

# Typestate Inference for Mungo: Algorithm and Implementation

---

Hans Hüttel – Department of Computer Science, Aalborg University, Denmark – with

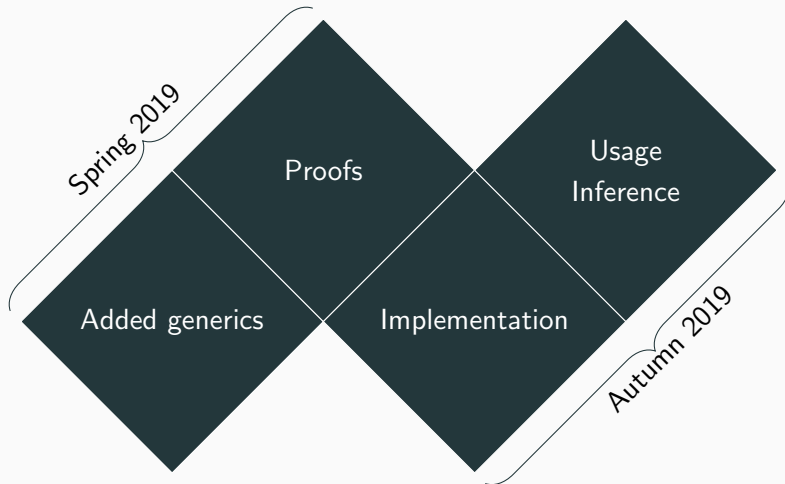
Iaroslav Golovanov   Mikkel Kettunen   Mathias Jakobsen

2 March 2020

Mungo is a typed Java-like language in the style of Featherweight Java of Igarashi et al. that contains usual object-oriented and imperative constructs.

Mungo is also the associated programming tool developed at Glasgow University by Dardha, Gay et al.

## Our contributions to Mungo



## Introduction: Typestates

Typestate definitions model dynamic program properties by letting types have content that can give an overapproximation of reachable program states.

Typestates allow us to specify the allowed order of operations.

In our setting, we can statically verify that

- Null-dereferencing does not occur
- Protocols are not violated

## Example: Protocol Error

```
1 public static void main(String[] args) {  
2     List list = new LinkedList();  
3     list.add(" List");  
4     list.add(" Modification");  
5     list.add(" Example");  
6  
7     ListIterator iterator = list.listIterator();  
8     iterator.next();  
9     iterator.remove();  
10    iterator.set(" Updating");  
11 }
```

Illegal state exception at line 10

Iterator should be invalidated after call to remove

## Introduction: Usages

Commonly, pre- and postconditions are used to define typestates

- On variables in Strom & Yemini
- On methods in most OOP contexts

Usages are another approach to typestate definitions.

A usage defines a global typestate instead of annotating each method or variable.

## Introduction: Usages cont.

Usages are similar to finite-state automata. We allow three behavioural constructs:

- Branching  $\{m_i; w_i\}_{i \in I}^{\vec{E}}$
- Choice  $\langle l_i : u_i \rangle_{i \in L}^{\vec{E}}$
- Recursion  $X^{\vec{E} \uplus \{X=u\}}$

$$\begin{aligned} \text{Usage} &= \{ \text{init} ; X \}^{\vec{E}} \\ \vec{E} &= \{ X = \{ \text{test} ; \left\langle \begin{array}{l} \text{T} : \{ \text{doStuff} ; X \} \\ \text{F} : X \end{array} \right\rangle \text{ stop} ; \text{end} \} \} \end{aligned}$$

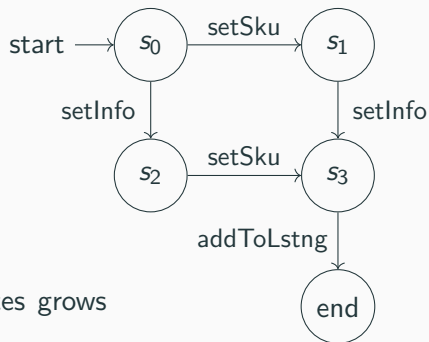
Each object is associated with a typestate consisting of a class and a usage

## Introduction: Example

Let Listing be a class with method addProduct with two parameters

Consider the following Product class:

```
1  class Product {  
2      String sku; PInfo info;  
3  
4      void setSku(String x) { sku = x; }  
5      void setInfo(PInfo x) { info = x; }  
6  
7      Listing addToLstng(Listing l) {  
8          l.addProduct(info, sku);  
9          return l;  
10     }  
11 }
```



Note that the number of reachable usage states grows exponentially with each additional linear field.



