

# A Dynamic Region System for Transforming Effect-dependent Programs

Zhixuan Yang

SOKENDAI and NII

May 2019

# Outline

## Background

Preliminaries: effect theories, a language and a logic

Dynamic region system

Separation guards

Assuming you are an optimising compiler, will you transform

$put\ p\ v; l \leftarrow get\ q; K$       (In C syntax  $*p = v; l = *q; K$ )

to

$l \leftarrow get\ q; put\ p\ v; K$       (In C syntax  $l = *q; *p = v; K$ )

( $put\ p\ v$  writes  $v$  to the cell  $p$ ;  $get\ q$  reads cell  $q$ )

Assuming you are an optimising compiler, will you transform

$put\ p\ v; l \leftarrow get\ q; K$       (In C syntax  $*p = v; l = *q; K$ )

to

$l \leftarrow get\ q; put\ p\ v; K$       (In C syntax  $l = *q; *p = v; K$ )

( $put\ p\ v$  writes  $v$  to the cell  $p$ ;  $get\ q$  reads cell  $q$ )

Only when  $put\ p$  and  $get\ q$  access different cells ( $p \neq q$ ) !

# Pointer Analyses

We need a way to statically track which memory cells a program may access.

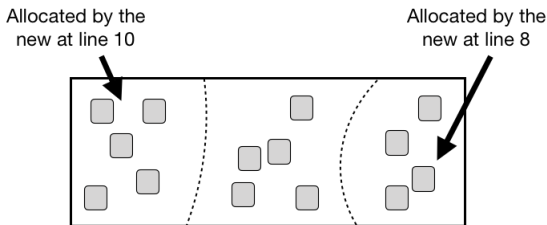
Two orthogonal problems:

1. How to “talk about” memory cells statically—a static abstraction of memory.  
E.g. “ $p$  points to a cell allocated by the *new* at line 10”.
2. How to track which memory cells a pointer variable may point to

# Pointer Analyses in Compilers

Core of static program analysis, intensively studied

- Usually with simple model to question 1: partitioning the memory cells by the *new* call in the code allocating it



- Automated solution to question 2: data-flow analyses, abstract interpretation ...

# Pointer Analyses Beyond Compilers

Program transformations are not limited to compilers.

- ▶ Equational proof: transforming a program to a simple and obviously correct program
- ▶ Done by human: allowing more refined memory models (that are not possible for automated analyses).

# Pointer Analyses Beyond Compilers

Region system (Lucassen and Gifford 1988): the memory is partitioned into *regions* (mentally, by the programmer)

- ▶ Reference (pointer) type  $Ref\ r\ D$  is indexed by a region
- ▶  $new : (r : Region) \rightarrow D \rightarrow Ref\ r\ D$  takes a region argument
- ▶ Use a type system to track which regions a program access

Compared to the line-number-based memory model, we can distinguish memory cells by passing different region argument to the same *new* call.



## Naive model vs. region model

```
1 makeList [] = return Nil
2 makeList (a : as) = { ls  $\leftarrow$  makeList as;
3                       p  $\leftarrow$  new (a, ls);
4                       return (Ptr p) }
5  $l_1 \leftarrow$  makeList a1
6  $l_2 \leftarrow$  makeList a2
```

In the naive model, both  $l_1$  and  $l_2$  point to cells allocated by the *new* at line #3.

# Naive model vs. region model

(Assuming  $r_1 \neq r_2$ )

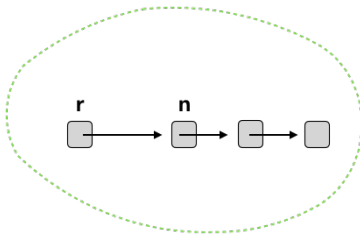
```
1 makeList  $r_1$  [] = return Nil
2 makeList  $r_1$  ( $a : ax$ ) = {  $ls \leftarrow \text{makeList } r \ ax;$ 
3                                $p \leftarrow \text{new } r \ (a, ls)$ 
4                               return (Ptr  $p$ ) }
5  $l_1 \leftarrow \text{makeList } r_1 \ a_1$ 
6  $l_2 \leftarrow \text{makeList } r_2 \ a_2$ 
```

In the region model,  $l_1$  and  $l_2$  point to cells in region  $r_1$  and  $r_2$  respectively.

