


# Class Analysis of Object-Oriented Programs through Abstract Interpretation

View metadata, citation and similar papers at [core.ac.uk](http://core.ac.uk)

brought to you by  CORE

provided by Catalogo dei prodotti della ricerca

Thomas Jensen and Fausto Spoto

IRISA, Campus de Beaulieu, 35042 Rennes Cedex, France  
{jensen,spoto}@irisa.fr

**Abstract.** We use abstract interpretation to define a uniform formalism for presenting and comparing class analyses for object-oriented languages. We consider three domains for class analysis derived from three techniques present in the literature, viz., rapid type analysis, a simple dataflow analysis and constraint-based 0-CFA analysis. We obtain three static analyses which are provably correct and whose abstract operations are provably optimal. Moreover, we prove that our formalisation of the 0-CFA analysis is more precise than that of the dataflow analysis.

**Keywords:** *Abstract interpretation, class analysis, object-oriented languages, domain theory, semantics.*

## 1 Introduction

Class analysis, one of the most important analyses for object-oriented languages, computes the set of classes that an expression can have at run-time [1,2,10,11,12]. It can serve to prove *type safety* by guaranteeing that methods are only invoked on objects that implement such a method. In certain cases it allows one to optimise virtual method invocations into direct calls to the code implementing the method. It is also used for building a precise call graph for a program which in turn enables other analyses. A variety of class analyses have been proposed (sometimes called receiver class analysis, type analysis, etc.), often using different analysis frameworks. This complicates the comparison of the analyses w.r.t. their precision. The complex question of proving the analyses correct is not always addressed.

Cousot [7] shows how *abstract interpretation* [8,9] can organise eleven different type systems for functional languages into a lattice that allows one to state their correctness and establish their relative precision. Here, we apply abstract interpretation to classify and prove correct class analyses for object-oriented languages. Abstract interpretation is a technique for the systematic construction of semantics and semantics-based static analyses for programming languages. Given a concrete semantics over a concrete domain and an *abstraction function* from the concrete to an abstract semantic domain, abstract interpretation shows how to define an abstract semantics over the abstract domain such that the result

of the analysis is provably *correct* w.r.t. the concrete one. The abstraction function can be seen as a specification of the program property of interest, and the abstract semantics as an analyser for that property. Furthermore, since abstract interpretation is a framework for relating semantic definitions, it can be used to compare the relative precision of two analyses. We present a framework for the static analysis of object-oriented programs. It gives semantics to a simple object-oriented language by specifying an algebra of states. The operations of this algebra are reminiscent of Java bytecode operations. Different static analyses are obtained by abstract interpretation of these operations. This approach allows one to make the following three contributions:

- We define three abstractions of the above algebra, namely, three domains for class analysis which use the same abstract information as [2], [10] and [11]. We obtain formal reconstructions of these analyses in a common framework.
- We use abstract interpretation to derive optimal (and hence correct) abstract operations for our three domains. Note that only the analysis in [11] has been provided with an (informal) correctness proof.
- We show that our reconstruction of the analysis in [11] is *provably* more precise than that of the analysis in [10].

We emphasise that we are not questioning the correctness of the original analyses. Nor are we aiming at defining new, more powerful class analyses (although the present framework could be used for this). Our goal is to set up a framework for studying the common structure of class analyses, and to facilitate their comparison and proof of correctness by means of abstract interpretation.

The paper is organised as follows. Section 2 presents notation and preliminaries. Section 3 shows our framework for the analysis of object-oriented languages, the states and their operations. Sections 4, 5 and 6 define the analyses in [2], [10] and [11], respectively, as abstract interpretations of the sets of states of Section 3. Section 7 compares the precision of our three analyses. Section 8 concludes. Due to space limitations, some proofs have been omitted.