# Introduction: Inference Motivation

Annotating each class with usages has some disadvantages:

- Usages can be large and trivial
- Costly to adopt
- Problems with maintainability
- Cannot check external modules

```
1  class Product{  
2      {setSku; {setInfo; X}  
3          setInfo; {setSku; X}}  
4      [X = {addToLstng; end}]  
5  
6      ...  
7  }
```

## Introduction: Inference Motivation cont.

It would be good to be able infer usages when protocol specification is not important:

- Statically ensures no null-dereferencing errors
- Reduces overhead
- Better maintainability
- Open source external modules can be checked



# Usage Inference: Problem

## Problem

Infer usages for classes when not explicitly specified

## Principal Usage

A usage that can simulate every other usage that well-types the class

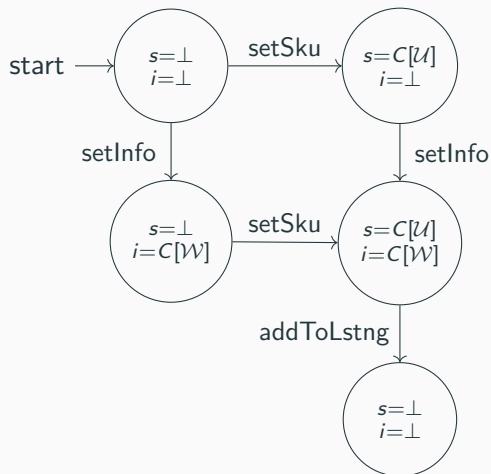
## Usage subtyping $\mathcal{U} \sqsubseteq \mathcal{U}'$ or $(\mathcal{U}, \mathcal{U}') \in R$

$R$  is a usage simulation iff for all  $(\mathcal{U}_1, \mathcal{U}_2) \in R$  we have that:

1. If  $\mathcal{U}_1 \xrightarrow{m} \mathcal{U}'_1$  then  $\mathcal{U}_2 \xrightarrow{m} \mathcal{U}'_2$  such that  $(\mathcal{U}'_1, \mathcal{U}'_2) \in R$
2. If  $\mathcal{U}_1 \xrightarrow{l} \mathcal{U}'_1$  then  $\mathcal{U}_2 \xrightarrow{l} \mathcal{U}'_2$  such that  $(\mathcal{U}'_1, \mathcal{U}'_2) \in R$

# Typesystem: Intuition

1. Start from the environment where all fields are null
2. Check that following the usage type-checks given the current environment
3. Finally, make sure that the environment is terminated when the usage is



Type judgements for classes are of the form

$$\Theta; env T_F \vdash_{\vec{D}} C[\mathcal{U}] \triangleright env T'_F$$

## Typesystem: Definition

(TCBR) makes sure that when we have a branching usage, all possible branches are well typed w.r.t. the field typing environment

(TCCH) checks that all labels leads to well-typed behaviour and same resulting environment

$$\begin{array}{c} \text{(TCBR)} \frac{\begin{array}{c} I \neq \emptyset \\ \text{terminated}(t_i'') \quad t_i \ m_i(t_i' \ x_i) \{e_i\} \in C.\text{methods}_{\vec{D}} \quad \Theta; \text{env}T_F'' \vdash_{\vec{D}} C[u_i^{\vec{E}}] \triangleright \text{env}T_F' \end{array}}{\Theta; \text{env}T_F \vdash_{\vec{D}} C[\{m_i; u_i\}_{i \in I}^{\vec{E}}] \triangleright \text{env}T_F'} \\ \text{(TCCH)} \frac{\forall l_i \in L . \Theta; \text{env}T_F \vdash_{\vec{D}} C[u_i^{\vec{E}}] \triangleright \text{env}T_F'}{\Theta; \text{env}T_F \vdash_{\vec{D}} C[\langle l_i : u_i \rangle_{l_i \in L}^{\vec{E}}] \triangleright \text{env}T_F'} \end{array}$$

