

Linear Regions Are All You Need

Matthew Fluet¹, Greg Morrisett², and Amal Ahmed²

¹ Cornell University
Ithaca, NY

`fluet@cs.cornell.edu`

² Harvard University
Cambridge, MA

`greg@eecs.harvard.edu`, `amal@eecs.harvard.edu`

Abstract. The type-and-effects system of the Tofte-Talpin region calculus makes it possible to safely reclaim objects without a garbage collector. However, it requires that regions have last-in-first-out (LIFO) lifetimes following the block structure of the language. We introduce λ^{rgnUL} , a core calculus that is powerful enough to encode Tofte-Talpin-like languages, and that eliminates the LIFO restriction. The target language has an extremely simple, substructural type system. To prove the power of the language, we sketch how Tofte-Talpin-style regions, as well as the first-class dynamic regions and unique pointers of the Cyclone programming language can be encoded in λ^{rgnUL} .

1 Introduction

Most type-safe languages rely upon a garbage collector to reclaim storage safely. But there are domains, such as device drivers and embedded systems, where today's garbage collection algorithms result in unacceptable space or latency overheads. In these settings, programmers have been forced to use languages, like C, where memory management can be tailored to the application, but where the lack of type-safety has led to numerous bugs. To address these concerns, we have been developing Cyclone [1], a type-safe dialect of C that is intended to give programmers as much control over memory management as possible while retaining strong, static typing.

The initial design of Cyclone was based upon the region type system of Tofte and Talpin [2]. Data are allocated within lexically-scoped *regions* and all of the objects in a region are deallocated at the end of the region's scope. Unfortunately, the last-in-first-out (LIFO) lifetimes of lexically-scoped regions place severe restrictions on when data can be effectively reclaimed, and we found many programs that resulted in (unbounded) leaks when compared to a garbage collected implementation.

To address these concerns, we added a number of new features to Cyclone, including *dynamic regions* and *unique pointers* that provide more control over memory management. Dynamic regions are not restricted to LIFO lifetimes and can be treated as first-class objects. They are particularly well suited for iterative

computations, CPS-based computations, and event-based servers where lexical regions do not suffice. Unique pointers are essentially lightweight, dynamic regions that hold exactly one object. To ensure soundness, both dynamic regions and unique pointers depend upon a notion of *linear capabilities* which must be carefully threaded through a program. To alleviate this tedium, Cyclone provides convenient mechanisms to temporarily “open” a dynamic region or unique pointer and treat it as if it were in a freshly allocated, lexically-scoped region.

The efficacy of these new memory management features was detailed in previous papers [3, 4], where we analyzed a range of applications, including a streaming media server, a space-conscious web server, and a Scheme runtime system with a copying garbage collector. And while the soundness of Cyclone’s lexical regions and type-and-effects system has been established [5, 6], a model that justifies the soundness of the new features has eluded our grasp, due to sheer complexity.

Therefore, the goal of this work is to provide a simple model where we can easily *encode* the key features of Cyclone in a uniform target language for which type soundness may be easily established. The first step of our encoding was detailed in a previous paper [6], where we gave a translation from a type-and-effects, region-based language to a monadic variant of System F called F^{RGN} . This calculus is summarized in Section 2. The meat of this paper picks up where this translation left off by further translating F^{RGN} to a substructural polymorphic lambda calculus where the internals of the indexed monad are exposed (Section 3). The target language and translation are extremely simple, yielding a relatively straightforward proof of soundness for lexically scoped regions. Then, in Section 5, we sketch how the features in the target language allow us to encode Cyclone’s dynamic regions and unique pointers, as well as their interactions with lexically-scoped regions. Throughout, it is our intention that the target calculus serve as a compiler intermediate language and as vehicle for formal reasoning, not as a high-level programming language.

2 Source Calculus: F^{RGN}

Launchbury and Peyton Jones introduced the **ST** monad to encapsulate stateful computations within the pure functional language Haskell [7]. Three key insights give rise to a safe and efficient implementation of stateful computations. First, a stateful computation is represented as a *store transformer*, a description of commands to be applied to an initial store to yield a final store. Second, the store can not be duplicated, because the state type is opaque and all primitive store transformers use the store in a single-threaded manner; hence, a stateful computation can update the store in place. Third, parametric polymorphism can be used to safely encapsulate and run a stateful computation.

All of these insights can be carried over to the region case, where we interpret stores as stacks of regions. We introduce the types and operations associated with the **rgn** monad:

$$\tau ::= \dots \mid \text{rgn } s \tau \mid \text{ref } s \tau \mid \text{hnd } s \mid \text{pf } (s_1 \leq s_2)$$

```

return :  $\forall \varsigma. \forall \alpha. \alpha \rightarrow \text{rgn } \varsigma \alpha$ 
then :  $\forall \varsigma. \forall \alpha, \beta. \text{rgn } \varsigma \alpha \rightarrow (\alpha \rightarrow \text{rgn } \varsigma \beta) \rightarrow \text{rgn } \varsigma \beta$ 
new :  $\forall \varsigma. \forall \alpha. \text{hnd } \varsigma \rightarrow \alpha \rightarrow \text{rgn } \varsigma (\text{ref } \varsigma \alpha)$ 
read :  $\forall \varsigma. \forall \alpha. \text{ref } \varsigma \alpha \rightarrow \text{rgn } \varsigma \alpha$ 
write :  $\forall \varsigma. \forall \alpha. \text{ref } \varsigma \alpha \rightarrow \alpha \rightarrow \text{rgn } \varsigma \mathbf{1}$ 
runRgn :  $\forall \alpha. (\forall \varsigma. \text{rgn } \varsigma \alpha) \rightarrow \alpha$ 
letRgn :  $\forall \varsigma_1. \forall \alpha. (\forall \varsigma_2. \text{pf } (\varsigma_1 \leq \varsigma_2) \rightarrow \text{hnd } \varsigma_2 \rightarrow \text{rgn } \varsigma_2 \alpha) \rightarrow \text{rgn } \varsigma_1 \alpha$ 
coerceRgn :  $\forall \varsigma_1, \varsigma_2. \forall \alpha. \text{pf } (\varsigma_1 \leq \varsigma_2) \rightarrow \text{rgn } \varsigma_1 \alpha \rightarrow \text{rgn } \varsigma_2 \alpha$ 
reflSub :  $\forall \varsigma. \text{pf } (\varsigma \leq \varsigma)$ 
transSub :  $\forall \varsigma_1, \varsigma_2, \varsigma_3. \text{pf } (\varsigma_1 \leq \varsigma_2) \rightarrow \text{pf } (\varsigma_2 \leq \varsigma_3) \rightarrow \text{pf } (\varsigma_1 \leq \varsigma_3)$ 

```

The type $\text{rgn } s \tau$ is the type of computations which transform a stack indexed by s and deliver a value of type τ . The type $\text{ref } s \tau$ is the type of mutable references allocated in the region at the top of the stack indexed by s and containing a value of type τ . The type $\text{hnd } s$ is the type of handles for the region at the top of the stack indexed by s ; we require a handle to allocate a reference in a region, but do not require a handle to read or write a reference.

The operations **return** and **then** are the *unit* and *bind* operations of the **rgn** monad, the former lifting a value to a computation and the latter sequencing computations. The next three operations are primitive stack transformers. **new** takes a region handle and an initial value and yields a stack transformer, which, when applied to a stack of regions, allocates and initializes a fresh reference in the appropriate region, and delivers the reference and the augmented stack. Similarly, **read** and **write** yield computations that respectively query and update the mappings of references to values in the current stack of regions. Note that all of these operations require the stack index ς of **rgn** and **ref** to be equal.