# Static region is not enough

- ▶ Existing region systems assume memory can be *statically* partitioned.
- ▶ For some pointer-manipulating programs, this assumption does not hold.

```
1 traverse (Ptr r)  
2 traverse (Ptr t)
```

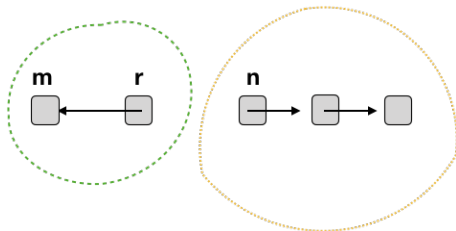


At some place in the code, we may want to put the list *r* in a region.

## Static region is not enough

At some other point in the program, we may want to put the first cell of list  $r$  and the rest of the list in different regions to reason about the program.

```
1  $p \text{ Nil} = \mathbf{return} ()$   
2  $p (Ptr\ r) = \{(a, n) \leftarrow get\ r$   
3            $put\ r\ (a, m)$   
4            $traverse\ n$   
5            $put\ r\ (a, n)\}$ 
```



Thus, a *static* region system is not sufficient.

# A dynamic region system

We wish a more *dynamic* region system such that

- ▶ the region to which a cell belongs can be determined dynamically (by the contents of memory cells)
- ▶ we can track the regions a program accesses

This is what this work wants to do.

# A dynamic region system

Currently, my dynamic region system is very simple. A region is

- ▶ a single memory cell or
- ▶ the *current* reachable closure from a memory cell.

For example, the judgement

$$l : ListPtr\ a \vdash t\ l : \mathbf{1} ! \{get_{r(l)}\}$$

asserts the program  $t\ l$  only reads the linked list starting from  $l$ .

# Separation guard

A complementary construct: *separation guard*

- ▶ it checks some pointers or their reachable closures are disjoint, otherwise terminates the program
- ▶ resembles the precondition of separation logic

For example,

$$l : \text{ListPtr } a, l_2 : \text{Ref } (a, \text{ListPtr } a) \vdash [r(l) * l_2] : \mathbf{FUnit}$$

checks the cell  $l_2$  is not any node of linked list  $l$ .

# Program transformation

Put dynamic region system and separation guard together, we can express and prove some transformations.

## Example

$$\frac{l : \text{ListPtr } a \vdash t \, l : \mathbf{1} ! \{get_{r(l)}\}}{[r(l) * l_2]; put \, l_2 \, v; t \, l \quad = \quad [r(l) * l_2]; t \, l; put \, l_2 \, v}$$



# Outline

Background

Preliminaries: effect theories, a language and a logic

Dynamic region system

Separation guards

# Language

For the purpose of discussion, we fix a small programming language based Levy's *call-by-push-value* calculus with algebraic effects (but no handlers)

Base types:  $\sigma ::= \text{Bool} \mid \text{Unit} \mid \text{Void} \mid \dots$

Value types:  $A ::= \sigma \mid \text{ListPtr } D \mid \text{Ref } D \mid A_1 \times A_2$   
 $\mid A_1 + A_2 \mid \mathbf{U}\underline{A}$

Storable types:  $D ::= \sigma \mid \text{ListPtr } D \mid \text{Ref } D \mid D_1 \times D_2$   
 $\mid D_1 + D_2$

Computation types:  $\underline{A} ::= \mathbf{F}A \mid A_1 \rightarrow \underline{A_2}$