## 2 Preliminaries

We recall some notions of abstract interpretation [8,9]. In the following, a total function is denoted by  $\mapsto$  and a partial function by  $\rightarrow$ . Let  $(C, \leq)$  and  $(A, \preceq)$  be two posets (the concrete and the abstract domain). A *Galois connection* is a pair of monotonic maps  $\alpha : C \mapsto A$  and  $\gamma : A \mapsto C$  such that for each  $x \in C$  we have  $x \leq (\gamma \circ \alpha)(x)$ , and for each  $y \in A$  we have  $(\alpha \circ \gamma)(y) \preceq y$ . A *Galois insertion* is a Galois connection where  $\alpha \circ \gamma$  is the identity map. This means that the abstract domain does not contain *useless* elements. The map  $\alpha$  ( $\gamma$ ) is called the *abstraction* (*concretisation*) function. The abstraction map uniquely determines the concretisation map and vice versa.

Let  $f : C^n \rightarrow C$  be an operator. We say that  $\hat{f} : A^n \rightarrow A$  is *correct* w.r.t.  $f$  if and only if for all  $y_1, \dots, y_n \in A$  we have  $\alpha(f(\gamma(y_1), \dots, \gamma(y_n))) \preceq$

$\hat{f}(y_1, \dots, y_n)$ . For each operator  $f$ , there exists an *optimal* (most precise) correct abstract operator  $\hat{f}$  defined as  $\hat{f}(y_1, \dots, y_n) = \alpha(f(\gamma(y_1), \dots, \gamma(y_n)))$ . This means that  $\hat{f}$  is the most precise abstraction of  $f$  which can be defined on  $A$ .

In the theory of abstract interpretation, the *semantics* of a program is the fixpoint of a continuous operator  $f : C \mapsto C$ , where  $C$  is the computational domain [6]. Its *collecting version* [8] works over *properties* of  $C$ , i.e., over  $\wp(C)$  and is the fixpoint of the powerset extension of  $f$ . If  $f$  is defined as composition of suboperations, their powerset extensions *and*  $\cup$  induce the extension of  $f$ . The  $\cup$  operation merges the semantics of the branches of a conditional. Note that the powerset extension of a predicate over  $C$  is a constant of  $\wp(C)$ . In Props. 2, 4 and 6, we compute the optimal abstractions of the powerset extension of every small operation  $g$  above as  $\alpha \circ g \circ \gamma$ , consistently with the previous paragraph.

Upper closure operators (*uco*'s) are monotonic, extensive and idempotent maps. They are isomorphic to Galois insertions [9]. The abstract domain induced by an uco  $\rho$  over the concrete domain  $L$  (for our purposes, the semantic domain  $\wp(C)$  of a collecting semantics) is its image  $\rho(L)$ , i.e., the set of its fixpoints. This image is a *Moore family* of  $L$ , i.e., a non-empty meet-closed subset of  $L$ . Conversely, any Moore family of  $L$  is the image of some uco on  $L$ . Thus, in order to define a Galois insertion on  $L$ , i.e., an abstraction of  $L$ , we can equivalently define a Moore family on  $L$ , whose induced uco is the abstraction map. In Props. 1, 3 and 5 we use this technique to prove that we deal with Galois insertions.

We denote by  $[v_1 \mapsto t_1, \dots, v_n \mapsto t_n]$  a function  $f$  whose domain is  $\{v_1, \dots, v_n\}$  and such that  $f(v_i) = t_i$  for  $i = 1, \dots, n$ . Its update is  $f[d_1/w_1, \dots, d_m/w_m]$ , where the domain of the function can be potentially enlarged. Its domain and codomain are  $\text{dom}(f)$  and  $\text{codom}(f)$ , respectively. By  $f|_s$  we denote the restriction of  $f$  to  $s \subseteq \text{dom}(f)$ , by  $f|_{-s}$  its restriction to  $\text{dom}(f) \setminus s$ . A definition like  $S = \langle a, b \rangle$ , with  $a$  and  $b$  meta-variables, defines the selectors  $s.a$  and  $s.b$  for  $s \in S$ . For instance, Def. 2 defines  $o.\kappa$  and  $o.\phi$ , for  $o \in \text{Obj}$ . If  $\langle D, \leq \rangle$  is a poset and  $d \in D$ , we define  $\downarrow(d) = \{d' \in D \mid d' \leq d\}$ .