Finally, the operation **runRgn** encapsulates a stateful computation. To do so, it takes a stack transformer as its argument, applies it to an initial empty stack of regions, and returns the result while discarding the final stack (which should be empty). Note that to apply **runRgn**, we instantiate α with the type of the result to be returned, and then supply a stack transformer, *which is polymorphic in the stack index* ς . The effect of this universal quantification is that the stack transformer makes no assumptions about the initial stack (e.g., the existence of pre-allocated regions or references). Furthermore, the instantiation of the type variable α occurs outside the scope of the stack variable ς ; this prevents the stack transformer from delivering a value whose type mentions ς . Thus, references or computations depending on the final stack cannot escape beyond the encapsulation of **runRgn**.

However, the above does not suffice to encode region-based languages. The difficulty is that, in a region-based language, it is critical to allocate variables in and read variables from an outer (older) region while in the scope of an inner (younger) region. To accommodate this essential idiom, we include a powerful **letRgn** operation that is similar to **runRgn** in the sense that it encapsulates a stateful computation. Operationally, **letRgn** transforms a stack by (1) creating a new region on the top of the stack, (2) applying a stack transformer to the augmented stack to yield a transformed stack, (3) destroying the region on the

<i>Kinds</i>	$\kappa ::= \text{STACK} \mid \star$
<i>Type-level Variables</i>	$\varepsilon, \varsigma, \alpha ::= T\text{Vars}$
<i>Stack Indices</i>	$s ::= \varsigma$
<i>Types</i>	$\tau ::= \alpha \mid \tau_1 \rightarrow \tau_2 \mid \mathbf{1} \mid \tau_1 \times \tau_2 \mid \forall \varepsilon:\kappa. \tau$ $\text{rgn } s \tau \mid \text{ref } s \tau \mid \text{hnd } s \mid \text{pf } (s_1 \leq s_2)$
<i>Type-level Terms</i>	$\epsilon ::= s \mid \tau$
<i>Type-level Contexts</i>	$\Delta ::= \bullet \mid \Delta, \varepsilon:\kappa$
<i>rgn Monad Operations</i>	$\text{ops} ::= \text{runRgn} \mid \text{coerceRgn} \mid \text{transSub} \mid$ $\text{return} \mid \text{then} \mid \text{letRgn} \mid \text{new} \mid \text{read} \mid \text{write}$
<i>Expressions</i>	$e ::= \text{ops} \mid x \mid \lambda x:\tau. e \mid e_1 e_2 \mid \langle \rangle \mid \text{let } \langle \rangle = e_1 \text{ in } e_2 \mid$ $\langle e_1, e_2 \rangle \mid \text{let } \langle x_1, x_2 \rangle = e_1 \text{ in } e_2 \mid \Lambda \varepsilon:\kappa. e \mid e[\epsilon]$
<i>Expression-level Contexts</i>	$\Gamma ::= \bullet \mid \Gamma, x:\tau$

Fig. 1. Syntax of F^{RGN}

top of the transformed stack and yielding the bottom of the transformed stack. Ignoring for the moment the argument of type $\text{pf}(s_1 \leq s_2)$, we see that we may apply exactly the same reasoning as applied to `runRgn`: the computation makes no assumptions about the newly augmented stack s_2 , nor can the newly augmented stack s_2 be leaked through the return value.

What, then, is the role of the $\text{pf}(s_1 \leq s_2)$? The answer lies in the fact that the stack index s_2 does not denote an arbitrary stack; rather, it should denote a stack that is related to s_1 by the addition of a newly created region (i.e., $s_2 \equiv r::s_1$). In fact, we may consider s_1 to be a subtype of s_2 , since every region in the stack s_1 is also in the stack s_2 ; values of type $\text{pf}(s_1 \leq s_2)$ are witnesses of this relationship. The operation `coerceRgn` applies a subtyping witness to a stack transformer for the substack to yield a stack transformer for the superstack; intuitively, the operation is sound as a stack transformer may simply ignore extra regions. The operations `reflSub` and `transSub` are combinators witnessing the reflexivity and transitivity of the subtyping relation.

Figure 1 gives the complete syntax for F^{RGN} , which is a natural extension of System F. We introduce a simple kind system to support abstraction over both types and stack indices. (In the text, we often omit kind annotations, using the convention that ς stands for a type-level variable of `STACK` kind, and α of \star .)

We adopt the standard type system for System F; the only typing judgement of interest is $\Delta; \Gamma \vdash e : \tau$ meaning that expression e has type τ , where Δ records the free type-level variables and their kinds and Γ records the free expression-level variables and their types. The types for the `rgn` monad operations are as given in the text above.

Our previous work [6] gave an operational semantics for F^{RGN} and proved the type soundness of F^{RGN} . However, the operational semantics of F^{RGN} is somewhat cumbersome, due to the intertwining of contexts for pure evaluation and monadic evaluation. Hence, in the present setting, we will *define* the operational behavior of F^{RGN} by its translation into the target language of Section 3.

3 Target Calculus: λ^{rgnUL}

In another line of work [8], we introduced λ^{URAL} , a core substructural polymorphic λ -calculus, and then extended it to λ^{refURAL} by adding a rich collection of mutable references. Providing four sorts of substructural qualifiers (unrestricted, relevant, affine, and linear) allowed us to encode and study the interactions of different forms of uniqueness that appear in other high-level programming languages. Notable features of λ^{refURAL} include: deallocation of references; strong (type-varying) updates; and storage of unique objects in shared references.

Here, we augment λ^{refURAL} by adding region primitives, and also simplify the language by removing features, such as the relevant and affine qualifiers, that do not play a part in the translation. We call the resulting language λ^{rgnUL} .

In contrast to the `letRgn` operation of \mathbf{F}^{RGN} , which encapsulates the creation and destruction of a region, the primitives of λ^{rgnUL} include `newrgn` and `freergn` for separately creating and destroying a region. All access to a region (for allocating, reading, and writing references) is mediated by a linear *capability* that is produced by `newrgn` and consumed by `freergn`.

As noted above, λ^{rgnUL} is a *substructural* polymorphic λ -calculus. A *substructural* type system provides the core mechanisms necessary to restrict the number and order of uses of data and operations. In our calculus, types and variables are qualified as unrestricted (U) or linear (L). Essentially, unrestricted variables are allowed to be used an arbitrary number of times, while linear variables are allowed to be used exactly once.

Figure 2 gives the syntax for λ^{rgnUL} , excluding intermediate terms that would appear in an operational semantics. Many of the types and expressions are based on a traditional polymorphic λ -calculus.

We structure our types τ as a qualifier q applied to a pre-type $\bar{\tau}$, yielding the two sorts of types noted above. The qualifier of a type dictates the number of uses of variables of the type, while the pre-type dictates the introduction and elimination forms. The pre-types $\mathbf{1}_{\otimes}$, $\tau_1 \otimes \dots \otimes \tau_n$, and $\tau_1 \multimap \tau_2$ correspond to the unit, product, and function types of the polymorphic λ -calculus. Quantification over qualifiers, region names, pre-types, and types is provided by the pre-types $\forall \varepsilon : \kappa. \tau$ and $\exists \varepsilon : \kappa. \tau$. (In the text, we often omit kind annotations, using the convention that ξ stands for a type-level variable of QUAL kind, ϱ of RGN, $\bar{\alpha}$ of $\bar{\star}$, and α of \star .)

The pre-types `ref` r τ and `hnd` r are similar to the corresponding types in \mathbf{F}^{RGN} ; the former is the type of mutable references allocated in the region r and the latter is the type of handles for the region r . The pre-type `cap` r is the type of capabilities for accessing the region named r . We shall shortly see how linear capabilities effectively mediate access to a region.