# Syntax

Value terms:	$\begin{aligned} v ::= & x \mid c \mid Nil \mid Ptr\ v \mid (v_1, v_2) \\ & \mid \mathbf{inj}_1^{A_1+A_2}\ v \mid \mathbf{inj}_2^{A_1+A_2}\ v \mid \mathbf{thunk}\ t \end{aligned}$
Computation terms:	$\begin{aligned} t ::= & \mathbf{return}\ v \mid \{x : A \leftarrow t_1; t_2\} \\ & \mid \mathbf{match}\ v\ \mathbf{as}\ \{(x_1, x_2) \rightarrow t\} \\ & \mid \mathbf{match}\ v\ \mathbf{as}\ \{Nil \rightarrow t_1, Ptr\ x \rightarrow t_2\} \\ & \mid \mathbf{match}\ v\ \mathbf{as}\ \{\mathbf{inj}_1\ x_1 \rightarrow t_1, \mathbf{inj}_2\ x_2 \rightarrow t_2\} \\ & \mid \lambda x : A. t \mid t\ v \mid \mathbf{force}\ v \mid op\ v \mid \\ & \mu x : \mathbf{U}\underline{A}. t \end{aligned}$
Operations:	$op ::= fail \mid get \mid put \mid new \mid \dots$

# Effect theories

For the operations in the language, we have equations on them:

- ▶  $l \leftarrow \text{get } p; \text{put } p \ l = \mathbf{return} \ ()$
- ▶  $\text{put } p \ v_1; \text{put } p \ v_2 = \text{put } p \ v_2$
- ▶ ...
- ▶ And this *separation axiom*

$$\begin{array}{lcl} 1 & \{ l_1 \leftarrow \text{new}_D \ v_1 & = & \{ l_1 \leftarrow \text{new}_D \ v_1 \\ 2 & l_2 \leftarrow \text{new}_D \ v_2 & & l_2 \leftarrow \text{new}_D \ v_2 \\ 3 & \mathbf{match} \ l_1 \equiv l_2 \ \mathbf{as} & & t_1 \} \\ 4 & \{ \text{False} \rightarrow t_1 & & \\ 5 & \text{True} \rightarrow t_2 \} \end{array}$$

# Type system and semantics

## Type system

Please refer to Plotkin and Pretnar's *A Logic for Algebraic Effects* paper.

## Semantics

$\llbracket A \rrbracket$  maps to some set.

$\llbracket \underline{A} \rrbracket$  maps to the set of tree whose internal nodes of operations and leaves are  $A$ -values.

$\llbracket \Gamma \vdash t : \tau \rrbracket$  maps to a function  $\llbracket \Gamma \rrbracket \rightarrow \llbracket \tau \rrbracket$ .

# Logic

We have a first-order equational logic for reasoning about programs of this language.

Judgements

$$\Gamma \mid \Psi \vdash \phi$$

$\Gamma$  is a context of the types of free variables and  $\Phi$  is a list of formulas that are the premises of  $\phi$ .

Terms

$$\begin{aligned} \phi ::= & t_1 =_{\text{CBPV}} t_2 \mid t_1 =_{\text{Alg}} t_2 \mid \forall x : A. \phi \mid \forall x : \underline{A}. \phi \\ & \mid \exists x : A. \phi \mid \exists x : \underline{A}. \phi \mid \phi_1 \wedge \phi_2 \mid \phi_1 \vee \phi_2 \\ & \mid \neg \phi \mid \phi_1 \rightarrow \phi_2 \mid \top \mid \perp \end{aligned}$$

# Inference rules

The logic system's inference rules are:

- ▶ Standard rules for classical first order connectives and structural rules for judgements,

# Inference rules

The logic system's inference rules are:

- ▶ standard  $\beta$ - and  $\eta$ - equivalence for CBPV language constructs, for example,

$$\overline{\{x \leftarrow \mathbf{return} \ v; t\} =_{\text{CBPV}} t[v/x]} \qquad \overline{(\lambda x. t) \ v =_{\text{CBPV}} t[v/x]}$$

$$\overline{\mathbf{case} \ True \ \mathbf{of} \ \{ \ True \rightarrow t_1; \ False \rightarrow t_2 \} =_{\text{CBPV}} t_1}$$



# Inference rules

The logic system's inference rules are:

- rules inherited from the effect theories, for example,

$$\frac{}{\{put(l_1, v_1); put(l_2, v_2); t\} =_{\mathbf{Alg}} \{put(l_2, v_2); t\}}$$

# Inference rules

The logic system's inference rules are:

- ▶ algebraicity of effect operations, an inductive principle over computations and a universal property of computation types.