## Typesystem: Definition

(TCALLF) type checks a method call of a field

It does so by making sure that a method call is available at the current type

$$\text{(TCALLF)} \frac{\begin{array}{c} \Lambda; \Delta \cdot (o, S) \vdash e : t \triangleright \Lambda' \{o.f \mapsto C[\mathcal{U}]\}; \Delta' \cdot (o, S') \\ t' \ m(t \ x)\{e'\} \in C.\text{methods}_{\vec{D}} \qquad \mathcal{U} \xrightarrow{m} \mathcal{W} \end{array}}{\Lambda; \Delta \cdot (o, S) \vdash f.m(e) : t' \triangleright \Lambda' \{o.f \mapsto C[\mathcal{W}]\}; \Delta' \cdot (o, S')}$$

## Inference Algorithm: Intuition

We need to do what (TCBR) and (TCC<sub>H</sub>) typing does in reverse

1. Start from the initial field typing environment
2. Typecheck all method bodies with the field typing environment and filter those from that leads to errors
3. Now continue to typecheck all new field typing environment
4. Build the usage based on the graph generated

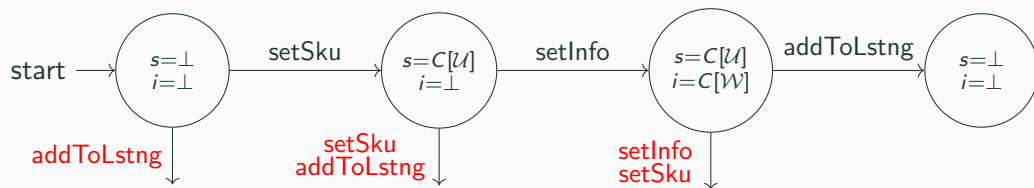


# Inference Algorithm: Intuition

Creating a method availability graph for the previous example

Graph shows a valid sequence of method calls

The red transitions are method calls that lead to errors



## Inference Algorithm: Rules

We define a transition relation between field typing environments.

A transition only exists if the method call is well typed.

$$(\text{CLASS}) \frac{}{envT_F \xrightarrow{m} envT'_F}$$

## Inference Algorithm: Rules

We define a transition relation between field typing environments.

A transition only exists if the method call is well typed.

$$\{\text{this} \mapsto \text{env}T_F\}; \emptyset \cdot (\text{this}, [x \mapsto t]) \vdash^{\emptyset} e : t' \triangleright \{\text{this} \mapsto \text{env}T'_F\}; \emptyset \cdot (\text{this}, [x \mapsto t''])$$

(CLASS)

---

$$\text{env}T_F \xrightarrow{m} \text{env}T'_F$$

## Inference Algorithm: Rules

We define a transition relation between field typing environments.

A transition only exists if the method call is well typed.

$$\text{(CLASS)} \frac{\begin{array}{c} \{ \text{this} \mapsto \text{env}T_F \}; \emptyset \cdot (\text{this}, [x \mapsto t]) \vdash^\emptyset e : t' \triangleright \{ \text{this} \mapsto \text{env}T'_F \}; \emptyset \cdot (\text{this}, [x \mapsto t'']) \\ t' \text{ } m(t \text{ } x) \{e\} \in C.\text{methods} \end{array}}{\text{env}T_F \xrightarrow{m} \text{env}T'_F} \neg \text{lin}(t'')$$

## Inference Algorithm: Rules

$S$  is a set of the environments which can lead to a terminated environment.

$$S = \{envT_F \mid envT_{F\perp} \rightarrow^* envT_F \wedge envT_F \rightarrow^* envT_{F\perp}\} \cup \{\text{end}\}$$

The usage graph allows us to reach the terminated usage end.

$$(\text{TRANS}) \frac{envT_F \xrightarrow{m} envT'_F}{envT_F \xRightarrow{m} envT'_F}$$

$$(\text{END}) \frac{envT_F \xrightarrow{m} envT_{F\perp}}{envT_F \xRightarrow{m} \text{end}}$$

## Inference Algorithm: Rules

The usage is created by creating a usage variable and usage for each field typing environment

---

```
1: function INFER( $S, A, \Rightarrow, envT_{F\perp}$ )
2:   function REACH( $envT_F$ )
3:     return  $\{envT'_F \mid \exists m \in A. envT_F \xRightarrow{m} envT'_F\}$ 
4:   if REACH( $envT_{F\perp}$ ) =  $\emptyset$  then
5:     return end $\emptyset$ 
6:    $E \leftarrow \emptyset$ 
7:   for all  $envT_F \in S \setminus \{\text{end}\}$  do
8:      $E \leftarrow E \cup \{X_{envT_F} = \text{CREATESTATE}(envT_F)\}$ 
9: return  $X_{envT_{F\perp}}^E$ 
```