Space precludes us from giving a detailed description of the type system for λ^{rgnUL} ; the major features are entirely standard for a substructural setting [9, 8]. First, in order to ensure the correct relationship between a data structure and its components, we extend the lattice ordering on constant qualifiers to arbitrary qualifiers ($\Delta \vdash q \preceq q'$), types ($\Delta \vdash \tau \preceq q'$), and contexts ($\Delta \vdash \Gamma \preceq q'$). Second, we introduce a judgement $\Delta \vdash \Gamma_1 \boxplus \Gamma_2 \rightsquigarrow \Gamma$ that splits the assumptions in Γ

<i>Kinds</i>	$\kappa ::= \text{QUAL} \mid \text{RGN} \mid \bar{\kappa} \mid \star$
<i>Type-level Variables</i>	$\varepsilon, \xi, \varrho, \bar{\alpha}, \alpha ::= \text{TVars}$
<i>Constant Qualifiers</i>	$\mathbf{q} \in \text{Quals} = \{\mathbf{U}, \mathbf{L}\} \quad \mathbf{U} \sqsubseteq \mathbf{L}$
<i>Qualifiers</i>	$q ::= \xi \mid \mathbf{q}$
<i>Constant Region Names</i>	$\mathbf{r} \in \text{RNames}$
<i>Region Names</i>	$r ::= \varrho \mid \mathbf{r}$
<i>PreTypes</i>	$\bar{\tau} ::= \bar{\alpha} \mid \tau_1 \multimap \tau_2 \mid \mathbf{1}_{\otimes} \mid \tau_1 \otimes \cdots \otimes \tau_n \mid \forall \varepsilon:\kappa. \tau \mid \exists \varepsilon:\kappa. \tau \mid \text{ref } r \tau \mid \text{hnd } r \mid \text{cap } r$
<i>Types</i>	$\tau ::= \alpha \mid {}^q \bar{\tau}$
<i>Type-level Terms</i>	$\varepsilon ::= q \mid r \mid \bar{\tau} \mid \tau$
<i>Type-level Contexts</i>	$\Delta ::= \bullet \mid \Delta, \varepsilon:\kappa$
<i>Region Primitives</i>	$\text{prims} ::= \text{newrgn} \mid \text{freergn} \mid \text{new} \mid \text{read} \mid \text{write}$
<i>Expressions</i>	$e ::= \text{prims} \mid x \mid {}^q \lambda x:\tau. e \mid e_1 e_2 \mid {}^q \langle \rangle \mid \text{let } \langle \rangle = e_1 \text{ in } e_2 \mid {}^q \langle e_1, \dots, e_n \rangle \mid \text{let } \langle x_1, \dots, x_n \rangle = e_1 \text{ in } e_2 \mid {}^q \Lambda \varepsilon:\kappa. e \mid e[\varepsilon] \mid {}^q \text{pack}(\varepsilon:\kappa, e) \mid \text{let pack}(\varepsilon:\kappa, x) = e_1 \text{ in } e_2$
<i>Expression-level Contexts</i>	$\Gamma ::= \bullet \mid \Gamma, x:\tau$

Fig. 2. Syntax of λ^{rgnUL}

between the contexts Γ_1 and Γ_2 . Splitting the context is necessary to ensure that variables are used appropriately by sub-expressions. Note that \square must ensure that an \mathbf{L} assumption appears in exactly one sub-context, while \mathbf{U} assumptions may appear in both sub-contexts.

The main typing judgement has the form $\Delta; \Gamma \vdash e : \tau$; Figure 3 gives typing rules for each of the expression forms of λ^{rgnUL} .

Finally, we assign types for each of the region primitives of λ^{rgnUL} :

$$\begin{aligned}
\text{newrgn} &: \mathbf{U}(\mathbf{L}\mathbf{1}_{\otimes} \multimap \mathbf{L}\exists \varrho. \mathbf{L}(\mathbf{L}\text{cap } \varrho \otimes \mathbf{U}\text{hnd } \varrho)) \\
\text{freergn} &: \mathbf{U}\forall \varrho. \mathbf{U}(\mathbf{L}(\mathbf{L}\text{cap } \varrho \otimes \mathbf{U}\text{hnd } \varrho) \multimap \mathbf{L}\mathbf{1}_{\otimes}) \\
\text{new} &: \mathbf{U}\forall \varrho. \mathbf{U}\forall \bar{\alpha}. \mathbf{U}(\mathbf{L}(\mathbf{L}\text{cap } \varrho \otimes \mathbf{U}\text{hnd } \varrho \otimes \mathbf{U}\bar{\alpha}) \multimap \mathbf{L}(\mathbf{L}\text{cap } \varrho \otimes \mathbf{U}(\text{ref } \varrho \mathbf{U}\bar{\alpha}))) \\
\text{read} &: \mathbf{U}\forall \varrho. \mathbf{U}\forall \bar{\alpha}. \mathbf{U}(\mathbf{L}(\mathbf{L}\text{cap } \varrho \otimes \mathbf{U}(\text{ref } \varrho \mathbf{U}\bar{\alpha})) \multimap \mathbf{L}(\mathbf{L}\text{cap } \varrho \otimes \mathbf{U}\bar{\alpha})) \\
\text{write} &: \mathbf{U}\forall \varrho. \mathbf{U}\forall \bar{\alpha}. \mathbf{U}(\mathbf{L}(\mathbf{L}\text{cap } \varrho \otimes \mathbf{U}(\text{ref } \varrho \mathbf{U}\bar{\alpha}) \otimes \mathbf{U}\bar{\alpha}) \multimap \mathbf{L}(\mathbf{L}\text{cap } \varrho \otimes \mathbf{U}\mathbf{1}_{\otimes}))
\end{aligned}$$

We have purposefully “streamlined” the type of the reference primitives in order to simplify the exposition. For example, note that we may only allocate, read, and write references whose contents are unrestricted. However, there is no fundamental difficulty in adopting a richer set of reference primitives (à la λ^{refURAL} [8]), which would allow references to contain arbitrary values.

Space again precludes us from giving a detailed description of the operational semantics for λ^{rgnUL} ; however, it is entirely standard for a region-based language. The small-step operational semantics is defined by a relation between configurations of the form (ψ, e) , where ψ is a global heap mapping region names to regions and regions are mappings from pointers to values.

The primitives **newrgn** and **freergn** perform the complementary actions of creating and destroying a region in the global heap. Note that the type of **newrgn**

$$\boxed{\Delta; \Gamma \vdash e : \tau}$$

$$\frac{}{\Delta; \bullet, x : \tau \vdash x : \tau}
\quad
\frac{\Delta; \Gamma_1 \boxplus \Gamma_2 \rightsquigarrow \Gamma \quad \Delta \vdash \Gamma_1 \preceq \mathbf{U} \quad \Delta; \Gamma_2 \vdash e : \tau}{\Delta; \Gamma \vdash e : \tau}$$

$$\frac{\Delta \vdash q : \text{QUAL}}{\Delta; \bullet \vdash \langle \rangle : {}^q \mathbf{1}_{\otimes}}
\quad
\frac{\Delta \vdash \Gamma_1 \boxplus \Gamma_2 \rightsquigarrow \Gamma \quad \Delta; \Gamma_1 \vdash e_1 : {}^q \mathbf{1}_{\otimes} \quad \Delta; \Gamma_2 \vdash e_2 : \tau}{\Delta; \Gamma \vdash \mathbf{let} \langle \rangle = e_1 \mathbf{in} e_2 : \tau}$$

$$\frac{\Delta \vdash \Gamma_1 \boxplus \dots \boxplus \Gamma_n \rightsquigarrow \Gamma \quad \Delta \vdash \tau_1 \preceq q \quad \dots \quad \Delta; \Gamma_n \vdash e_n : \tau_n \quad \Delta \vdash \tau_n \preceq q}{\Delta; \Gamma \vdash \langle e_1, \dots, e_n \rangle : {}^q (\tau_1 \otimes \dots \otimes \tau_n)}
\quad
\frac{\Delta \vdash \Gamma_1 \boxplus \Gamma_2 \rightsquigarrow \Gamma \quad \Delta; \Gamma_1 \vdash e_1 : {}^q (\tau_1 \otimes \dots \otimes \tau_n) \quad \Delta; \Gamma_2, x_1 : \tau_1, \dots, x_n : \tau_n \vdash e_2 : \tau}{\Delta; \Gamma \vdash \mathbf{let} \langle x_1, \dots, x_n \rangle = e_1 \mathbf{in} e_2 : \tau}$$