# Sum up of preliminaries

Plotkin's algebraic effects have these layers of concepts:

1. Basic types (e.g. *Nat*, *Bool*, ...)
2. Effect theories (operations and equations)
3. A programming language with effect operations, sequential composition, and (possibly) handlers
4. An equational logic for reasoning about programs of this language.

# Sum up of preliminaries

Plotkin's algebraic effects have these layers of concepts:

1. Basic types (e.g. *Nat*, *Bool*, ...)
2. Effect theories (operations and equations)
3. A programming language with effect operations, sequential composition, and (possibly) handlers
4. An equational logic for reasoning about programs of this language.

The logic is rich and powerful—we will express our region system in this logic.

# Outline

Background

Preliminaries: effect theories, a language and a logic

Dynamic region system

Separation guards

# Region system as logic predicates

Our dynamic region systems are defined as logic predicates on computation terms in the logic.

Let  $op$  range over possible effect operations in the language. We extend the term of the logic:

$$\phi ::= \dots \mid t ! \epsilon$$

$$\epsilon ::= \emptyset \mid \epsilon, op \mid \epsilon, get_{r(v)} \mid \epsilon, put_{r(v)} \mid \epsilon, get_v \mid \epsilon, put_v$$

# Well-formedness

The new term is well-formed when

$$\frac{\Gamma \vdash t : \mathbf{FA}}{\Gamma \vdash t ! \cdot : \mathbf{form}}$$

$$\frac{\Gamma \vdash t ! \epsilon : \mathbf{form}}{\Gamma \vdash t ! \epsilon, op : \mathbf{form}}$$

$$\frac{\Gamma \vdash t ! \epsilon : \mathbf{form} \quad \Gamma \vdash v : Ref\ D}{\Gamma \vdash t ! \epsilon, o_v : \mathbf{form}} \quad (o \in \{get, put\})$$

$$\frac{\Gamma \vdash t ! \epsilon : \mathbf{form} \quad \Gamma \vdash v : ListPtr\ D}{\Gamma \vdash t ! \epsilon, o_{r(v)} : \mathbf{form}} \quad (o \in \{get, put\})$$

# Semantics

We intend the semantics of  $\llbracket \Gamma \vdash t ! \epsilon : \mathbf{form} \rrbracket$  to be

$$\{\gamma \in \llbracket \Gamma \rrbracket \mid \llbracket t \rrbracket(\gamma) \text{ is a computation only invokes operations in } \epsilon\}$$

I am hesitating about the definition—should we interpret *get/put* in  $t$ ?



# Inference rules

$$\frac{}{\Gamma \mid \Psi \vdash \mathbf{return} \ x ! \emptyset} \text{R-PURE} \qquad \frac{\Gamma \mid \Psi \vdash t ! \epsilon \quad \epsilon \subseteq \epsilon'}{\Gamma \mid \Psi \vdash t ! \epsilon'} \text{R-SUB}$$

$$\frac{\Gamma \mid \Psi \vdash t =_{\text{CBPV}} t' \wedge t' ! \epsilon}{\Gamma \mid \Psi \vdash t ! \epsilon} \text{R-EQ}$$

And for any  $op : A \rightarrow \underline{B} \in \epsilon$  that is not  $get/put_{r(v)}$ ,

$$\frac{\Gamma, a : B \mid \Psi \vdash k ! \epsilon}{\Gamma \mid \Psi \vdash (a \leftarrow op \ v; k) ! \epsilon} \text{R-OP}$$

# Inference rules

Intuitively,

if  $\epsilon$  contains  $get_{r(v)}$  or  $put_{r(v)}$ , the program can read or write the cells linked from  $v : ListPtr\ D$ .