### 3 Semantic Domains and Operations

We describe here the semantic domain of program states and their operations. A transition trace semantics (where the state is a stack of activation frames) can be defined either in the style of Bertelsen's small-step operational semantics [3] or the Abstract State Machines of Börger and Schulte [4]. Our set of operations is not sufficiently complete to give semantics for a concrete language like Java, (neither static fields and methods nor interfaces are considered). However it is representative of the set of operations that would be required for such a task. The correctness of the whole abstract interpretation follows by fairly standard means from the correctness of the abstraction of the operations [8,9].

For simplicity, integers and booleans are the only basic types. There are no static fields or methods. Methods are uniquely identified by a *Method\_ref* (for instance, their fully qualified name, like in Java), and classes by a *Class\_ref*. We define  $\text{Type} = \{\text{int}, \text{bool}\} \cup \text{Class\_ref}$ ,  $\text{Int} = \mathbb{Z}$  and  $\text{Bool} = \{\text{true}, \text{false}\}$ .

In the following, we assume that the class structure of a program is specified by some  $d \in Decl$ ,  $p \in Pars$  and  $n \in Name$  (see below), which give the classes, the parameters of the methods and their names, respectively, together with a subtype relation  $\leq$  on *Type*. For simplicity, we do not allow basic types to have subtypes. By *instance variables* we refer to fields. A *program variable* is a local variable, a parameter of a method or **this**. By *class variables* we denote the variables and fields which can reference objects.

**Definition 1.** Let  $Id$  be an infinite set of identifiers, with **this**  $\in Id$ . We define

$$Typing = \left\{ \tau : Id \rightarrow Type \left| \begin{array}{l} \text{dom}(\tau) \text{ is finite} \\ \text{if } \mathbf{this} \in \text{dom}(\tau) \text{ then } \tau(\mathbf{this}) \in \text{Class\_ref} \end{array} \right. \right\}$$

$$Class = Typing \times (Id \rightarrow \text{Method\_ref}) \quad Decl = (\text{Class\_ref} \mapsto \text{Class})$$

$$Pars = (\text{Method\_ref} \mapsto \text{seq}(Id) \times Typing) \quad Name = (\text{Method\_ref} \mapsto Id) .$$

If  $cl \in Class$ , we name its components as  $cl = \langle \tau, \nu \rangle$ . If  $p \in Pars$ ,  $n \in Name$  and  $mr \in \text{Method\_ref}$  then  $p(mr) = \langle s, \tau \rangle$ , where **this**,  $n(mr) \in s$  and  $\text{dom}(\tau) = s$ .

The hypothesis that  $n(mr) \in s$  means that the name of a method is among its local variables. Namely, it is used to store its result, like in Pascal procedures. Typings map variables to types, but **this** can be bound only to objects. The indirect definition of classes through *Decl* allows one to define two classes with the same structure. This structure is a pair whose first component is the typing of the fields defined or inherited by the class, and the second the table of the methods defined or inherited by the class. The map *Pars* binds every method reference to the typing (its signature) and its list of parameters. This list, whose elements are the domain of the typing, is useful since it yields their order in the definition of the method. The map *Name* binds a method reference to its name.

*Example 1.* Assume that we have two classes **a** and **b**, subclass of **a**, and that **a** (and, hence, **b**) has a field **n** of class **a**. We model this situation as  $\text{Class\_ref} = \{\kappa_a, \kappa_b\}$ ,  $[\kappa_a \mapsto a, \kappa_b \mapsto b] \in Decl$ ,  $a = b = \langle [\mathbf{n} \mapsto \kappa_a], [] \rangle$  and  $\kappa_b \leq \kappa_a$ . Note how *a* and *b* have the same structure but have different class references.

