# Simplifying Regions

## Greg Morrisett

### Harvard University

Collaborators:  Matthew Fluet & Amal Ahmed

# The Cyclone Safe-C Project

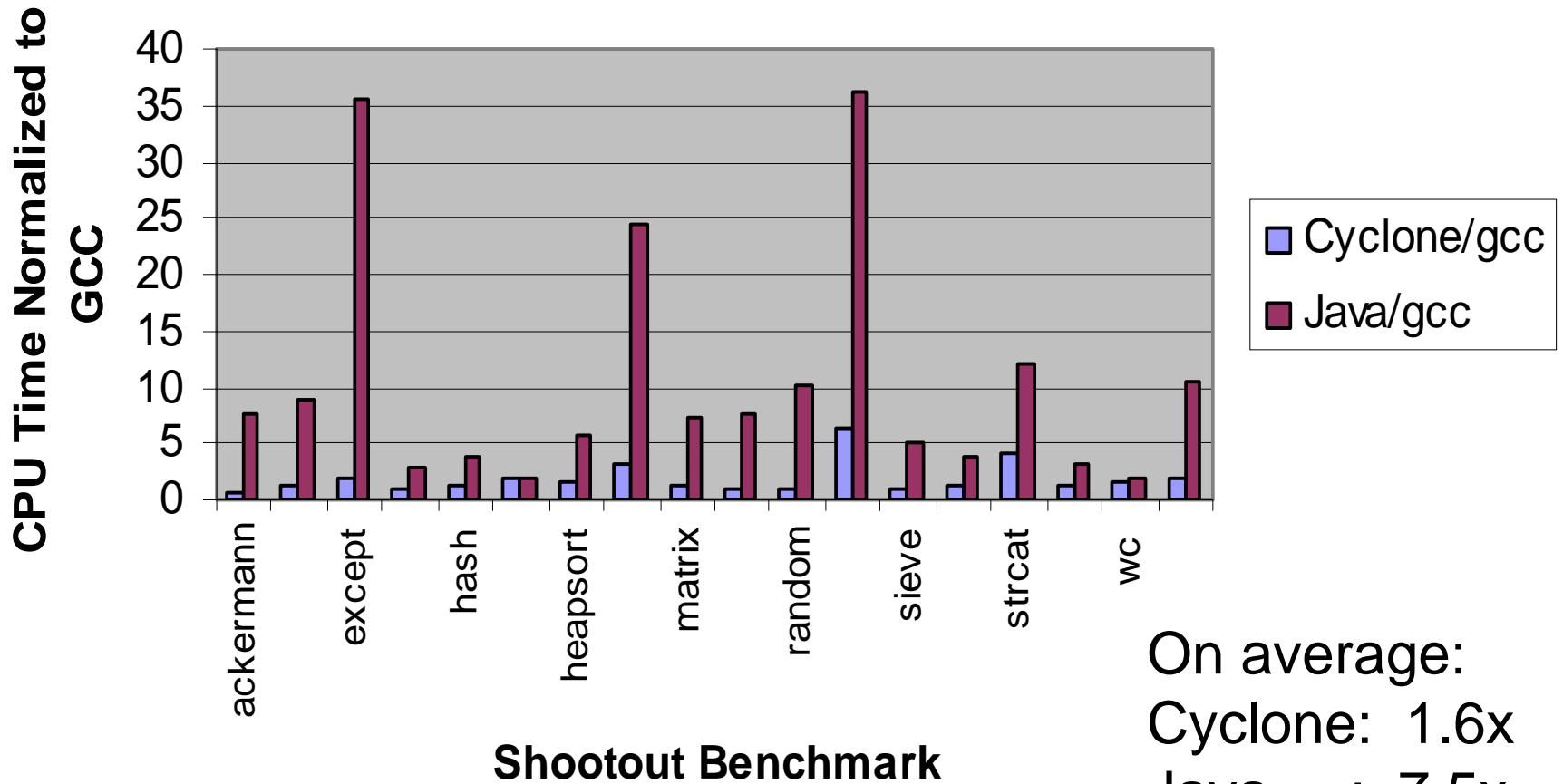Primary goal: type-safety

Secondary goal: retain virtues of C

- C programmers should feel comfortable.