- ▶ When  $v = Nil$ , the program get no cells to access from  $r(v)$ .
- ▶ When  $v = Ptr\ v'$ , the program can read or write the cell  $v'$ ,
  - ▶ and if it reads it by  $(a, n) \leftarrow get\ v'$ , its allowed operation on  $r(v)$  is inherited by  $r(n)$  and  $v'$

# Inference rules

If  $get_{r(v)} \in \epsilon$

$$\frac{\begin{array}{l} \Gamma \mid \Psi \vdash t \text{ Nil} ! \epsilon \setminus \{get_{r(v)}, put_{r(v)}\} \\ \Gamma, v' \mid \Psi \vdash t (Ptr \ v') =_{\text{CBPV}} \{(a, n) \leftarrow get \ v'; k\} \\ \Gamma, v', a, n \mid \Psi, \ t \ n ! \epsilon[r(n)/r(v)] \vdash k ! \epsilon[r(n)/r(v)] \cup \epsilon[v'/r(v)] \end{array}}{\Gamma \mid \Psi \vdash t \ v ! \epsilon}$$

If  $put_{r(v)} \in \epsilon$

$$\frac{\begin{array}{l} \Gamma \mid \Psi \vdash t \text{ Nil} ! \epsilon \setminus \{get_{r(v)}, put_{r(v)}\} \\ \Gamma \mid \Psi \vdash t (Ptr \ v') =_{\text{CBPV}} \{put \ v' \ c; k\} \quad \Gamma \mid \Psi \vdash k ! \epsilon[v'/r(v)] \end{array}}{\Gamma \mid \Psi \vdash t \ v ! \epsilon}$$

## Theorem (Soundness)

*If  $\Gamma \mid \Psi \vdash t ! \epsilon$ , then  $\llbracket \Psi \rrbracket \subseteq \llbracket t ! \epsilon \rrbracket$ .*

Proof.

To do



# Outline

Background

Preliminaries: effect theories, a language and a logic

Dynamic region system

Separation guards

# Separation guards

- ▶ Our dynamic region systems proves a program only operates on certain memory cells

$$t_1 ! \{get_{r(l_1)}, put_{r(l_1)}\} \quad \text{and} \quad t_2 ! \{get_{r(l_2)}, put_{r(l_2)}\}$$

- ▶ This information is useful only when we can also shows the cells that two programs respectively operates on are disjoint

$$r(l_1) \cap r(l_2) = \emptyset \implies t_1; t_2 = t_2; t_1$$

- ▶ Ultimately, disjointness comes from the separation axiom of *new*, but it is too primitive for practical use.

# Separation guards

- ▶ Our dynamic region systems proves a program only operates on certain memory cells

$$t_1 ! \{get_{r(l_1)}, put_{r(l_1)}\} \quad \text{and} \quad t_2 ! \{get_{r(l_2)}, put_{r(l_2)}\}$$

- ▶ This information is useful only when we can also shows the cells that two programs respectively operates on are disjoint

$$r(l_1) \cap r(l_2) = \emptyset \implies t_1; t_2 = t_2; t_1$$

- ▶ Ultimately, disjointness comes from the separation axiom of *new*, but it is too primitive for practical use.

We introduce *separation guards* for tracking disjointness more easily at a higher level.

Following separation logic, we write  $\phi = l_1 * l_2 * \dots * l_n$  to denote that cells described by  $l_i$  are disjoint.

Each  $l_i$  can be either a value  $v$  of type  $\text{Ref } D$ , or  $r(v)$  for a value  $v$  of type  $\text{ListPtr } D$ .

A separation guard  $[\phi]$  is a computation of type  $\mathbf{FUnit}$ .



# Definition

$[\phi]$  is defined as

```
1  $[\phi] = sepChk \ \phi \ \emptyset$ 
2  $sepChk \ [] \ s = \mathbf{return} \ ()$ 
3  $sepChk \ (v * \phi) \ x = \mathbf{if} \ l \in x \ \mathbf{then} \ fail \ \mathbf{else} \ sepChk \ \phi \ (x \cup l)$ 
4  $sepChk \ (r(v) * \phi) \ x = \{x' \leftarrow chkList \ v \ x; sepChk \ \phi \ x'\}$ 
5  $chkList \ Nil \ x = \mathbf{return} \ x$ 
6  $chkList \ (Ptr \ p) \ x = \mathbf{if} \ p \in x$ 
7  $\qquad \qquad \qquad \mathbf{then} \ \{fail; \mathbf{return} \ ()\}$ 
8  $\qquad \qquad \qquad \mathbf{else} \ \{(-, n) \leftarrow get \ p; chkList \ n \ (x \cup p)\}$ 
```

$[\phi]$  checks regions  $l_i$  of  $\phi$  are distinct.