*Frames* yield values to the variables addressable at a given program point. The special variable **this** cannot be unbound. *Memories* map locations to objects. *Objects* contain their *Class\_ref* and the frame of their instance variables.

**Definition 2.** Let  $Loc$  be an infinite set of locations. Let  $Value = Int + Bool + Loc + \{nil\}$ . Given  $\tau \in Typing$ , we define frames, memories and objects as in Fig. 1. Given  $\mu_1, \mu_2 \in Memory$ , we say that  $\mu_2$  is an update of  $\mu_1$ , written  $\mu_1 \triangleleft \mu_2$ , if  $\text{dom}(\mu_1) \subseteq \text{dom}(\mu_2)$  and for every  $l \in \text{dom}(\mu_1)$  we have  $\mu_1(l).\kappa = \mu_2(l).\kappa$ .

Def. 2 states that the types of the variables in a frame are consistent with its typing, but class variables are just required to be bound to *nil* or a location (only a location for **this**). In order to guarantee that they contain objects of a class compatible with the typing, we define states, i.e., pairs of frame and memory.

$$\begin{array}{l}
\text{Frame}_\tau = \left\{ \phi \left| \begin{array}{l} \phi \in \text{dom}(\tau) \mapsto \text{Value and for every } v \in \text{dom}(\tau) \\ \text{if } \tau(v) = \text{int then } \phi(v) \in \text{Int} \\ \text{if } \tau(v) = \text{bool then } \phi(v) \in \text{Bool} \\ \text{if } \tau(v) \in \text{Class\_ref then } \phi(v) \in \{\text{nil}\} \cup \text{Loc} \\ \text{if } \text{this} \in \text{dom}(\tau) \text{ then } \phi(\text{this}) \in \text{Loc} \end{array} \right. \right\} \\
\text{Memory} = \left\{ \mu \in \text{Loc} \rightarrow \text{Obj} \right\} \quad \text{Obj} = \left\{ \langle \kappa, \phi \rangle \left| \begin{array}{l} \kappa \in \text{Class\_ref} \\ \phi \in \text{Frame}_{d(\kappa).\tau} \end{array} \right. \right\}
\end{array}$$

Fig. 1. Frames, memories and objects.

