



Garbage-Collection Safety for Region-Based Type-Polymorphic Programs

MARTIN ELSMAN, University of Copenhagen, Denmark

Region inference offers a mechanism to reduce (and sometimes entirely remove) the need for reference-tracing garbage collection by inferring where to insert allocation and deallocation instructions in a program at compile time. When the mechanism is combined with techniques for reference-tracing garbage collection, which is helpful in general to support programs with very dynamic memory behaviours, it turns out that region-inference is complementary to adding generations to a reference-tracing collector. However, region-inference and the associated region-representation analyses that make such a memory management strategy perform well in practice are complex, both from a theoretical point-of-view and from an implementation point-of-view.

In this paper, we demonstrate a soundness problem with existing theoretical developments, which have to do with ensuring that, even for higher-order polymorphic programs, no dangling-pointers appear during a reference-tracing collection. This problem has materialised as a practical soundness problem in a real implementation based on region inference. As a solution, we present a modified, yet simple, region type-system that captures garbage-collection effects, even for polymorphic higher-order code, and outline how region inference and region-representation analyses are adapted to the new type system. The new type system allows for associating simpler region type-schemes with functions, compared to original work, makes it possible to combine region-based memory management with partly tag-free reference-tracing (and generational) garbage-collection, and repairs previously derived work that is based on the erroneous published results.

CCS Concepts: • **Software and its engineering** → **Functional languages; Runtime environments.**

Additional Key Words and Phrases: region-inference, garbage-collection, Standard ML

ACM Reference Format:

Martin Elsmann. 2023. Garbage-Collection Safety for Region-Based Type-Polymorphic Programs. *Proc. ACM Program. Lang.* 7, PLDI, Article 115 (June 2023), 23 pages. <https://doi.org/10.1145/3591229>

1 INTRODUCTION

Region-based memory management allows programmers to associate life-times of objects with so-called regions and to reason about how and when such regions are allocated and deallocated. Region-based memory management, as it is implemented for instance in Rust [Aldrich et al. 2002], can be a valuable tool for constructing certain kinds of critical systems, such as real-time embedded systems [Salagnac et al. 2006]. Region inference differs from explicit region-based memory management by taking a non-annotated program as input and producing a region-annotated program, including directives for allocating and deallocating regions [Tofte et al. 2004]. The result is a programming paradigm where programmers can learn to write region-friendly code (by following certain patterns [Tofte et al. 2022]) to obtain good space and time performance for critical parts of the program.

The region-based memory management scheme that we consider here is based on the stack discipline. Whenever e is some expression, region inference may decide to replace e with the term $\text{letregion } \rho \text{ in } e'$, where e' is the result of transforming the expression e , which includes

Author's address: Martin Elsmann, Department of Computer Science, University of Copenhagen, Universitetsparken 5, Copenhagen, DK-2100, Denmark, mael@di.ku.dk.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2023 Copyright held by the owner/author(s).

2475-1421/2023/6-ART115

<https://doi.org/10.1145/3591229>

annotating allocating expressions with particular region variables (e.g., ρ) specifying the region each value should be stored in. The semantics of the `letregion` term is first to allocate a region (initially an empty list of pages) on the region stack, bind the region to the region variable ρ , evaluate e' , and, finally, deallocate the region bound to ρ (and its pages). The region type system allows regions to be passed to functions at run time (i.e., functions can be region-polymorphic) and to be captured in closures. The soundness of region inference ensures that a region is not deallocated as long as a value within it may be used by the remainder of the computation.

To remedy the problem that region inference does not always capture precisely the lifetime properties of objects, previous work has augmented the static inference scheme with more dynamic lifetime-based reference-tracing copying garbage collectors [Elsman and Hallenberg 2020, 2021; Hallenberg et al. 2002]. For such integrations, care must be taken to rule out the possibility of deallocating regions with incoming pointers from live objects.

It turns out, however, that region inference (and the accompanying region typing rules) allows for so-called *dangling pointers*, which are pointers to objects that region inference has determined will not be needed by the remainder of the computation, yet are captured in objects (e.g., in closures) that escape a `letregion` construct and are live from a reference-tracing point-of-view.

Previous work attempt to rule out the possibility of dangling pointers by adjusting the region typing rules (and region inference) in such a way that the type of an object will mention all regions that the object may live in [Elsman 2003] (by enlarging the latent effect sets of certain function types). As a consequence, such an adjustment will capture the effect of a reference-tracing garbage collection appearing when control enters a function, for instance. Unfortunately, the previous attempts at ruling out dangling pointers fail for certain programs that involve a combination of higher-order dead values and type polymorphism. From a theoretical point-of-view, the problem is that the region-typing rules, which form the basis of region inference, are not closed under type substitution, which is erroneously claimed by previous work [Elsman 2003; Elsmann and Hallenberg 2021]. As we shall see, this problem is not straightforward to overcome.

The erroneous theoretical results are exposed through the MLKit Standard ML compiler [Tofte et al. 2022]. The MLKit compiles programs to native code for Linux and macOS [Elsman and Hallenberg 1995] and implements techniques for dividing regions into those that are bounded and may be stack allocated (called *finite regions*) and those that must be heap allocated (called *infinite regions*) [Birkedal et al. 1996; Tofte et al. 2004], which are subject to reference-tracing collections. Based on the theoretical insights described above, dangling pointers may occur at runtime, which may cause programs to fail (e.g., segfault) during a reference-tracing collection. Fortunately, it is possible to adjust the region type system to mitigate the problem and provide guarantees, also for higher-order type-polymorphic programs, that dangling pointers do not appear at runtime.

The contributions of this paper are the following:

- (1) We identify a safety problem with existing techniques for abandoning dangling-pointers at runtime, which serves as an assumption for combining region-inference with reference-tracing garbage collection.
- (2) We present a non-trivial and novel modification to an existing region-based type system that rules out dangling pointers and allows for combining region-based memory management with reference-tracing (and even generational) garbage collection.
- (3) We describe how the modified type system affects region inference and the region representation analyses that form the basis for a practical compiler infrastructure.
- (4) We demonstrate that, in practice, the necessary modifications have little effect on performance, affect only a small set of functions, and resolve problems with running large programs, such as MLKit and MLton, two different Standard ML compilers.

The paper is organised as follows. In the following section, we first give an informal example demonstrating a program for which dangling pointers will occur at runtime unless the region typing rules that form the basis of region inference are adjusted beyond previous suggestions. In this section, we also demonstrate, informally, how we may adjust the region typing rules further to eliminate completely the presence of dangling pointers.

In Section 3, we present a simplified, but formal, region type system for a language that serves as a target language for region inference. We present a number of properties of the type system, including region type soundness and the property that no dangling pointers are introduced during evaluation. In Section 4, we describe various aspects of the implementation, including how region inference is implemented for the system. We also give examples demonstrating some non-trivial aspects of the system. In Section 5, we present experimental results and evaluate the work. In Section 6, we describe related work, and in Section 7, we conclude.

2 THE PROBLEM

We now demonstrate the unsoundness problem that may occur when reference-tracing garbage collection is combined with higher-order functions and type-polymorphism. We present the problem in the context of a slightly modified region-based type system, compared to the original Tofte-Talpin region-type system, but emphasise that the unsoundness can be demonstrated also for the original system even if the typing rules are modified as described in [Tofte and Talpin 1993, page 50] and [Elsmann 2003], which aim at abandoning dangling pointers (but fail).

Consider the function-composition function, which has the following ML type-scheme:

$$\mathbf{val} \ o : (\gamma \rightarrow \beta) \times (\alpha \rightarrow \gamma) \rightarrow \alpha \rightarrow \beta$$

Here α , β , and γ are (implicitly quantified) type variables and o is an infix function that takes a pair of two functions as argument and returns a function as the result.

The region annotated version of the function o has the following region (and effect) type-scheme:

$$\forall \epsilon \epsilon_0 \epsilon_1 \epsilon_2 \rho_0 \rho_1 \rho_2 \rho_3 \alpha \beta \gamma. \quad (1)$$

$$((\gamma \xrightarrow{\epsilon_2, \emptyset} \beta, \rho_2) \times (\alpha \xrightarrow{\epsilon_1, \emptyset} \gamma, \rho_1), \rho_0) \xrightarrow{\epsilon_0, \{\rho_0, \rho_3\}} (\alpha \xrightarrow{\epsilon, \{\epsilon_1, \epsilon_2, \rho_1, \rho_2\}} \beta, \rho_3)$$

Here ϵ , ϵ_0 , ϵ_1 , and ϵ_2 are effect variables and ρ_0 , ρ_1 , ρ_2 , and ρ_3 are region variables. We see that function type constructors are annotated with so-called arrow effects, each of which is a set of atomic effects (effect variables and region variables) identified by an effect variable.¹ Moreover, type constructors for products (\times) and functions are annotated with region variables that indicate in which region a particular constructed value resides. The arrow effect $\epsilon_0, \{\rho_0, \rho_3\}$ expresses that when the function o is applied to a pair of functions, the pair, which resides in ρ_0 is deconstructed and a new closure is stored in region ρ_3 . The arrow effect $\epsilon, \{\epsilon_1, \epsilon_2, \rho_1, \rho_2\}$, which appears on the arrow of the type of the resulting function, expresses that, when the function is applied, the two argument functions are accessed (ρ_2, ρ_1) and evaluated (ϵ_2, ϵ_1).

When a function such as o is applied, a particular instantiation of the function's type-scheme is described by a particular substitution that maps generic effect variables to arrow effects, generic region variables to region variables, and generic type variables to region-annotated types.

Consider now the problematic function run in Figure 1, which first creates a function h , thereby capturing a dead value in a closure, calls a function $work$ (for the sake of triggering a reference-tracing collection), and finally calls the function h . Notice that the argument to the function o

¹As described in details later, allowing arrow effects to be identified by effect variables (so-called *effect-handles*) enables the possibility that effects may grow by applying effect substitutions (which map effect variables to arrow effects).

```

fun run () : unit =
  let val h : unit -> unit = (op o) let val x = "oh" ^ "no"
                                in (fn x => (), fn () => x)
                                end
      val _ = work () (* trigger gc *)
  in h ()
  end

```

Fig. 1. Problematic source program involving higher-order functions, type-polymorphism, and dead values.

<pre> fun run () : unit = letregion ρ_1, ρ_2, ρ_3 in let val h : (unit $\xrightarrow{\epsilon:\{\rho_1, \rho_2\}}$ unit, ρ_3) = letregion ρ, ρ_0 in (op o [ρ_3]) let val x = op ^ [ρ] ("oh", "no") in (fn at ρ_1 x => (), fn at ρ_2 () => x) at ρ_0 end end val _ = work () (* trigger gc *) in h () end end </pre> <p style="text-align: center;">(a)</p>	<pre> fun run () : unit = letregion $\rho, \rho_1, \rho_2, \rho_3$ in let val h : (unit $\xrightarrow{\epsilon:\{\rho_1, \rho_2, \rho\}}$ unit, ρ_3) = letregion ρ_0 in (op o [ρ_3]) let val x = op ^ [ρ] ("oh", "no") in (fn at ρ_1 x => (), fn at ρ_2 () => x) at ρ_0 end end val _ = work () (* trigger gc *) in h () end end </pre> <p style="text-align: center;">(b)</p>
--	--

Fig. 2. An unsound region-annotated program (a) and an alternative sound region-annotated program (b).

evaluates to a pair of functions for which the second function returns a pointer to an already allocated value ("ohno") and the first function will silently discard its argument.

Next, consider the region-annotated version of the function run, given in Figure 2(a). We see that region inference has determined that the closure bound to h will reside in the region ρ_3 and that the value bound to x will reside in the region ρ , which is deallocated after the function h is constructed.² The effect is that when the function work is called, which may perhaps trigger a reference-tracing collection, the value bound to h, which is live (and therefore part of the garbage-collection root set), will contain a pointer to an object that no longer exists.

Whereas the appearance of such dangling pointers is perfectly ok for a region-based memory management scheme that does not integrate with reference-tracing garbage collection (as long as the program itself does not dereference dangling pointers), a reference-tracing garbage collector will stumble over dangling pointers. An alternative region-annotated version of the program appears in Figure 2(b). This version of the program, which is the result of region inference based on our modified type system, does not introduce dangling pointers at runtime as the region ρ is live at least as long as the function h, which is enforced by ensuring that the type of the function h mentions the region ρ (in the arrow effect of the function type).