$$\frac{\Delta \vdash \Gamma \preceq q \quad \Delta; \Gamma, x : \tau_x \vdash e : \tau}{\Delta; \Gamma \vdash {}^q \lambda x : \tau_x. e : {}^q (\tau_x \multimap \tau)}
\quad
\frac{\Delta \vdash \Gamma_1 \boxplus \Gamma_2 \rightsquigarrow \Gamma \quad \Delta; \Gamma_1 \vdash e_1 : {}^q (\tau_x \multimap \tau) \quad \Delta; \Gamma_2 \vdash e_2 : \tau_x}{\Delta; \Gamma \vdash e_1 e_2 : \tau}$$

$$\frac{\Delta \vdash \Gamma \preceq q \quad \Delta, \varepsilon : \kappa; \Gamma \vdash e : \tau}{\Delta; \Gamma \vdash {}^q \Lambda \varepsilon : \kappa. e : {}^q (\forall \varepsilon : \kappa. \tau)}
\quad
\frac{\Delta; \Gamma \vdash e_1 : {}^q (\forall \varepsilon : \kappa. \tau) \quad \Delta \vdash e_2 : \kappa}{\Delta; \Gamma \vdash e_1 [e_2] : \tau [e_2 / \varepsilon]}$$

$$\frac{\Delta \vdash \varepsilon_1 : \kappa \quad \Delta; \Gamma \vdash e_2 : \tau [\varepsilon_1 / \varepsilon] \quad \Delta \vdash \tau [\varepsilon_1 / \varepsilon] \preceq q}{\Delta; \Gamma; \Sigma \vdash {}^q \mathbf{pack}(\varepsilon_1 : \kappa, e_2) : {}^q (\exists \varepsilon : \kappa. \tau)}
\quad
\frac{\Delta \vdash \Gamma_1 \boxplus \Gamma_2 \rightsquigarrow \Gamma \quad \Delta \vdash \tau' : \star \quad \Delta; \Gamma_1 \vdash e_1 : {}^q (\exists \varepsilon : \kappa. \tau) \quad \Delta, \varepsilon : \kappa; \Gamma_2, x : \tau \vdash e_2 : \tau'}{\Delta; \Gamma \vdash \mathbf{let} \mathbf{pack}(\varepsilon : \kappa, x) = e_1 \mathbf{in} e_2 : \tau'}$$

Fig. 3. Static Semantics of λ^{rgnUL}

specifies that it returns an existential package, hiding the name of the fresh region. The primitives `new`, `read`, and `write` behave precisely as their counterparts in any region-based language. Additionally, their types specify that they thread $\text{Lcap } \varrho$ values through the evaluation; the capability is simply presented at each access of a region and returned to allow future access. In the semantics, the capability is represented as a dummy token, which has no run-time significance.

As expected, the type system for λ^{rgnUL} is sound with respect to its operational semantics:

Theorem 1 (λ^{rgnUL} **Safety**). *If $\bullet, \bullet \vdash e_1 : \tau$ and $(\{\}, e_1) \mapsto^* (\psi_2, e_2)$, then either there exists v such that $e_2 \equiv v$ or there exists ψ_3 and e_3 such that $(\psi_2, e_2) \mapsto (\psi_3, e_3)$.*

We have formally verified this result (for a rich superset of λ^{rgnUL}) in the Twelf system [10] using its metatheorem checker [11]. The mechanized proof can be obtained at <http://www.cs.cornell.edu/People/fluett/research/substruct-regions/>.

4 Translation: F^{RGN} to λ^{rgnUL}

Having introduced both our source and target calculi, we are in a position to consider a (type-preserving) translation from F^{RGN} to λ^{rgnUL} . Before giving the details, we discuss a few of the high-level issues.

First, we note that F^{RGN} has no notion of linearity in the syntax or type system. Rather, all variables and types are implicitly considered unrestricted. Hence, we can expect that the translation of all F^{RGN} expressions will yield λ^{rgnUL} expressions with a U qualified type.

On the other hand, we claimed that a stateful region computation could be interpreted as a stack transformer. Recall that the type $\text{rgn } s \tau$ is the type of computations which transform a stack indexed by s and deliver a value of type τ . A key characteristic of F^{RGN} is that all primitive stack transformers are meant to use the stack in a single-threaded manner; hence, a stateful computation can update the stack in place. This single-threaded behavior is precisely the sort of resource management that may be captured by a substructural type system. Hence, we can expect that the representation of a stack of regions in λ^{rgnUL} will be a value with L qualified type. In particular, we will represent a stack of regions as a sequence of linear capabilities, formed out of nested linear tuples.

Third, we must be mindful of a slight mismatch between the hnd and ref types in F^{RGN} and the corresponding types in λ^{rgnUL} . Recall that, in F^{RGN} , $\text{hnd } s$ and $\text{ref } s \tau$ are handles for and references allocated in the region at the top of the stack indexed by s . Whereas, in λ^{rgnUL} , $\text{hnd } r$ and $\text{ref } r \tau$ explicitly name the region of the handle or reference. This subtle distinction (whether the region is implicit or explicit) will need to be handled by the translation.

Bearing these issues in mind, we turn our attention to the translation of F^{RGN} type-level terms given in Figure 4. $\mathcal{S}_\star \llbracket s \rrbracket$ translates a F^{RGN} term of STACK kind to a λ^{rgnUL} term of \star kind. As the STACK kind of F^{RGN} is inhabited only by variables, the translation is trivial: in λ^{rgnUL} , ς is considered a variable of \star kind.

$\mathcal{T}_\star \llbracket \tau \rrbracket$ and $\mathcal{T}_{\bar{\star}} \llbracket \tau \rrbracket$ translate a F^{RGN} term of \star kind to λ^{rgnUL} terms of \star and $\bar{\star}$ kinds, respectively. As we observed above, when we translate a F^{RGN} type to a λ^{rgnUL} type, we ensure that the result is a U qualified type. The $\mathcal{T}_{\bar{\star}} \llbracket \tau \rrbracket$ translation is straightforward on the functional types. (However, note that a F^{RGN} variable α of \star kind is translated to a λ^{rgnUL} variable $\bar{\alpha}$ of $\bar{\star}$ kind; this ensures that *every* type corresponding to a F^{RGN} type is manifestly qualified with U .)

More interesting are the translations of the types associated with the rgn monad. In the translation of the $\text{rgn } s \tau$ type, we see the familiar store (stack) passing interpretation of computations. Since the representation of a stack of regions is linear, the resulting store/value pair is qualified with L . Next, consider the translation of the $\text{pf } (s_1 \leq s_2)$ type. Recall that it is the type of witnesses to the fact that the stack indexed by s_1 is a subtype of the stack indexed by s_2 . Hence, we translate to a type that expresses the isomorphism between $\mathcal{S}_\star \llbracket s_2 \rrbracket$ and ${}^L(\mathcal{S}_\star \llbracket s_1 \rrbracket \otimes \beta)$, for some “slack” β . Note that while the types $\mathcal{S}_\star \llbracket s_2 \rrbracket$, $\mathcal{S}_\star \llbracket s_1 \rrbracket$, and β may be linear, the pair of functions witnessing the isomorphism is unrestricted. This corresponds to the fact that the *proof* that s_1 is a subtype of s_2 is *persistent*, while the *existence* of the stacks s_1 and s_2 are *ephemeral*.