$$\begin{array}{l}
\text{nop}_\tau : \text{State}_\tau \mapsto \text{State}_\tau \\
\text{get\_int}_\tau^i : \text{State}_\tau \mapsto \text{State}_{\tau[\text{int}/\text{res}]} \quad \text{with } \text{res} \notin \text{dom}(\tau), i \in \text{Int} \\
\text{get\_nil}_\tau^{\text{cr}} : \text{State}_\tau \mapsto \text{State}_{\tau[\text{cr}/\text{res}]} \quad \text{with } \text{res} \notin \text{dom}(\tau), \text{cr} \in \text{Class\_ref} \\
\text{get\_bool}_\tau^b : \text{State}_\tau \mapsto \text{State}_{\tau[\text{bool}/\text{res}]} \quad \text{with } \text{res} \notin \text{dom}(\tau), b \in \text{Bool} \\
\text{get\_var}_\tau^v : \text{State}_\tau \mapsto \text{State}_{\tau[v]/\text{res}} \quad \text{with } v \in \text{dom}(\tau), \text{res} \notin \text{dom}(\tau) \\
\text{get\_field}_\tau^f : \text{State}_\tau \mapsto \text{State}_{\tau[d(\tau(\text{res})).\tau(f)]/\text{res}} \\
\quad \text{with } \text{res} \in \text{dom}(\tau), \tau(\text{res}) \in \text{Class\_ref}, f \in \text{dom}(d(\tau(\text{res})).\tau) \\
\text{put\_var}_\tau^v : \text{State}_\tau \mapsto \text{State}_{\tau|_{-\text{res}}} \quad \text{with } \text{res} \in \text{dom}(\tau), v \in \text{dom}(\tau), v \neq \text{res}, \tau(\text{res}) \leq \tau(v) \\
\text{put\_field}_\tau^f : \text{State}_\tau \mapsto (\text{State}_{\tau'} \mapsto \text{State}_{\tau|_{-\text{res}}}) \\
\quad \text{with } \text{res} \in \text{dom}(\tau), \tau(\text{res}) \in \text{Class\_ref}, f \in \text{dom}(d(\tau(\text{res})).\tau), \\
\quad \tau' = \tau[t/\text{res}] \text{ with } t \leq d(\tau(\text{res})).\tau(f) \\
=_\tau : \text{State}_\tau \mapsto (\text{State}_\tau \mapsto \text{State}_{\tau[\text{bool}/\text{res}]}) \quad \text{with } \text{res} \in \text{dom}(\tau) \\
\text{scope}_\tau^{mr, v_1, \dots, v_{\#p(mr).s-2}} : \text{State}_\tau \mapsto \text{State}_{p(mr).\tau|_{-n(mr)}} \\
\quad \text{with } \text{res} \in \text{dom}(\tau), \tau(\text{res}) \in \text{Class\_ref}, \\
\quad p(mr).s \setminus \{n(mr), \text{this}\} = \langle \iota_1, \dots, \iota_n \rangle, \tau(\text{res}) \leq p(mr).\tau(\text{this}), \\
\quad v_i \in \text{dom}(\tau) \text{ and } \tau(v_i) \leq p(mr).\tau(\iota_i) \text{ for every } i = 1, \dots, \#p(mr).s - 2 \\
\text{unscope}_\tau^{mr} : \text{State}_\tau \mapsto (\text{State}_{p(mr).\tau|_{-n(mr)}} \mapsto \text{State}_{\tau[p(mr).\tau(n(mr))/\text{res}]}) \quad \text{with } \text{res} \notin \text{dom}(\tau) \\
\text{restrict}_\tau^{vs} : \text{State}_\tau \mapsto \text{State}_{\tau|_{-vs}} \quad \text{with } vs \subseteq \text{dom}(\tau) \\
\text{expand}_\tau^{v:t} : \text{State}_\tau \mapsto \text{State}_{\tau[t/v]} \quad \text{with } v \notin \text{dom}(\tau), t \in \text{Type} \\
\text{lookup}_\tau^{id, mr} \subseteq \wp(\text{State}_\tau) \\
\quad \text{with } \text{res} \in \text{dom}(\tau), \tau(\text{res}) \in \text{Class\_ref}, id \in \text{dom}(d(\tau(\text{res})).\nu), mr \in \text{Method\_ref} \\
\text{new}_\tau^{\text{cr}} : \text{State}_\tau \mapsto \text{State}_{\tau[\text{cr}/\text{res}]} \quad \text{with } \text{res} \notin \text{dom}(\tau), \text{cr} \in \text{Class\_ref} \\
\text{is\_true}_\tau, \text{is\_false} \subseteq \wp(\text{State}_\tau) \quad \text{with } \tau(\text{res}) = \text{bool}
\end{array}$$

Fig. 2. The signatures of the operations over the states.

**Definition 3.** Let  $\tau \in \text{Typing}$ . Recall that  $\leq$  is the subtype relation. We say that  $\phi \in \text{Frame}_\tau$  is  $\tau$ -correct in a memory  $\mu$ , written  $\phi \propto_\tau \mu$ , if it binds the class variables in its domain to objects of classes compatible with  $\tau$ . Namely,  $\phi \propto_\tau \mu$  if and only if for every  $v \in \text{dom}(\phi)$  such that  $\phi(v) \in \text{Loc}$  we have  $\phi(v) \in \text{dom}(\mu)$  and  $(\mu(\phi(v))).\kappa \leq \tau(v)$ . We define

$$\text{State}_\tau = \left\{ \langle \phi, \mu \rangle \left| \begin{array}{l} \phi \in \text{Frame}_\tau, \mu \in \text{Memory}, \phi \propto_\tau \mu \\ \text{and for every } \langle \kappa, \phi' \rangle \in \text{codom}(\mu) \text{ we have } \phi' \propto_{d(\kappa).\tau} \mu \end{array} \right. \right\}.$$

