

# Verifying Pointer and String Analyses with Region Type Systems

Lennart Beringer<sup>a,\*</sup>, Robert Grabowski<sup>b</sup>, Martin Hofmann<sup>b,\*</sup>

<sup>a</sup>*Princeton University, 35 Olden Street, Princeton 08540, New Jersey, USA*

<sup>b</sup>*Ludwig-Maximilians-Universität, Oettingenstrasse 67, 80538 München, Germany*

---

## Abstract

Pointer analysis statically approximates the heap pointer structure during a program execution in order to track heap objects or to establish alias relations between references, and usually contributes to other analyses or code optimizations. In recent years, a number of algorithms have been presented that provide an efficient, scalable, and yet precise pointer analysis. However, it is unclear how the results of these algorithms compare to each other semantically.

In this paper, we present a general region type system for a Java-like language and give a formal soundness proof. The system is subsequently specialized to obtain a platform for embedding the results of various existing context-sensitive pointer analysis algorithms, thereby equipping the computed relations with a common interpretation and verification. We illustrate our system by outlining an extension to a string value analysis that builds on pointer information.

*Key words:* pointer alias analysis, formal verification, region type system, object-oriented languages, string analysis

---

## 1. Introduction

Pointer (or points-to) analysis is a static program analysis technique that determines an over-approximation of possible points-to relations that occur during the execution of a program. More precisely, it chooses an abstraction of the pointers and references, and computes which pointers may possibly point to which data. A conservative approximation of this structure is also an alias analysis, as the computed points-to relation directly includes the information which pointers may point to the same object. A pointer analysis is often used for compiler optimizations, but may also serve as the basis for other analyses, such as the computation of possible string values in order to prevent string-based security holes.

---

\*Corresponding author

*Email addresses:* `eberinge@cs.princeton.edu` (Lennart Beringer),  
`robert.grabowski@ifi.lmu.de` (Robert Grabowski), `hofmann@ifi.lmu.de` (Martin Hofmann)

There exists a large number of pointer analysis algorithms [1, 2, 3] for different languages. Each algorithm faces the trade-off between precision and efficiency of the analysis: it should choose the right abstractions in order to produce as much useful information as possible while at the same time being able to process large code bases in a reasonable time. These analyses have different techniques, implementations, and complexities. Especially BDD-based algorithms [4, 5] have been shown to be very efficient and precise at the same time.

While several of these analyses also consider soundness, it appears that there does not yet exist a uniformly agreed-upon formal framework that encompasses the interpretations of at least a substantial subset of the analyses. We argue that such a unifying treatment is important, for theoretical as well as pragmatic and practical reasons. First, it is a basis for fair comparisons regarding the precision, flexibility or expressivity of different analyses, the theoretical complexity of the associated algorithms, and their experimental evaluation on common benchmarks. Second, once the analyses agree on the formal property they guarantee, we can safely replace one analysis by another in a compiler or verification tool. Third, a uniform guarantee provides the basis for the formal verification of security-relevant properties that rely on pointer analysis results, as is required in proof-carrying code scenarios.

The first purpose of the present paper is to provide such a framework for Java-like languages, given by a hierarchy of region-based type systems for a language in the style of Featherweight Java [6]. Uniformity (i.e. semantic agreement) is guaranteed by equipping the bottom layer in the hierarchy with a formal interpretation and soundness result. We derive the higher levels by specializing this bottom layer to move towards calculi representing concrete analyses.

Second, we demonstrate that a number of existing pointer analyses for object-oriented programs are based on abstraction disciplines that arise as specializations of a generic parametrized refinement of our base-level type system. We focus on disciplines that specify the abstraction of references and execution points [7], and therefore enable different forms of field-sensitive and context-sensitive analyses. For example, objects may be abstracted to their allocation sites and their class, and execution points may be abstracted by receiver-object or call-site stack contexts.

At the more applied levels of our hierarchy, the selection of a particular analysis is as much influenced by algorithmic aspects such as decidability and efficiency of type inference or checking as by the formal expressivity and soundness of the underlying abstraction. Our third contribution consists of showing how the parametrized type system can be reformulated algorithmically to yield a type checking algorithm. Thanks to the hierarchical structure of our framework, we thus immediately obtain algorithms for automatically validating the correctness of the results of concrete analyses, as long as the results have been interpreted in the framework by instantiating its parameters. As we only consider the results, our verification is also independent of implementation details of concrete analysis algorithms.

Finally, we apply our framework to the analysis of string values, in order

to lay the foundations for eliminating SQL injections and cross-site scripting attacks. We extend the language to a simple string model, and outline how data flow analyses for possible string values that build on pointer analyses [8, 9] can be verified with the correspondingly extended type system.

*Related work.* We include concepts of context-sensitivity to distinguish different analyses of the same method implementation. In particular,  $k$ -CFA [10] is a higher-order control flow analysis where contexts are call stack abstractions of finite length  $k$ . The  $k$ -CFA mainly addresses control flows in functional languages, where functions are first-class values. It requires a combination of value (data flow) analysis and control flow analysis to approximate the possible lambda abstractions that an expression may evaluate to. The similar dynamic dispatch problem for methods in object-oriented languages is easier to solve, as the possible implementation targets of a method invocation can be retrieved from the class information.  $k$ -CFA has been extended with polyvariant types for functions [11], such that different types can be used for the function at different application sites. In our type system, we borrow this concept of polyvariance.

Hardekopf and Lin [12] transform variables that are not address-taken into SSA form, such that running a flow-insensitive analysis on these converted variables has the effect of running a flow-sensitive analysis on the original variables. Our framework assumes SSA-transformed code input, presented in the form of a functional program. Identifying opportunities for SSA transformation, which is a central concern of [12], can thus be seen as a preprocessing phase for our framework to apply.

This paper uses the notion of “regions” in the sense of Lucassen and Gifford [13], i.e. as representations of disjoint sets of memory locations. Equivalently, regions partition or color the memory. In the literature, this disjointness is used to show that certain memory manipulations do not influence other parts of a program, in order to e.g. show semantic equivalences [14], enable safe garbage collection [15], or infer properties for single objects by tracking unique objects in a region [16, 17]. In the pointer analysis setting presented here, regions are simply seen as abstract memory locations that summarize one or more concrete locations, thereby helping to discover may-alias relations. We do not consider uniqueness or must-alias information, and do not aim to justify garbage collection.

Paddle [18] and Alias Analysis Library [19] are implementation frameworks that embed concrete pointer analyzers by factoring out different parameters, as is done here. However, the works do not aim at a formal soundness result. Indeed, they embed the *algorithms* into a common *implementation* framework, whereas our type system embeds the *results* of such algorithms into a common *semantic* framework.

An extended abstract of this work has been published earlier in the proceedings of the 16th LPAR conference (Spring 2010) [20]. The present paper does not have new results but complements the short version with detailed rule systems, proofs, and definitions.

*Synopsis.* The next section introduces the FJEU language and its semantics. We introduce the base-level region type system and the soundness proof in Section 3. In Section 4, we specialize the type system to a parametrized version, such that abstraction principles found in pointer analyses can be modeled explicitly as instantiations of the parameters, and present a type-checking algorithm. Section 5 extends the system, such that results of a string analysis based on pointer analysis can also be verified. We give a variant of the type system with stronger proof rules for precise conditionals in Section 6, and conclude the work with a discussion in Section 7. Formalizations of the soundness proofs in Isabelle/HOL and other accompanying material can be found on our website [21].

## 2. Featherweight Java with Updates

We examine programs of the language FJEU [22], a simplified formal model of the sequential fragment of Java that is relevant for pointer analysis. FJEU extends Featherweight Java (FJ) [6] with attribute updates, such that programs may have side effects on a heap. Object constructors do not take arguments, but initialize all fields with *null*, as they can be updated later. Also, the language adds **let** constructs and conditionals to FJ.

### 2.1. Preliminaries

We write  $\mathcal{P}^{fin}(X)$  for the set of all finite subsets of  $X$ , and  $\mathcal{P}(X)$  for the set of all finite and infinite subsets of  $X$ . The notations  $A \rightarrow B$  and  $A \dashrightarrow B$  stand for the set of total and partial functions from  $A$  to  $B$ , respectively. We write  $[x \mapsto v]$  for the partial function that maps  $x$  to  $v$  and is else undefined. The function  $f[x \mapsto v]$  is equal to  $f$ , except that it returns  $v$  for  $x$ . Both function notations may be used in an indexed fashion, e.g.  $f[x_i \mapsto v_i]_{\{1, \dots, n\}}$ , to define multiple values. In typing rules, we sometimes write  $x : v$  and  $f, x : v$  corresponding to the above notation. Finally, we write  $\bar{a}$  for a sequence of entities  $a$ .

### 2.2. Syntax

The following table summarizes the (infinite) abstract identifier sets in the language, the meta-variables we use to range over them, and the syntax of FJEU expressions:

$$\begin{array}{ll}
 \text{variables: } x, y \in \mathcal{X} & \text{classes: } C, D, E \in \mathcal{C} \\
 \text{fields: } f \in \mathcal{F} & \text{methods: } m \in \mathcal{M}
 \end{array}$$

$$\mathcal{E} \ni e ::= \text{null} \mid x \mid \text{new } C \mid \text{let } x = e \text{ in } e \mid x.f \mid x.f := y \mid x.m(\bar{y}) \mid \\
 \text{if } x \text{ instanceof } C \text{ then } e \text{ else } e \mid \text{if } x = y \text{ then } e \text{ else } e$$

To keep our calculus minimal and focused on pointer alias analysis, we omit primitive data types such as integers or booleans. However, such data types and their operations can easily be added to the language. We also omit type casts found in FJ, as they do not provide new insights to pointer analysis. Nevertheless, casts can be included in a straight-forward manner. In order to simplify the

proofs, we require programs to be in let normal form. The somewhat unusual conditional constructs for dynamic class tests and value equality are included to have reasonable if-then-else expressions in the language while avoiding the introduction of booleans.