$$\begin{array}{l}
\text{F}^{\text{RGN}} \text{ STACK to } \lambda^{\text{rgnUL}} \star \\
\mathcal{S}_\star \llbracket \varsigma \rrbracket = \varsigma \\
\text{F}^{\text{RGN}} \star \text{ to } \lambda^{\text{rgnUL}} \star \\
\mathcal{T}_\star \llbracket \tau \rrbracket = \text{U} \mathcal{T}_\star \llbracket \tau \rrbracket
\end{array}
\qquad
\begin{array}{l}
\text{F}^{\text{RGN}} \star \text{ to } \lambda^{\text{rgnUL}} \bar{\star} \text{ (functional types)} \\
\mathcal{T}_\star \llbracket \alpha \rrbracket = \bar{\alpha} \\
\mathcal{T}_\star \llbracket \tau_1 \rightarrow \tau_2 \rrbracket = \mathcal{T}_\star \llbracket \tau_1 \rrbracket \multimap \mathcal{T}_\star \llbracket \tau_2 \rrbracket \\
\mathcal{T}_\star \llbracket \mathbf{1} \rrbracket = \mathbf{1}_\otimes \\
\mathcal{T}_\star \llbracket \tau_1 \times \tau_2 \rrbracket = \mathcal{T}_\star \llbracket \tau_1 \rrbracket \otimes \mathcal{T}_\star \llbracket \tau_2 \rrbracket \\
\mathcal{T}_\star \llbracket \forall \alpha: \star. \tau \rrbracket = \forall \bar{\alpha}: \bar{\star}. \mathcal{T}_\star \llbracket \tau \rrbracket \\
\mathcal{T}_\star \llbracket \forall \varsigma: \text{STACK}. \tau \rrbracket = \forall \varsigma: \star. \mathcal{T}_\star \llbracket \tau \rrbracket
\end{array}$$

$$\begin{array}{l}
\text{F}^{\text{RGN}} \star \text{ to } \lambda^{\text{rgnUL}} \bar{\star} \text{ (rgn monad types)} \\
\mathcal{T}_\star \llbracket \text{rgn } s \tau \rrbracket = \mathcal{S}_\star \llbracket s \rrbracket \multimap \text{L}(\mathcal{S}_\star \llbracket s \rrbracket \otimes \mathcal{T}_\star \llbracket \tau \rrbracket) \\
\mathcal{T}_\star \llbracket \text{pf } (s_1 \leq s_2) \rrbracket = \exists \beta: \star. \mathbf{Iso}(\mathcal{S}_\star \llbracket s_2 \rrbracket, \text{L}(\mathcal{S}_\star \llbracket s_1 \rrbracket \otimes \beta)) \\
\mathcal{T}_\star \llbracket \text{hnd } s \rrbracket = \exists \varrho: \text{RGN}. \text{U}(\text{U} \exists \beta: \star. \mathbf{Iso}(\mathcal{S}_\star \llbracket s \rrbracket, \text{L}(\beta \otimes \text{L}(\text{cap } \varrho))) \otimes \text{U}(\text{hnd } \varrho)) \\
\mathcal{T}_\star \llbracket \text{ref } s \tau \rrbracket = \exists \varrho: \text{RGN}. \text{U}(\text{U} \exists \beta: \star. \mathbf{Iso}(\mathcal{S}_\star \llbracket s \rrbracket, \text{L}(\beta \otimes \text{L}(\text{cap } \varrho))) \otimes \text{U}(\text{ref } \varrho \mathcal{T}_\star \llbracket \tau \rrbracket))
\end{array}$$

$$\begin{array}{l}
\lambda^{\text{rgnUL}} \text{ Type-level Macros} \\
\mathbf{Iso}(\tau_1, \tau_2) = \text{U}(\text{U}(\tau_1 \multimap \tau_2) \otimes \text{U}(\tau_2 \multimap \tau_1))
\end{array}$$

Fig. 4. Type-level Translation

The translation of the `hnd s` and `ref s τ` types are similar. An existentially bound region name ϱ fixes the region for the λ^{rgnUL} handle or reference, while an isomorphism witnesses the fact that ϱ may be found within the stack $\mathcal{S}_\star \llbracket s \rrbracket$.

With the translation of F^{RGN} type-level terms in place, the translation of F^{RGN} expressions follows almost directly. We elide the translation of the introduction and elimination forms for the functional types in F^{RGN} (it is simply the homomorphic mapping of the given expression translations) and focus on the translation of the `rgn` monad operations. For readability, we give translations for fully applied region primitives only, assuming that partially applied primitives have been eta-expanded. The translation of `return` and `then` follow directly from our store (stack) passing interpretation of `rgn s τ` types:

$$\begin{array}{l}
\mathcal{E} \llbracket \text{return } [s] [\tau_\alpha] e \rrbracket = \\
\text{let } \text{res}: \mathcal{T}_\star \llbracket \tau_\alpha \rrbracket = \mathcal{E} \llbracket e \rrbracket \text{ in} \\
\text{U} \lambda \text{stk}: \mathcal{S}_\star \llbracket s \rrbracket. \text{L} \langle \text{stk}, \text{res} \rangle
\end{array}
\qquad
\begin{array}{l}
\mathcal{E} \llbracket \text{then } [s] [\tau_\alpha] [\tau_\beta] e_1 e_2 \rrbracket = \\
\text{let } f: \mathcal{T}_\star \llbracket \text{rgn } s \tau_\alpha \rrbracket = \mathcal{E} \llbracket e_1 \rrbracket \text{ in} \\
\text{let } g: \mathcal{T}_\star \llbracket \tau_\alpha \rightarrow \text{rgn } s \tau_\beta \rrbracket = \mathcal{E} \llbracket e_2 \rrbracket \text{ in} \\
\text{U} \lambda \text{stk}: \mathcal{S}_\star \llbracket s \rrbracket. \text{let } \langle \text{stk}, \text{res} \rangle = f \text{ stk in} \\
g \text{ res stk}
\end{array}$$

The translation of `letRgn` is the most complicated, but breaks down into conceptually simple components. We bracket the execution of the inner computation with a `newrgn/freergn` pair, creating and destroying a new region. We construct the representation of the new stack stk_2 for the inner computation by pairing the old stack stk_1 with the new region capability cap . Finally, we construct isomorphisms witnessing the relationships between the new region capability and the new stack and between the old stack and the new stack. We carefully chose the isomorphism types so that the identity function suffices as a witness. Putting all of these pieces together, we have the following:

$$\begin{aligned}
\mathcal{E} \llbracket \text{letRgn } [s_1] [\tau_\alpha] e \rrbracket &= \\
&\text{let } f: \mathcal{T}_* \llbracket \forall \varsigma_2: \text{STACK. pf } (s_1 \leq \varsigma_2) \rightarrow \text{hnd } \varsigma_2 \rightarrow \text{rgn } \varsigma_2 \tau_\alpha \rrbracket = \mathcal{E} \llbracket e \rrbracket \text{ in} \\
&\quad \text{let } \text{pack}(\varrho: \text{RGN}, \langle \text{cap}, \text{hnd} \rangle) = \text{newrgn}^{\text{L}} \langle \rangle \text{ in} \\
&\quad \text{let } id = \text{let } \text{stk}_1: \mathcal{S}_* \llbracket s_1 \rrbracket . \text{let } \text{pack}(\beta: \star, \langle \text{spl}, \text{cmb} \rangle) = \text{ppf} \text{ in} \\
&\quad \quad \text{let } \langle \text{stk}_1, \text{stk}_\beta \rangle = \text{spl } \text{stk}_2 \text{ in} \\
&\quad \quad \text{let } \langle \text{stk}_1, \text{res} \rangle = f \text{ } \text{stk}_1 \text{ in} \\
&\quad \quad \text{let } \text{stk}_2 = \text{cmb}^{\text{L}} \langle \text{stk}_1, \text{stk}_\beta \rangle \text{ in} \\
&\quad \quad \quad \text{L} \langle \text{stk}_2, \text{res} \rangle \\
&\quad \text{let } id = \text{let } \text{stk}_1: \mathcal{S}_* \llbracket s_1 \rrbracket \otimes^{\text{L}} \text{cap } \varrho . \text{stk} \text{ in} \\
&\quad \text{let } \text{ppf} = \text{pack}^{\text{L}}(\text{cap } \varrho) : \star, \text{let } \langle id, id \rangle \text{ in} \\
&\quad \text{let } \text{phnd} = \text{pack}(\varrho: \text{RGN}, \text{let } \langle id, id \rangle \text{ in } \text{pack}(\mathcal{S}_* \llbracket s_1 \rrbracket : \star, \text{let } \langle id, id \rangle \text{ in } \text{hnd})) \text{ in} \\
&\quad \text{let } \text{stk}_2 = \text{L} \langle \text{stk}_1, \text{cap} \rangle \text{ in} \\
&\quad \text{let } \langle \text{stk}_2, \text{res} \rangle = f \text{ } [\text{L}(\mathcal{S}_* \llbracket s_1 \rrbracket \otimes^{\text{L}} \text{cap } \varrho)] \text{ppf } \text{phnd } \text{stk}_2 \text{ in} \\
&\quad \text{let } \langle \text{stk}_1, \text{cap} \rangle = \text{stk}_2 \text{ in} \\
&\quad \text{let } \langle \rangle = \text{freergn}[\varrho]^{\text{L}} \langle \text{cap}, \text{hnd} \rangle \text{ in} \\
&\quad \quad \text{L} \langle \text{stk}_1, \text{res} \rangle
\end{aligned}$$