$\text{nop}_\tau((\phi, \mu)) = \langle \phi, \mu \rangle$	$\text{get\_int}_\tau^i((\phi, \mu)) = \langle \phi[\text{res}], \mu \rangle$
$\text{get\_nil}_\tau^{cr}((\phi, \mu)) = \langle \phi[\text{nil}/\text{res}], \mu \rangle$	$\text{get\_bool}_\tau^b((\phi, \mu)) = \langle \phi[\text{b}/\text{res}], \mu \rangle$
$\text{get\_var}_\tau^v((\phi, \mu)) = \langle \phi[\phi(v)/\text{res}], \mu \rangle$	$\text{put\_var}_\tau^v((\phi, \mu)) = \langle \phi[\phi(\text{res})/v]_{-\text{res}}, \mu \rangle$
$\text{get\_field}_\tau^f((\phi', \mu)) = \begin{cases} \langle \phi'[(\mu(\phi'(\text{res})).\phi)(f)/\text{res}], \mu \rangle & \text{if } \phi'(\text{res}) \neq \text{nil} \\ \text{undefined} & \text{otherwise} \end{cases}$	
$\text{put\_field}_\tau^f((\phi_1, \mu_1))((\phi_2, \mu_2)) = \begin{cases} \langle \phi_2 _{-\text{res}}, \mu_2[\langle \mu_2(l).\kappa, \mu_2(l).\phi[\phi_2(\text{res})/f] \rangle / l] \rangle & \text{if } \mu_1 \triangleleft \mu_2 \text{ and } (l = \phi_1(\text{res})) \neq \text{nil} \\ \text{undefined} & \text{otherwise} \end{cases}$	
$=_\tau((\phi_1, \mu_1))((\phi_2, \mu_2)) = \begin{cases} \langle \phi_2[\text{true}/\text{res}], \mu_2 \rangle & \text{if } \mu_1 \triangleleft \mu_2 \text{ and } \phi_1(\text{res}) = \phi_2(\text{res}) \\ \langle \phi_2[\text{false}/\text{res}], \mu_2 \rangle & \text{if } \mu_1 \triangleleft \mu_2 \text{ and } \phi_1(\text{res}) \neq \phi_2(\text{res}) \\ \text{undefined} & \text{otherwise} \end{cases}$	
$\text{scope}_\tau^{mr, v_1, \dots, v_n}((\phi, \mu)) = \langle [l_1 \mapsto \phi(v_1), \dots, l_n \mapsto \phi(v_n), \text{this} \mapsto \phi(\text{res})], \mu \rangle$ where $\langle l_1, \dots, l_n \rangle = p(mr).s \setminus \{n(mr), \text{this}\}$	
$\text{unscope}_\tau^{mr}((\phi_1, \mu_1))((\phi_2, \mu_2)) = \begin{cases} \langle \phi_1[\phi_2(n(mr))/\text{res}], \mu_2 \rangle & \text{if } \mu_1 \triangleleft \mu_2 \\ \text{undefined} & \text{otherwise} \end{cases}$	
$\text{restrict}_\tau^{vs}((\phi, \mu)) = \langle \phi _{-vs}, \mu \rangle$	$\text{expand}_\tau^{vt}((\phi, \mu)) = \langle \phi[\text{init}(t)/v], \mu \rangle$
$\text{lookup}_\tau^{id, mr}((\phi, \mu))$ if and only if $\phi(\text{res}) \neq \text{nil}$ and $d(\mu(\phi(\text{res})).\kappa).\nu(id) = mr$	
$\text{new}_\tau^{cr}((\phi, \mu)) = \langle \phi[l/\text{res}], \mu[\langle cr, \text{init}(d(cr).\tau) \rangle / l] \rangle$ $l \in \text{Loc}$ fresh	
$\text{is\_true}_\tau((\phi, \mu))$ if and only if $\phi(\text{res}) = \text{true}$	$\text{is\_false}_\tau((\phi, \mu))$ if and only if $\phi(\text{res}) = \text{false}$ ,