- It should be easy to interoperate with legacy C.

- Most importantly, costs should be manifest:

  - Programmers can understand the physical layout of data structures by looking at the types.

  - Programmers can avoid overheads of run-time tags and checks by programming with certain idioms.

  - Want this to be suitable for real-time and embedded settings where space and time may be scarce.

# Some Cyclone Users

- In-kernel Network Monitoring [Penn]
- MediaNet [Maryland & Cornell]
- Open Kernel Environment [Leiden]
- RBClick Router [Utah]
- xTCP [Utah & Washington]
- Lego Mindstorm on BrickOS [Utah]
- Cyclone on Nintendo DS [AT&T]
  - Scheme run-time & interpreter
- Cyclone compiler, tools, & libraries
  - Over 100 KLOC
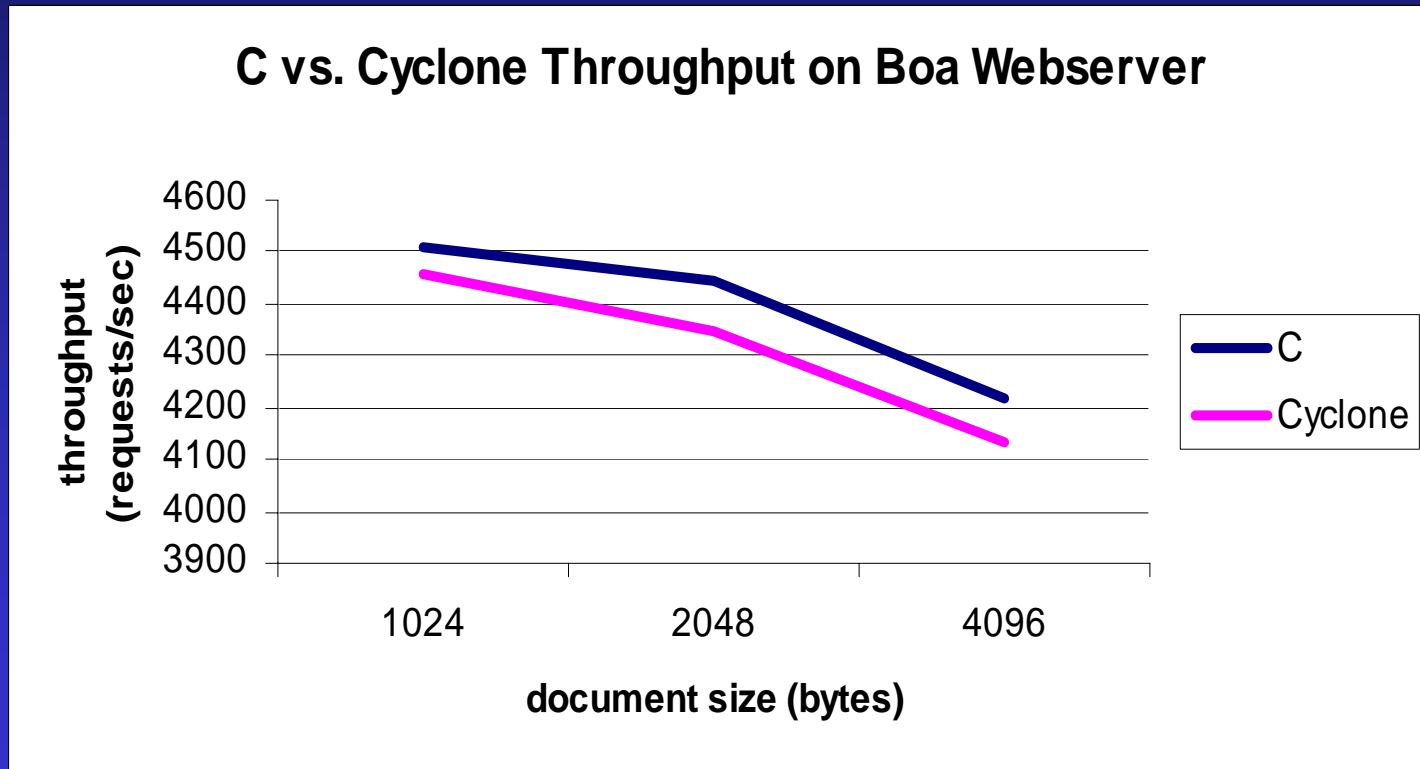  - Plus many sample apps, benchmarks, etc.