An FJEU program is defined by the following relations and functions:

$$\begin{array}{lll}
\text{subclass relation:} & \prec & \in \mathcal{P}^{fin}(\mathcal{C} \times \mathcal{C}) \\
\text{field list:} & \mathit{fields} & \in \mathcal{C} \rightarrow \mathcal{P}^{fin}(\mathcal{F}) \\
\text{method list:} & \mathit{methods} & \in \mathcal{C} \rightarrow \mathcal{P}^{fin}(\mathcal{M}) \\
\text{method table:} & \mathit{mtable} & \in \mathcal{C} \times \mathcal{M} \rightarrow \mathcal{E} \\
\text{FJEU program:} & P & = (\prec, \mathit{fields}, \mathit{methods}, \mathit{mtable})
\end{array}$$

FJEU is a language with nominal subtyping:  $D \prec C$  means  $D$  is an immediate subclass of  $C$ . The relation is well-formed if it is a tree successor relation; multiple inheritance is not allowed. We write  $\preceq$  for the reflexive and transitive hull of  $\prec$ . The functions  $\mathit{fields}$  and  $\mathit{methods}$  describe for each class  $C$  the fields and methods of objects of class  $C$ . The functions are well-formed if for all classes  $C$  and  $D$  such that  $D \preceq C$ ,  $\mathit{fields}(C) \subseteq \mathit{fields}(D)$  and  $\mathit{methods}(C) \subseteq \mathit{methods}(D)$ , i.e. classes inherit fields and methods from their superclasses. A method table  $\mathit{mtable}$  gives for each class and each method identifier its implementation, i.e. the FJEU expression that forms the body of the method. To simplify the presentation, we assume that formal argument variables in the body of a method  $m$  are named  $x_1^m, x_2^m$ , etc., abbreviated to  $\overline{x^m}$ , besides the implicit and reserved variable  $\mathit{this}$ . All free variables of an implementation of  $m$  must be from the set  $\{\mathit{this}, x_1^m, x_2^m, \dots\}$ . A method table is well-formed if  $\mathit{mtable}(C, m)$  is defined whenever  $m \in \mathit{methods}(C)$ . In other words, all methods declared by  $\mathit{methods}$  must be implemented, though the implementation may be overridden in subclasses for the same number of formal parameters. In the following, we assume a fixed FJEU program  $P$  whose components are all well-formed.

### 2.3. Semantics

A state consists of a store (variable environment or stack) and a heap (memory). Stores map variables to values, while heaps map locations to objects. An object consists of a class identifier and a valuation of its fields. The only kinds of values in FJEU are locations and *null* references.

$$\begin{array}{lll}
\text{locations:} & l & \in \mathcal{L} \\
\text{values:} & v & \in \mathcal{V} = \mathcal{L} \cup \{\mathit{null}\} \\
\text{objects:} & (C, F) & \in \mathcal{O} = \mathcal{C} \times (\mathcal{F} \rightarrow \mathcal{V}) \\
\text{stores:} & s & \in \mathcal{X} \rightarrow \mathcal{V} \\
\text{heaps:} & h, k & \in \mathcal{L} \rightarrow \mathcal{O}
\end{array}$$

The semantics of an FJEU program is defined as a standard big-step relation  $(s, h) \vdash e \Downarrow v, h'$ , which means that an FJEU expression  $e$  evaluates in store  $s$  and heap  $h$  to the value  $v$  and modifies the heap to  $h'$ . Figure 1 shows the defining rules of the operational semantics. A premise involving a partial function, like  $s(x) = l$ , implies the side condition  $x \in \text{dom}(s)$ .

$$\boxed{(s, h) \vdash e \Downarrow v, h'}$$

$$\frac{}{(s, h) \vdash \mathbf{null} \Downarrow \mathit{null}, h} \qquad \frac{}{(s, h) \vdash x \Downarrow s(x), h}$$

$$\frac{l \notin \mathit{dom}(h) \quad F = [f \mapsto \mathit{null}]_{f \in \mathit{fields}(C)}}{(s, h) \vdash \mathbf{new} C \Downarrow l, h[l \mapsto (C, F)]}$$

$$\frac{(s, h) \vdash e_1 \Downarrow v_1, h_1 \quad (s[x \mapsto v_1], h_1) \vdash e_2 \Downarrow v_2, h_2}{(s, h) \vdash \mathbf{let} x = e_1 \mathbf{in} e_2 \Downarrow v_2, h_2}$$

$$\frac{s(x) = l \quad h(l) = (D, -) \quad D \preceq C \quad (s, h) \vdash e_1 \Downarrow v, h'}{(s, h) \vdash \mathbf{if} x \mathbf{instanceof} C \mathbf{then} e_1 \mathbf{else} e_2 \Downarrow v, h'}$$

$$\frac{s(x) = l \quad h(l) = (D, -) \quad D \not\preceq C \quad (s, h) \vdash e_2 \Downarrow v, h'}{(s, h) \vdash \mathbf{if} x \mathbf{instanceof} C \mathbf{then} e_1 \mathbf{else} e_2 \Downarrow v, h'}$$

$$\frac{s(x) = s(y) \quad (s, h) \vdash e_1 \Downarrow v, h'}{(s, h) \vdash \mathbf{if} x = y \mathbf{then} e_1 \mathbf{else} e_2 \Downarrow v, h'}$$

$$\frac{s(x) \neq s(y) \quad (s, h) \vdash e_2 \Downarrow v, h'}{(s, h) \vdash \mathbf{if} x = y \mathbf{then} e_1 \mathbf{else} e_2 \Downarrow v, h'} \qquad \frac{s(x) = l \quad h(l) = (C, F)}{(s, h) \vdash x.f \Downarrow F(f), h}$$

$$\frac{s(x) = l \quad h(l) = (C, F) \quad h' = h[l \mapsto (C, F[f \mapsto s(y)])]}{(s, h) \vdash x.f := y \Downarrow s(y), h'} \qquad \frac{s(x) = l \quad h(l) = (C, -) \quad |\overline{x^m}| = |\overline{y}| = n \quad s' = [\mathit{this} \mapsto s(x)] \cup [x_i^m \mapsto s(y_i)]_{i \in \{1, \dots, n\}} \quad (s', h) \vdash \mathit{mtable}(C, m) \Downarrow v, h'}{(s, h) \vdash x.m(\overline{y}) \Downarrow v, h'}$$

Figure 1: Operational semantics of FJEU

#### 2.4. Class Tables

A class table  $\mathfrak{C}_0 = (A_0, M_0)$  models FJ's standard type system, where types are simply classes. The *field typing*  $A_0 : (\mathcal{C} \times \mathcal{F}) \rightarrow \mathcal{C}$  assigns to each class  $C$  and each field  $f \in \text{fields}(C)$  the class of the field, which is required to be invariant with respect to subclasses of  $C$ . The *method typing*  $M_0 : (\mathcal{C} \times \mathcal{M}) \rightarrow \overline{\mathcal{C}} \times \mathcal{C}$  assigns to each class  $C$  and each method  $m \in \text{methods}(C)$  a *method type*, which specifies the classes of the formal argument variables and of the result value. It is required to be contravariant in the argument classes and covariant in the result class with respect to subclasses of  $C$ . Our discipline is thus slightly more permissive than the one of Java where arguments types are required to be invariant.

#### 2.5. Example Program

As a very small example, look at the Java program in Figure 2 that implements a list via `Nil` and `Cons` classes, with a `copy` method that duplicates the list structure, but not the actual elements. The same program in FJEU and its class table is shown in Figure 3, where  $\epsilon$  denotes the empty sequence.

### 3. Region Type System

In this section we define the base region type system, which serves as a main unifying calculus for pointer analysis and is given an interpretation and soundness proof. We assume an infinite set  $\mathcal{R}$  of *regions*  $r$ , which are abstract memory locations. Each region stands for zero or more concrete locations. Different regions represent disjoint sets of concrete locations, hence they partition or *color* the memory. Two pointers to different regions can therefore never alias.

#### 3.1. Refined Types and Subtyping

The region type system is a refinement of the plain type system: we equip classes with (possibly infinite) subsets from  $\mathcal{R}$ . For example, a location  $l$  is typed with the *refined type*  $C_{\{r,s\}}$  if it points to an object of class  $C$  (or a subclass of  $C$ ), and if  $l$  is abstracted to either  $r$  or  $s$ , but no other region. The *null* value can be given any type, while the type of locations must have a non-empty region set. The following table summarizes the definitions:

$$\begin{array}{ll} \text{regions:} & r, s, t \in \mathcal{R} \\ \text{region sets:} & R, S, T \in \mathcal{P}(\mathcal{R}) \\ \text{refined types:} & \sigma, \tau \in \mathcal{T} = \mathcal{C} \times \mathcal{P}(\mathcal{R}) \end{array}$$

In the following, we use the notation  $C_R$  instead of  $(C, R)$  for types. Though the region identifier  $s$  is already used for variable stores, the difference should be clear from the context. Since region sets are an over-approximation of the possible locations where an object resides, we can easily define a subtyping relation  $<$ : based on set inclusion:

$$C_R <: D_S \iff R \subseteq S \wedge C \preceq D$$

```

class Elem {}
class List { List copy() { return new List(); } }
class Nil extends List { Nil copy() { return new Nil(); } }
class Cons extends List {
  Elem elem; List next;
  Cons copy() { Cons res = new Cons(); res.elem = elem;
               res.next = next.copy(); return res; } }

```

Figure 2: List copy example in Java

$$\begin{aligned}
P &= (\preceq, \text{fields}, \text{methods}, \text{mtable}) \\
\preceq &= \{(Nil, List), (Cons, List)\} \\
\text{fields}(Elem) &= \emptyset \\
\text{fields}(List) &= \emptyset \\
\text{fields}(Nil) &= \emptyset \\
\text{fields}(Cons) &= \{elem, next\} \\
A_0(Cons, elem) &= Elem \\
A_0(Cons, next) &= List \\
M_0(List, copy) &= (\epsilon, List) \\
M_0(Nil, copy) &= (\epsilon, Nil) \\
M_0(Cons, copy) &= (\epsilon, Cons) \\
\\
\text{methods}(Elem) &= \emptyset \\
\text{methods}(List) &= \{copy\} \\
\text{methods}(Nil) &= \{copy\} \\
\text{methods}(Cons) &= \{copy\} \\
\text{mtable}(List, copy) &= \text{new List} \\
\text{mtable}(Nil, copy) &= \text{new Nil} \\
\text{mtable}(Cons, copy) &= \text{let } res = \text{new Cons in} \\
&\quad \text{let } copyTail = next.copy() \text{ in} \\
&\quad \text{let } _ = res.elem := elem \text{ in} \\
&\quad \text{let } _ = res.next := copyTail \text{ in } res
\end{aligned}$$

Figure 3: List copy example in FJEU

We also extend subtyping to method types:

$$\begin{aligned}\bar{\sigma} <: \bar{\tau} &\iff |\bar{\sigma}| = |\bar{\tau}| \wedge \forall i \in 1, \dots, |\bar{\sigma}|. \sigma_i <: \tau_i \\ (\bar{\sigma}, \tau) <: (\bar{\sigma}', \tau') &\iff \bar{\sigma}' <: \bar{\sigma} \wedge \tau <: \tau'\end{aligned}$$

### 3.2. Annotated Class Tables

We extend class tables  $\mathfrak{C}_0$  to *annotated class tables*  $\mathfrak{C} = (A^{get}, A^{set}, M)$ .