We now describe, informally, the mechanism that enforces region inference to assign the type $\text{unit} \xrightarrow{\epsilon:\{\rho_1, \rho_2, \rho\}} \text{unit}$ to the function h. First, notice that the type of h is the result type of an instance of the type scheme for the function o and that we must somehow capture, in the type scheme for o, that the type instance for the type variable γ specifies values that live in the region ρ .

²Notice that string concatenation (^) takes, besides the two argument strings, the region (ρ) into which the result is allocated.

We can capture this property by giving \circ the following type scheme:

$$\begin{aligned} & \forall \epsilon_0 \epsilon_1 \epsilon_2 \epsilon' \rho_0 \rho_1 \rho_2 \rho_3 \alpha \beta (\gamma : \epsilon'. \emptyset). \quad (2) \\ & ((\gamma \xrightarrow{\epsilon_2. \emptyset} \beta, \rho_2) \times (\alpha \xrightarrow{\epsilon_1. \emptyset} \gamma, \rho_1), \rho_0) \xrightarrow{\epsilon_0. \{\rho_0, \rho_3\}} (\alpha \xrightarrow{\epsilon. \{\epsilon_1, \epsilon_2, \epsilon', \rho_1, \rho_2\}} \beta, \rho_3) \end{aligned}$$

Compared to (1), the modified type scheme expresses a relationship between the type variable γ , through the *type variable descriptor* $\gamma : \epsilon'. \emptyset$, and the effect of the resulting function, which can be used to establish that regions appearing in the type of the instantiated type for γ must appear in the effect identified by the effect variable ϵ' . In the theoretical development presented in the following sections, this establishment will be implemented as part of the instance-of relation between region type-schemes and region-annotated types. Moreover, the typing rule for functions will ensure that type variables that appear in the type of a free variable occurring in the body of the function are associated with effect variables that are added to the arrow effect of the function type.

An alternative sound type scheme for \circ is the following:

$$\begin{aligned} & \forall \epsilon_0 \epsilon_1 \epsilon_2 \rho_0 \rho_1 \rho_2 \rho_3 \alpha \beta (\gamma : \epsilon. \{\epsilon_1, \epsilon_2, \rho_1, \rho_2\}). \quad (3) \\ & ((\gamma \xrightarrow{\epsilon_2. \emptyset} \beta, \rho_2) \times (\alpha \xrightarrow{\epsilon_1. \emptyset} \gamma, \rho_1), \rho_0) \xrightarrow{\epsilon_0. \{\rho_0, \rho_3\}} (\alpha \xrightarrow{\epsilon. \{\epsilon_1, \epsilon_2, \rho_1, \rho_2\}} \beta, \rho_3) \end{aligned}$$

Compared to (2), the alternative type scheme identifies the arrow effects associated with the result function type and the type variable γ , which is fine for the function \circ . Such an identification, which is perfectly sound, can be problematic (i.e., cause larger live ranges of regions), however, for type schemes with multiple type variables occurring free in the types of free identifiers of a function. On the positive side, however, the alternative type scheme can be expressed without introducing new *secondary* effect variables,³ which can be problematic for region inference.

Both of the above type schemes are sound candidates for the composition function \circ and both type schemes are accepted by the GC-safe region type system that we present in the next section. Distinguishing between the type system and the inference algorithm is vital here as it provides us with important implementation flexibility. For the theoretical development, we associate all bound type variables with arrow effects, whereas, for implementation purposes, we first identify the so-called *spurious type variables* for which we need to associate arrow effects. We return to the details of the inference algorithm and the concept of spurious type variables in Section 4.

3 A GC-SAFE REGION TYPE SYSTEM

In this section, we present a type system that provides us with the necessary guarantees for integrating region inference and reference-tracing garbage collection. Compared to the Tofte-Talpin type system [Tofte and Talpin 1997], the type system that we present ensures that no dangling pointers are introduced during evaluation even for programs that involve higher-order type-polymorphic functions. In the remainder of this section, we present a formal treatment for a small ML-like intermediate language extended with region annotations.

3.1 Regions and Effects

We assume a denumerably infinite set of *program variables*, ranged over by x and f . We also assume a denumerably infinite set of *region variables*, ranged over by ρ . Moreover, we assume a denumerably infinite set of *effect variables*, ranged over by ϵ . An *atomic effect*, ranged over by η , is either a region variable or an effect variable, and an *effect*, ranged over by φ , is a set of atomic effects. An *arrow effect*, written $\epsilon. \varphi$, and ranged over by ν , is a pair of an effect variable and an

³A secondary effect variable is an effect variable that does not appear syntactically as a handle on an arrow type constructor anywhere in the type-annotated version of the program.

effect. For simplicity, we do not distinguish between put- and get-effects. Finally, we assume a denumerably infinite set of *type variables*, ranged over by α .

A *type variable context*, ranged over by Ω (or Δ), is a finite map from type variables to arrow effects. When M and M' are two finite maps, we write $M + M'$ to denote the map with domain $\text{dom}(M) \cup \text{dom}(M')$ and values $(M + M')(x) = M'(x)$, if $x \in \text{dom}(M')$ and $M(x)$, otherwise.

3.2 Types and Type Schemes

The grammars for *types* (τ), *type and places* (μ), *type schemes* (σ), and *type schemes and places* (π) are as follows:

$$\begin{array}{ll} \mu ::= (\tau, \rho) \mid \alpha \mid \text{int} & \tau ::= \mu_1 \times \mu_2 \mid \mu_1 \xrightarrow{\epsilon \cdot \varphi} \mu_2 \\ \sigma ::= \forall \vec{\rho} \vec{\epsilon}. \sigma \mid \forall \Delta. \tau & \pi ::= (\sigma, \rho) \mid \mu \end{array}$$

For type schemes of the form $\forall \vec{\rho} \vec{\epsilon}. \sigma$, the region variables $\vec{\rho}$ and the effect variables $\vec{\epsilon}$ are considered *bound* in σ . In type schemes of the form $\forall \Delta. \tau$, which are novel, the type variables in $\text{dom}(\Delta)$ are considered bound in τ . Type schemes are considered identical up to renaming of bound variables.

Following the usual definition of bound variables, we define, for any kind of object o , the *free region variables* and the *free region and effect variables* of o , written $\text{frv}(o)$ and $\text{frev}(o)$, respectively. We write $\text{fv}(o)$ to denote the *free type, region, and effect variables* of o .

A type and place μ (or type τ) is *well-formed* with respect to a type variable context Ω , if the sentence $\Omega \vdash \mu$ (or $\Omega \vdash \tau$) can be derived from the following rules:

A type and place μ (or type scheme and place π) is *well-formed* with respect to a type variable context Ω , if the sentence $\Omega \vdash \mu$ (or $\Omega \vdash \pi$) can be derived from the following rules:

Well-formed types and type scheme and places

$$\begin{array}{c} \boxed{\Omega \vdash \mu \text{ and } \Omega \vdash \pi} \\ \frac{\alpha \in \text{dom}(\Omega)}{\Omega \vdash \alpha} \quad \frac{}{\Omega \vdash \text{int}} \quad \frac{\Omega \vdash \mu_1 \quad \Omega \vdash \mu_2}{\Omega \vdash (\mu_1 \times \mu_2, \rho)} \quad \frac{\Omega \vdash \mu_1 \quad \Omega \vdash \mu_2}{\Omega \vdash (\mu_1 \xrightarrow{\epsilon \cdot \varphi} \mu_2, \rho)} \\ \frac{\Omega \vdash (\sigma, \rho)}{\Omega \vdash (\forall \vec{\rho} \vec{\epsilon}. \sigma, \rho)} \quad \frac{\Omega + \Delta \vdash (\tau, \rho) \quad \text{dom}(\Delta) \cap \text{dom}(\Omega) = \emptyset}{\Omega \vdash (\forall \Delta. \tau, \rho)} \end{array}$$

Before we define the notion of substitution, we define a notion of *containment*, which expresses that a type and place μ (or type scheme and place π) is contained in an effect φ , under the assumption of a type variable context Ω . The relation is written $\Omega \vdash \mu : \varphi$ (or $\Omega \vdash \pi : \varphi$) and is defined according to the following rules:

Type (scheme) and place containment

$$\begin{array}{c} \boxed{\Omega \vdash \mu : \varphi \text{ and } \Omega \vdash \pi : \varphi} \\ \frac{\Omega \vdash \mu_1 : \varphi \quad \Omega \vdash \mu_2 : \varphi \quad \rho \in \varphi}{\Omega \vdash (\mu_1 \times \mu_2, \rho) : \varphi} \quad \frac{\varphi_0 \subseteq \varphi \quad \{\rho, \epsilon\} \subseteq \varphi}{\Omega \vdash \mu_1 : \varphi \quad \Omega \vdash \mu_2 : \varphi} \quad \frac{}{\Omega \vdash \text{int} : \varphi} \quad \frac{\text{frev}(\Omega(\alpha)) \subseteq \varphi}{\Omega \vdash \alpha : \varphi} \\ \frac{\Omega \vdash \sigma : \varphi \quad \rho \in \varphi \quad \{\vec{\rho} \vec{\epsilon}\} \cap \text{frev}(\Omega, \rho) = \emptyset}{\Omega \vdash (\forall \vec{\rho} \vec{\epsilon}. \sigma, \rho) : \varphi \setminus \{\vec{\rho} \vec{\epsilon}\}} \quad \frac{\Omega + \Delta \vdash (\tau, \rho) : \varphi \quad \text{dom}(\Delta) \cap \text{dom}(\Omega) = \emptyset}{\Omega \vdash (\forall \Delta. \tau, \rho) : \varphi} \end{array}$$

Containment implies well-formedness, which can be demonstrated by simple structural induction:

PROPOSITION 1 (CONTAINMENT IMPLIES WELL-FORMEDNESS). *Assume o is one of μ or π . If $\Omega \vdash o : \varphi$ then $\Omega \vdash o$.*

Both well-formedness and containability features context extensibility properties. Assume o is one of μ or π and $\text{dom}(\Omega) \cap \text{dom}(\Delta) = \emptyset$. If $\Omega \vdash o : \varphi$ then $\Omega + \Delta \vdash o : \varphi$. Moreover, if $\Omega \vdash o$ then $\Omega + \Delta \vdash o$. It is also straightforward to demonstrate an effect extensibility property for type

containment stating that, when o is one of μ or π , if $\Omega \vdash o : \varphi$ and $\varphi \subseteq \varphi'$ then $\Omega \vdash o : \varphi'$. Finally, the following effect containment property can be shown by simple structural induction:

PROPOSITION 2 (CONTAINMENT). *Assume o is one of μ or π . If $\Omega \vdash o : \varphi$ then $\text{frev}(o) \subseteq \varphi$.*

3.3 Substitutions

A *substitution* (S) is a triple (S^t, S^r, S^e) , where S^t is a *type substitution*, a finite map from type variables to type and places, S^r is a *region substitution*, a finite map from region variables to region variables, and S^e is an *effect substitution*, a finite map from effect variables to arrow effects. The effect of applying a substitution to a particular object is to carry out the three substitutions simultaneously to the three kinds of variables in the object (possibly by renaming of bound variables within the object to avoid capture) and acting as the identity outside of its domain. For effect sets and arrow effects, substitution is defined as follows [Tofte and Birkedal 2000], assuming $S = (S^t, S^r, S^e)$:

$$\begin{aligned} S(\varphi) &= \{S^t(\rho) \mid \rho \in \varphi\} \cup \{\eta \mid \exists \epsilon. \epsilon \in \varphi \wedge \eta \in \text{frev}(S^e(\epsilon))\} \\ S(\epsilon. \varphi) &= \epsilon'. (\varphi' \cup S(\varphi)), \text{ where } S^e(\epsilon) = \epsilon'. \varphi' \end{aligned}$$

Notice in particular, that when a substitution is applied to an effect φ , the result is also an effect.