# C vs. Cyclone vs. Java

**Cyclone vs. Java**



On average:
Cyclone:  1.6x
Java     :  7.5x

# Macro-benchmarks:

We have also ported a variety of security-critical applications where we see little overhead (e.g., 3% throughput for the Boa Webserver.)

**C vs. Cyclone Throughput on Boa Webserver**

# Memory Management

A range of options:

- Heap allocation with conservative GC
- Lexical Regions
  - Stack allocation
  - Lexical arena allocation
  - Tofte & Talpin + region subtyping
- 1st class Regions
  - Enables "tail-calls" -- can code copying GC
- Unique pointers
  - Enables reclamation of individual objects

Each has different tradeoffs.

# The Flexibility Pays: MediaNET

TTCP benchmark (packet forwarding):

Cyclone v.0.1 (lexical regions & BDW GC)

- High water mark: 840 KB
- 130 collections
- Basic throughput:  50 MB/s

Cyclone v.0.5 (unique ptrs + dynamic regions)

- High water mark:  8 KB
- 0 collections
- Basic throughput: 74MB/s

# A Model?

The combination of lexical regions, unique pointers, region subtyping, etc. makes the meta-theory of Cyclone a nightmare.

- Gave up on usual syntactic proof.

At the heart of the problem:

- Certain types are "ephemeral".
- The interaction between persistent and ephemeral types is extremely subtle.
- Polymorphism really complicates things.
- Same issue arises in many other settings: TAL(T), Vault, Cqual, Haskell's runST, …

# Outline

Core Cyclone $\rightarrow$ F+RGN   [ICFP'04]

- Effects map to an indexed store monad
- Coercion-based interpretation of subtyping