# Inference rules

Define  $c \langle a =_{\text{Alg}} b \rangle$  to be  $(c; a) =_{\text{Alg}} (c; b)$ .

We have the following inference rules for separation guards

$$\frac{}{[\phi] \langle \text{new } t =_{\text{Alg}} (l \leftarrow \text{new } t; [\phi * l]; \text{return } l) \rangle}$$

$$\frac{\Gamma \mid \Psi, l_1 \neq l_2 \vdash t_1 =_{\text{Alg}} t_2}{\Gamma \mid \Psi \vdash [l_1 * l_2] \langle t_1 =_{\text{Alg}} t_2 \rangle}$$

$$\frac{}{\text{return } Nil =_{\text{Alg}} (l \leftarrow \text{return } Nil; [r(l)]; \text{return } l)}$$

$$\frac{}{[r(l)] \langle \text{new } (Cell \ a \ l) =_{\text{Alg}} (l' \leftarrow \text{new } (Cell \ a \ l); [r(l')]; \text{return } l') \rangle}$$

# Inference rules

The rule for list induction

$$\frac{\text{base case} \quad \text{inductive case}}{\Gamma \mid \Psi \vdash [r(l) * \phi] \langle t_1 =_{\text{Alg}} t_2 \rangle} \quad (\text{LISTIND})$$

- ▶ base case is  $\Gamma \mid \Psi, l =_{\text{Alg}} \text{Nil} \vdash [\phi] \langle t_1 =_{\text{Alg}} t_2 \rangle$
- ▶ inductive case is

$$\Gamma \mid l = \text{Ptr } l', \text{hyp} \vdash \\ ((\text{Cell } \_, n) \leftarrow \text{get } l'; [l' * r(n) * \phi]) \langle t_1 =_{\text{Alg}} t_2 \rangle$$

$$\text{and hyp} =_{\text{def}} [r(n) * \phi] \langle t_1 =_{\text{Alg}} t_2 \rangle$$

# Soundness

## Theorem

*The inference rules for separation guards are sound.*

## Proof.

To do. It'll be a large verifying proof.



# Program transformations

Now we have enough ingredients to express the program transformations I wanted:

A frame rule:

$$\frac{\Gamma \mid \Psi \vdash t_1 ! \overline{\phi_1} \quad \Gamma \mid \Psi \vdash [\phi_1] \langle t_1 =_{\text{Alg}} t_2 \rangle}{\Gamma \mid \Psi \vdash [\phi_1 * \phi_2] \langle t_1 =_{\text{Alg}} (t_2; [\phi_2]) \rangle}$$

Commutativity lemma

$$\frac{\Gamma \mid \Psi \vdash t_i ! \overline{\phi_i} \ (i = 1, 2)}{\Gamma \mid \Psi \vdash [\phi_1 * \phi_2] \langle (t_1; t_2) =_{\text{Alg}} (t_2; t_1) \rangle}$$

# Program transformations

Now we have enough ingredients to express the program transformations I wanted:

A frame rule:

$$\frac{\Gamma \mid \Psi \vdash t_1 ! \overline{\phi_1} \quad \Gamma \mid \Psi \vdash [\phi_1] \langle t_1 =_{\text{Alg}} t_2 \rangle}{\Gamma \mid \Psi \vdash [\phi_1 * \phi_2] \langle t_1 =_{\text{Alg}} (t_2; [\phi_2]) \rangle}$$

Commutativity lemma

$$\frac{\Gamma \mid \Psi \vdash t_i ! \overline{\phi_i} \ (i = 1, 2)}{\Gamma \mid \Psi \vdash [\phi_1 * \phi_2] \langle (t_1; t_2) =_{\text{Alg}} (t_2; t_1) \rangle}$$

Proof: induction on the derivation of region predicates and using the inference rules of separation guards.

# Conclusion

- ▶ We proposed static methods to track
  - ▶ which memory cells a program access, and
  - ▶ if these cells are disjoint.
- ▶ Technical parts are remained to be completed and polished.