Applying a substitution S to a type variable context Δ is defined only if $\text{dom}(S) \cap \text{dom}(\Delta) = \emptyset$:

$$S(\{\alpha_1 : v_1, \dots, \alpha_n : v_n\}) = \{\alpha_1 : S(v_1), \dots, \alpha_n : S(v_n)\}$$

For type schemes and for type-schemes and places, substitution is defined as follows, assuming that bound variables in type schemes have been renamed to avoid capture:

$$S(\forall \vec{\rho} \vec{\epsilon}. \sigma) = \forall \vec{\rho} \vec{\epsilon}. S(\sigma) \quad S(\forall \Delta. \tau) = \forall S(\Delta). S(\tau) \quad S(\sigma, \rho) = (S(\sigma), S(\rho))$$

It turns out that substitution is a monotone operation with respect to effects, which follows immediately from the definition of substitution on effects:

PROPOSITION 3 (SUBSTITUTION EFFECT MONOTONICITY). *If $\varphi \subseteq \varphi'$ then $S(\varphi) \subseteq S(\varphi')$, for any substitution S and effects φ and φ' .*

Another property that holds, which we shall call the *arrow-effect-substitution interchange property*, is that for any substitution S and arrow effect $\epsilon. \varphi$, we have $\text{frev}(S(\epsilon. \varphi)) = S(\{\epsilon\} \cup \varphi)$.

If $S = (S^t, S^r, S^e)$, we call S a *region-effect substitution* if $\text{dom}(S^t) = \emptyset$. Type containment is closed under region-effect substitutions, which also follows by straightforward induction:

PROPOSITION 4 (REGION-EFFECT SUBSTITUTION CLOSEDNESS). *Assume o is one of μ or π . If $\Omega \vdash o : \varphi$ and S is a region-effect substitution then $S(\Omega) \vdash S(o) : S(\varphi)$.*

For type containment to be closed under type substitutions, on the other hand, a substitution coverage requirement is needed. A type substitution S^t is *covered* by a type variable context Ω , through another type variable context Δ , written $\Omega \vdash S^t : \Delta$, if $\text{dom}(S^t) = \text{dom}(\Delta)$ and, for all $\alpha \in \text{dom}(S^t)$, we have $\Omega \vdash S^t(\alpha) : \text{frev}(\Delta(\alpha))$.

In connection with the notion of instantiation, which we shall define shortly, it is the notion of substitution coverage that ensures that the arrow effect associated with a bound type variable captures the free region and effect variables of the types instantiated for the type variable (which also holds transitively via the type containment relation.)

Provided the type substitution is properly covered, type containment is closed under type substitution, for which a proof by structural induction appears in Appendix A (auxiliary material):

PROPOSITION 5 (TYPE SUBSTITUTION CLOSEDNESS). *Assume o is one of μ or π . If $\Omega + \Delta \vdash o : \varphi$ and $\Omega \vdash S : \Delta$ then $\Omega \vdash S(o) : \varphi$.*

3.4 Instantiation

Given a type variable context Ω and a type scheme $\sigma = \forall \Delta. \tau'$ such that $\Omega \vdash \sigma$, a type τ is an *instance of σ* via a type substitution S^t , written $\Omega \vdash \sigma \geq \tau$ via S^t , if (1) $\Omega \vdash S^t : \Delta$ and (2) $S^t(\tau') = \tau$.

Given a type variable context Ω and a type scheme $\sigma = \forall \vec{\rho} \vec{e}. \sigma'$ such that $\Omega \vdash \sigma$, a type τ is an *instance of σ* via a substitution $S = (S^t, S^r, S^e)$, written $\Omega \vdash \sigma \geq \tau$ via S , if (1) $\text{dom}(S^r) = \{\vec{\rho}\}$ and $\text{dom}(S^e) = \{\vec{e}\}$ and (2) $\Omega \vdash S^e(S^r(\sigma')) \geq \tau$ via S^t .

When we are interested in only the region instance list, we write $\Omega \vdash \sigma \geq \tau$ via $\vec{\rho}$ to mean there exists a substitution $S = (S^t, S^r, S^e)$ such that $\Omega \vdash \sigma \geq \tau$ via S and $\text{rng}(S^r) = \{\vec{\rho}\}$.

It holds that if $\Omega \vdash \sigma \geq \tau$ via $\vec{\rho}$, for some σ, τ , and $\vec{\rho}$, and S is a region-effect substitution, then $S(\Omega) \vdash S(\sigma) \geq S(\tau)$ via $S(\vec{\rho})$. Moreover, if $\Omega + \Delta \vdash \sigma \geq \tau$ via $\vec{\rho}$, for some $\Omega, \Delta, \sigma, \tau$, and $\vec{\rho}$, if $\Omega \vdash S : \Delta$, then $\Omega \vdash S(\sigma) \geq S(\tau)$ via $\vec{\rho}$. These properties are corollaries of the following, more general, propositions, for which proofs appear in Appendix A:

PROPOSITION 6 (INSTANTIATION CLOSED UNDER REGION-EFFECT SUBSTITUTION). *If S is a region-effect substitution and $\Omega \vdash \sigma \geq \tau$ via S' then $S(\Omega) \vdash S(\sigma) \geq S(\tau)$ via S'' , where $S'' = (S \circ S') \downarrow \text{dom}(S')$.*

PROPOSITION 7 (INSTANTIATION CLOSED UNDER TYPE SUBSTITUTION). *If $\Omega + \Delta \vdash \sigma \geq \tau$ via S' and $\Omega \vdash S : \Delta$ then $\Omega \vdash S(\sigma) \geq S(\tau)$ via S'' , where $S'' = (S \circ S') \downarrow \text{dom}(S')$.*

3.5 The Role of Arrow Effects

We emphasise that function types are annotated with arrow effects $\epsilon. \varphi$ and not only with effects φ ; with arrow effects, we can allow for effects to grow (by applying substitutions) and we can make sure that if a non-region-annotated type is given two distinct region-annotations, then there exists a substitution, a *unifier*, that, when applied to the two types, will make the two resulting region-annotated types equal. This property is essential for the applied unification-based region inference algorithm [Tofte and Birkedal 1998], which we shall discuss further later.

Moreover, for each object that we deal with, when an effect variable appears free in the object, it is made explicit what effect it denotes, except when an effect variable appears free in an effect; in this case, however, we know that, due to an assumed transitivity of effects, the effect already includes the effect denoted by the included effect variable. It is for this reason that we annotate quantified type variables with arrow effects and not only with effect variables; we often (e.g., in the typing rules) need to know what effect the effect variable denotes. An alternative would be, in the typing rules, to keep track of the denotation of effect variables in an external *effect basis*, similarly to how effects are treated in the description of region inference [Tofte and Birkedal 1998, 2000]; making the effect basis explicit makes it straightforward to formulate certain well-formedness and consistency constraints on the effect variables and their denotations in the rules. For instance, if $\epsilon. \varphi$ and $\epsilon'. \varphi'$ are two arrow effects appearing in the derivation of some judgement, then $\epsilon = \epsilon'$ implies $\varphi = \varphi'$ (the basis is *functional*) and $\epsilon' \in \varphi$ implies $\varphi' \subseteq \varphi$ (the basis is *transitive*).

3.6 Terms

The grammars for *expressions* (e) and *values* (v) are as follows:

$$\begin{aligned} v & ::= d \mid \langle v_1, v_2 \rangle^\rho \mid \langle \lambda x. e \rangle^\rho \mid \langle \text{fun } f \ [\vec{\rho}] \ x = e \rangle^\rho \\ e & ::= v \mid x \mid \text{let } x = e_1 \text{ in } e_2 \mid e_1 \ e_2 \mid \lambda x. e \text{ at } \rho \mid \text{letregion } \rho \text{ in } e \\ & \quad \mid \text{fun } f \ [\vec{\rho}] \ x = e \text{ at } \rho \mid e \ [\vec{\rho}] \text{ at } \rho \mid (e_1, e_2) \text{ at } \rho \mid \#i \ e \end{aligned}$$

Values include unboxed integers (d), pairs, ordinary closures, and recursive function closures (which may also take regions as parameters). All values, except integers, are boxed and associated with distinguished regions. An expression can be a value, a variable, a let-expression, a function

Values

$$\boxed{\varphi \models v}$$

$$\varphi \models d \quad \frac{\varphi \models_v e \quad \rho \in \varphi}{\varphi \models \langle \lambda x.e \rangle^\rho} \quad \frac{\rho \in \varphi \quad \varphi \models v_1 \quad \varphi \models v_2}{\varphi \models \langle v_1, v_2 \rangle^\rho} \quad \frac{\rho \in \varphi \quad \varphi \models_v e \quad \{\vec{\rho}\} \cap \varphi = \emptyset}{\varphi \models \langle \text{fun } f [\vec{\rho}] x = e \rangle^\rho}$$

Expressions

$$\boxed{\varphi \models_v e}$$

$$\frac{\varphi \models v}{\varphi \models_v v} \quad \varphi \models_v x \quad \frac{\varphi \models_v e_1 \quad \varphi \models_v e_2}{\varphi \models_v (e_1, e_2) \text{ at } \rho} \quad \frac{\varphi \models_v e}{\varphi \models_v \#i e} \quad \frac{\varphi \models_v e}{\varphi \models_v \lambda x.e \text{ at } \rho}$$

$$\frac{\varphi \models_v e_1 \quad \varphi \models_v e_2}{\varphi \models_v e_1 e_2} \quad \frac{\varphi \models_v e \quad \{\vec{\rho}\} \cap \varphi = \emptyset}{\varphi \models_v \text{fun } f [\vec{\rho}] x = e \text{ at } \rho} \quad \frac{\varphi \models_v e}{\varphi \models_v e [\vec{\rho}] \text{ at } \rho}$$

$$\frac{\varphi \models_v e_1 \quad \varphi \models_v e_2}{\varphi \models_v \text{let } x = e_1 \text{ in } e_2} \quad \frac{\rho \notin \varphi \quad \varphi \models_v e}{\varphi \models_v \text{letregion } \rho \text{ in } e}$$

Fig. 3. Value containment.

application, a lambda-expression, a `letregion`-construct, a recursive function binding, an application of a recursive function to a list of region parameters, a pair-construct, and a pair-projection expression. Notice that allocating expressions are annotated with an `at`-specifier, which specifies in which region the value should be allocated. Notice also that expressions may contain values. A program does not contain values initially. During evaluation, however, variables in the program may be substituted with values, which is captured precisely by the small-step dynamic semantics that we shall define later.

In expressions of the forms `let $x = e_1$ in e` and `$\lambda x.e$ at ρ` , the variable x is *bound* in e . In expressions of the form `fun $f [\vec{\rho}] x = e$ at ρ` , the variables f , $\vec{\rho}$, and x are *bound* in e . Similarly for values. In expressions `letregion ρ in e` , the variable ρ is *bound* in e . We identify terms up to renaming of bound variables. The *free (program) variables* of an expression e is written $\text{fpv}(e)$.

A *type environment* (Γ) maps program variables to type schemes and places. It is *well-formed* in a type variable context Ω , written $\Omega \vdash \Gamma$, if $\Omega \vdash \pi$, for all type schemes and places $\pi \in \text{rng}(\Gamma)$.

3.7 Value Containment and GC Safety

To guarantee safety of garbage collection, we must ensure that no dangling pointers are introduced during evaluation, which is not guaranteed by the Tofte-Talpin region type system [Tofte and Talpin 1997]. The solution that we apply here is to add additional side conditions to the typing rules for functions that guarantee the absence of dangling pointers [Elsman 2003].

First, we define a notion of *value containment*; all values in an expression e are contained in a set of regions φ , if the sentence $\varphi \models_v e$ is derivable from the rules in Figure 3. It is straightforward to demonstrate that if $\varphi \models_v e$ and $\varphi \subseteq \varphi'$ then $\varphi' \models_v e$ (*value containment extensibility*). Moreover, for any substitution S , it follows that $S(\varphi) \models_v S(e)$. Finally, if $\varphi \models_v e$ and $\varphi \models v$ then $\varphi \models_v e[v/x]$ (*value containment substitution*).