F+RGN $\rightarrow$ Linear F+Stores

- Monad abandoned in favor of linearity.
- Regions become 1st-class, unique pointers fall out as a special case.
- Developing a semantic model of the target.
- Believe it serves as foundation for Cqual, Vault, etc.
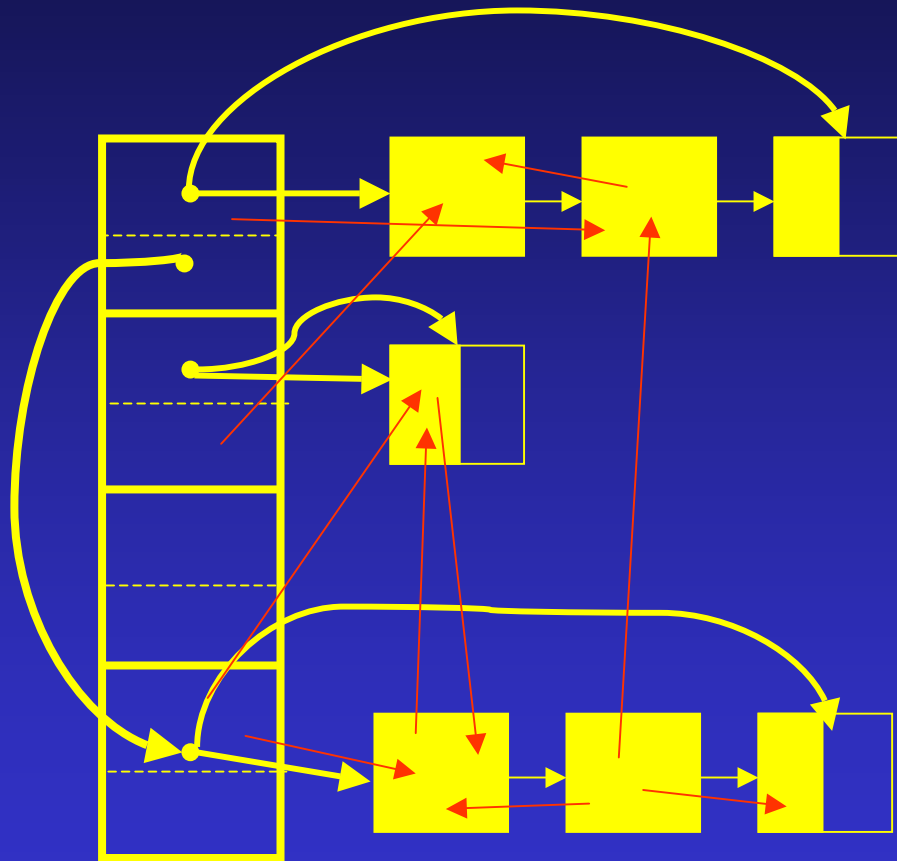
# The Tofte-Talpin Region Calculus

Operationally:

- Memory is divided into *regions* ($\rho$)

- Objects are allocated in a region: $(3,2)@\rho$

- Regions are created and destroyed with a lexically-scoped construct:

$$\textbf{letregion } \rho \textbf{ in } e$$

- All objects allocated in $\rho$ are deallocated at the end of $\rho$'s scope.

- Region names can be passed into functions to support a "callee-allocates in caller's region idiom."

# Runtime Organization

Regions are linked lists of pages.

Arbitrary inter-region references.

Similar to arena-style allocators.

runtime stack

# Typing

- Pointer types indicate referent's region: $(int, int) @ \rho$

- The type system tracks the set $\varphi$ of regions that are accessed when a computation is run: $\Gamma \blacktriangleright e : T, \varphi$

- Function types include a latent effect:

$$T_1 \xrightarrow{\varphi} T_2$$

- The role of $\varphi$ is to tell us when it's *not* safe to deallocate a region.

# Letregion

The typing for `letregion` is subtle:

$$\frac{\Gamma \blacktriangleright e : \tau, \varphi \qquad \rho \notin FRV(\Gamma, \tau)}{\Gamma \blacktriangleright \text{letregion } \rho \text{ in } e : \tau, \varphi \backslash \rho}$$

In particular, pointers into $\rho$ can escape the scope of the letregion.

# Example:

```
letregion ρ in
```
  $\texttt{let } x = (1,2)@ \ \rho \ \texttt{in}$

  $\texttt{let } z = (3,4)@ \ \rho' \ \texttt{in}$

  $\texttt{let } w = (x,z)@ \ \rho' \ \texttt{in}$

  $\lambda y.\#1(\#2 \ w) + y \quad : \quad \text{int} \xrightarrow{\{\rho'\}} \text{int}, \{\rho'\}$