- The *annotated field typings*  $A^{get}, A^{set} : (\mathcal{C} \times \mathcal{R} \times \mathcal{F}) \rightarrow \mathcal{T}$  assign to each class  $C$ , region  $r$ , and field  $f \in \text{fields}(C)$  the refined type of the field for all objects of class  $C$  in region  $r$ . The field type is split into a covariant get-type  $A^{get}$  for the data read from the field, and a contravariant set-type  $A^{set}$  that is needed for data to be written to the field. This technique improves precision and is borrowed from Hofmann and Jost [22]. More formally, annotated field typings are *well-formed* if for all classes  $C$ , subclasses  $D \preceq C$ , regions  $r$  and fields  $f \in \text{fields}(C)$ ,

- $A^{set}(C, r, f) <: A^{get}(C, r, f)$ , and
- $A^{get}(D, r, f) <: A^{get}(C, r, f)$  and  $A^{set}(C, r, f) <: A^{set}(D, r, f)$ .

Also, the class component of  $A^{get}(C, r, f)$  and  $A^{set}(C, r, f)$  must be the class specified in the unannotated table, i.e.  $A_0(C, f)$ .

- The *annotated method typing*  $M : (\mathcal{C} \times \mathcal{R} \times \mathcal{M}) \rightarrow \mathcal{P}(\bar{\mathcal{T}} \times \mathcal{T})$  assigns to each class  $C$ , region  $r$ , and method  $m \in \text{methods}(C)$  an unbounded number of refined method types for objects of class  $C$  in region  $r$ , enabling infinite polymorphic method types. This makes it possible to use a different type at different invocation sites (program points) of the same method. Even more importantly, the same invocation site can be checked in different type derivations with different method types. For every *well-formed* annotated method type, there must be an improved method type in each subclass: for all classes  $C$ , subclasses  $D \preceq C$ , regions  $r$ , and methods  $m \in \text{methods}(C)$ , we require

- $\forall (\bar{\sigma}, \tau) \in M(C, r, m). \exists (\bar{\sigma}', \tau') \in M(D, r, m). (\bar{\sigma}', \tau') <: (\bar{\sigma}, \tau)$ .

Again, the class components of the refined types  $M(C, r, m)$  have to match the classes of the underlying unannotated method type  $M_0(C, m)$ .

In the following, we assume a fixed annotated class table  $\mathfrak{C}$  with well-formed field and method typings.

### 3.3. Proof Rules

The type system (see Figure 4) derives judgments  $\Gamma \vdash e : \tau$ , meaning FJEU expression  $e$  has type  $\tau$  with respect to a *variable context* (store typing)  $\Gamma : \mathcal{X} \rightarrow \mathcal{T}$  that maps variables to types.

$$\begin{array}{c}
\text{T-SUB} \frac{\Gamma \vdash e : \sigma \quad \sigma <: \tau}{\Gamma \vdash e : \tau} \quad \text{T-LET} \frac{\Gamma \vdash e_1 : \sigma \quad \Gamma, x : \sigma \vdash e_2 : \tau}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \tau} \\
\\
\text{T-VAR} \frac{}{\Gamma, x : \tau \vdash x : \tau} \quad \text{T-NULL} \frac{}{\Gamma \vdash \text{null} : \tau} \\
\\
\text{T-IFINST} \frac{x \in \text{dom}(\Gamma) \quad \Gamma \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash \text{if } x \text{ instance of } C \text{ then } e_1 \text{ else } e_2 : \tau} \\
\\
\text{T-IFEQ} \frac{\Gamma, x : C_{R \cap S}, y : D_{R \cap S} \vdash e_1 : \tau \quad \Gamma, x : C_R, y : D_S \vdash e_2 : \tau}{\Gamma, x : C_R, y : D_S \vdash \text{if } x = y \text{ then } e_1 \text{ else } e_2 : \tau} \\
\\
\text{T-NEW} \frac{}{\Gamma \vdash \text{new } C : C_{\{r\}}} \\
\\
\text{T-INVOKE} \frac{\forall r \in R. \exists (\bar{\sigma}', \tau') \in M(C, r, m). (\bar{\sigma}', \tau') <: (\bar{\sigma}, \tau)}{\Gamma, x : C_R, \bar{y} : \bar{\sigma} \vdash x.m(\bar{y}) : \tau} \\
\\
\text{T-GETF} \frac{\forall r \in R. A^{\text{get}}(C, r, f) <: \tau}{\Gamma, x : C_R \vdash x.f : \tau} \quad \text{T-SETF} \frac{\forall r \in R. \tau <: A^{\text{set}}(C, r, f)}{\Gamma, x : C_R, y : \tau \vdash x.f := y : \tau}
\end{array}$$

Figure 4: Region type system

The rule T-SUB is used to obtain weaker types for the expression. The rules T-LET, T-VAR, and T-IFINST are standard. The rule T-IFEQ exploits the fact that the two variables must point to the same object (or be *null*) in the **then** branch, therefore the intersection of the region sets can be assumed. In T-NULL, the *null* value may have any type (any class and any region set including the empty one). In the rule T-NEW, we may choose any region  $r$ , which is an abstract location that includes (possibly among others) the concrete location of the object allocated by this expression.

When a field is read (T-GETF), we look up the type of the field in the  $A^{\text{get}}$  table. As the variable  $x$  may point to a number of regions, we need to ensure that  $\tau$  is an upper bound of the get-types of  $f$  over all  $r \in R$ . In contrast, when a field is written (T-SETF), the written value must have a subtype of the types allowed for that field by the  $A^{\text{set}}$  table with respect to each possible region  $r \in R$ . Finally, the rule T-INVOKE requires that for all regions  $r \in R$  where the receiver object  $x$  may reside, there must exist a method typing that is suitable for the argument and result types.

An FJEU program  $P = (\prec, \text{fields}, \text{methods}, \text{mtable})$  is *well-typed* if for all classes  $C$ , regions  $r$ , methods  $m$  and method types  $(\bar{\sigma}, \tau)$  such that  $(\bar{\sigma}, \tau) \in M(C, r, m)$ , the following judgment is derivable:

$$[\text{this} \mapsto C] \cup [x_i^m \mapsto \sigma_i]_{i \in \{1, \dots, |\overline{x^m}|\}} \vdash \text{mtable}(C, m) : \tau$$

The polymorphic method types make the region type system very expressive in terms of possible analyses of a given program. Each method may have many types, each corresponding to a derivation of the respective typing judgment. In different derivations, different regions may be chosen for new objects, and different types may be chosen for called methods. This flexibility provides the basis for our embedding of external pointer analyses. Moreover, since there may be infinitely many method types for each method, this system is equivalent to one which allows infinite unfoldings of method calls.

### 3.4. Interpretation

We now give a formal interpretation of the typing judgment in form of a soundness theorem. Later, we shall restrict the expressivity of the system and reformulate the rules in order to move towards an actual type checking algorithm. The underlying idea is that we only have to prove soundness once for the general system; for later systems, it suffices to show that they are special cases of the general system, such that the soundness theorem applies to these systems as well.

A *heap typing*  $\Sigma, \Pi : \mathcal{L} \rightarrow (\mathcal{C} \times \mathcal{R})$  assigns to heap locations a static class (an upper bound of the actual class found at that location) and a region. Heap typings, a standard practice in type systems for languages with dynamic memory allocation [23, chapter 13.4], separate the well-typedness definitions of locations in stores and objects from the actual heap, thereby avoiding the need for a co-inductive definition for well-typed heaps in the presence of cyclic structures. Heap typings map locations to very specific types, namely those where the region set is a singleton. A heap typing thus partitions a heap into (disjoint) regions.

We define a typing judgment for values  $\Sigma \vdash v : \tau$ , which means that according to heap typing  $\Sigma$ , the value  $v$  may be typed with  $\tau$ . In particular, the information in  $\Sigma(l)$  specifies the type of  $l$ . Also, the typing judgment of locations is lifted to stores and variable contexts.

$$\frac{}{\Sigma \vdash \text{null} : \tau} \qquad \frac{\Sigma(l) = (C, r)}{\Sigma \vdash l : C_{\{r\}}} \qquad \frac{\Sigma \vdash v : \sigma \quad \sigma <: \tau}{\Sigma \vdash v : \tau}$$

$$\Sigma \vdash s : \Gamma \iff \forall x \in \text{dom}(\Gamma). \Sigma \vdash s(x) : \Gamma(x)$$

A heap  $h$  is *well-typed* with respect to a heap typing  $\Sigma$  and implicitly a field typing  $A^{get}$ , written  $h \models \Sigma$ , if the types for all locations given by  $\Sigma$  are actually “valid” with respect to the classes of the objects, and if the field values are well-typed with respect to  $A^{get}$  and  $\Sigma$ :

$$h \models \Sigma \iff \forall l \in \text{dom}(\Sigma). l \in \text{dom}(h) \wedge \Sigma \models h(l) : \Sigma(l)$$

where

$$\Sigma \models (C, F) : (D, r) \iff C \preceq D \wedge \text{dom}(F) = \text{fields}(C) \wedge \forall f \in \text{fields}(C). \Sigma \vdash F(f) : A^{get}(C, r, f)$$

As the memory locations are determined at runtime, the heap typings cannot be derived statically. Instead, our interpretation of the typing judgment  $\Gamma \vdash e : \tau$  states that whenever a well-typed program is executed on a heap that is well-typed with respect to some typing  $\Sigma$ , then the final heap after the execution is well-typed with respect to some possibly larger heap typing  $\Pi$ . The typing  $\Pi$  may be larger to account for new objects that may have been allocated during execution, but the type of locations that already existed in  $\Sigma$  may not change. More formally, a heap typing  $\Pi$  *extends* a heap typing  $\Sigma$ , written  $\Pi \sqsupseteq \Sigma$ , if  $\text{dom}(\Sigma) \subseteq \text{dom}(\Pi)$  and  $\forall l \in \text{dom}(\Sigma). \Sigma(l) = \Pi(l)$ .

**Theorem 1** (Soundness Theorem). *Fix a well-typed FJEU program  $P$ . For all  $\Sigma, \Gamma, \tau, s, h, e, v, k$  with*

$$\Gamma \vdash e : \tau \quad \text{and} \quad \Sigma \vdash s : \Gamma \quad \text{and} \quad (s, h) \vdash e \Downarrow v, k \quad \text{and} \quad h \models \Sigma$$

*there exists some  $\Pi \sqsupseteq \Sigma$  such that*

$$\Pi \vdash v : \tau \quad \text{and} \quad k \models \Pi.$$

*Proof.* By induction over derivation of the operational semantics and the typing judgment. We have also formalized the proof in Isabelle/HOL [21].

We first consider the case where  $\Gamma \vdash e : \tau$  has been derived by the subtyping rule. By rule inversion, we get  $\Gamma \vdash e : \sigma$  and  $\sigma <: \tau$ . With this (smaller) derivation of the typing judgment, we can apply the theorem inductively, and get  $\Pi \vdash v : \sigma$  and  $k \models \Pi$  for some  $\Pi \sqsupseteq \Sigma$ . As  $\sigma <: \tau$ , we can deduce  $\Pi \vdash v : \tau$ .

In the following, we assume that the typing judgment has not been derived by a subtyping rule, and continue with a case distinction over the possible forms of the big-step semantics relation.