We now introduce a *GC-Safety relation* G , which we shall use to strengthen the typing rules for functions to avoid dangling pointers during evaluation. The relation is derived from the side condition for functions suggested by Tofte and Talpin in [Tofte and Talpin 1993, page 50] and is parameterised over a type variable context Ω , an environment Γ , a function body e , a set of function parameters X , and the type scheme and place π of the function:

$$G(\Omega, \Gamma, e, X, \pi) = \text{frv}(\pi) \models_v e \wedge \forall y \in \text{fpv}(e) \setminus X. \Omega \vdash \Gamma(y) : \text{frev}(\pi) \quad (4)$$

The relation combines value containment for a function body and type containment for free variables of the function. It is closed under region-effect substitution, which follows immediately from the definition of garbage-collection safety, the property that value-containment is closed under substitution, and Proposition 4:

PROPOSITION 8 (GC-SAFETY RELATION CLOSED UNDER REGION-EFFECT SUBSTITUTION). *If $G(\Omega, \Gamma, e, X, \pi)$ and S is a region-effect substitution then $G(S(\Omega), S(\Gamma), S(e), X, S(\pi))$.*

Assuming that the type substitution is properly covered, the garbage-collection safety relation is also closed under type substitution, which follows from properties of value containment, type-containment effect-extensibility, and Proposition 5, with a detailed proof appearing in Appendix A:

PROPOSITION 9 (GC-SAFETY RELATION CLOSED UNDER TYPE SUBSTITUTION). *Assume $\Omega \vdash S : \Delta$. If $G(\Omega + \Delta, \Gamma, e, X, \pi)$ then $G(\Omega, S(\Gamma), e, X, S(\pi))$.*

Finally, we can establish that the garbage-collection safety relation is closed under value substitution, which follows from the definition of garbage-collection and Proposition 2, with a detailed proof appearing in Appendix A:

PROPOSITION 10 (GC-SAFETY RELATION CLOSED UNDER VALUE SUBSTITUTION). *If $x \notin X$ and $G(\Omega, \Gamma + \{x : \pi\}, e, X, \pi')$ and $\text{frv}(\pi) \models v$ and $\text{fpv}(v) = \emptyset$ then $G(\Omega, \Gamma, e[v/x], X, \pi')$.*

3.8 Typing Rules

The typing rules for values and expressions are mutually dependent and are shown in Figure 4. The typing rules for values allow inference of sentences of the form $\vdash v : \pi$, which states that “the value v has type scheme and place π ”. The typing rules for expressions allow inference of sentences of the form $\Omega, \Gamma \vdash e : \pi, \varphi$, which states that “in the type variable context Ω and in the type environment Γ , the expression e has type scheme and place π and effect φ ”.

There are a number of observations to be made about the typing rules. First, notice that the typing of values is specified without a variable environment, which, implicitly, specifies that well-typed values must be closed with respect to program variables. Moreover, values have no effect. Notice also that the typing rules for closures and for region- and effect-polymorphic function values specify that values within function bodies are contained in regions that appear in the type schemes for the functions (ensured using the value-containment judgement).

For lambda-expressions and region- and effect-polymorphic function expressions, gc-safety properties are specified using the gc-safety relation, which generalises the containment conditions specified in the corresponding value typing rules. Moreover, notice that there are two rules for typing region- and effect-polymorphic function expressions (and values), one that supports recursion (and even region- and effect-polymorphic recursion) and one that supports parameterisation of effects that are associated with quantified type variables. The duplication ensures that polymorphic recursion only quantify over region and effect variables that do not appear in type variable contexts that specify quantified type variables in the type scheme of the function.

For simplicity, the typing rule for let-bindings does not allow for generalisation.

3.9 Typing Properties

The typing rules are closed under region-effect substitution, which can be demonstrated by straightforward simultaneous induction over the typing judgments for values and expressions:

PROPOSITION 11 (TYPING CLOSED UNDER REGION-EFFECT SUBSTITUTION). *Assume S is a region-effect substitution.*

- (1) *If $\Omega, \Gamma \vdash e : \pi, \varphi$ then $S(\Omega), S(\Gamma) \vdash S(e) : S(\pi), S(\varphi)$.*

Values

$$\begin{array}{c}
\boxed{\vdash v : \pi} \\
\frac{}{\vdash d : \text{int}} \quad \frac{\{x : \mu_1\} \vdash e : \mu_2, \varphi \quad \vdash \mu}{\mu = (\mu_1 \xrightarrow{\epsilon, \varphi} \mu_2, \rho) \quad \text{frv}(\mu) \models_v e} \text{ [TVLAM]} \quad \frac{\mu = \mu_1 \times \mu_2 \quad \vdash v_1 : \mu_1 \quad \vdash v_2 : \mu_2}{\vdash \langle v_1, v_2 \rangle^\rho : (\mu, \rho)} \text{ [TVPAIR]} \\
\frac{\Delta, \{f : (\forall \vec{\rho} \vec{\epsilon}. \mu_1 \xrightarrow{\epsilon, \varphi} \mu_2, \rho), x : \mu_1\} \vdash e : \mu_2, \varphi \quad \vdash \pi \quad \text{frev}(\vec{\rho} \vec{\epsilon}) \cap \text{frev}(\Delta) = \emptyset \quad \text{frev}(\vec{\rho} \vec{\epsilon}) \cap \{\rho\} = \emptyset \quad \pi = (\forall \vec{\rho} \vec{\epsilon} \Delta. \mu_1 \xrightarrow{\epsilon, \varphi} \mu_2, \rho) \quad \text{frv}(\pi) \models_v e}{\vdash \langle \text{fun } f [\vec{\rho}] x = e \rangle^\rho : \pi} \text{ [TVREC]} \\
\frac{\Delta, \{x : \mu_1\} \vdash e : \mu_2, \varphi \quad \vdash \pi \quad \text{frv}(\pi) \models_v e \quad \text{frev}(\vec{\rho} \vec{\epsilon}) \cap \{\rho\} = \emptyset \quad \pi = (\forall \vec{\rho} \vec{\epsilon} \Delta. \mu_1 \xrightarrow{\epsilon, \varphi} \mu_2, \rho)}{\vdash \langle \text{fun } f [\vec{\rho}] x = e \rangle^\rho : \pi} \text{ [TVFUN]}
\end{array}$$

Expressions

$$\begin{array}{c}
\boxed{\Omega, \Gamma \vdash e : \pi, \varphi} \\
\frac{\vdash v : \pi}{\Omega, \Gamma \vdash v : \pi, \emptyset} \text{ [TEVAL]} \quad \frac{\varphi' \supseteq \varphi \quad \Omega, \Gamma \vdash e : \pi, \varphi}{\Omega, \Gamma \vdash e : \pi, \varphi'} \text{ [TESUB]} \quad \frac{\Gamma(x) = \pi}{\Omega, \Gamma \vdash x : \pi, \emptyset} \text{ [TEVAR]} \\
\frac{\Omega, \Gamma + \{x : \mu_1\} \vdash e : \mu_2, \varphi \quad \mu = (\mu_1 \xrightarrow{\epsilon, \varphi} \mu_2, \rho) \quad \Omega \vdash \mu \quad G(\Omega, \Gamma, e, \{x\}, \mu)}{\Omega, \Gamma \vdash \lambda x. e \text{ at } \rho : \mu, \{\rho\}} \text{ [TELAM]} \\
\frac{\Omega, \Gamma \vdash e : (\sigma, \rho'), \varphi \quad \Omega \vdash \sigma \geq \tau \text{ via } \vec{\rho} \quad \Omega \vdash \tau}{\Omega, \Gamma \vdash e [\vec{\rho}] \text{ at } \rho : (\tau, \rho), \varphi \cup \{\rho, \rho'\}} \text{ [TERAPP]} \quad \frac{\Omega, \Gamma \vdash e_1 : (\mu' \xrightarrow{\epsilon, \varphi_0} \mu, \rho), \varphi_1 \quad \Omega, \Gamma \vdash e_2 : \mu', \varphi_2 \quad \varphi = \{\epsilon, \rho\}}{\Omega, \Gamma \vdash e_1 e_2 : \mu, \varphi_0 \cup \varphi_1 \cup \varphi_2 \cup \varphi} \text{ [TEAPP]} \\
\frac{\Omega, \Gamma \vdash e_1 : \mu_1, \varphi_1 \quad \Omega, \Gamma \vdash e_2 : \mu_2, \varphi_2 \quad \varphi = \varphi_1 \cup \varphi_2 \cup \{\rho\}}{\Omega, \Gamma \vdash (e_1, e_2) \text{ at } \rho : (\mu_1 \times \mu_2, \rho), \varphi} \text{ [TEPAIR]} \quad \frac{i \in \{1, 2\} \quad \Omega, \Gamma \vdash e : (\mu_1 \times \mu_2, \rho), \varphi}{\Omega, \Gamma \vdash \#i e : \mu_i, \varphi \cup \{\rho\}} \text{ [TESEL]} \\
\frac{\Omega, \Gamma \vdash e : \mu, \varphi' \quad \varphi = \varphi' \setminus \{\rho, \vec{\epsilon}\} \quad \{\rho, \vec{\epsilon}\} \cap \text{frev}(\Omega, \Gamma, \mu) = \emptyset}{\Omega, \Gamma \vdash \text{letregion } \rho \text{ in } e : \mu, \varphi} \text{ [TEREG]} \quad \frac{\Omega, \Gamma \vdash e_1 : \pi, \varphi_1 \quad \Omega, \Gamma + \{x : \pi\} \vdash e_2 : \mu, \varphi_2}{\Omega, \Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \mu, \varphi_1 \cup \varphi_2} \text{ [TELET]} \\
\frac{\Omega + \Delta, \Gamma + \{x : \mu_1\} \vdash e : \mu_2, \varphi \quad \Omega \vdash \pi \quad G(\Omega, \Gamma, e, \{f, x\}, \pi) \quad (\text{dom}(\Delta) \cup \text{frev}(\vec{\rho} \vec{\epsilon})) \cap \text{fv}(\Omega, \Gamma, \rho) = \emptyset \quad \pi = (\forall \vec{\rho} \vec{\epsilon} \Delta. \mu_1 \xrightarrow{\epsilon, \varphi} \mu_2, \rho)}{\Omega, \Gamma \vdash \text{fun } f [\vec{\rho}] x = e \text{ at } \rho : \pi, \{\rho\}} \text{ [TEFUN]} \\
\frac{\Omega + \Delta, \Gamma + \{f : (\forall \vec{\rho} \vec{\epsilon}. \mu_1 \xrightarrow{\epsilon, \varphi} \mu_2, \rho), x : \mu_1\} \vdash e : \mu_2, \varphi \quad \Omega \vdash \pi \quad G(\Omega, \Gamma, e, \{f, x\}, \pi) \quad \text{frev}(\vec{\rho} \vec{\epsilon}) \cap \text{frev}(\Delta) = \emptyset \quad (\text{dom}(\Delta) \cup \text{frev}(\vec{\rho} \vec{\epsilon})) \cap \text{fv}(\Omega, \Gamma, \rho) = \emptyset \quad \pi = (\forall \vec{\rho} \vec{\epsilon} \Delta. \mu_1 \xrightarrow{\epsilon, \varphi} \mu_2, \rho)}{\Omega, \Gamma \vdash \text{fun } f [\vec{\rho}] x = e \text{ at } \rho : \pi, \{\rho\}}
\end{array}$$

Fig. 4. Typing rules for values and expressions.

(2) If $\vdash v : \pi$ then $\vdash S(v) : S(\pi)$.

The typing rules are also closed under type substitution provided the substitution is properly covered, which is demonstrated by induction on the typing derivation, with a detailed proof appearing in Appendix A:

PROPOSITION 12 (TYPING CLOSED UNDER TYPE SUBSTITUTION). *If $\Omega + \Delta, \Gamma \vdash e : \pi, \varphi$ and $\Omega \vdash S : \Delta$ then $\Omega, S(\Gamma) \vdash S(e) : S(\pi), S(\varphi)$.*

$$\begin{array}{ll}
E_\varphi & ::= [\cdot] & (\varphi = \emptyset) \\
& | \text{letregion } \rho \text{ in } E_{\varphi \setminus \{\rho\}} & (\rho \in \varphi) \\
& | E_\varphi e \mid v E_\varphi \mid E_\varphi [\vec{\rho}] \text{ at } \rho \mid \text{let } x = E_\varphi \text{ in } e \\
& | (E_\varphi, e) \text{ at } \rho \mid (v, E_\varphi) \text{ at } \rho \mid \#i E_\varphi \\
\iota & ::= d \mid \lambda x. e \text{ at } \rho \mid (v_1, v_2) \text{ at } \rho \mid \#1 \langle v_1, v_2 \rangle^\rho \mid \#2 \langle v_1, v_2 \rangle^\rho \\
& | \langle \lambda x. e \rangle^\rho v \mid \langle \text{fun } f [\vec{\rho}] x = e \rangle^\rho [\vec{\rho}'] \text{ at } \rho'
\end{array}$$

Fig. 5. The grammars for *evaluation contexts* (E) and *instructions* (ι).