# Example:

**letregion** $\rho$ **in**

  **let** x = (1,2)@ $\rho$ **in**

  **let** z = (3,4)@ $\rho'$ **in**

  **let** w = (x,z)@ $\rho'$ **in**

  $\lambda$y.#1(#2 w) + y    :    int $\xrightarrow{\{\rho'\}}$ int, {$\rho'$}

$\rho'$

# Example:

**letregion** $\rho$ **in**

  **let** $x = (1,2)@\ \rho$ **in**

  **let** $z = (3,4)@\ \rho'$ **in**

  **let** $w = (x,z)@\ \rho'$ **in**

  $\lambda y.\#1(\#2\ w) + y$   :   $int \xrightarrow{\{\rho'\}} int, \{\rho'\}$

# Example:

**letregion** $\rho$ **in**

  <u>**let** $x = (1,2)$ @ $\rho$ **in**</u>

  **let** $z = (3,4)$ @ $\rho'$ **in**

  **let** $w = (x,z)$ @ $\rho'$ **in**

  $\lambda y.\#1(\#2\ w) + y$    :    $\text{int} \xrightarrow{\{\rho'\}} \text{int}, \{\rho'\}$

$\rho'$ → [ ]

$\rho$ → [ (1,2) ]

# Example:

**letregion** $\rho$ **in**

  **let** x = (1,2)@ $\rho$ **in**

  <u>**let** z = (3,4)@ $\rho$' **in**</u>

  **let** w = (x,z)@ $\rho$' **in**

  $\lambda$y.#1(#2 w) + y    :   int $\xrightarrow{\{\rho'\}}$ int, $\{\rho'\}$

$\rho'$ → (3,4)

$\rho$ → (1,2)

# Example:

**letregion** $\rho$ **in**

  **let** $x = (1,2)@ \rho$ **in**

  **let** $z = (3,4)@ \rho'$ **in**

  <u>**let** $w = (x,z)@ \rho'$ **in**</u>

$\lambda y.\#1(\#2\ w) + y \quad : \quad \text{int} \xrightarrow{\{\rho'\}} \text{int}, \{\rho'\}$

# Example:

**letregion** $\rho$ **in**

  **let** $x = (1,2)@\ \rho$ **in**

  **let** $z = (3,4)@\ \rho'$ **in**

  **let** $w = (x,z)@\ \rho'$ **in**

  $\underline{\lambda y.\#1(\#2\ w) + y}$    :    $int \xrightarrow{\{\rho'\}} int, \{\rho'\}$



$\rho'$

(3,4)

$( \bullet , \bullet )$

$\rho$

(1,2)

closure

# Example:

**letregion** $\rho$ **in**

  **let** $x = (1,2)@\ \rho$ **in**

  **let** $z = (3,4)@\ \rho'$ **in**

  **let** $w = (x,z)@\ \rho'$ **in**

  $\lambda y.\#1(\#2\ w) + y \quad : \quad int \xrightarrow{\{\rho'\}} int, \{\rho'\}$

Pointers are persistent, regions aren't…



$\rho'$   (3,4)   $(\bullet, \bullet)$

closure

# Subtyping

Tofte & Talpin's effect weakening:

$$\frac{\Gamma \blacktriangleright e : \tau, \varphi \qquad \varphi \subseteq \varphi'}{\Gamma \blacktriangleright e : \tau, \varphi'}$$

Cyclone's region "outlives":

$$\frac{\Gamma \blacktriangleright \rho \leq \rho'}{\Gamma \blacktriangleright \tau @ \rho \leq \tau @ \rho'}$$

$$\frac{\Gamma, \text{FRV}(\Gamma) \leq \rho \blacktriangleright e : \tau, \varphi \qquad \rho \notin \text{FRV}(\Gamma, \tau)}{\Gamma \blacktriangleright \texttt{letregion } \rho \texttt{ in } e : \tau, \varphi \backslash \rho}$$

# Core Cyclone to F+RGN

The source language is complicated by:

- Effects: sets of regions
- Subtyping, letregion, polymorphism.

Choose as intermediate language:

- CBV System-F plus…
- An indexed monad family: RGN $\sigma$ $\tau$
  - Inspired by Haskell's ST monad.
  - Key: run can be provided in the language.
- Eliminate subtyping via coercions

# Type Constructors

RGN $\sigma\ \tau$

computation running in store $\sigma$ producing a $\tau$.

ptr $\rho\ \tau$

pointer into region $\rho$ holding a $\tau$ value.