where  $\text{init}(int) = 0$ ,  $\text{init}(bool) = \text{false}$ ,  $\text{init}(cr) = \text{nil}$  for  $cr \in \text{Class\_ref}$  and  $\text{init}(\tau) = \lambda v \in \text{dom}(\tau).\text{init}(\tau(v))$  for  $\tau \in \text{Typing}$ .

**Fig. 3.** The operations over the states.

This set forms an algebraic structure with the operations whose signatures are given in Figure 2 and which are explicitly defined in Figure 3.

In these operations the variable  $\text{res}$  stands for the place where intermediate results are stored and retrieved. In a semantics for a stack-based language such as the Java bytecode,  $\text{res}$  would be the top of the operand stack. The binary operations of Figure 3 are undefined when the memory of the second argument is not an update of that of the first one, to guarantee that the  $\alpha_\tau$  relation between frame and memory of a state (Definition 3) is not broken by the operation.

The  $\text{nop}$  operation does nothing. The  $\text{get\_int}$  operation loads an integer into  $\text{res}$ . Similarly for  $\text{get\_nil}$  and  $\text{get\_bool}$ . The  $\text{get\_var}$  operation fetches the value of a variable  $v$  and loads it into  $\text{res}$ . The  $\text{get\_field}$  operation fetches from  $\text{res}$  a non- $\text{nil}$  reference to the object containing the field. The content of this field is loaded into  $\text{res}$ , whose type is updated with the declared type of the field. Note how this content is found. Since  $\phi'(\text{res})$  points to the object containing the field,  $\mu(\phi'(\text{res}))$  is that object. Its fields are in  $\mu(\phi'(\text{res})).\phi$  and the field  $f$  is then  $\mu(\phi'(\text{res})).\phi(f)$ . The  $\text{put\_var}$  operation copies the content of  $\text{res}$  into  $v$ . The declared type of  $\text{res}$  must be a subtype of that of  $v$ . There is no resulting value. Thus, the variable  $\text{res}$  is removed from the typing of the result. In the  $\text{put\_field}$

operation the *res* variable of the first argument points to the target object, while that of the second argument contains the new value for the field, whose declared type must be a subtype of that of the field. By using two states instead of an object and a value, we deal with just one semantic domain. In Fig. 3 we see its implementation. Since  $l = \phi_1(res)$  points to the target object, it must not be *nil* and  $\mu_2(l)$  is that object (since  $\mu_1 < \mu_2$ ). We write  $\phi_2(res)$  in the variable *f* of the frame  $\mu_2(l).\phi$  of that object. For every binary operation of the language there is a corresponding binary operation on states, like the shown case of =.

Four operations (*scope*, *unscope*, *restrict* and *expand*) modify the structure of frames and states. The operations *scope* and *unscope* are used before and after a method call, respectively. The former creates a new frame in which the invoked method *mr* executes. Its typing is  $p(mr).\tau|_{-n(mr)}$  since the variable  $n(mr)$ , i.e., the name of the method, is not among its inputs. Note that  $p(mr).\tau|_{-n(mr)}$  contains exactly the parameters of the method and the implicit **this** parameter. In this operation *res* points to the object over which the method is called and becomes the new **this** variable, which justifies the subtype check in Fig. 2. The first argument of the *unscope* operation is the state before the method call. The second is the state at the end of the execution of the method. The frame of this second state contains just a variable with the name of the method, and holds its result like in Pascal procedures (in a complete operational semantics, the frames of these two states would be the top two elements of the stack). The operation restores the frame at the beginning of the execution of the method, by copying into *res* its result, and yields the memory at the end of its execution. The *restrict* operation removes some variables from a state and *expand* adds a new initialised variable to a state.