Two other essential properties, which are demonstrated by induction over the typing judgments, is an environment extensibility property and a type-variable context extensibility property:

PROPOSITION 13 (ENVIRONMENT EXTENSIBILITY). *If $\Omega, \Gamma \vdash e : \pi, \varphi$ and $\text{dom}(\Gamma) \cap \text{dom}(\Gamma') = \emptyset$ then $\Omega, \Gamma + \Gamma' \vdash e : \pi, \varphi$.*

PROPOSITION 14 (TYPE-VARIABLE CONTEXT EXTENSIBILITY). *If $\Omega, \Gamma \vdash e : \pi, \varphi$ and $\text{dom}(\Omega) \cap \text{dom}(\Omega') = \emptyset$ then $\Omega + \Omega', \Gamma \vdash e : \pi, \varphi$.*

Central to the proofs is that the gc-safety relation is closed under environment extensibility and under type-variable context extensibility.

Typed values contain no free program variables and the free variables of typed expressions are captured by the environment, which are also demonstrated by induction over typing derivations:

PROPOSITION 15 (FREE VARIABLES). *If $\vdash v : \pi$ then $\text{fpv}(v) = \emptyset$. Moreover, if $\Omega, \Gamma \vdash e : \pi, \varphi$ then $\text{fpv}(e) \subseteq \text{dom}(\Gamma)$.*

A final, but essential, property is the following value substitution property:

PROPOSITION 16 (VALUE SUBSTITUTION). *If $\Omega, \Gamma + \{x : \pi\} \vdash e : \pi', \varphi$ and $\vdash v : \pi$ then $\Omega, \Gamma \vdash e[v/x] : \pi', \varphi$.*

PROOF. By induction on the derivation of $\Omega, \Gamma + \{x : \pi\} \vdash e : \pi', \varphi$. For the cases that involve the gc-safety relation, Proposition 15 and Proposition 10 are applied. See Appendix A for details. \square

3.10 A Small Step Dynamic Semantics

The dynamic semantics that we present is in the style of a contextual dynamic semantics [Morrisett 1995] and is similar to the semantics given by Helsen and Thiemann [Calcagno et al. 2002; Helsen and Thiemann 2001], although it differs in the way that inaccessibility to values in deallocated regions is modeled. Whereas Helsen and Thiemann “null out” references to deallocated regions (to avoid future access), our semantics keep track of a set of currently allocated regions and disallow access to regions that are not in this set.

The grammars for *evaluation contexts* (E) and *instructions* (ι) are shown in Figure 5. Contexts E_φ make explicit the set of regions φ bound by `letregion` constructs encapsulating context holes.

The evaluation rules are given in Figure 6 and consist of *allocation and deallocation rules*, *reduction rules*, and a *context rule*. The rules are of the form $e \xrightarrow{\varphi} e'$, which says that, given a set of allocated regions φ , the expression e reduces (in one step) to the expression e' . Next, the *evaluation relation* $\xrightarrow{\varphi^*}$ is defined as the least relation formed by the reflexive transitive closure of the relation $\xrightarrow{\varphi}$. We further define $e \Downarrow_\varphi v$ to mean $e \xrightarrow{\varphi^*} v$, and $e \Uparrow_\varphi$ to mean that there exists an infinite sequence, $e \xrightarrow{\varphi} e_1 \xrightarrow{\varphi} e_2 \xrightarrow{\varphi} \dots$.

Allocation and Deallocation

$$e \xrightarrow{\varphi} v$$

$$\begin{aligned} \lambda x. e \text{ at } \rho &\xrightarrow{\varphi \cup \{\rho\}} \langle \lambda x. e \rangle^\rho \quad [\text{LAM}] & (v_1, v_2) \text{ at } \rho &\xrightarrow{\varphi \cup \{\rho\}} \langle v_1, v_2 \rangle^\rho \quad [\text{PAIR}] \\ \text{fun } f \text{ [}\vec{\rho}\text{]} x = e \text{ at } \rho &\xrightarrow{\varphi \cup \{\rho\}} \langle \text{fun } f \text{ [}\vec{\rho}\text{]} x = e \rangle^\rho \quad [\text{FUN}] & \text{letregion } \rho \text{ in } v &\xrightarrow{\varphi} v \quad [\text{REG}] \end{aligned}$$

Reduction and Context

$$e \xrightarrow{\varphi} e'$$

$$\begin{aligned} \langle \lambda x. e \rangle^\rho v &\xrightarrow{\varphi \cup \{\rho\}} e[v/x] \quad [\text{APP}] & \text{let } x = v \text{ in } e &\xrightarrow{\varphi} e[v/x] \quad [\text{LET}] \\ \langle \text{fun } f \text{ [}\vec{\rho}\text{]} x = e \rangle^\rho [\vec{\rho}'] \text{ at } \rho' &\xrightarrow{\varphi \cup \{\rho\}} \lambda x. e[\vec{\rho}'/\vec{\rho}] [(\langle \text{fun } f \text{ [}\vec{\rho}\text{]} x = e \rangle^\rho) / f] \text{ at } \rho' \quad [\text{RAPP}] \\ \#1 \langle v_1, v_2 \rangle^\rho &\xrightarrow{\varphi \cup \{\rho\}} v_1 \quad [\text{SEL1}] & \frac{e \xrightarrow{\varphi' \cup \varphi} e' \quad \varphi \cap \varphi' = \emptyset \quad E_\varphi \neq [\cdot]}{E_\varphi[e] \xrightarrow{\varphi'} E_\varphi[e']} \quad [\text{CTX}] \\ \#2 \langle v_1, v_2 \rangle^\rho &\xrightarrow{\varphi \cup \{\rho\}} v_2 \quad [\text{SEL2}] \end{aligned}$$

Fig. 6. Evaluation rules.

3.11 Type Safety

The proof of type safety is based on well-known techniques for proving type safety for statically typed languages [Morrisett 1995; Wright and Felleisen 1994].

We first state a property saying that a well-typed expression is either a value or can be separated into an evaluation context and an instruction (shown by induction on the structure of e):

PROPOSITION 17 (UNIQUE DECOMPOSITION). *If $\vdash e : \pi, \varphi$, then either (1) e is a value, or (2) there exist a unique $E_{\varphi'}$, e' , and π' such that $e = E_{\varphi'}[e']$ and $\vdash e' : \pi', \varphi \cup \varphi'$ and e' is an instruction.*

A type preservation property (i.e., subject reduction) for the language, as well as progress and type soundness, can be stated as follows:

PROPOSITION 18 (TYPE PRESERVATION). *If $\vdash e : \pi, \varphi$ and $e \xrightarrow{\varphi} e'$ then $\vdash e' : \pi, \varphi$.*

PROOF. By induction on the derivation $e \xrightarrow{\varphi} e'$. Details are provided in Appendix A. \square

PROPOSITION 19 (PROGRESS). *If $\vdash e : \pi, \varphi$ then either e is a value or $e \xrightarrow{\varphi} e'$, for some e' .*

PROOF. If e is not a value, then by Proposition 17 there exist a unique $E_{\varphi'}$, ι , and π' such that $e = E_{\varphi'}[\iota]$ and $\vdash \iota : \pi', \varphi \cup \varphi'$. The remainder of the proof argues that $\iota \xrightarrow{\varphi \cup \varphi'} e_2$, for some e_2 , so that $E_{\varphi'}[\iota] \xrightarrow{\varphi} E_{\varphi'}[e_2]$ follows from [CTX] in Figure 6. Details are provided in Appendix A. \square

THEOREM 1 (TYPE SOUNDNESS). *If $\vdash e : \pi, \varphi$, then either $e \uparrow_\varphi$ or $e \Downarrow_\varphi v$ and $\vdash v : \pi, \varphi$, for some v .*

PROOF. By induction on the number of rewriting steps, using Propositions 18 and 19. \square

We now introduce the notion of *context containment*, written $\varphi \models_c e$, which expresses that when e can be written in the form $E_{\varphi'}[e']$, values in e' must be contained in the regions in the set $\varphi \cup \varphi'$, where φ' are regions on the stack represented by the evaluation context $E_{\varphi'}$. The definition of context containment is given in Figure 7. For well-typed programs, containment is preserved under evaluation, which is demonstrated by straightforward structural induction:

THEOREM 2 (CONTAINMENT). *If $\vdash e : \pi, \varphi$ and $\varphi \models_c e$ and $e \xrightarrow{\varphi} e'$ then $\varphi \models_c e'$.*

$$\begin{array}{c}
\frac{\varphi \models_c x \quad \frac{\varphi \models v}{\varphi \models_c v} \quad \frac{\rho \notin \varphi \quad \varphi \cup \{\rho\} \models_c e}{\varphi \models_c \text{letregion } \rho \text{ in } e} \quad \frac{\varphi \models_c e \quad \varphi \models_v e'}{\varphi \models_c \text{let } x = e \text{ in } e'} \quad \frac{\varphi \models_c e \quad \varphi \models_v e'}{\varphi \models_c e e'} \\
\frac{\varphi \models v \quad \varphi \models_c e}{\varphi \models_c v e} \quad \frac{\varphi \models_c e}{\varphi \models_c e [\vec{\rho}] \text{ at } \rho} \quad \frac{\varphi \models_c e \quad \varphi \models_v e'}{\varphi \models_c (e, e') \text{ at } \rho} \quad \frac{\varphi \models v \quad \varphi \models_c e}{\varphi \models_c (v, e) \text{ at } \rho} \quad \frac{\varphi \models_c e}{\varphi \models_c \#i e}
\end{array}$$

Fig. 7. Context containment.

Essentially, the containment theorem states that evaluation allocates only in regions that are either global or present on the region stack, represented by the evaluation context. Moreover, at any time during evaluation, live reachable values are stored in regions that are either global or present on the region stack. This last property is essential for reference-tracing garbage collection, which relies on the safety of dereferencing live reachable values [Morrisett et al. 1995]. In particular, the containment theorem allows for a reference-tracing garbage collector to be interleaved with evaluation (as captured by the small-step evaluation semantics).

4 IMPLEMENTATION

For practical purposes, it is desirable to identify a quantified type variable to be *spurious* if it either appears free in the type of identifiers occurring free in a function expression (but not in the type of the function) or occurs free in a type that is instantiated for another spurious type variable. In particular, it turns out that only spurious type variables need to be associated with arrow effects in type variable contexts, which, in general, leads to simpler region type schemes, while limiting the computational overhead of applying effect substitutions. We shall refer to a *spurious function* as one with spurious type variables in its inferred type scheme. It turns out that spurious functions occur only rarely in real programs. For example, the MLKit implementation of the Standard ML Basis Library [Gansner and Reppy 2004] contains only three spurious functions, which include the top-level composition function `o` and the functions `Option.compose` and `Option.mapPartial`.

The region type system presented in the previous section extends to other ML-language features, including references, exceptions, and algebraic datatypes.

4.1 Region Inference

Region inference takes as input a well-typed source program and returns a region annotated version of the program that is well-typed according to the region typing rules. A simple region inference algorithm stores all values in the global region ρ and associates all function arrows and quantified occurrences of spurious type variables with the arrow effect $\epsilon.\{\rho\}$, where ϵ is a global effect variable. It is straightforward to prove that this trivial region inference algorithm leads to well-typed region-annotated programs and works for all source programs that are well-typed according to a classic Hindley-Milner style type system.

A proper region-inference algorithm introduces regions locally and seek to quantify over region variables and effect variables in order to pass regions to functions at runtime and to make it possible to use functions in different contexts without necessarily having to keep function arguments and results alive as long as the function is alive. In order to guarantee an upper limit to the number of introduced region variables and effect variables (to ensure termination), region inference can be divided into two phases. Here, the first phase, called the *spreading* phase, adds distinct fresh region variables to all allocation points and distinct fresh effect variables to all function type arrows. The second phase, called the *fix-point* phase, runs repeatedly until a fix-point is found by unifying region types according to the requirement of the region type system and by abstracting over region

and effect variables, when possible, either by inserting `letregion` expressions or by abstracting over region variables and effect variables in `fun` expressions. The result is a well-typed region-annotated program. Implementing a proper region-inference algorithm for the region type system presented in the previous section differs from earlier work [Tofte and Birkedal 1998] by having to deal with spurious type variables and their associated arrow effects.