$\rho \in \sigma$

a proof that $\sigma$ includes the region $\rho$

$\sigma_1 \le \sigma_2 \qquad [= \mathbf{8}\rho.(\rho \in \sigma_1)\ \mathbf{!}\ (\rho \in \sigma_2)]$

a proof of store inclusion

# Translation Essence:

$$\langle\!\langle \text{int} @ \rho_1 \xrightarrow{\{\rho_1, \rho_2, \rho_3\}} \text{int} @ \rho_3 \rangle\!\rangle \frac{1}{4}$$

$$\forall \sigma. (\rho_1 \in \sigma) \to (\rho_2 \in \sigma) \to (\rho_3 \in \sigma) \to$$
$$(\text{ptr } \rho_1 \text{ int}) \to \text{RGN } \sigma (\text{ptr } \rho_3 \text{ int})$$

# Monadic Operations

**`return`** : $\mathbf{8}\alpha, \sigma.\ \alpha\ !\ \mathrm{RGN}\ \sigma\ \alpha$

**`then`** : $\mathbf{8}\alpha, \beta, \sigma.\ \mathrm{RGN}\ \sigma\ \alpha\ !$

$$(\alpha\ !\ \mathrm{RGN}\ \sigma\ \beta)\ !\ \mathrm{RGN}\ \sigma\ \beta$$

- Can only sequence in *same* store.
- Need some way to lift computations in sub-stores

**`run`** : $\mathbf{8}\alpha.\ (\mathbf{8}\sigma.\ \mathrm{RGN}\ \sigma\ \alpha)\ !\ \alpha$

- Note that $\alpha$ cannot mention $\sigma$!
- Quite similar to letregion.

# Primitives:

**`new`**:

$$8\alpha,\sigma,\rho.\ \alpha \rightarrow (\rho \in \sigma) \rightarrow \text{RGN } \sigma\ (\text{ptr } \rho\ \alpha)$$

**`read`**:

$$8\alpha,\sigma,\rho.\ \text{ptr } \rho\ \alpha \rightarrow (\rho \in \sigma) \rightarrow \text{RGN } \sigma\ \alpha$$

**`letRGN`** :

$$8\alpha,\sigma_1.\ (8\sigma_2.\ (\sigma_1 \leq \sigma_2) \rightarrow (\rho \in \sigma_2) \rightarrow \text{RGN } \sigma_2\ \alpha)$$
$$\rightarrow \text{RGN } \sigma_1\ \alpha$$

**`subRGN`** :

$$8\alpha,\sigma_1,\sigma_2.\ (\sigma_1 \leq \sigma_2) \rightarrow \text{RGN } \sigma_1\ \alpha \rightarrow \text{RGN } \sigma_2\ \alpha$$

# Notes:

We constructed an operational model and proved a soundness result at this level, as well as the correctness of the translation.

In practice, you need to phase-split the evidence (e.g., $\rho \in \sigma$) and coercions.

F+RGN is somewhat simpler than T.T. and sheds light on regions and Haskell's ST, but not 1st class regions or unique pointers.

# New Target:  Linear F + regions

- We'll use a *linear* version of F similar to Walker & Watkins.

- We'll eliminate the RGN monad in favor of explicit store-passing but use linearity to ensure store remains single-threaded.

- Unique pointers & 1st class regions pop out for free…

# Types:

$T ::= \alpha \mid \text{int}$

$\mid \text{ptr } \rho \; T$      (pointer into region $\rho$)

$\mid \text{cap } \rho$      (capability for region $\rho$)

$\mid 1 \mid T_1 \otimes T_2$

$\mid T_1 \multimap T_2$

$\mid !T$

$\mid 8\alpha.T \mid 8\rho.T$

$\mid \exists\alpha.T \mid \exists\rho.T$

# Primitives:

newrgn :  $1 \multimap \exists\rho.\text{cap }\rho$

freergn :  $\forall\rho.\text{cap }\rho \multimap 1$

new : $\forall\alpha,\rho.!\alpha \multimap \text{cap }\rho \multimap \text{cap }\rho \otimes !\text{ptr }\rho\ !\alpha$

read : $\forall\alpha,\rho.\text{ptr }\rho\ !\alpha \multimap \text{cap }\rho \multimap \text{cap }\rho \otimes !\alpha$

# Dynamics

Mostly just CBV lambda calculus.