We can see the isomorphisms in action in the translation of `coerceRgn`:

$$\begin{aligned}
\mathcal{E} \llbracket \text{coerceRgn } [s_1] [s_2] [\tau_\alpha] e_1 e_2 \rrbracket &= \\
&\text{let } \text{ppf}: \mathcal{T}_* \llbracket \text{pf } (s_1 \leq s_2) \rrbracket = \mathcal{E} \llbracket e_1 \rrbracket \text{ in} \\
&\text{let } f: \mathcal{T}_* \llbracket \text{rgn } s_1 \tau_\alpha \rrbracket = \mathcal{E} \llbracket e_2 \rrbracket \text{ in} \\
&\quad \text{let } \text{pack}(\beta: \star, \langle \text{spl}, \text{cmb} \rangle) = \text{ppf} \text{ in} \\
&\quad \quad \text{let } \langle \text{stk}_1, \text{stk}_\beta \rangle = \text{spl } \text{stk}_2 \text{ in} \\
&\quad \quad \text{let } \langle \text{stk}_1, \text{res} \rangle = f \text{ } \text{stk}_1 \text{ in} \\
&\quad \quad \text{let } \text{stk}_2 = \text{cmb}^{\text{L}} \langle \text{stk}_1, \text{stk}_\beta \rangle \text{ in} \\
&\quad \quad \quad \text{L} \langle \text{stk}_2, \text{res} \rangle
\end{aligned}$$

Note how the the stack “slack” stk_β is split out and then combined in, bracketing the execution of the `rgn` $s_1 \tau_\alpha$ computation.

As a final example, we can see an “empty” stack (represented by a $\text{L} \mathbf{1}_\otimes$ value) being provided as the initial stack in the translation of `runRgn`:

$$\begin{aligned}
\mathcal{E} \llbracket \text{runRgn } [\tau_\alpha] e \rrbracket &= \\
&\text{let } f: \mathcal{T}_* \llbracket \forall \varsigma: \text{STACK. rgn } \varsigma \tau_\alpha \rrbracket = \mathcal{E} \llbracket e \rrbracket \text{ in} \\
&\text{let } \langle \rangle, \text{res} = f \text{ } [\text{L} \mathbf{1}_\otimes]^{\text{L}} \langle \rangle \text{ in } \text{res}
\end{aligned}$$

The translations of the remaining `rgn` monad primitives are given in Figure 5. We strongly believe, but have not mechanically verified, that the translation is type preserving.

5 Extensions

The primary advantage of working at the target level is that we can expose the capabilities for regions as first-class objects instead of indirectly manipulating a stack of regions. In turn, this allows us to avoid the last-in-first-out lifetimes dictated by a lexically-scoped `letRgn`. For example, we can now explain the semantics for Cyclone’s *dynamic regions* and *unique pointers* using the concepts in the target language.

Dynamic Regions In Cyclone, a dynamic region r is represented by a *key* (`key r`) which is treated linearly by the type system. At the target level, a key can be represented by a pair of the capability for the region and its handle:

$$\text{key } r = \text{L}(\text{L} \text{cap } r \otimes^{\text{U}} \text{hnd } r)$$

Then creating a new key is accomplished by calling `newrgn`, and destroying the key is accomplished by calling `freergn`.