4.2 The MLKit

The MLKit is a Standard ML compiler that compiles programs to efficient native machine code for Linux and macOS [Elsman and Hallenberg 1995] and implements a number of techniques for refining the representations of regions [Birkedal et al. 1996; Tofte et al. 2004], including dividing regions into stack allocated (bounded) regions (also called *finite regions*) and heap allocated regions (also called *infinite regions*), which are subject to reference-tracing garbage collections.

The region type system presented in Section 3 is implemented in the MLKit in terms of a region-inference algorithm that deals properly with spurious type variables. The changes to the region inference algorithm are orthogonal to many of the later region-representation phases of the MLKit, including dropping of quantified parameter regions that are not stored into by a function and distinguishing between regions holding different types of values (for supporting tag-free representations of values of certain types).

We emphasise that the implementation changes enforced by the modified region-type system are of mandatory importance for ensuring soundness of integrating region-inference and reference-tracing garbage collection. The MLKit compiles all of Standard ML, including itself and the MLton compiler. With the implementation changes, the generated executable, resulting from compiling MLKit with MLKit, no longer fails during a minor collection. Moreover, compiling MLton with MLKit no longer results in an executable that fails under a major collection.

MLton and the MLKit are two very different compilers with different characteristics. Whereas MLton generates very compact (and often very efficient) executables, by featuring aggressive inlining and optimisation strategies, the MLKit features efficient recompilation and relative fast compilation for large programs. For instance, compiling the MLKit from scratch with MLKit itself takes 201 seconds (real time) whereas the same task takes 1039 seconds with MLton.⁴ Moreover, upon changes of source code, recompiling the MLKit with the MLKit compiler often takes less than 10 seconds.

In some cases, a spurious function can be rewritten as a non-spurious function (function composition cannot, unless it is given a less polymorphic type). Consider first the function `List.app` from the Standard ML Basis Library. This function has type scheme $\forall\alpha.(\alpha \rightarrow \text{unit}) \rightarrow \alpha \text{ list} \rightarrow \text{unit}$ with the following possible implementation:

```
fun app f = let fun loop nil = ()
                | loop (x::xs) = (f x ; loop xs)
in loop
end
```

Unfortunately, a Standard ML compiler based on algorithm W [Milner 1978] will give `app` the type scheme $\forall\alpha\beta.(\alpha \rightarrow \beta) \rightarrow \alpha \text{ list} \rightarrow \text{unit}$ and `loop` the type $\alpha \text{ list} \rightarrow \text{unit}$ (a module signature constraint may later constrain the type scheme of `app` to be less generic). Because `f` has type $\alpha \rightarrow \beta$ and occurs free in `loop`, β is inferred to be a spurious type variable. In general, the number

⁴All measurements are performed on a MacBook Pro (15-inch, 2016), 16GB RAM, 2.7 GHz Quad-Core Intel Core i7, running MLKit v4.7.2 and MLton 20210117 with option `-disable-pass deepFlatten`.

```

fun g (f : unit->'a)
  : unit->unit =
  op o
  let val x = f()
  in (fn x => (),
      fn () => x)
  end
val h =
  g (fn () => "oh" ^ "no")

```

(a)

```

fun g [ρ5, ρ6] f =
  letregion ρ3
  in op o[ρ6] let val x = f ()
                in (fn at ρ5 x => (),
                    fn at ρ5 v110 => x) at ρ3
                end
  end
val h =
  letregion ρ4
  in g[ρ1, ρ1] (fn at ρ4 () => op ^[ρ2] ("oh", "no"))
  end

```

(b)

Fig. 8. A problematic source code program featuring a dependency between two spurious type variables (a) and a sound region-annotated version of the program (b).

of inferred spurious type variables may be decreased by applying a type minimization algorithm [Bjørner 1994]. For the example with `app`, it suffices to constrain `f` to have type $\alpha \rightarrow \text{unit}$.

As a second example, the Standard ML Basis Library contains a function `Array.copy`, which copies elements from a generic source array (of type α array) into locations in a target array (also of type α array). One possible implementation of this function uses a local utility function `loop` that, assuming no overlaps, loops through the indexes of the source and at each index, fetches the corresponding value from the source array and updates the appropriate location in the target array. This local function will have type $\text{int} \rightarrow \text{unit}$, which means that the type variable α will be inferred to be spurious. In practice, the inference of α to be spurious will likely have little influence on region-inference for programs that use the `Array.copy` function. It is possible, however, also in this case, to modify the code slightly, by passing the source array as an additional argument to the `loop` function, in order to ensure that α is not considered spurious.

4.3 Tracking Spurious Type-Variable Dependencies

It may be enlightening to see how the type system tracks spurious type variable dependencies. Consider the program in Figure 8(a). This program is much similar to the program presented in the introduction, except that the spurious type variable bound by the composition function `o` is here not instantiated to a ground type immediately. Instead, it is instantiated to a new spurious type variable, which is bound by the function `g` with the type scheme $\forall \alpha. (\text{unit} \rightarrow \alpha) \rightarrow \text{unit} \rightarrow \text{unit}$.

There are a couple of interesting aspects about the region-annotated program in Figure 8(b). First, notice how the inference algorithm has arranged for the two intermediate functions (passed to the composition function `o`) to be stored in the same region ρ_5 , which is bound by (and passed to) the function `g`. This unification is due to the region inference algorithm unifying secondary quantified region and effect variables in type schemes, which is a central part of ensuring termination of region inference. Second, consider the region type scheme for the function `g`:

$$\forall \rho_5 \rho_6 \rho_7 \epsilon_1 \epsilon_2 \in \epsilon_4 (\alpha : \epsilon. \emptyset). (\text{unit} \xrightarrow{\epsilon_1. \emptyset} \alpha, \rho_7) \xrightarrow{\epsilon_2. \{\epsilon_1, \rho_7, \rho_5, \rho_6\}} (\text{unit} \xrightarrow{\epsilon_4. \{\epsilon, \rho_5\}} \text{unit}, \rho_6)$$

We see that α is inferred to be a spurious type variable and that it is associated with the arrow effect $\epsilon. \emptyset$. The reason α is inferred to be spurious is not because it appears free in the type of a variable captured in a closure but because it appears free in a type instantiated for another spurious type variable, namely that bound in the type scheme for the function `o`. Notice that the type scheme for `g` captures that the argument function is applied immediately (ϵ_1 appears in the

effect $\epsilon_2.\{\epsilon_1, \rho_7, \rho_5, \rho_6\}$. Here is the type instance of the type scheme, with $\epsilon'_1, \epsilon'_2, \epsilon', \epsilon'_4$ being fresh:

$$(\text{unit} \xrightarrow{\epsilon'_1.\emptyset} (\text{string}, \rho_2), \rho_4) \xrightarrow{\epsilon'_2.\{\epsilon'_1, \rho_4, \rho_1, \rho_2\}} (\text{unit} \xrightarrow{\epsilon'_4.\{\epsilon', \rho_1, \rho_2\}} \text{unit}, \rho_1)$$

We see that region ρ_4 does not occur in the type of the function returned by g , which perfectly aligns with the fact that the function passed to g , which is stored in ρ_4 , is applied immediately and not accessed again. For this reason, region inference can surround the call to g with a `letregion ρ_4 in ... end` construct. Contrary, the region type scheme for g captures the relationship between the spurious quantified type variable α and the capture of a value of type α in the returned closure through the associated effect variable ϵ , which occurs in the effect of the resulting function. As a consequence, the string "ohno" is rightfully forced into a global region (i.e., ρ_2).

5 BENCHMARKS

In this section, we report on the consequences of the type system changes for a variety of benchmark programs. We compare the benchmark programs using three different compilation strategies using the MLKit (v4.7.2) and a single compilation strategy using MLton 20210117 [Weeks 2006], a whole-program optimising Standard ML compiler, which serves to relate the performance of the code generated by the MLKit with the performance of a state-of-the-art compiler. We emphasise again that the MLKit and MLton are two very different compilers.

All benchmark programs are executed on a MacBook Pro (15-inch, 2016) with a 2.7GHz Intel Core i7 processor and 16GB of memory. Times reported are wall clock times and memory usage is measured using the macOS `/usr/bin/time` program. The benchmark programs span from micro-benchmarks such as `fib37` and `tak` (7 and 12 lines), which use only the runtime stack for allocation, to larger programs, such as `vliw` and `mlyacc` (3681 and 7385 lines), that solve real-world problems.

The three MLKit compilation strategies have identical compilation times and include the **rg** compilation strategy, which is based on the region type-system presented in this paper and which combines region-inference and reference-tracing garbage collection, the **rg-** compilation strategy, which is like **rg** but without taking spurious type variables into account (and which is therefore unsound), and, finally, the **r** compilation strategy, which is based alone on region-inference.

Figure 9 lists the benchmark programs and reports on how the type system changes influence the generated code for each of the benchmarks. Measurements are averages over 10 runs. The **real time** columns list the average execution time in seconds, annotated with relative standard deviations. For the **rss** and **gc #** columns, the relative standard deviations are less than two percent. There are a number of observations to make. First notice that, for many of the programs, the type system has no influence on the generated code (and thus on region live ranges) even in cases where many of the functions are spurious and when boxed types are instantiated for spurious type variables (column ∂). We also see that for programs that contain no spurious functions (column **fcns**), the type system changes have no influence on the generated code (column ∂). However, for certain programs containing spurious functions, even when there are no instantiations of boxed types for spurious type variables (column **inst**), the type system changes may have resulted in different generated code in terms of longer region live ranges (programs `b-hut`, `kbc`, `simple`, `zebra`, and `zern`). There are two reasons why generated code may be different in these cases. The first reason may be that the implementation identifies the effect variable associated with a spurious type variable with the effect variable associated with the function type for which the type variable appears free in the type of a free variable, as illustrated by (3). The second reason may be that the implementation unifies secondary effect variables, which may lead to unifying of unrelated effects.

Concerning execution times (the **real time** columns), we see that there are no significant differences between the execution times for the **rg** and **rg-** strategies, even for cases where the generated