- $(s, h) \vdash x \Downarrow s(x), h$ .

Then we have  $\Gamma, x : \tau \vdash x : \tau$  and  $k = h \ v = s(x)$ ,  $\Gamma = [x \mapsto \tau]$ . With  $\Pi = \Sigma$ , we get  $k \models \Pi$ . As  $\Pi \vdash s : (\Gamma, x : \tau)$ , we can deduce  $\Pi \vdash v : \tau$ .

- $(s, h) \vdash \text{null} \Downarrow \text{null}, h$ .

Then we have  $\Gamma \vdash \text{null} : \tau$  and  $v = \text{null}$  and  $k = h$ . With  $\Pi = \Sigma$ , we get  $\Pi \vdash v : \tau$  by definition of well-typed values, and  $k \models \Pi$  follows directly from the assumption.

- $(s, h) \vdash \text{new } C \Downarrow l, h[l \mapsto (C, F)]$ .

Then  $\Gamma \vdash \text{new } C : C_{\{r\}}$  and  $v = l$  and  $k = h[l \mapsto (C, F)]$  where  $l \notin \text{dom}(h)$  and  $F = [f \mapsto \text{null}]_{f \in \text{fields}(C)}$ . Since  $\text{dom}(\Sigma) \subseteq \text{dom}(h)$ , we have  $l \notin \text{dom}(\Sigma)$ . We choose  $\Pi = \Sigma[l \mapsto (C, r)]$ , and thus have  $\Pi \vdash l : C_{\{r\}}$ . To show  $k \models \Pi$ , it suffices to show  $\Pi \models k(l) : C_{\{r\}}$ . This holds trivially, as  $k(l) = (C, F)$  and  $C \preceq C$  and  $F(f) = \text{null}$  for all  $f \in \text{dom}(F)$ .

- $(s, h) \vdash \text{let } x = e_1 \text{ in } e_2 \Downarrow v_2, h_2$ .

Then the typing judgment is  $\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \tau$ . Then by rule inversion of the typing and semantics rule, we get  $\Gamma \vdash e_1 : \sigma$  and  $\Gamma, x : \sigma \vdash e_2 : \tau$ , as well as  $(s, h) \vdash e_1 \Downarrow v_1, h_1$  and  $(s[x \mapsto v_1], h_1) \vdash e_2 \Downarrow v_2, h_2$ . By induction, we get  $\Sigma_1 \vdash v_1 : \sigma$  and  $h_1 \models \Sigma_1$  for some  $\Sigma_1 \sqsupseteq \Sigma$ . As  $s$  can include at most locations in  $\text{dom}(\Sigma)$  and since  $\Sigma_1$  is a disjoint extension of  $\Sigma$ , we also get  $\Sigma_1 \vdash s : \Gamma$ , and thereby  $\Sigma_1 \vdash s[x \mapsto v_1] : \Gamma[x \mapsto \sigma]$ . Again by induction, we get  $\Sigma_2 \vdash v_2 : \tau$  and  $h_2 \models \Sigma_2$  for some  $\Sigma_2 \sqsupseteq \Sigma_1$ . We can transitively deduce  $\Sigma_2 \sqsupseteq \Sigma$ , and have proved the case for  $\Pi = \Sigma_2$ .

- $(s, h) \vdash \text{if } x \text{ instanceof } C \text{ then } e_1 \text{ else } e_2 \Downarrow v, h'$ .

The typing judgment is  $\Gamma \vdash \text{if } x \text{ instanceof } C \text{ then } e_1 \text{ else } e_2 : \tau$ . As there are two rules for the semantics of if-then-else constructs, we either get  $(s, h) \vdash e_1 \Downarrow v, h'$  or  $(s, h) \vdash e_2 \Downarrow v, h'$  by rule inversion. Without loss of generality, we only treat the first case here. By inversion of the typing rule, we get  $\Gamma \vdash e_1 : \tau$ . By induction, there exists a  $\Pi \sqsupseteq \Sigma$  such that  $\Pi \vdash v : \tau$  and  $h' \models \Pi$ .

- $(s, h) \vdash \text{if } x = y \text{ then } e_1 \text{ else } e_2 \Downarrow v, h'$ .

The typing judgment is  $\Gamma, x : C_R, y : D_S \vdash \text{if } x = y \text{ then } e_1 \text{ else } e_2 : \tau$ . Let  $s(x) = s(y)$ . Then we get from the operational semantics  $(s, h) \vdash e_1 \Downarrow v, h'$ . We have to show that the variable context for the “then” branch is well-formed, i.e. that  $\Sigma \vdash s(x) : C_{\{R \cap S\}}$  and  $\Sigma \vdash s(y) : D_{\{R \cap S\}}$ . If  $s(x) = s(y) = \text{null}$ , this holds trivially. If  $s(x) = s(y) = l$ , then we know  $\Sigma(s(x)) = \Sigma(s(y)) = \Sigma(l) = (E, r)$ , and  $r \in R$  and  $r \in S$ , hence  $r \in R \cap S$ . By induction, there exists a  $\Pi \sqsupseteq \Sigma$  such that  $\Pi \vdash v : \tau$  and  $h' \models \Pi$ . If  $s(x) \neq s(y)$ , then variable context remains the same, therefore the induction can be applied immediately.

- $(s, h) \vdash x.f \Downarrow F(f), h$ .

Then  $e = x.f$  and  $k = h$  and there is some  $l \in \text{dom}(h)$  and some value  $v$  such that  $s(x) = l$  and  $h(l) = (D, F)$  and  $F(f) = v$ . The typing judgment is  $\Gamma, x : C_R \vdash x.f : \tau$ . With  $\Pi = \Sigma$ , we get  $\Pi \vdash s : (\Gamma, x : C_R)$ , which implies  $\Pi \vdash s(x) : (C_R)$ . We have  $k \models \Pi$  and need to show  $\Pi \vdash F(f) : \tau$ . With  $k \models \Pi$ , we have  $\Pi \models (D, F) : (C_R)$  and therefore by definition  $D \preceq C$  and there is some  $r \in R$  such that  $\Pi \vdash F(f) : A^{get}(D, r, f)$ . By well-typedness of the class table and premise of the typing rule, we get  $A^{get}(D, r, f) <: A^{get}(C, r, f) <: \tau$ , hence  $\Pi \vdash F(f) : \tau$ .

- $(s, h) \vdash x.f := y \Downarrow s(y), h'$ .

Then  $e = x.f := y$  and  $v = s(y)$ . The typing judgment is  $\Gamma, x : C_R, y : \tau \vdash x.f := y : \tau$ . We choose  $\Pi = \Sigma$ . With  $\Sigma \vdash s : (\Gamma, x : C_R, y : \tau)$ , we have  $\Pi \vdash v : \tau$ . By  $h \models \Sigma$ , we know  $D \preceq C$ . With the rule of the operational semantics, there is some  $l \in \text{dom}(h)$  such that  $s(x) = l$  and  $h(l) = (D, F)$  and  $k = h[l \mapsto (D, F[f \mapsto s(y)])]$ . We still need to show

$k \models \Pi$ . Because  $k$  and  $h$  are identical with the exception of  $k(l)(f)$ , we need to show that the semantic heap relation is still preserved for the field value, i.e. that there is a region  $r \in R$  such that  $\Pi \vdash s(y) : A^{get}(C, r, f)$ . With well-formedness of the class table and the premise of the rule, we have  $\tau <: A^{set}(C, r, f) <: A^{get}(C, r, f)$ , hence  $\Pi \vdash s(y) : \tau$ .

- $(s, h) \vdash x.m(\bar{y}) \Downarrow v, h'$ .

By inversion of the rule for the operational semantics, we know there is a location  $l \in \text{dom}(h)$  such that  $s(x) = l$  and  $h(l) = (D, \_)$ , and  $(s', h) \vdash \text{mtable}(D, m) \Downarrow v, k$  where  $s' = [\text{this} \mapsto l] \cup [x_i^m \mapsto s(y_i)]_{i \in \{1, \dots, |\bar{x}^m\}}$ .

The typing judgment is  $\Gamma, x : C_R, \bar{y} : \bar{\sigma} \vdash x.m(\bar{y}) : \tau$ . As  $\Sigma \vdash s : (\Gamma, x : C_R, \bar{y} : \bar{\sigma})$ , we have  $\Sigma \vdash l : C_R$ . By definition of well-typed values, there must exist some class  $E$  and some class  $r$  with  $\Sigma(l) = (E, r)$  such that  $r \in R$  and  $E \preceq C$ . We can derive  $\Sigma \vdash l : C_{\{r\}}$ . With  $h \models \Sigma$ , we can also infer  $D \preceq C$ .

By inversion of the typing rule, we know there is some method typing  $(\bar{\sigma}', \tau') \in M(C, r, m)$  such that  $\bar{\sigma} <: \bar{\sigma}'$  and  $\tau' <: \tau$ . As  $D \preceq C$  and the class table is well-formed, there is a method typing  $(\bar{\sigma}'', \tau'') \in M(D, r, m)$  such that  $\bar{\sigma}' <: \bar{\sigma}''$  and  $\tau'' <: \tau'$ . From  $\Sigma \vdash s : (\Gamma, x : C_R, \bar{y} : \bar{\sigma})$ , it follows  $\Sigma \vdash s(y_i) : \sigma_i''$  for  $i \in \{1, \dots, |\bar{x}^m|\}$ .

As the FJEU program is well-typed, we get  $\Gamma' \vdash \text{mtable}(D, m) : \tau''$  where  $\Gamma' = [\text{this} \mapsto C_{\{r\}}] \cup [x_i^m \mapsto \sigma_i'']_{i \in \{1, \dots, |\bar{x}^m\}}$ . With the facts from above, we get  $\Sigma \vdash s' : \Gamma'$ , so we can finally apply the theorem inductively on the derivation of the semantics and get  $\Pi \vdash v : \tau''$  and  $k \models \Pi$  for some  $\Pi \sqsupseteq \Sigma$ . From  $\tau'' <: \tau$  follows  $\Pi \vdash v : \tau$ .

□

#### 4. Parametrized Region Type System

Our next goal is to use the type system for the automatic algorithmic verification of results of external analyses. In this section, we first show how to *interpret* given analysis results in the type system, and then present an algorithmic type checking algorithm for the automatic verification.

The interpretation of given results requires the reformulation of the above type system to explicitly model abstraction principles that are fundamental for a number of different pointer analysis techniques [7]. We focus on two general classes of abstractions: the abstraction of the call graph using contexts, and the abstraction of objects on the heap. The region type system will be equipped with parameters that can be instantiated to specific abstraction principles. We show that the parametrized version of the type system arises as a specialization of the general region type system.

In the following, the notion of program points is made explicit by annotating expressions with *expression labels*  $i \in \mathcal{I}$ : we write  $[e]^i$  for FJEU expressions, where  $e$  is defined as before. An FJEU program is well-formed if each expression

label  $i$  appears at most once in it. In the following, we only consider well-formed programs, and simply write  $e$  instead  $[e]^i$  if the expression label  $i$  is not important.