$$\begin{aligned}
\mathcal{E} \llbracket \text{new } [s] [\tau_\alpha] e_1 e_2 \rrbracket &= \\
&\text{let } phnd : \mathcal{T}_* \llbracket \text{hnd } s \rrbracket = \mathcal{E} \llbracket e_1 \rrbracket \text{ in} \\
&\text{let } x : \mathcal{T}_* \llbracket \tau_\alpha \rrbracket = \mathcal{E} \llbracket e_2 \rrbracket \text{ in} \\
&\text{let } \lambda stk : \mathcal{S}_* \llbracket [s] \rrbracket . \text{let } \text{pack}(\varrho : \text{RGN}, \langle \text{pack}(\beta : \star, \langle prj, inj \rangle), hnd \rangle) = phnd \text{ in} \\
&\quad \text{let } \langle stk_\beta, cap \rangle = prj \text{ stk in} \\
&\quad \text{let } \langle cap, ref \rangle = \text{new } [\varrho] [\mathcal{T}_* \llbracket \tau_\alpha \rrbracket] \llcorner \langle cap, hnd, x \rangle \text{ in} \\
&\quad \text{let } pref = \text{pack}(\rho : \text{RGN}, \text{pack}(\beta : \star, \langle prj, inj \rangle), ref) \text{ in} \\
&\quad \text{let } stk = inj \llcorner \langle stk_\beta, cap \rangle \text{ in} \\
&\quad \llcorner \langle stk, pref \rangle \\
\mathcal{E} \llbracket \text{read } [s] [\tau_\alpha] e \rrbracket &= \\
&\text{let } pref : \mathcal{T}_* \llbracket \text{ref } s \tau_\alpha \rrbracket = \mathcal{E} \llbracket e \rrbracket \text{ in} \\
&\text{let } \lambda stk : \mathcal{S}_* \llbracket [s] \rrbracket . \text{let } \text{pack}(\varrho : \text{RGN}, \langle \text{pack}(\beta : \star, \langle prj, inj \rangle), ref \rangle) = pref \text{ in} \\
&\quad \text{let } \langle stk_\beta, cap \rangle = prj \text{ stk in} \\
&\quad \text{let } \langle cap, res \rangle = \text{read } [\varrho] [\mathcal{T}_* \llbracket \tau_\alpha \rrbracket] \llcorner \langle cap, ref \rangle \text{ in} \\
&\quad \text{let } stk = inj \llcorner \langle stk_\beta, cap \rangle \text{ in} \\
&\quad \llcorner \langle stk, res \rangle \\
\mathcal{E} \llbracket \text{write } [s] [\tau_\alpha] e_1 e_2 \rrbracket &= \\
&\text{let } pref : \mathcal{T}_* \llbracket \text{ref } s \tau_\alpha \rrbracket = \mathcal{E} \llbracket e_1 \rrbracket \text{ in} \\
&\text{let } x : \mathcal{T}_* \llbracket \tau_\alpha \rrbracket = \mathcal{E} \llbracket e_2 \rrbracket \text{ in} \\
&\text{let } \lambda stk : \mathcal{S}_* \llbracket [s] \rrbracket . \text{let } \text{pack}(\varrho : \text{RGN}, \langle \text{pack}(\beta : \star, \langle prj, inj \rangle), ref \rangle) = pref \text{ in} \\
&\quad \text{let } \langle stk_\beta, cap \rangle = prj \text{ stk in} \\
&\quad \text{let } \langle cap, res \rangle = \text{write } [\varrho] [\mathcal{T}_* \llbracket \tau_\alpha \rrbracket] \llcorner \langle cap, ref, x \rangle \text{ in} \\
&\quad \text{let } stk = inj \llcorner \langle stk_\beta, cap \rangle \text{ in} \\
&\quad \llcorner \langle stk, res \rangle \\
\mathcal{E} \llbracket \text{reflSub } [s] \rrbracket &= \\
&\text{let } spl = \text{let } \lambda stk : \mathcal{S}_* \llbracket [s] \rrbracket . \text{let } su = \llcorner \langle stk, \llcorner \rangle \text{ in } su \text{ in} \\
&\text{let } cmb = \text{let } \lambda su : \llcorner (\mathcal{S}_* \llbracket [s] \rrbracket \otimes \llcorner \mathbf{1}_\otimes) . \text{let } \langle stk, \llcorner \rangle = su \text{ in } stk \text{ in} \\
&\text{pack}(\llcorner \mathbf{1}_\otimes : \star, \llcorner \langle spl, cmb \rangle) \\
\mathcal{E} \llbracket \text{transSub } [s_1] [s_2] [s_3] e_1 \rightsquigarrow_2 e_2 \rightsquigarrow_3 \rrbracket &= \\
&\text{let } ppf_{1 \rightsquigarrow_2} : \mathcal{T}_* \llbracket \text{pf } (s_1 \leq s_2) \rrbracket = \mathcal{E} \llbracket e_1 \rightsquigarrow_2 \rrbracket \text{ in} \\
&\text{let } ppf_{2 \rightsquigarrow_3} : \mathcal{T}_* \llbracket \text{pf } (s_2 \leq s_3) \rrbracket = \mathcal{E} \llbracket e_2 \rightsquigarrow_3 \rrbracket \text{ in} \\
&\text{let } \text{pack}(\alpha : \star, \langle spl_{2 \rightsquigarrow_1 \otimes \alpha}, cmp_{1 \otimes \alpha \rightsquigarrow_2} \rangle) = ppf_{1 \rightsquigarrow_2} \text{ in} \\
&\text{let } \text{pack}(\beta : \star, \langle spl_{3 \rightsquigarrow_2 \otimes \beta}, cmp_{2 \otimes \beta \rightsquigarrow_3} \rangle) = ppf_{2 \rightsquigarrow_3} \text{ in} \\
&\text{let } spl = \text{let } \lambda stk_3 : \mathcal{S}_* \llbracket [s_3] \rrbracket . \text{let } \langle stk_2, stk_\beta \rangle = spl_{3 \rightsquigarrow_2 \otimes \beta} stk_3 \text{ in} \\
&\quad \text{let } \langle stk_1, stk_\alpha \rangle = spl_{2 \rightsquigarrow_1 \otimes \alpha} stk_2 \text{ in} \\
&\quad \text{let } sss = \llcorner \langle stk_1, \llcorner \langle stk_\beta, stk_\alpha \rangle \rangle \text{ in} \\
&\quad \quad sss \quad \text{in} \\
&\text{let } cmb = \text{let } \lambda sss : \llcorner (\mathcal{S}_* \llbracket [s_1] \rrbracket \otimes \llcorner (\beta \otimes \alpha)) . \text{let } \langle stk_1, \langle stk_\beta, stk_\alpha \rangle \rangle = sss \text{ in} \\
&\quad \text{let } stk_2 = cmb_{1 \otimes \alpha \rightsquigarrow_2} \llcorner \langle stk_1, stk_\alpha \rangle \text{ in} \\
&\quad \text{let } stk_3 = cmb_{2 \otimes \beta \rightsquigarrow_3} \llcorner \langle stk_2, stk_\beta \rangle \text{ in} \\
&\quad \quad stk_3 \quad \text{in} \\
&\text{pack}(\llcorner (\beta \otimes \alpha) : \star, \llcorner \langle spl, cmb \rangle)
\end{aligned}$$

Fig. 5. Translation of rgn Monad Operations

To access a value allocated in a dynamic region, or to allocate a value in a dynamic region, Cyclone requires that the region be *opened* by presenting its key. The `openDRgn` is similar to a `letRgn` in that it conceptually pushes the dynamic region onto the stack of regions, executes the body, and then pops the region off the stack. During execution of the `openDRgn`'s body, the key becomes inaccessible, ensuring that the region cannot be deallocated. At the end of the `openDRgn` scope, the key is given back. The programmer is then able to destroy the region or later re-open it.

The `openDRgn` primitive can be implemented as a higher-order function with a signature like this (eliding the L and U qualifiers, and using source-level `rgn` to abbreviate the store passing translation):

$$\text{openDRgn} : \forall \varrho, \varsigma, \alpha. \text{key } \varrho \multimap (\text{hnd } \varrho \multimap \text{rgn } (\varsigma \otimes \text{cap } \varrho) \alpha) \multimap \text{rgn } \varsigma (\alpha \otimes \text{key } \varrho)$$

The function takes the key for ϱ and a computation, which, when given the handle for ϱ , expects to run on a stack of the form $\varsigma \otimes \text{cap } \varrho$ for some ς . Once applied, `openDRgn` returns a computation, which, when run on a stack ς , opens up the key to get the capability and handle, pushes the capability for ϱ on the stack, passes the handle to the computation and runs it in the extended stack to produce an α value. Then, it pops the capability, and returns a (linear) pair of the result and the re-packaged key. (We leave the definition of `openDRgn` as an exercise for the reader.)

Furthermore, since keys are first-class objects, they can be placed in data structures. For example, in our space-conscious web server [3], we use a list of dynamic regions, each of which holds data corresponding to a particular connection. When we receive data from a connection, we find the corresponding key, open it up, and then place the data in the region. When a connection is terminated, we pull the corresponding key out of the queue, perform `freeDRgn` on it, and thus deallocate all of the data associated with the connection. The price paid for this flexibility is that the list of keys must be treated linearly to avoid creating multiple aliases to the keys.

Unique Pointers Cyclone's unique pointers are anonymous dynamic regions, without the handle. Like the keys of dynamic regions, they can be destroyed at any time and the type system treats them linearly. At the target level, a unique pointer to a τ object can be represented as a term with type:

$$\text{L}\exists\varrho. \text{L}(\text{Lcap } \varrho \otimes \text{U}(\text{Lcap } \varrho \multimap \text{U}\mathbf{1}_{\otimes}) \otimes \text{U}(\text{ref } \varrho \tau))$$

Note that the actual reference is unrestricted, whereas the capability is linear; the handle is not available for further allocations, but is caught up in the function closure, which may be applied to free the region. This encoding allows us to "open" a unique pointer, just as we do dynamic regions, and for a limited scope, freely access (and duplicate) the underlying reference. Of course, during the scope of the open, we temporarily lose the capability to deallocate the object, but regain the capability upon exit from the scope.

In practice, the ability to open dynamic regions and unique pointers has proven crucial for integrating these facilities into the language. They make it relatively easy to access a data structure and mitigate some of the pain of threading linear resources through the program. Furthermore, they make it possible

to write re-usable libraries that manipulate data allocated in lexically-scoped regions, dynamic regions, or as unique objects.

Phase-Splitting In our target language, we represented capabilities and proof witnesses as explicit terms. But we have also crafted the language and translation so that these values are never actually needed at run-time. For instance, our witnesses only manipulate (products of) capabilities, which are themselves operationally irrelevant. These objects are only used to make the desired safety properties easy to check statically. So in principle, we should be able to erase the capabilities and witnesses before running a program.