---

# Inference Algorithm: Rules

---

```
1: function CREATESTATE( $envT_F$ )
2:    $u \leftarrow \emptyset$ 
3:   for all  $m \in A$  do
4:     for all  $s \in S$  do
5:       if  $envT_F \xrightarrow{m} s$  then
6:         if  $s = \text{end}$  then
7:            $u \leftarrow u \cup \{m; \text{end}\}$ 
8:           if  $L\ m(-\ x)\{-\} \in C.\text{methods}$  then
9:              $u \leftarrow u \cup \{m; \langle l_i : \text{end} \rangle_{l_i \in L}\}$ 
10:        else
11:           $u \leftarrow u \cup \{m; X_s\}$ 
12:          if  $L\ m(-\ x)\{-\} \in C.\text{methods}$  then
13:             $u \leftarrow u \cup \{m; \langle l_i : X_s \rangle_{l_i \in L}\}$ 
14:   return  $u$ 
```

---

1. Find all transitions from  $envT_F$
2. Create a branch for a transition
3. If the method returns a label, add a choice usage

Haskell implementation of type system and inference module in `mungoi`

Implemented for only the formalised subset of Java

Motivation: Verify complexity in real-life situations



## Implementation: Example

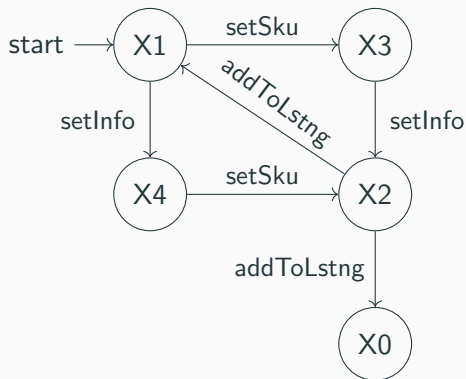
We introduce the special usage variable `infer` to indicate inference

```
1  class Product {
2      infer []
3
4      String sku
5      PInfo info
6
7      void setSku(String[initialised] x) { sku = x }
8      void setInfo(PInfo[initialised] x) { info = x }
9      Listing[initialised] addToLstng(Listing[uninitialised] l) {
10         l.addPinfo(info); l.addSku(sku);
11         l
12     }
13 }
```

## Implementation

The example can be inferred, and is also well-typed

```
1 $ stack run Product.mg
2 X1[ X0 = end
3     X1 = { setSku; X3
4           setInfo; X4 }
5     X2 = { addToLstng; X0
6           addToLstng; X1 }
7     X3 = { setInfo; X2 }
8     X4 = { setSku; X2 } ]
9 new Right ()
```



### **Theorem (Principal Usage Inference)**

*Let  $C$  be a class and  $\mathcal{U}_I$  be an inferred usage, then  $\mathcal{U}_I$  is a principal usage for  $C$ .*

### **Proof.**

By showing that the inferred usage is the largest and that it makes the class well typed □

Size of usages are bounded by

$$O((|e| + |f| \cdot |\mathcal{U}|)^2 + 2^{2 \cdot |f| \cdot |\mathcal{U}|} \cdot |m|) \approx O(2^{2 \cdot |f| \cdot |\mathcal{U}|} \cdot |m|)$$

## **Worst case assumption**

Every usage state can transition to every other usage state, with every method call

While technically possible, it is a huge overestimation for actual programs

Large usages are plausible for usual programs

Classes with  $n$  unrelated linear fields will have usages of size  $O(2^n)$

Non-determinism does not seem to be useful often

## Practical Complexity

```
1  class SmartHomeController {
2      DoorController dc;
3      TemperatureController tc;
4      ElectronicsController ec;
5      [...]
6      void initTempController(TemperatureController c) {
7          this.tc = c;
8          this.tc.initialise();
9      }
10     [...]
11 }
```

Large usages are plausible for usual programs

Classes with  $n$  unrelated linear fields will have usages of size  $O(2^n)$

Non-determinism does not seem to be useful often

Inference makes usages more practical

We infer the largest (*principal*) usage

- Infer once, use everywhere

Is the extra effort and complexity worth it?