#### 4.1. Abstraction Principles

Context-insensitive pointer analyses examine each method exactly once. Their result provides for each method a single pointer structure specification which is used for every call of that method. *Context-sensitive* algorithms improve the precision: each method may be analyzed multiple times under different *contexts*, so that different specifications can be used for different calls to the same method. A context-sensitive algorithm settles on a specific (finite) set of contexts, and produces for each method one pointer structure specification per context. Pointer structure specifications correspond to our method types. We therefore model the concept of contexts by introducing a finite abstract set of contexts  $\mathcal{Z}$ , and by parametrizing the method typing function  $M$  to associate one method type per context in  $\mathcal{Z}$ .

The choice of contexts and specifications for each method call depends on the analysis in question. For example, call-site sensitive algorithms [24] process each method once for each program point where the method is called. Receiver-object sensitive analyses [25] differentiate pointer structure specifications of a method according to the abstraction of the invoking object. More powerful analyses use *call stacks* as contexts to differentiate method calls. For example, a method  $m$  may be analyzed for each possible call-site stack that may occur at the time when  $m$  is called. Similarly, receiver-object stacks can be used. In other words, one considers the *call graph* of the program. A method  $m$  is then represented by a node in this graph, and a context corresponds to a possible path that leads to the node. As there may be infinitely many paths in recursive programs, the number of paths needs to be restricted by some mechanism. A common way is to consider only the last  $k$  entries on the stack ( $k$ -CFA [10]), or to collapse each strongly connected component into a node, thereby eliminating recursive cycles from the set of possible paths [1, 26]. Following this observation, we employ a general *context transfer function*  $\phi$  which represents the edges in the abstract call graph. The function selects a context for the callee based on the caller’s context, the class of the receiver object, its region, the method name, and the call site.

Another abstraction principle is the object abstraction, i.e. the abstract location assigned to allocated objects. This corresponds to our concepts of regions. As pointer analysis algorithms differentiate only finitely many abstract objects, we can restrict the set of regions  $\mathcal{R}$  to a finite size.

A common abstraction is to distinguish objects according to their allocation site and/or their class. More precise analyses also take into account the context under which allocation takes place. For example, in object-sensitive analysis by Milanova et al. [25], objects are distinguished by their own allocation site and the allocation site of the `this` object. Objects may also be distinguished by the call site of the invoked method, a technique called *heap specialization* [27]. We model these concepts by an *object abstraction function*  $\psi$  that assigns the region

Regions (finite):	$r, s, t \in \mathcal{R}$
Contexts (finite):	$z \in \mathcal{Z}$
Context transfer function:	$\phi \in \mathcal{Z} \times \mathcal{C} \times \mathcal{R} \times \mathcal{M} \times \mathcal{I} \rightarrow \mathcal{Z}$
Object abstraction function:	$\psi \in \mathcal{Z} \times \mathcal{I} \rightarrow \mathcal{R}$

<i>set of regions</i>	<i>object abstraction function</i>	<i>principle</i>
$\mathcal{R} = \mathcal{I}$	$\psi(z, i) = i$	allocation site abstraction
$\mathcal{R} = \mathcal{Z} = \mathcal{I} \times \mathcal{I}$	$\psi((i_1, i_2), i_0) = (i_0, i_1)$	object-sensitive allocation site abstraction
$\mathcal{R} = \mathcal{Z} = \mathcal{I}$	$\psi(i_c, i) = i_c$	heap specialization

  

<i>set of contexts</i>	<i>context transfer function</i>	<i>principle</i>
$\mathcal{Z} = \{z_0\}$	$\phi(z, C, r, m, i) = z_0$	context-insensitivity
$\mathcal{Z} = \mathcal{R}$	$\phi(z, C, r, m, i) = r$	object-sensitive 1-CFA
$\mathcal{Z} = \bigcup_{n \in \{1, \dots, k\}} \mathcal{M}^n$	$\phi(z, C, r, m, i) = (m :: z) _k$	method identifier $k$ -CFA
$\mathcal{Z} = \{\bar{i} \in \sigma(\mathcal{I}') \mid \mathcal{I}' \subseteq \mathcal{I}\}$	$\phi(z_1, C, r, m, i) = z_2$ s.th. $IE_c(z_1, i, z_2, m)$	CFA with eliminated recursive cycles

where

- $\sigma(X)$  is the set of all permutations of  $X$
- $L|_k$  is the truncation of list  $L$  after the first  $k$  elements
- $IE_c$  is the call graph relation of [26] where recursive cycles have been replaced by single nodes

Figure 5: The four type system parameters, and possible instantiations to common abstraction principles

for the new object, given the allocation site and the current method context. Our system is by design class-sensitive, as the class information is part of the type.

Figure 5 summarizes the four parameters of our system: contexts, context transfer function, regions, and object abstraction function. Also, it shows how to instantiate the parameters to obtain various standard abstraction principles.

The parametrized method typing  $\hat{M} : (\mathcal{C} \times \mathcal{R} \times \mathcal{Z} \times \mathcal{M}) \rightarrow \bar{\mathcal{T}} \times \mathcal{T}$  replaces the annotated polymorphic method typing  $M$  in the annotated class table. It is well-formed if for all classes  $C$  and subclasses  $D \preceq C$ , regions  $r \in \mathcal{R}$ , methods  $m \in \text{methods}(C)$ , and contexts  $z \in \mathcal{Z}$ , it holds  $\hat{M}(D, r, z, m) < \hat{M}(C, r, z, m)$ .

#### 4.2. The Parametrized Type System

The parametrized type system extends the typing judgment of the general region type system by a context component  $z$ . Except for the following two rules, all rules remain as presented in Section 3 (with the addition of the context  $z$  in each judgment):

$$\text{TP-NEW} \frac{r = \psi(z, i)}{\Gamma ; z \vdash [\text{new } C]^i : C_{\{r\}}}$$

$$\text{TP-INVOKE} \frac{\forall r \in R. \hat{M}(C, r, \phi(z, C, r, m, i), m) <: (\bar{\sigma}, \tau)}{\Gamma, x : C_R, \bar{y} : \bar{\sigma} ; z \vdash [x.m(\bar{y})]^i : \tau}$$

While in the previous system any region  $r$  could be chosen for new objects, we have restricted this flexibility in the TP-NEW rule to the region specified by  $\psi$ . Moreover, we previously allowed arbitrarily many types per method. For the invocation, any type could be selected for each of the possibly infinitely many regions  $r \in R$  given by the type of the receiver object  $x$ . In contrast, the type of  $x$  in TP-INVOKE can only have a finite set of regions, and for each region  $r \in R$ , the method type is determined by the context selected by  $\phi$ .

For a given parametrized method typing  $\hat{M}$ , we define the corresponding polymorphic method typing  $M(C, r, m) := \bigcup_{z \in \mathcal{Z}} \{\hat{M}(C, r, z, m)\}$ . The rules of the parametrized system are derivable in the previous system: The only rules that differ from the previous system, TP-NEW and TP-INVOKE, have more restrictive premises than their counterparts T-NEW and T-INVOKE. Hence if  $\Gamma ; z \vdash e : \tau$  can be derived from some  $\hat{M}$  in the parametrized system, then  $\Gamma \vdash e : \tau$  can be derived in the previous system with respect to the corresponding method typing  $M$ .

A method table is *well-typed* for  $\hat{M}$  if for all classes  $C$ , contexts  $z$ , regions  $r$ , and methods  $m$  such that  $\hat{M}(C, r, z, m) = (\bar{\sigma}, \tau)$ , the judgment  $\Gamma ; z \vdash \text{mtable}(C, m) : \tau$  can be derived with  $\Gamma = [\text{this} \mapsto C] \cup [x_i^m \mapsto \sigma_i]_{i \in \{1, \dots, |\bar{x}^m|\}}$ . It is easy to see that if a method table is well-typed with respect to  $\hat{M}$  in the parametrized system, then it is also well-typed with respect to the corresponding method typing  $M$  in the general system. Therefore, the soundness theorem is applicable to the parametrized region type system.

#### 4.3. Algorithmic Type System

We now present a syntax-directed form of the parametrized type system from which one can directly read off an algorithm  $A(\Gamma, e) = \tau$  that computes “from left to right” the type  $\tau$  of an expression  $e$  based on a store typing  $\Gamma$ . For this, we eliminate the subtyping rule, and instead specify the most precise resulting type  $\tau$  for each expression, similarly to the approach taken by Pierce [23]. The full algorithmic type system is shown in Figure 6.

For the TA-NULL rule, the type system needs some guidance for the choice of the type. We solve this problem by assuming that the syntax is annotated with

$$\begin{array}{c}
\text{TA-VAR} \frac{}{\Gamma, x : \tau ; z \vdash x : \tau} \\
\text{TA-LET} \frac{\Gamma ; z \vdash e_1 : \sigma \quad \Gamma, x : \sigma ; z \vdash e_2 : \tau}{\Gamma ; z \vdash \text{let } x = e_1 \text{ in } e_2 : \tau} \\
\text{TA-IFEQ} \frac{\Gamma, x : C_{R \cap S}, y : D_{R \cap S} ; z \vdash e_1 : \tau_1 \quad \Gamma, x : C_R, y : D_S ; z \vdash e_2 : \tau_2}{\Gamma, x : C_R, y : D_S ; z \vdash \text{if } x = y \text{ then } e_1 \text{ else } e_2 : \tau_1 \vee \tau_2} \\
\text{TA-IFINST} \frac{\Gamma ; z \vdash e_1 : \tau_1 \quad \Gamma ; z \vdash e_2 : \tau_2}{\Gamma ; z \vdash \text{if } x \text{ instanceof } C \text{ then } e_1 \text{ else } e_2 : \tau_1 \vee \tau_2} \\
\text{TA-NULL} \frac{}{\Gamma ; z \vdash \text{null}_C : C_\emptyset} \quad \text{TA-NEW} \frac{r = \psi(z, i)}{\Gamma ; z \vdash \text{new } C : C_{\{r\}}} \\
\text{TA-GETF} \frac{\bigvee_{r \in R} A^{get}(C, r, f) = \tau}{\Gamma, x : C_R ; z \vdash x.f : \tau} \\
\text{TA-PUTF} \frac{\forall r \in R. \tau <: A^{set}(C, r, f)}{\Gamma, x : C_R, y : \tau ; z \vdash x.f := y : \tau} \\
\text{TA-INVOKE} \frac{\forall r \in R. \text{argtypes}(\hat{M}(C, r, \phi(z, C, r, m, i), m)) <: \bar{\sigma} \quad \bigvee_{r \in R} \text{restype}(\hat{M}(C, r, \phi(z, C, r, m, i), m)) = \tau}{\Gamma, x : C_R, \bar{y} : \bar{\sigma} ; z \vdash x.m(\bar{y}) : \tau}
\end{array}$$

Figure 6: Algorithmic type system

a class: the expression  $\text{null}_C$  means that the type system shall use the class  $C$  for that particular *null* value. The reason why we have not introduced this annotated syntax from the beginning is that the inference of class information is orthogonal to pointer analysis and the verification of region sets.