To realize this goal, we should introduce a phase distinction via another modality, where we treat capabilities and proof witnesses as static objects, and all other terms as dynamic. The modality would demand that, as usual, static computations cannot depend upon dynamic values. Furthermore, we must be sure that witness functions (i.e., proof objects) are in fact total, effect-free functions so that they and their applications to capabilities may be safely erased. This sort of phase-splitting is effectively used in other settings that mix programming languages and logics, such as Xi *et al.*'s Applied Type System [12] and Sheard's Omega [13]. Perhaps the most promising approach is suggested by Mandelbaum, Walker, and Harper's work [14], where they developed a two-level language for reasoning about effectful programs.

The primary reason we did not introduce phase splitting here is that it complicates the translation and the target language, and thus obscures what is actually a relatively simple and straightforward encoding. A secondary reason is that, as demonstrated by the cited work above, there are many domains that would benefit from a general solution to the problem of type relevant, but operationally irrelevant, values.

6 Related Work and Open Issues

There has been much prior work aimed at relaxing the stack discipline imposed on region lifetimes by the Tofte-Talpin (TT) approach. The ML Kit [15] uses a storage-mode analysis to determine when it is safe to deallocate data in a region (known as region resetting) prior to the deallocation of the region itself. The safety of the storage-mode analysis has not been established formally.

Aiken *et al.* [16] eliminate the requirement that region allocation and deallocation should coincide with the beginning and end of the scope of region variables introduced by the `letregion` construct. They use a *late allocation/early deallocation* approach that delays the allocation of a region until just before its first access, and deallocates the region just after its last access. We believe that the results of their analysis can be encoded explicitly in our target language.

Unlike the previous two approaches which build on TT, Henglein *et al.* [17] present a region system that (like ours) replaces the `letregion` primitive with explicit commands to create and free a region. To ensure safety, they use a Hoare-logic-based region type system and consequently have no support for higher-order

functions. While they provide an inference algorithm to annotate programs with region manipulation commands, we intend for our system to serve as a target language for programs annotated using TT region inference, or those written in languages like Cyclone. The Calculus of Capabilities [18] is also intended as a target for TT-annotated programs, but unlike λ^{rgnUL} , it is defined in terms of a continuation-passing style language and does not support first-class regions.

The region system presented by Walker and Watkins [19] is perhaps the most closely related work. Like our target, they require a linear capability to be presented upon each access to a region. However, they provide a primitive, similar to `letregion`, that allows a capability to be temporarily treated as unrestricted for convenience’s sake. We have shown that no such primitive is needed. Rather, we use a combination of monadic encapsulation (to thread capabilities) coupled with *unrestricted witnesses* to achieve the same flexibility. In particular, our open construct for dynamic regions (and unique pointers) achieves the same effect as the Walker-Watkin’s primitive.

A related body of work has used regions as a low-level primitive on which to build type-safe garbage collectors [20–22]. Each of these approaches requires non-lexical regions, since, in a copying collector, the from- and to-spaces have non-nested lifetimes. Hawblitzel *et al.* [22] introduce a very low-level language in which they begin with a single linear array of words, construct lists and arrays out of the basic linear memory primitives, introduce *type sequences* for building regions of nonlinear data. Such a foundational approach is admirable, but there is a large semantic gap between a high-level language and such a target. Hence, λ^{rgnUL} serves as a useful intermediate point, and we may envision further translation from λ^{rgnUL} to such a low-level language.

The Vault language [23, 24] includes many of the features described in our target, including linear capabilities for accessing resources and a mechanism, called adoption, for temporarily transferring ownership of a capability to another capability, for a limited scope. But Vault also includes support for strong (i.e., type-changing) updates on linear resources, as well as features for temporarily treating an unrestricted resource as if it were linear. On the other hand, to the best of our knowledge, there exists no formal model that justifies the soundness of all of these mechanisms. We believe that it may be possible to combine λ^{rgnUL} with our previous work on strong updates [25, 26] to achieve this.

References

1. Cyclone, version 0.9. (2005) <http://www.eecs.harvard.edu/~greg/cyclone/>.
2. Tofte, M., Talpin, J.P.: Region-based memory management. *Information and Computation* **132**(2) (1997) 109–176
3. Hicks, M., Morrisett, G., Grossman, D., Jim, T.: Experience with safe manual memory-management in Cyclone. In: *Proc. International Symposium on Memory Management*. (2004) 73–84
4. Fluet, M., Wang, D.: Implementation and performance evaluation of a safe runtime system in Cyclone. In: *Proc. SPACE Workshop*. (2004)

5. Grossman, D., Morrisett, G., Jim, T., Hicks, M., Wang, Y., Cheney, J.: Region-based memory management in Cyclone. In: Proc. Programming Language Design and Implementation. (2002) 282–293
6. Fluet, M., Morrisett, G.: Monadic regions. In: Proc. International Conference on Functional Programming. (2004) 103–114
7. Launchbury, J., Peyton Jones, S.: State in Haskell. *Lisp and Symbolic Computation* **8**(4) (1995) 293–341
8. Ahmed, A., Fluet, M., Morrisett, G.: A step-indexed model of substructural state. In: Proc. International Conference on Functional Programming. (2005) 78–91
9. Walker, D.: Substructural type systems. In Pierce, B., ed.: *Advanced Topics in Types and Programming Languages*. MIT Press, Cambridge, MA (2005) 3–43
10. Pfenning, F., Schürmann, C.: Twelf – a meta-logic framework for deductive systems. In: Proc. Conference on Automated Deduction. (1999) 202–206
11. Schürmann, C., Pfenning, F.: A coverage checking algorithm for LF. In: Proc. Theorem Proving in Higher Order Logics. (2003) 120–135 LNCS 2758.
12. Chen, C., Xi, H.: Combining programming with theorem proving. In: Proc. International Conference on Functional Programming. (2005) 66–77
13. Sheard, T., Pasalic, E.: Meta-programming with built-in type equality (extended abstract). In: International Workshop on Logical Frameworks and Meta-Languages. (2004)
14. Mandelbaum, Y., Walker, D., Harper, R.: An effective theory of type refinements. In: Proc. International Conference on Functional Programming. (2003) 213–225
15. Tofte, M., Birkedal, L., Elsmann, M., Hallenberg, N., Olesen, T.H., Sestoft, P.: Programming with regions in the ML Kit (for version 4). Technical report, IT University of Copenhagen (2002)
16. Aiken, A., Fähndrich, M., Levien, R.: Better static memory management: Improving region-based analysis of higher-order languages. In: Proc. Programming Language Design and Implementation. (1995) 174–185
17. Henglein, F., Makhholm, H., Niss, H.: A direct approach to control-flow sensitive region-based memory management. In: Proc. Principles and Practice of Declarative Programming. (2001) 175–186
18. Walker, D., Crary, K., Morrisett, G.: Typed memory management in a calculus of capabilities. *ACM Transactions on Programming Languages and Systems* **24**(4) (2000) 701–771
19. Walker, D., Watkins, K.: On regions and linear types. In: Proc. International Conference on Functional Programming. (2001) 181–192
20. Wang, D., Appel, A.: Type-preserving garbage collectors. In: Proc. Principles of Programming Languages. (2001) 166–178
21. Monnier, S., Saha, B., Shao, Z.: Principled scavenging. In: Proc. Programming Language Design and Implementation. (2001) 81–91
22. Hawblitzel, C., Wei, E., Huang, H., Krupski, E., Wittie, L.: Low-level linear memory management. In: Proc. SPACE Workshop. (2004)
23. DeLine, R., Fähndrich, M.: Enforcing high-level protocols in low-level software. In: Proc. Programming Language Design and Implementation. (2001) 59–69
24. Fähndrich, M., DeLine, R.: Adoption and focus: Practical linear types for imperative programming. In: Proc. Programming Language Design and Implementation. (2002) 13–24
25. Morrisett, G., Ahmed, A., Fluet, M.: \mathbf{L}^3 : A linear language with locations. In: Proc. Typed Lambda Calculi and Applications. (2005) 293–307
26. Ahmed, A., Fluet, M., Morrisett, G.: \mathbf{L}^3 : A linear language with locations. Technical Report TR-24-04, Harvard University (2004)