Program	loc	fcns	inst	∂	real time (s)				rss (Mb \pm 1.2%)				gc #	
					rg	rg-	r	MLKit	rg	rg-	r	MLKit	rg	rg-
DLX	2841	2/149	0/690	✓	0.14 \pm 0%	0.15 \pm 4%	<u>0.12 \pm 6%</u>	0.40 \pm 5%	8	8	7	31	3	3
b-hut	1245	2/140	0/459	✓	0.67 \pm 1%	0.70 \pm 4%	0.63 \pm 1%	<u>0.14 \pm 3%</u>	5	5	169	3	471	471
fft	71	0/9	0/34		0.64 \pm 1%	0.66 \pm 2%	0.51 \pm 1%	<u>0.26 \pm 2%</u>	69	69	56	125	10	10
fib37	7	0/1	0/0		0.98 \pm 1%	0.98 \pm 1%	<u>0.21 \pm 2%</u>	0.85 \pm 1%	3	3	3	1	1	1
kbc	679	1/90	0/249	✓	0.21 \pm 2%	0.21 \pm 0%	0.19 \pm 0%	<u>0.07 \pm 0%</u>	10	10	10	2	10	10
lexgen	1322	0/108	0/531		0.74 \pm 3%	0.81 \pm 5%	0.62 \pm 2%	<u>0.43 \pm 3%</u>	15	15	67	18	109	109
life	202	0/35	0/146		0.44 \pm 5%	0.46 \pm 2%	0.43 \pm 3%	0.43 \pm 1%	3	3	14	2	58	58
logic	351	0/22	0/806		0.63 \pm 2%	0.64 \pm 1%	0.43 \pm 1%	<u>0.13 \pm 0%</u>	4	4	251	2	1843	1843
mandel	62	0/5	0/0		0.53 \pm 2%	0.53 \pm 2%	0.38 \pm 1%	<u>0.31 \pm 1%</u>	3	3	3	1	1	1
mlyacc	7385	10/966	5/3256	✓	0.36 \pm 7%	0.34 \pm 0%	<u>0.30 \pm 1%</u>	0.33 \pm 3%	18	15	115	10	29	28
mpuz	124	0/13	0/44		0.68 \pm 1%	0.66 \pm 1%	0.46 \pm 2%	<u>0.27 \pm 2%</u>	3	3	3	1	2	2
msort-r	119	0/14	0/27		0.69 \pm 1%	0.67 \pm 2%	<u>0.47 \pm 2%</u>	0.93 \pm 3%	116	116	97	577	16	16
msort	113	0/13	0/22		0.89 \pm 1%	0.89 \pm 1%	<u>0.54 \pm 2%</u>	0.93 \pm 1%	131	131	375	421	26	26
nucleic	3215	1/40	475/1273		0.34 \pm 2%	0.34 \pm 1%	0.33 \pm 2%	<u>0.17 \pm 0%</u>	5	5	231	3	645	645
prof	282	0/57	0/99		0.49 \pm 2%	0.48 \pm 1%	<u>0.38 \pm 1%</u>	0.42 \pm 2%	4	4	11	1	263	263
ratio	620	0/54	0/848		1.71 \pm 2%	1.67 \pm 0%	1.33 \pm 1%	<u>0.48 \pm 1%</u>	16	16	36	46	14	14
ray	529	1/48	0/120		0.69 \pm 2%	0.67 \pm 1%	0.64 \pm 1%	<u>0.25 \pm 2%</u>	14	14	14	15	12	12
simple	1053	15/327	0/448	✓	0.19 \pm 6%	0.16 \pm 7%	<u>0.13 \pm 9%</u>	0.28 \pm 4%	5	5	4	8	4	4
tak	12	0/2	0/0		2.09 \pm 1%	2.04 \pm 1%	<u>0.55 \pm 1%</u>	2.12 \pm 0%	3	3	3	1	1	1
tsp	493	0/26	0/19		0.14 \pm 5%	0.14 \pm 3%	<u>0.11 \pm 0%</u>	0.14 \pm 0%	11	11	6	11	7	7
vliw	3681	5/563	4/2133	✓	0.78 \pm 3%	0.73 \pm 2%	0.56 \pm 2%	<u>0.30 \pm 3%</u>	14	14	44	9	15	15
zebra	313	2/50	0/288	✓	1.58 \pm 2%	1.58 \pm 3%	1.15 \pm 0%	<u>0.45 \pm 1%</u>	3	3	122	1	1066	1504
zern	605	3/103	0/34	✓	0.80 \pm 1%	0.80 \pm 2%	0.52 \pm 1%	<u>0.30 \pm 2%</u>	6	5	5	11	4503	4503

Fig. 9. Benchmark programs. The second column (**loc**) lists the size of the benchmark in terms of lines of code, excluding Basis Library code. The third column (**fcns**) lists the number of spurious functions, relative to the total number of functions. The fourth column (**inst**) lists the number of times a spurious type variable is instantiated with a boxed type, relative to the total number of type variable instantiations. The fifth column (∂) indicates if the notion of spurious type variables made a difference to the generated target program. The next four columns (**real time**) list execution times in seconds for different benchmark compilation strategies. The next three columns (**rss**) list memory usage (in Mb) for the compilation strategies. Finally, the last two columns list the number of reference tracing garbage collections for the strategies **rg** and **rg-**.

code differs (due to different region live ranges). Notice also that for none of the benchmarks do we experience failures due to the possibility of dangling-pointers in the **rg-** compilation strategy. We also see that the **r** compilation strategy performs better than the **rg** and **rg-** strategies. Sometimes MLKit generates faster code than MLton, which is the case for DLX, msort-r, msort, simple, and tak), but, for most benchmarks, MLton outperforms the MLKit.

With respect to memory usage (the **rss** columns), we see that the **rg** and **rg-** compilation strategies have similar behavior. We also see that the **r** compilation strategy sometimes perform better (e.g., fft), which is due to its more compact (tag-free) value representation and the less-restrictive region type system (dangling pointers are permitted). Sometimes, however, reference-tracing garbage collection is essential, which is exemplified by the benchmarks b-hut, logic, nucleic, and zebra. We also see that the memory usage of MLton generated executables is often higher than the memory usage of the **rg** compilation strategy (we have not explored MLKit's and MLton's runtime flags for adjusting heap-to-live ranges, etc.)

Finally, from the `gc #` columns, we see that, across benchmarks, except zebra, the `rg` and `rg-` compilation strategies lead to executables that trigger similar numbers of garbage collections.

6 RELATED WORK

Most related to this work is the previous work on combining region inference and garbage collection in the MLKit [Hallenberg et al. 2002], the work on integrating region-based memory management and generational garbage collection [Elsman and Hallenberg 2020], and the previous work on guaranteeing the absence of dangling pointers for region-based memory management [Elsman 2003]. Compared to previous work, the present work does not aim at distinguishing between regions containing different types of values, but is concerned purely about establishing a sound foundation for integrating region inference and reference-tracing garbage collection. The region type system (and the region inference algorithm) presented in this paper integrates well with the techniques for typing regions. These techniques allow for a tag-free representation of pairs, triples, and references, which provides dramatic savings on allocated memory and execution time.

Another strand of related work is the large body of related work concerning general garbage collection techniques [Jones et al. 2011] and garbage collection techniques for functional languages, including [Doligez and Leroy 1993; Huelsbergen and Winterbottom 1998; Reppy 1994; Ueno and Ogori 2016]. Incremental, concurrent, and real-time garbage collection techniques for functional languages have recently obtained much attention. In particular, the presence of generations has been shown useful for collecting parts of the heap incrementally and in a concurrent and parallel fashion [Anderson 2010; Marlow and Peyton Jones 2011; Marlow et al. 2009]. We leave it to future work to investigate the use of regions and generations for supporting concurrency and parallelism.

There is also a series of proposals for tag-free garbage collection schemes [Aditya et al. 1994; Appel 1989; Goldberg 1991; Goldberg and Gloger 1992; Tolmach 1994] and nearly tag-free garbage collection schemes [Morrisett et al. 1996; Tarditi et al. 1996]. The partly tag-free garbage collection scheme supported by the region type system does not involve untagging of all values. In particular, unboxed objects (e.g., integers and booleans) are tagged in our system, which makes it possible to distinguish pointers from unboxed objects at runtime. However, the scheme allows for commonly used data structures, such as tuples, reals, and reference cells, to be untagged, which, as mentioned, can lead to significant time and memory savings, in particular because pairs and triples are used for the implementation of many dynamic data structures, including lists and trees.⁵

As other techniques that support full untagging, our technique does not involve traversing the runtime stack to determine types during garbage collection [Appel 1989; Goldberg 1991; Goldberg and Gloger 1992] or require special type information to be passed to functions at runtime [Tolmach 1994]. By requiring values in certain regions to be of the same kind, our approach has much in common with BIBOP (Big Bag Of Pages) systems, with regions as pages [Hanson 1980].

Another body of related work investigates the notion of escape analysis for improving stack allocation in garbage collected systems [Blanchet 1998; Salagnac et al. 2005]. Region inference and MLKit's polymorphic multiplicity analysis [Birkedal et al. 1996] allow more objects to be stack allocated than traditional escape analyses, which allow only local, non-escaping values to be stack allocated. Other work investigates the use of static prediction techniques and linear typing for inferring heap space usage [Jost et al. 2010].

Cyclone [Swamy et al. 2006] is a region-based type-safe C dialect, for which, the programmer can decide if an object should reside in the GC heap or in a region. Cyclone is constructed to disallow program code to dereference dangling pointers. For the GC heap, Cyclone uses a conservative

⁵The scheme works well together with support for unboxed data constructors, such as `cons (: :)`, which, for instance, leads to a compact representation of linked lists [Elsman 1998].

reference-tracing collector and no guarantee is given that it does not trace dangling pointers (safety is ensured by the collector being conservative). Another region-based language is Gay and Aiken's RC system, which features limited explicit regions for C, combined with reference counting of regions [Gay and Aiken 2001]. A modern language for system programming is Rust [Jung et al. 2017], which is based on ownership types for controlling the use of resources, including memory [Aldrich et al. 2002]. Ownership types are also used for real-time implementations of Java [Boyapati et al. 2003]. None of the above systems are combined with techniques for reference-tracing garbage collection of each individual region (Cyclone allows values to be stored in the global garbage collected heap region, but other regions are not collected using reference-tracing collection). Ownership types also lead to problems with constructing cyclic data structures, which are straightforward to work with in effect-based systems.

Also, Aiken et al. [Aiken et al. 1995] show how region inference may be improved for some programs by removing the constraints of the stack discipline, which may cause a garbage collector to run less often. Other work in this area includes [Fluet et al. 2006], which removes the constraints of the region stack discipline for an intermediate language using a linear type system.

7 CONCLUSION AND FUTURE WORK

We have identified and fixed a soundness problem with combining region inference and reference-tracing garbage collection. The solution involves associating so-called spurious type variables with effect sets and tracking effect dependencies to ensure that no dangling pointers appear during evaluation of a program. The work thus justifies earlier work by (1) suggesting how the unsafe type system is modified into a sound type system and (2) demonstrating that the necessary modifications to the region type system have little influence on the generated code (and thus on previous reported results on combining region inference and reference-tracing garbage collection), although they are essential for sound execution (MLKit now compiles itself and MLton without the generated compilers failing during reference-tracing garbage collections).

There are multiple paths of relevant future work. Whereas the type system presented in this paper has been proven sound on paper, we do not have a mechanised version of the proof, which would be a major engineering task. We consider efforts in this direction as possible future work. Another possibility for future work is on allowing programmers to interfere with region inference by being explicit about regions and effects in types and expressions. Finally, a possibility for future work would be to improve instruction selection and optimisations of MLKit programs to match the performance of MLton executables in more cases.

From a sustainability point-of-view, region inference may limit the memory footprint of programs as garbage collections can occur less frequently if a part of the heap is managed by explicit memory allocation and deallocation. Future work may investigate this path in more details.

DATA AVAILABILITY STATEMENT

An artifact demonstrating the results of Figure 9 is archived in Zenodo [Elsman 2023].

REFERENCES

- Shail Aditya, Christine H. Flood, and James E. Hicks. 1994. Garbage Collection for Strongly-Typed Languages Using Run-Time Type Reconstruction. In *Proceedings of the 1994 ACM Conference on LISP and Functional Programming* (Orlando, Florida, USA) (LFP '94). Association for Computing Machinery, New York, NY, USA, 12–23. <https://doi.org/10.1145/182409.182414>
- Alexander Aiken, Manuel Fähndrich, and Raph Levien. 1995. Better Static Memory Management: Improving Region-Based Analysis of Higher-Order Languages. In *Proceedings of the ACM SIGPLAN 1995 Conference on Programming Language Design and Implementation* (La Jolla, California, USA) (PLDI '95). Association for Computing Machinery, New York, NY, USA, 174–185. <https://doi.org/10.1145/207110.207137>