For the rules TA-IF, TA-GETF, and TA-INVOKE, we compute the least upper bound of two or more types. We define the join  $C_R \vee D_S = E_{R \cup S}$  where  $E$  is the most specific common superclass of  $C$  and  $D$ . The notation  $\bigvee X$  denotes join of a set of types  $X$ . The functions *argtypes* and *restype* used in TA-INVOKE return the argument and result types of a method type, respectively.

The soundness proof shows that the internalization of the subtyping rule is correct, i.e. that the judgments derived with the algorithmic type system can also be derived with the parametrized type system. We only outline the proof: The rule TA-IFEQ is derivable by a combination of the rules TP-IFEQ and TP-SUB; similarly for TA-IFINST. The rule TA-NULL is a specialization

of TP-NUL. In TA-GETF, the computed type  $\tau$  satisfies the premise of the TP-GETF. The same holds for the type  $\tau$  in the TA-INVOKE rule, which is a candidate for  $\tau$  in TP-INVOKE. All other rules as well as the definition of annotated class tables remain unchanged. We can therefore apply the soundness theorem to the algorithmic type system.

#### 4.4. FJEU Pointer Analysis in Datalog

The algorithmic type system verifies results of pointer analysis algorithms for FJEU programs. In this section, we present a concrete context-sensitive FJEU pointer analysis to demonstrate the verification capabilities of the algorithmic type system.

Building on the work by Whaley and Lam [26], we specify the algorithm declaratively as a set of recursive Datalog rules, which can be translated into an efficient BDD-based algorithm. To simplify matters, we omit class types, and assume the program in question has been verified by some other (class) type checking facility before. Methods are supposed to have only one argument variable  $x^m$ . There is no class type analysis, i.e. all implementations of a method are regarded as possible targets of a method invocation.

Figure 7 shows the Datalog rules that define the algorithm. The underlined sections determine the abstraction principles of the analysis: just as in the algorithm by Whaley and Lam, we choose the  $IE_c$  call graph abstraction relation for the contexts, and the program points (expression labels) as the object abstraction.

The algorithm computes the fixed point of the relations that express which variables, fields, and method arguments may point to which regions, depending on the context. The resulting points-to relations  $hP$  and  $mT$  can then be interpreted as region annotations for the class table. The algorithmic type system can check the correctness of the annotated class table (and thus of the analysis result) if the abstraction functions  $\psi$  and  $\phi$  are instantiated correctly such that they model the abstraction principles used by our Datalog algorithm:

$$\psi(z, i) = i \quad \text{and} \quad \phi(C, z, r, m, i) = z' \iff IE_c(z, i, z', m).$$

If the program is typable using the annotated class table and these instantiations of  $\phi$  and  $\psi$ , then the result of the pointer analysis is indeed correct.

## 5. String Analysis

String analysis is a dataflow analysis technique to determine possible string values (character sequences) that may occur during the execution of a program. Since strings appear as objects in Java, it is natural to implement a string analysis by building on pointer analysis: string objects are identified and tracked by the pointer analysis, while their possible values are determined by the string analysis.

We now equip the FJEU language with special string objects and operations to give a simplified formalization of Java’s `String` class, and extend the region

$$\begin{aligned}
eP(z, i, r) &: - \text{var}(i, x), vP(z, x, r) \\
eP(z, i, r_2), vP(z, x, r_1) &: - \text{let}(i, x, i_1, i_2), eP(z, i_1, r_1), eP(z, i_2, r_2) \\
eP(z, i, r_1), eP(z, i, r_2) &: - \text{if}(i, i_1, i_2), eP(z, i_1, r_1), eP(z, i_2, r_2) \\
eP(z, i, r) &: - \text{new}(i), r = i \\
eP(z, i, r_2) &: - \text{getf}(i, x, f), vP(z, x, r_1), hP(z, r_1, f, r_2) \\
eP(z, i, r_1), hP(z, r_1, f, r_2) &: - \text{putf}(i, x, f, z), vP(z, x, r_1), vP(z, y, r_2) \\
eP(z, i, r), mT(r_1, m, z', r_2, r) &: - \text{call}(i, m, x, y), vP(z, x, r_1), vP(z, y, r_2), \\
&\quad \text{body}(m, i_{\text{body}}), eP(z', i_{\text{body}}, r), \\
&\quad \underline{IE_c(z, i, z', m)} \\
vP(z, x^m, r) &: - mT(r, m, z, r_v, r_r)
\end{aligned}$$

where

- $i, r, z$  range over finite sets of expression labels, regions, and contexts
- $\text{var}(i, x), \text{let}(i, x, i_1, i_2), \text{new}(i), \dots$ : the expression at  $i$  is a variable/let expression/new expression/...
- $eP(z, i, r)$ : in context  $z$ , the evaluation of the expression at  $i$  may point to  $r$
- $vP(z, x, r)$ : in context  $z$ , variable  $x$  may point to region  $r$
- $hP(z, r_1, f, r_2)$ : in context  $z$ , field  $f$  of objects in region  $r_1$  may point to region  $r_2$ ; this corresponds to  $A^{\text{get}}$  and  $A^{\text{set}}$
- $mT(r, m, z, r_v, r_r)$ : the type for method  $m$  for context  $z$  of objects in region  $r$  is  $(r_v, r_r)$ ; this corresponds to  $M$
- $\text{body}(m, i)$ : the expression with label  $i$  is an implementation of method  $m$

Figure 7: FJEU Pointer Analysis defined as Datalog rules

type system to enable the verification of pointer analyses of string objects. Afterwards, we show how to use the region information to interpret the results of a specific string analysis.

### 5.1. FJEU with Strings

The language FJEUS is an extension of FJEU with operations to create and concatenate strings. While the `String` class in Java is just another class in the Java class hierarchy, it is regarded as a separate type in FJEUS. This allows us to treat string objects differently: An object of class `String` in FJEUS is simply a string value (character sequence) on the heap. The meta-variable  $w$  ranges over character sequences  $\mathcal{W}$ , and  $W$  ranges over sets of string values. We rely on a given sequence concatenation function  $+$ .

$$\begin{array}{lll} \text{character sequences:} & w & \in \mathcal{W} \\ \text{sets of character sequences:} & W & \in \mathcal{P}(\mathcal{W}) \\ \text{extended expressions:} & \mathcal{E} \ni e & ::= \dots \mid \text{"}w\text{"} \mid x.\text{concat}(y) \end{array}$$

$$\begin{array}{lll} \text{character sequence concatenation:} & + & \in \mathcal{W} \times \mathcal{W} \rightarrow \mathcal{W} \\ \text{heaps:} & h, k & \in \mathcal{L} \rightarrow \mathcal{O} \cup \mathcal{W} \end{array}$$

The literal expression `"w"` allocates a new string object with the character sequence  $w$  on the heap. The string operation  $x.\text{concat}(y)$  has its own special semantics and is implemented with the  $+$  operator. As we only model strings with non-mutable values in the language, a string concatenation always creates a new string object on the heap. The operational semantics is extended as follows:

$$\frac{l \notin \text{dom}(h)}{(s, h) \vdash \text{"}w\text{"} \Downarrow l, h[l \mapsto w]}$$

$$\frac{\begin{array}{llll} s(x) = l_1 & s(y) = l_2 & h(l_1) = w_1 & h(l_2) = w_2 \\ l \notin \text{dom}(h) & w = w_1 + w_2 & & \end{array}}{(s, h) \vdash x.\text{concat}(y) \Downarrow l, h[l \mapsto w]}$$

Note that this rather modest extension is intended to keep the formalization simple. Other extensions could include more string operations, or mutable string objects that model Java's `StringBuffer` class.

### 5.2. Pointer Analysis for String Objects

We extend the region type system from Section 3 to accommodate the new string objects. We distinguish references to "proper" objects and to string objects: a type is either a class with a region set ( $C_R$ ), or the special `String` class with a region set ( $\text{String}_R$ ). The `String` class is independent from other classes in the class hierarchy.

$$\begin{array}{c} \text{types: } \sigma, \tau \in \mathcal{T} = (\mathcal{C} \cup \{\mathbf{String}\}) \times \mathcal{P}(\mathcal{R}) \\ \frac{C \preceq D \quad R \subseteq S}{C_R <: D_S} \qquad \frac{R \subseteq S}{\mathbf{String}_R <: \mathbf{String}_S} \end{array}$$

Fields and method typings may now include the  $\mathbf{String}_R$  type. Also, all existing typing rules from Section 3 remain unchanged. In particular, the *null* value may be assigned a  $\mathbf{String}$  type. A field  $x.f$  may only be accessed if  $x$  is a (non- $\mathbf{String}$ ) class  $C$ ; similarly for method calls  $x.m$ . These are the two additional typing rules for " $w$ " and  $x.concat(y)$ :

$$\frac{}{\Gamma \vdash "w" : \mathbf{String}_{\{r\}}}$$

$$\frac{}{\Gamma, x : \mathbf{String}_R, y : \mathbf{String}_S \vdash x.concat(y) : \mathbf{String}_{\{t\}}}$$

The heap typings  $\Sigma, \Pi : \mathcal{L} \rightarrow (\mathcal{C} \cup \{\mathbf{String}\}) \times \mathcal{R}$  may now also map locations to the  $\mathbf{String}$  class. We extend the well-typed value relation  $\Sigma \vdash v : \tau$  accordingly:

$$\frac{}{\Sigma \vdash \mathbf{null} : \tau} \qquad \frac{\Sigma(l) = (C, r)}{\Sigma \vdash l : C_{\{r\}}} \qquad \frac{\Sigma(l) = (\mathbf{String}, r)}{\Sigma \vdash l : \mathbf{String}_{\{r\}}}$$

$$\frac{\Sigma \vdash v : \sigma \quad \sigma <: \tau}{\Sigma \vdash v : \tau}$$

The definition of well-typed heaps additionally requires well-typed character sequences (last line):

$$\begin{array}{l} h \models \Sigma \iff \forall l \in \text{dom}(\Sigma). l \in \text{dom}(h) \wedge \Sigma \models h(l) : \Sigma(l) \\ \Sigma \models (C, F) : (D, r) \iff C \preceq D \wedge \dots \text{ (as before)} \\ \Sigma \models w : (\mathbf{String}, r) \iff \text{PROP}(w, r) \end{array}$$

In other words, the property that a heap  $h$  is well-typed with respect to  $\Sigma$  now includes the condition that for all locations  $l$  such that  $\Sigma(l) = (\mathbf{String}, r)$ ,  $h(l)$  contains a string value  $w$  that satisfies a certain property  $\text{PROP}$  with respect to  $r$ . For the moment, assume  $\text{PROP}$  is simply *True*. The proof of the soundness theorem is extended in a straight-forward way for the extensions described above. Moreover, the type system can be parametrized in the same fashion as described in Section 4: as both string operations create new objects, we use the  $\psi$  function to determine the region of these objects.