Semantic values:

- ptr $\rho$ $\tau$ $\approx$ $\text{Loc}_\rho$

- cap $\rho$ $\approx$ $\text{Loc}_\rho \to \text{Val}$

- NB: $!(\text{cap } \rho) \approx \varnothing$

We actually use a step-indexed model a la Appel & McAllester to avoid problems with recursive types.

# Encoding F+RGN Types

$\langle\!\langle \text{int} \urcorner = !\langle\!\langle \text{int} \urcorner$

$\langle\!\langle \text{ptr } \sigma\ \tau \urcorner = !\text{ptr } \sigma\ \langle\!\langle \tau \urcorner$

$\langle\!\langle \tau_1 \ !\ \tau_2 \urcorner = !(\langle\!\langle \tau_1 \urcorner \multimap \langle\!\langle \tau_2 \urcorner)$

$\langle\!\langle \text{RGN } \sigma\ \tau \urcorner = \sigma \multimap \sigma \otimes \langle\!\langle \tau \urcorner$

$\langle\!\langle \rho \in \sigma \urcorner = !\ \exists \sigma'. (\sigma \multimap \sigma' \otimes \text{cap } \rho) \otimes$
$(\sigma' \otimes \text{cap } \rho \multimap \sigma)$

$\langle\!\langle \sigma_1 \leq \sigma_2 \urcorner = !\ \exists \sigma'. (\sigma_2 \multimap \sigma_1 \otimes \sigma') \otimes$
$(\sigma_1 \otimes \sigma' \multimap \sigma_2)$

# Encoding Monadic Primitives:

Just store-passing:

$«\texttt{return}¬ = \Lambda\alpha,\sigma.\ \lambda x{:}!\alpha.\ \lambda s{:}\sigma.\ (s,x)$

$«\texttt{then}¬ = \Lambda\alpha,\beta,\sigma.$
$\qquad \lambda f{:}«RGN\ \sigma\ \alpha¬.$
$\qquad \lambda g{:}!(!\alpha \multimap «RGN\ \sigma\ \beta¬).$
$\qquad \lambda s{:}\sigma.\ \texttt{let}\ (s',y) = f\ s\ \texttt{in}\ g\ y\ s'$

# Encoding Let-region

$«\texttt{letRGN}¬ =$

$\Lambda\alpha,\sigma_1.\lambda f: «8\sigma_2.\ \sigma_1 \leq \sigma_2\ !\ \rho \in \sigma_2\ !\ RGN\ \sigma_2\ \alpha\ ¬.$

$\quad \lambda s:\sigma_1.$

$\quad \texttt{unpack}\ [\rho,c] = newrgn\ ()\ \texttt{in}$

$\quad \texttt{let}\ w_2 = pack[\sigma_1,(id,id)]: «\rho \in (\sigma_1 \otimes cap\ \rho)¬\ \texttt{in}$

$\quad \texttt{let}\ w_1 = pack[cap\ \rho,(id,id)]: «\sigma_1 \leq (\sigma_1 \otimes cap\ \rho)¬$

$\quad \texttt{in}\ \texttt{let}\ ((s,c),x) = f\ [\sigma_1 \otimes cap\ \rho]\ w_1\ w_2\ (s,c)\ \texttt{in}$

$\quad freergn\ c;$

$\quad (s,x)$

Key:  new store is $\sigma_1 \otimes cap\ \rho$

# Encoding New and Read:

Use witnesses to get capability from store:

