/10.1145/3453483.3454032>.
- Sven-Bodo Scholz. Nov. 2003. Single Assignment C: Efficient Support for High-Level Array Operations in a Functional Setting. *Journal of Functional Programming* 13 (6). Cambridge University Press, USA: 1005–1059. <https://doi.org/10.1017/S0956796802004458>.
- H. Schorr, and W. M. Waite. Aug. 1967. An Efficient Machine-Independent Procedure for Garbage Collection in Various List Structures. *Comm. ACM* 10 (8). ACM: 501–506. <https://doi.org/10.1145/363534.363554>.
- KC Sivaramakrishnan, Stephen Dolan, Leo White, Sadiq Jaffer, Tom Kelly, Anmol Sahoo, Sudha Parimala, Atul Dhiman, and Anil Madhavapeddy. Aug. 2020. Retrofitting Parallelism onto OCaml. *Proc. ACM Program. Lang.* 4 (ICFP). ACM. <https://doi.org/10.1145/3408995>.
- David N. Turner, and Phillip Wadler. 1999. Operational Interpretations of Linear Logic, number 227: 231–248. [https://doi.org/10.1016/S0304-3975\(99\)00054-7](https://doi.org/10.1016/S0304-3975(99)00054-7).
- Sebastian Ullrich, and Leonardo de Moura. Sep. 2019. Counting Immutable Beans – Reference Counting Optimized for Purely Functional Programming. In *Proceedings of the 31st Symposium on Implementation and Application of Functional Languages (IFL'19)*. Singapore. <https://doi.org/10.1145/3412932.3412935>.
- David Ungar, David Grove, and Hubertus Franke. 2017. Dynamic Atomicity: Optimizing Swift Memory Management. In *Proceedings of the 13th ACM SIGPLAN International Symposium on on Dynamic Languages*, 15–26. DLS 2017. Vancouver, BC, Canada. <https://doi.org/10.1145/3133841.3133843>.
- Edsko de Vries, Rinus Plasmeijer, and David M. Abrahamson. 2008. Uniqueness Typing Simplified. In *Implementation and Application of Functional Languages (IFL'08)*, edited by Olaf Chitil, Zoltán Horváth, and Viktória Zsók, 201–218. Springer. https://doi.org/10.1007/978-3-540-85373-2_12.
- Phillip Wadler. 1990. Linear Types Can Change the World! In *Programming Concepts and Methods*.
- Andrew K. Wright, and Matthias Felleisen. Nov. 1994. A Syntactic Approach to Type Soundness. *Inf. Comput.* 115 (1): 38–94. <https://doi.org/10.1006/inco.1994.1093>.
- Ningning Xie, and Daan Leijen. Aug. 2021. Generalized Evidence Passing for Effect Handlers – Efficient Compilation of Effect Handlers to C. In *Proceedings of the 26th ACM SIGPLAN International Conference on Functional Programming (ICFP'2021)*. ICFP '21. Virtual. <https://doi.org/10.1145/3473576>.
- Danil Yarantsev. Oct. 2020. ORC - Nim's Cycle Collector. <https://nim-lang.org/blog/2020/10/15/introduction-to-arc-orc-in-nim.html>.