In the following subsection, we present an analysis that can help to prevent cross-site scripting attacks, and give a semantic formalization by instantiating the string property  $\text{PROP}$ .

### 5.3. String Analysis with Operation Contexts

In a typical cross-site scripting scenario, a user input is embedded into a string that is executed or interpreted. To prevent the injection of malicious code, one wants to track how a string is constructed, and ensure that its executable parts originate from trusted sources, like string literals in the program code.

We therefore add a global signature of string operations  $\Omega$ , which expresses the possible string operations that objects in specific regions may be the result of. The typing rules are extended with constraints on  $\Omega$ .

string operations:  $\mathcal{S} \ni \omega ::= \text{lit } w \ r \mid \text{concat } r \ s \ t \mid \text{unknown } W \ r$   
string operation context:  $\Omega \in \mathcal{P}(\mathcal{S})$

$$\frac{\text{lit } w \ r \in \Omega}{\Gamma \vdash \text{"}w\text{"} : \mathbf{String}_{\{r\}}}$$

$$\frac{\forall r \in R, s \in S. \text{concat } r \ s \ t \in \Omega}{\Gamma, x : \mathbf{String}_R, y : \mathbf{String}_S \vdash x.\text{concat}(y) : \mathbf{String}_{\{t\}}}$$

Informally, an operation  $\text{lit } w \ r \in \Omega$  means that region  $r$  is permitted to include the string  $w$ .  $\text{concat } r \ s \ t$  means that region  $t$  is permitted to include a string object that is obtained by concatenating some strings from regions  $r$  and  $s$ . The operation  $\text{unknown } W \ r$  means that all strings in  $W$  are permitted to appear in region  $r$ . This is useful for external methods whose types are given to but not verified by the type system. Note that the inclusion of  $\text{unknown } W \ r$  in  $\Omega$  is equivalent to the inclusion of  $\text{lit } w \ r$ , for all  $w \in W$ . Whenever additional primitive string operations are included to the language, the set  $\mathcal{S}$  may be extended accordingly.

Formally, we define a semantic interpretation  $\Omega[[r]]$  that gives the possible values for string objects in region  $r$ . It is defined as the smallest set satisfying the following conditions:

$$\begin{aligned} \text{lit } w \ r \in \Omega &\Rightarrow w \in \Omega[[r]] \\ \text{concat } r \ s \ t \in \Omega &\Rightarrow \forall w_1 \in \Omega[[r], w_2 \in \Omega[[s]]. w_1 + w_2 \in \Omega[[t]] \\ \text{unknown } W \ r \in \Omega &\Rightarrow W \subseteq \Omega[[r]] \end{aligned}$$

After instantiating the string property as  $\text{PROP}(w, r) \equiv w \in \Omega[[r]]$ , the relation  $h \models \Sigma$  ensures that string values on the heap are indeed in the interpretation of the string operation context.

Apart from this extensional interpretation, the string operation context and the typing of external methods also contain intensional information about the origin and the possible constructions of strings in a specific region, which enables the verification of more complex string policies.

For example, consider the following string-manipulating program that relies on external functions  $\text{getUserInput}()$  to retrieve data from the user,  $\text{escapeHTML}(s)$  that quotes all HTML tags in string  $s$  and returns the result as a new string, and  $\text{output}(s)$  that outputs the string  $s$ .

```

let firstPart = new String("<someTag>")
in let contents = getUserInput()
    in let escContents = escapeHTML(contents)
        in output(firstPart.concat(escContents))

```

The security policy is that the output may only be the result of a concatenation of a string literal with a string that does not contain HTML tags. The policy is expressed using the following types for external methods and the string operation context:

$$\begin{aligned}
\textit{getUserInput} & : \text{unit} \longrightarrow \text{String}_{\{q\}} \\
\textit{escapeHTML} & : \text{String}_{\{q\}} \longrightarrow \text{String}_{\{s\}} \\
\textit{output} & : \text{String}_{\{t\}} \longrightarrow \text{unit}
\end{aligned}$$

$$\begin{aligned}
\Omega = \{ & \textit{lit} \text{ "<someTag>" } l, \\
& \textit{unknown } \mathcal{W} q, \textit{ unknown } \hat{\mathcal{W}} s, \\
& \textit{concat } l \textit{ s } t \}
\end{aligned}$$

( $\hat{\mathcal{W}}$  is the set of all strings that do not contain HTML tags, and `unit` is a unit type, which could be modeled in FJEUS as `String∅`, containing only the null value.) The external function `getUserInput` returns strings with arbitrary values of the set  $\mathcal{W}$  in region  $q$  (“questionable”), which are converted by `escapeHTML` into strings of region  $s$  (“sanitized”), which are assumed to not contain any HTML tags (set  $\hat{\mathcal{W}}$ ). For the literal `firstPart`, the type checker can assign the region  $l$  (“literals”), as the literal value in  $\Omega$  matches. The `output` function only accepts strings from region  $t$  (“trusted”), which must be, according to  $\Omega$ , a concatenation of strings from region  $l$  and (HTML tag-free) region  $s$ . Therefore, the typability of the program proves that the security policy is indeed fulfilled. The example demonstrates that handling trusted sanitizing functions is actually a strength of type-based presentations: simply assign an appropriate type to a function if you believe the associated semantic property of the function.

The approach is related to the work by Christensen et al. [8] on the analysis of string values using context-free grammars with operation productions. The symbolic string operations in the context correspond to nodes in their annotated flow graph, and their semantics of the flow graph resembles our interpretation of  $\Omega[[r]]$ . Similarly, Crégut and Alvarado [9] have presented an algorithm that tracks string objects with pointer analysis, and collects intensional information about the string operations applied to them. We thus expect that aspects of the results of these algorithms are verifiable in our system. Our approach is also related to taint analysis [28], as the region identifiers can convey information about the trustworthiness of strings, which is preserved throughout assignments and method invocations.

In the approach shown above, the regions used in `String` types can be seen as “type identifiers”, which are referenced in the string operation context  $\Omega$  to

$$\begin{array}{c}
\text{TI-IFINSTT} \frac{\Gamma(x) = C_R \quad C \preceq E \quad \Gamma \vdash e_1 : \tau}{\Gamma \vdash \text{if } x \text{ instanceof } E \text{ then } e_1 \text{ else } e_2 : \tau} \\
\text{TI-IFINSTF} \frac{\Gamma(x) = C_R \quad \neg E \preceq C \quad \neg C \preceq E \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash \text{if } x \text{ instanceof } E \text{ then } e_1 \text{ else } e_2 : \tau} \\
\text{TI-IFINST} \frac{\Gamma(x) = C_R \quad E \preceq C \quad \Gamma, x : E_R \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash \text{if } x \text{ instanceof } E \text{ then } e_1 \text{ else } e_2 : \tau} \\
\text{TI-IFEQC} \frac{\Gamma(x) = C_R \quad \Gamma(y) = D_S \quad C \preceq D \quad \Gamma, x : C_{R \cap S}, y : C_{R \cap S} \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash \text{if } x = y \text{ then } e_1 \text{ else } e_2 : \tau} \\
\text{TI-IFEQD} \frac{\Gamma(x) = C_R \quad \Gamma(y) = D_S \quad D \preceq C \quad \Gamma, x : D_{R \cap S}, y : D_{R \cap S} \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash \text{if } x = y \text{ then } e_1 \text{ else } e_2 : \tau} \\
\text{TI-IFEQ} \frac{\Gamma(x) = C_R \quad \Gamma(y) = D_S \quad \neg C \preceq D \quad \neg D \preceq C \quad \Gamma, x : C_\emptyset, y : D_\emptyset \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash \text{if } x = y \text{ then } e_1 \text{ else } e_2 : \tau}
\end{array}$$

Figure 8: More precise rules for conditionals

provide further information on these string types. In a recent paper [29], we have instead directly included string information in the type annotation for the `String` class. This annotation represents the string history and value in the type in a more concise way than the string operation contexts presented here. Even there, however, regions are useful e.g. to distinguish string objects that are referenced in a field of two different container objects of the same class.

## 6. Stronger Proof Rules for Precise Conditionals

Following a suggestion by a reviewer, we also investigated stronger proof rules for conditionals. These allow the branches to be typed under additional assumptions regarding the outcome of the branch conditions. Figure 8 shows the more precise rules for the two forms of conditionals, which are meant to replace the original rules `T-IFINST` and `T-IFEQ`.

Rule `TI-IFINSTT` relies on the fact that whenever a subclass test can statically be shown to always succeed, no typing hypothesis for the negative branch is required. In contrast, rule `TI-IFINSTF` applies if the classes  $C$  and  $E$  are incomparable – in this case, the condition will always evaluate negatively, due to the tree shape of the subclass relationship. Finally, rule `TI-IFINST` applies if

the branch outcome cannot be predicted statically, but exploits the additional knowledge  $E \preceq C$  in the positive branch.

In the case of pointer equality tests, the first two rules refine the class and region information depending on positive subclass information, while the final rule exploits the fact that incompatibility of the static types necessitates that pointers can only be equal if they evaluate to *null*.

Both original rules for conditionals given in Figure 4, T-IFINST and T-IFEQ, are derivable from the rules shown in Figure 8.

The soundness proof for the system including the more precise rules for conditionals proceeds similar to the proof in Section 3.4, but requires the definition of heap typing extensions to be altered: instead of defining  $\Pi \sqsupseteq \Sigma$  to hold whenever  $\text{dom}(\Sigma) \subseteq \text{dom}(\Pi)$  and  $\forall l \in \text{dom}(\Sigma). \Sigma(l) = \Pi(l)$ , we now require  $\text{dom}(\Sigma) \subseteq \text{dom}(\Pi)$  and  $\forall l \in \text{dom}(\Sigma). \Pi(l) <: \Sigma(l)$ .

The formalized proof of the appropriately altered soundness result is included in our Isabelle development.

## 7. Discussion

We presented a framework for classifying alias analyses for Java-like languages, given by a hierarchy of region-based type systems. We demonstrated how existing disciplines arise as instantiations of our framework and may be given a uniform interpretation by embedding their results in a single type system. We also gave an algorithmic variant of the type system, thus enabling syntax-directed type-checking and hence validation of analyses results. Finally, we showed how our framework may be extended to string analyses. In the following, we briefly discuss specific design decisions, and outline future work.

To our knowledge, most existing pointer analyses express their results at a coarse-grain level of syntactic structure such as methods. In accordance with this, we employed a phrase-based formulation of type systems and interpreted judgments with respect to a big-step evaluation semantics. An extension of the interpretation to include non-terminating executions appears possible, using techniques as in [30].