**«new¬** $= \Lambda\alpha,\sigma,\rho.\ \lambda x{:}!\alpha.\ \lambda w{:}«\rho \in \sigma¬.\lambda s{:}\sigma.$

        **unpack** $[\sigma',(f,g)] = w$ **in**

        **let** $(s',c) = f\ s$ **in**

        **let** $(c,r) = $ new $x\ c$ **in**

        **let** $s = g\ (s',c)$ **in** $(s,r)$

**«read¬** $= \Lambda\alpha,\sigma,\rho.\lambda x{:}\text{ptr}\ \rho\ !\alpha.\ \lambda w{:}«\rho \in \sigma¬.\lambda s{:}\sigma.$

        **unpack** $[\sigma',(f,g)] = w$ **in**

        **let** $(s',c) = f\ s$ **in**

        **let** $(c,x) = $ read $r\ c$ **in**

        **let** $s = g\ (s',c)$ **in** $(s,r)$

# Subrgn

Use witness to get sub-store:

**«subRGN¬ =**

$\Lambda\alpha,\sigma_1,\sigma_2.\ \lambda w{:}\text{«}\sigma_1 \leq \sigma_2\text{¬}.\ \lambda k{:}\text{«}RGN\ \sigma_1\ \alpha\text{¬}.$

$\lambda s_2{:}\sigma_2.$

  **unpack** $[\sigma',(f,g)] = w$ **in**

  **let** $(s_1,s') = f\ s_2$ **in**

  **let** $(s_1,x) = k\ s_1$ **in**

  **let** $s_2 = g\ (s_1,s')$ **in**  $(s_2,x)$

# 1st Class Regions

At the target level, regions are 1st class!

- Can export newrgn & freergn to the source.
- No LIFO constraints needed!
- Source-level 1st class region: $\exists\rho.(\text{cap }\rho \otimes !T[\rho])$

We can *open* such a region to regain the convenience of the monadic threading:

$8\rho.\text{cap }\rho \multimap$

$8\alpha,\sigma_1. \; (8\sigma_2. «\sigma_1 \leq \sigma_2¬ \multimap «\rho \in \sigma_2¬ \multimap «\text{RGN }\sigma_2\ \alpha¬)$
$\multimap \text{RGN }\sigma_1 \; (\text{cap }\rho \otimes \alpha)$

- So the monad is purely a convenience.

# Unique Pointers

These are just a degenerate case of 1$^{st}$ class regions: $\exists\rho.(cap\ \rho \otimes !ptr\ \rho\ \tau)$

We can deallocate these at will!

- In practice, we split cap $\rho$ into two capabilities.
- One (<u>access $\rho$</u>) lets us access $\rho$.
- The other (alloc $\rho$) lets us allocate in $\rho$.
- Only the alloc capability is needed at run-time.
- So a unique pointer is: $\exists\rho.(\underline{access\ \rho} \otimes !ptr\ \rho\ \tau)$
- Can "open" a unique pointer to again regain convenience of monadic abstraction.

# Recap:

- At source-level, we seem to have a variety of memory mgmt. facilities:
  - Stack allocation, lexical regions, 1st class regions, unique pointers, …
  - They're all useful in practice.
- The target exposes the commonalities:
  - Linear capabilities for access control ensure state is single-threaded *and* eventually reclaimed.
  - Monadic encapsulation is purely a convenience (implicit threading of capabilities).
  - That convenience has a price:  LIFO.
  - Fortunately, we don't *have* to encapsulate.

# Future Work:

- Need to fill in all of the details.
- Need to phase-split capabilities.
- In practice, need affine, linear, and unrestricted types to model Cyclone.
- Modeling other languages:
  - Alias types, Cqual: require only a slight refinement where we have two kinds of pointers (ephemeral vs. persistent).
  - Vault: still need to account for adoption and suspect that relevant types play role.