- Jonathan Aldrich, Valentin Kostadinov, and Craig Chambers. 2002. Alias Annotations for Program Understanding. In *Proceedings of the 17th ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications* (Seattle, Washington, USA) (*OOPSLA '02*). Association for Computing Machinery, New York, NY, USA, 311–330. <https://doi.org/10.1145/582419.582448>
- Todd A. Anderson. 2010. Optimizations in a Private Nursery-Based Garbage Collector. In *Proceedings of the 2010 International Symposium on Memory Management* (Toronto, Ontario, Canada) (*ISMM '10*). Association for Computing Machinery, New York, NY, USA, 21–30. <https://doi.org/10.1145/1806651.1806655>
- Andrew W. Appel. 1989. Runtime tags aren't necessary. *Lisp and Symbolic Computation* 2 (1989), 153–162. <https://doi.org/10.1007/BF01811537>
- Lars Birkedal, Mads Tofte, and Magnus Vejlstrup. 1996. From Region Inference to von Neumann Machines via Region Representation Inference. In *Proceedings of the 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (St. Petersburg Beach, Florida, USA) (*POPL '96*). Association for Computing Machinery, New York, NY, USA, 171–183. <https://doi.org/10.1145/237721.237771>
- Nikolaj Skallerud Bjørner. 1994. Minimal Typing Derivations. In *ACM SIGPLAN Workshop on ML and its Applications*. 120–126.
- Bruno Blanchet. 1998. Escape Analysis: Correctness Proof, Implementation and Experimental Results. In *Proceedings of the 25th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (San Diego, California, USA) (*POPL '98*). Association for Computing Machinery, New York, NY, USA, 25–37. <https://doi.org/10.1145/268946.268949>
- Chandrasekhar Boyapati, Alexandru Salcianu, William Beebe, and Martin Rinard. 2003. Ownership Types for Safe Region-Based Memory Management in Real-Time Java. In *Proceedings of the ACM SIGPLAN 2003 Conference on Programming Language Design and Implementation* (San Diego, California, USA) (*PLDI '03*). Association for Computing Machinery, New York, NY, USA, 324–337. <https://doi.org/10.1145/781131.781168>
- Christiano Calcagno, Simon Helsen, and Peter Thiermann. 2002. Syntactic Type Soundness Results for the Region Calculus. *Inf. Comput.* 173, 2 (mar 2002), 199–221. <https://doi.org/10.1006/inco.2001.3112>
- Damien Doligez and Xavier Leroy. 1993. A Concurrent, Generational Garbage Collector for a Multithreaded Implementation of ML. In *Proceedings of the 20th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Charleston, South Carolina, USA) (*POPL '93*). Association for Computing Machinery, New York, NY, USA, 113–123. <https://doi.org/10.1145/158511.158611>
- Martin Elsmann. 1998. Polymorphic Equality—No Tags Required. In *Second International Workshop on Types in Compilation*. <https://doi.org/10.1007/BFb0055516>
- Martin Elsmann. 2003. Garbage Collection Safety for Region-Based Memory Management. In *Proceedings of the 2003 ACM SIGPLAN International Workshop on Types in Languages Design and Implementation* (New Orleans, Louisiana, USA) (*TLDI '03*). Association for Computing Machinery, New York, NY, USA, 123–134. <https://doi.org/10.1145/604174.604190>
- Martin Elsmann. 2023. Artifact for the PLDI 2023 paper: Garbage-Collection Safety for Region-Based Type-Polymorphic Programs. Zenodo. <https://doi.org/10.5281/zenodo.7803910>
- Martin Elsmann and Niels Hallenberg. 1995. An Optimizing Backend for the ML Kit Using a Stack of Regions. Student Project 95-7-8, University of Copenhagen (DIKU).
- Martin Elsmann and Niels Hallenberg. 2020. On the Effects of Integrating Region-Based Memory Management and Generational Garbage Collection in ML. In *Practical Aspects of Declarative Languages (PADL '20)*. Springer International Publishing, 95–112. https://doi.org/10.1007/978-3-030-39197-3_7
- Martin Elsmann and Niels Hallenberg. 2021. Integrating region memory management and tag-free generational garbage collection. *Journal of Functional Programming* 31 (2021), e4. <https://doi.org/10.1017/S0956796821000010>
- Matthew Fluet, Greg Morrisett, and Amal Ahmed. 2006. Linear Regions Are All You Need. In *Proceedings of the 15th European Conference on Programming Languages and Systems* (Vienna, Austria) (*ESOP'06*). Springer-Verlag, Berlin, Heidelberg, 7–21. https://doi.org/10.1007/11693024_2
- Emden R. Gansner and John H. Reppy. 2004. *The Standard ML Basis Library*. Cambridge University Press. <https://doi.org/10.1017/CBO9780511546846>
- David Gay and Alex Aiken. 2001. Language Support for Regions. In *Proceedings of the ACM SIGPLAN 2001 Conference on Programming Language Design and Implementation* (Snowbird, Utah, USA) (*PLDI '01*). Association for Computing Machinery, New York, NY, USA, 70–80. <https://doi.org/10.1145/378795.378815>
- Benjamin Goldberg. 1991. Tag-Free Garbage Collection for Strongly Typed Programming Languages. In *Proceedings of the ACM SIGPLAN 1991 Conference on Programming Language Design and Implementation* (Toronto, Ontario, Canada) (*PLDI '91*). Association for Computing Machinery, New York, NY, USA, 165–176. <https://doi.org/10.1145/113445.113460>
- Benjamin Goldberg and Michael Gloger. 1992. Polymorphic Type Reconstruction for Garbage Collection without Tags. In *Proceedings of the 1992 ACM Conference on LISP and Functional Programming* (San Francisco, California, USA) (*LFP '92*). Association for Computing Machinery, New York, NY, USA, 53–65. <https://doi.org/10.1145/141471.141504>

- Niels Hallenberg, Martin Elsman, and Mads Tofte. 2002. Combining Region Inference and Garbage Collection. In *Proceedings of the ACM SIGPLAN 2002 Conference on Programming Language Design and Implementation* (Berlin, Germany) (PLDI '02). Association for Computing Machinery, New York, NY, USA, 141–152. <https://doi.org/10.1145/512529.512547>
- David R. Hanson. 1980. A portable storage management system for the Icon programming language. *Software—Practice and Experience* 10 (1980), 489–500. <https://doi.org/10.1002/spe.4380100607>
- Simon Helsen and Peter Thiemann. 2001. Syntactic Type Soundness for the Region Calculus. *Electronic Notes in Theoretical Computer Science* 41, 3 (2001), 1–19. [https://doi.org/10.1016/S1571-0661\(04\)80870-3](https://doi.org/10.1016/S1571-0661(04)80870-3) HOOTS 2000, 4th International Workshop on Higher Order Operational Techniques in Semantics (Satellite to PLI 2000).
- Lorenz Huelshberger and Phil Winterbottom. 1998. Very Concurrent Mark-&-Sweep Garbage Collection without Fine-Grain Synchronization. In *Proceedings of the 1st International Symposium on Memory Management* (Vancouver, British Columbia, Canada) (ISMM '98). Association for Computing Machinery, New York, NY, USA, 166–175. <https://doi.org/10.1145/286860.286878>
- Richard Jones, Antony Hosking, and Eliot Moss. 2011. *The Garbage Collection Handbook: The Art of Automatic Memory Management*. Chapman & Hall/CRC.
- Steffen Jost, Kevin Hammond, Hans-Wolfgang Loidl, and Martin Hofmann. 2010. Static Determination of Quantitative Resource Usage for Higher-Order Programs. In *Proceedings of the 37th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Madrid, Spain) (POPL '10). Association for Computing Machinery, New York, NY, USA, 223–236. <https://doi.org/10.1145/1706299.1706327>
- Ralf Jung, Jacques-Henri Jourdan, Robbert Krebbers, and Derek Dreyer. 2017. RustBelt: Securing the Foundations of the Rust Programming Language. *Proc. ACM Program. Lang.* 2, POPL, Article 66 (dec 2017), 34 pages. <https://doi.org/10.1145/3158154>
- Simon Marlow and Simon Peyton Jones. 2011. Multicore Garbage Collection with Local Heaps. In *Proceedings of the International Symposium on Memory Management* (San Jose, California, USA) (ISMM '11). Association for Computing Machinery, New York, NY, USA, 21–32. <https://doi.org/10.1145/1993478.1993482>
- Simon Marlow, Simon Peyton Jones, and Satnam Singh. 2009. Runtime Support for Multicore Haskell. In *Proceedings of the 14th ACM SIGPLAN International Conference on Functional Programming* (Edinburgh, Scotland) (ICFP '09). Association for Computing Machinery, New York, NY, USA, 65–78. <https://doi.org/10.1145/1596550.1596563>
- Robin Milner. 1978. A theory of type polymorphism in programming. *J. Comput. System Sci.* 17, 3 (1978), 348–375. [https://doi.org/10.1016/0022-0000\(78\)90014-4](https://doi.org/10.1016/0022-0000(78)90014-4)
- Greg Morrisett. 1995. *Compiling with Types*. Ph.D. Dissertation. School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213.
- Greg Morrisett, Matthias Felleisen, and Robert Harper. 1995. Abstract Models of Memory Management. In *Proceedings of the Seventh International Conference on Functional Programming Languages and Computer Architecture* (La Jolla, California, USA) (FPCA '95). Association for Computing Machinery, New York, NY, USA, 66–77. <https://doi.org/10.1145/224164.224182>
- Greg Morrisett, David Tarditi, Perry Cheng, Chris Stone, Robert Harper, and Peter Lee. 1996. The TIL/ML Compiler: Performance and Safety through Types. In *Workshop on Compiler Support for Systems Software*, Tucson..
- John H. Reppy. 1994. *A High-performance Garbage Collector for Standard ML*. Technical Report. AT&T Bell Laboratories.
- Guillaume Salagnac, Chaker Nakhli, Christophe Rippert, and Sergio Yovine. 2006. Efficient Region-Based Memory Management for Resource-limited Real-Time Embedded Systems.. In *Workshop on Implementation, Compilation, Optimization of Object-Oriented Languages, Programs and Systems*.
- G. Salagnac, S. Yovine, and D. Garbervetsky. 2005. Fast Escape Analysis for Region-Based Memory Management. *Electron. Notes Theor. Comput. Sci.* 131 (may 2005), 99–110. <https://doi.org/10.1016/j.entcs.2005.01.026>
- Nikhil Swamy, Michael Hicks, Greg Morrisett, Dan Grossman, and Trevor Jim. 2006. Safe manual memory management in Cyclone. *Science of Computer Programming* 62, 2 (2006), 122–144. <https://doi.org/10.1016/j.scico.2006.02.003> Special Issue: Five perspectives on modern memory management - Systems, hardware and theory.
- D. Tarditi, G. Morrisett, P. Cheng, C. Stone, R. Harper, and P. Lee. 1996. TIL: A Type-Directed Optimizing Compiler for ML. In *Proceedings of the ACM SIGPLAN 1996 Conference on Programming Language Design and Implementation* (Philadelphia, Pennsylvania, USA) (PLDI '96). Association for Computing Machinery, New York, NY, USA, 181–192. <https://doi.org/10.1145/231379.231414>
- Mads Tofte and Lars Birkedal. 1998. A Region Inference Algorithm. *ACM Trans. Program. Lang. Syst.* 20, 4 (jul 1998), 724–767. <https://doi.org/10.1145/291891.291894>
- Mads Tofte and Lars Birkedal. 2000. Unification and Polymorphism in Region Inference. *Proof, Language, and Interaction. Essays in Honour of Robin Milner* (May 2000). (25 pages).
- Mads Tofte, Lars Birkedal, Martin Elsman, and Niels Hallenberg. 2004. A Retrospective on Region-Based Memory Management. *Higher-Order and Symbolic Computation* 17, 3 (01 Sep 2004), 245–265. <https://doi.org/10.1023/B:LISP.0000029446.78563.a4>

- Mads Tofte, Lars Birkedal, Martin Elsman, Niels Hallenberg, Tommy Højfeldt Olesen, and Peter Sestoft. 2022. *Programming with Regions in the MLKit (Revised for Version 4.7.2)*. Technical Report. Department of Computer Science, University of Copenhagen, Denmark.
- Mads Tofte and Jean-Pierre Talpin. 1993. *A Theory of Stack Allocation in Polymorphically Typed Languages*. Technical Report DIKU-report 93/15. Department of Computer Science, University of Copenhagen.
- Mads Tofte and Jean-Pierre Talpin. 1997. Region-Based Memory Management. *Information and Computation* 132, 2 (1997), 109–176.
- Andrew Tolmach. 1994. Tag-Free Garbage Collection Using Explicit Type Parameters. In *Proceedings of the 1994 ACM Conference on LISP and Functional Programming (Orlando, Florida, USA) (LFP '94)*. Association for Computing Machinery, New York, NY, USA, 1–11. <https://doi.org/10.1145/182409.182411>
- Katsuhiko Ueno and Atsushi Ohori. 2016. A Fully Concurrent Garbage Collector for Functional Programs on Multicore Processors. In *Proceedings of the 21st ACM SIGPLAN International Conference on Functional Programming (Nara, Japan) (ICFP 2016)*. Association for Computing Machinery, New York, NY, USA, 421–433. <https://doi.org/10.1145/2951913.2951944>
- Stephen Weeks. 2006. Whole-Program Compilation in MLton. In *Proceedings of the 2006 Workshop on ML (Portland, Oregon, USA) (ML '06)*. Association for Computing Machinery, New York, NY, USA, 1. <https://doi.org/10.1145/1159876.1159877>
- A.K. Wright and M. Felleisen. 1994. A Syntactic Approach to Type Soundness. *Inf. Comput.* 115, 1 (nov 1994), 38–94. <https://doi.org/10.1006/inco.1994.1093>

Received 2022-11-10; accepted 2023-03-31