Our framework does not aim to be flow- or path-sensitive. We see these concepts as orthogonal to the central idea of our paper, namely interpreting context-sensitivity using polyvariant types. Nevertheless, we acknowledge the increasing relevance of flow and path sensitivity in recent work on pointer analysis. The T-IFEQ rule illustrates a possible extension of the type system with path-sensitive capabilities: the information from the branching expression is used to refine the analysis for the “then” branch of the conditional. Moreover, subderivations of a judgment contain implicitly more fine-grained (non-)alias relationships applicable at intermediate program points, and include aspects of flow-sensitivity as any variable may be associated with different types in a derivation. Arguably, local alias assertions could be made more explicit by moving to a small-step operational regime and/or to formulations of the type systems that are inspired by abstract interpretation and yield a global specification table

with entries for all program points [31]. However, the use of evaluation-style judgments greatly simplifies soundness proofs at least at the level of methods, as recursive calls follow the type derivation structure.

The Doop framework [5] enables the definition of highly precise and efficient pointer analyses declaratively using Datalog rules. While the authors do not seem to aim at a fully formal correctness proof that interprets the Datalog rules and relations with respect to the semantics of the Java language, they take great care to separate the essential aspects of analysis techniques from implementation details of the algorithm. We expect this separation of concerns to enable a future integration of Doop-style analyses into our framework.

In addition to type systems for pointer alias analysis, the literature contains the terminology “type-based alias analysis” [32]. The latter term appears to mean that datatype or static class information is used to *improve the precision* of an alias analysis, while region type systems directly *include* the points-to relations in the types. However, as our system extends the ordinary type system of Java, it arguably also encompasses type-based alias analyses.

Having been developed in order to embed static analysis results, it is not surprising that the type systems over-approximate their semantic guarantee. Thus, failure to validate the *legitimate* result of a specific analysis may be rooted in either an incorrect interpretation of the analysis into our framework or the fact that the analysis is more precise than the host type system. A particular direction in which our type system (and some analyses) might be generalized is shape analysis [33]. Another interesting recent development that our work does not accommodate is the analysis of programs that use reflection [28]; this would require a fundamental understanding of the semantic analysis of reflection.

Although we have only examined the *results* of pointer analysis algorithms, the algorithms can be seen as an external means of type inference. It seems promising to further investigate the *implementations* of these algorithms, and to recreate their logic in the type system in order to obtain a parametric type system with a fully automatic (internal) type inference. Alternatively, the identification of abstraction principles could also propose a way to parametrize existing pointer analysis implementations.

Regarding the string analysis, we have concentrated on the previously noted observation that the precision of the analysis benefits from the availability of (non-)aliasing information [9]. In principle, the benefits may be mutual. For example, the method call  $x.concat(y)$  on a `String` in Java actually returns the reference  $x$  if  $y$  has a length of zero. If the length of  $y$  can be obtained from a string analysis, this information helps to improve the region set for the result type in the rule for concatenation. Mutual dependencies between aliasing and string analyses may thus be an interesting topic for future work.

Also, we have only outlined the use of type systems for string security policies. In collaboration with SAP’s Sophia-Antipolis-based research lab on security, we have developed a type system to enforce secure programming guidelines which aim to prevent cross-site scripting attacks in concrete application scenarios [29].

*Acknowledgements.* We would like to thank Pierre Crégut for making us aware of the interplay between pointer and string analyses and providing us with several pointers to the literature, as well as Keqin Li (SAP Research) for giving us insights to concrete cross-site scripting scenarios and the corresponding string policies. This work was supported by the DFG-funded project PolyNI. We thank the reviewers for their valuable comments, in particular for suggesting the improved typing rules in Section 6.

## References

- [1] M. Emami, R. Ghiya, L. J. Hendren, Context-sensitive interprocedural points-to analysis in the presence of function pointers, in: 1994 Conference on Programming language design and implementation (PLDI'94), ACM, New York, NY, USA, 1994, pp. 242–256.
- [2] L. O. Andersen, Program analysis and specialization for the c programming language, Ph.D. thesis, DIKU, University of Copenhagen (1994).
- [3] B. Steensgaard, Points-to analysis in almost linear time, in: 23rd Symposium on Principles of programming languages (POPL '96), ACM, New York, NY, USA, 1996, pp. 32–41.
- [4] M. Berndt, O. Lhoták, F. Qian, L. Hendren, N. Umanee, Points-to analysis using BDDs, in: 2003 Conference on Programming Language Design and Implementation (PLDI '03), ACM, New York, NY, USA, 2003, pp. 103–114.
- [5] M. Bravenboer, Y. Smaragdakis, Strictly declarative specification of sophisticated points-to analyses, in: 24th Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA 2009), ACM, 2009, pp. 243–262.
- [6] A. Igarashi, B. Pierce, P. Wadler, Featherweight Java: A minimal core calculus for Java and GJ, in: 1999 Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA 1999), ACM, 1999, pp. 132–146.
- [7] B. G. Ryder, Dimensions of precision in reference analysis of object-oriented programming languages, in: 12th International Conference on Compiler Construction, Vol. 2262, Springer-Verlag, 2003, pp. 126–137.
- [8] A. S. Christensen, A. Møller, M. I. Schwartzbach, Precise analysis of string expressions, in: R. Cousot (Ed.), Static Analysis: 10th International Symposium (SAS 2003), Vol. 2694 of Lecture Notes in Computer Science, Springer-Verlag, 2003, pp. 1–18.
- [9] P. Crégut, C. Alvarado, Improving the Security of Downloadable Java Applications With Static Analysis, Electronic Notes in Theoretical Computer Science 141 (1) (2005) 129–144.

- [10] O. Shivers, Control-flow analysis of higher-order languages, or taming lambda, Ph.D. thesis, Carnegie Mellon University, Pittsburgh, PA, USA, technical Report CMU-CS-91-145 (1991).
- [11] A. Banerjee, T. Jensen, Modular control-flow analysis with rank 2 intersection types, *Mathematical Structures in Computer Science* 13 (1) (2003) 87–124.
- [12] B. Hardekopf, C. Lin, Semi-sparse flow-sensitive pointer analysis, *SIGPLAN Not.* 44 (1) (2009) 226–238.
- [13] J. M. Lucassen, D. K. Gifford, Polymorphic effect systems, in: 15th Symposium on Principles of programming languages (POPL '88), ACM, New York, NY, USA, 1988, pp. 47–57.
- [14] N. Benton, A. Kennedy, L. Beringer, M. Hofmann, Relational semantics for effect-based program transformations with dynamic allocation, in: *Principles and Practice of Decl. Prog. (PPDP '07)*, ACM, 2007, pp. 87–96.
- [15] M. Tofte, J.-P. Talpin, Region-based memory management, *Information and Computation* 132 (2) (1997) 109–176.
- [16] K. Cray, D. Walker, G. Morrisett, Typed memory management in a calculus of capabilities, in: 26th Symposium on Principles of programming languages (POPL '99), ACM, 1999, pp. 262–275.
- [17] J. S. Foster, T. Terauchi, A. Aiken, Flow-sensitive type qualifiers, in: 2002 Conference on Programming language design and implementation (PLDI '02), ACM, New York, NY, USA, 2002, pp. 1–12.
- [18] O. Lhoták, Program analysis using binary decision diagrams, Ph.D. thesis, McGill University (Jan. 2006).
- [19] T. Lenherr, Taxonomy and Applications of Alias Analysis, Master's thesis, ETH Zürich (2008).
- [20] L. Beringer, R. Grabowski, M. Hofmann, Verifying pointer and string analyses with region type systems, in: E. M. Clarke, A. Voronkov (Eds.), *LPAR (Dakar)*, Vol. 6355 of *Lecture Notes in Computer Science*, Springer-Verlag, 2010, pp. 82–102.
- [21] L. Beringer, R. Grabowski, M. Hofmann, Verifying Pointer and String Analyses with Region Type Systems: Soundness proofs and other material, <http://www.tcs.ifi.lmu.de/~grabow/regions> (2010).
- [22] M. Hofmann, S. Jost, Type-based amortised heap-space analysis, in: 15th European Symposium on Programming (ESOP 2006), Vol. 3924 of *Lecture Notes in Computer Science*, Springer-Verlag, Vienna, Austria, 2006, pp. 22–37.

- [23] B. C. Pierce, *Types and Programming Languages*, MIT Press, 2002.
- [24] M. Sharir, A. Pnueli, Two approaches to interprocedural data flow analysis, in: Muchnick, Jones (Eds.), *Program Flow Analysis: Theory and Applications*, Prentice Hall International, 1981.
- [25] A. Milanova, A. Rountev, B. G. Ryder, Parameterized object sensitivity for points-to analysis for Java, *ACM Trans. Softw. Eng. Methodol.* 14 (1) (2005) 1–41.
- [26] J. Whaley, M. S. Lam, Cloning-based context-sensitive pointer alias analysis using binary decision diagrams, *SIGPLAN Not.* 39 (6) (2004) 131–144.
- [27] E. M. Nystrom, H.-S. Kim, W. W. Hwu, Importance of heap specialization in pointer analysis, in: C. Flanagan, A. Zeller (Eds.), *5th Workshop on Program Analysis for Software Tools and Engineering (PASTE '04)*, ACM, New York, USA, 2004, pp. 43–48.
- [28] O. Tripp, M. Pistoia, S. J. Fink, M. Sridharan, O. Weisman, Taj: effective taint analysis of web applications, in: *2009 Conference on Programming Language Design and Implementation (PLDI '09)*, ACM, New York, NY, USA, 2009, pp. 87–97.
- [29] R. Grabowski, M. Hofmann, K. Li, Type-Based Enforcement of Secure Programming Guidelines – Code Injection Prevention at SAP, in: *8th International Workshop on Formal Aspects of Security & Trust (FAST 2011)*, Vol. 7140 of *Lecture Notes in Computer Science*, Springer-Verlag, 2011, pp. 182–197.
- [30] L. Beringer, M. Hofmann, M. Pavlova, Certification Using the Mobius Base Logic, in: *Formal Methods for Components and Objects: 6th International Symposium (FMCO 2007)*, Springer-Verlag, Berlin, Heidelberg, 2008, pp. 25–51.
- [31] G. Barthe, D. Pichardie, T. Rezk, A certified lightweight non-interference Java bytecode verifier, in: *European Symposium on Programming (ESOP 2007)*, Vol. 4421 of *Lecture Notes in Computer Science*, Springer-Verlag, 2007, pp. 125 – 140.
- [32] A. Diwan, K. S. McKinley, J. E. B. Moss, Type-based alias analysis, in: *1998 Conference on Programming Language Design and Implementation (PLDI '98)*, ACM, New York, NY, USA, 1998, pp. 106–117.
- [33] S. Loncaric, A survey of shape analysis techniques, *Pattern Recognition* 31 (1998) 983–1001.