The *lookup* predicate for dynamic method lookup holds if by invoking the method *id* on the object referenced by *res* the method referenced by *mr* is selected. The object on which the method is invoked is  $\mu(\phi(res))$ , its class reference  $\mu(\phi(res)).\kappa$  and the list of its methods  $d(\mu(\phi(res)).\kappa).\nu$ . If  $d(\mu(\phi(res)).\kappa).\nu(id) = mr$  means that by calling *id* we select the method *mr*. The new operation creates a new initialised object and loads it into *res*. The *is.true* (*is.false*) predicate contains those states whose *res* variable contains *true* (*false*).

*Example 2.* In the situation of Ex. 1, consider how the following sequence of operations change the state. Reading downwards, we start from a state containing just the variable  $v_1$ , we introduce a new variable  $v_2$ , we create and store in *res* a newly initialised object, we read its **n** field into *res* and we store it into  $v_2$ .

Operation	State
	$\langle [v_1 \mapsto l_1], [l_1 \mapsto \langle \kappa_a, [n \mapsto nil]] \rangle \rangle$ (initial state)
$\text{expand}_{[v_1 \mapsto \kappa_a]}^{v_2: \kappa_a}$	$\langle [v_1 \mapsto l_1, v_2 \mapsto nil], [l_1 \mapsto \langle \kappa_a, [n \mapsto nil]] \rangle \rangle$
$\text{new}_{[v_1, v_2 \mapsto \kappa_a]}^{\kappa_b}$	$\langle [v_1 \mapsto l_1, v_2 \mapsto nil, res \mapsto l_2], [l_1 \mapsto \langle \kappa_a, [n \mapsto nil]] \rangle, l_2 \mapsto \langle \kappa_b, [n \mapsto nil] \rangle \rangle$
$\text{get\_field}_{[v_1, v_2 \mapsto \kappa_a, res \mapsto \kappa_b]}^a$	$\langle [v_1 \mapsto l_1, v_2 \mapsto nil, res \mapsto nil], [l_1 \mapsto \langle \kappa_a, [n \mapsto nil]] \rangle, l_2 \mapsto \langle \kappa_b, [n \mapsto nil] \rangle \rangle$
$\text{put\_var}_{[v_1, v_2, res \mapsto \kappa_a]}^{v_2}$	$\langle [v_1 \mapsto l_1, v_2 \mapsto nil], [l_1 \mapsto \langle \kappa_a, [n \mapsto nil]] \rangle, l_2 \mapsto \langle \kappa_b, [n \mapsto nil] \rangle \rangle$

$$\begin{aligned}
& \text{nop}_\tau^{rta}(s) = (\text{get.int}_\tau^i)^{rta}(s) = (\text{get.nil}_\tau^{cr})^{rta}(s) = \text{is.true}_\tau^{rta}(s) = s \\
& (\text{get.bool}_\tau^b)^{rta}(s) = (\text{get.var}_\tau^v)^{rta}(s) = (\text{put.var}_\tau^v)^{rta}(s) = \text{is.false}_\tau^{rta}(s) = s \\
& (\text{get.field}_\tau^f)^{rta}(s) = \begin{cases} \emptyset & \text{if } s \cap \downarrow\tau(\text{res}) = \emptyset \\ s & \text{otherwise} \end{cases} \\
& (\text{put.field}_{\tau,\tau'}^f)^{rta}(s_1)(s_2) = \begin{cases} \emptyset & \text{if } s_1 \cap s_2 \cap \downarrow\tau(\text{res}) = \emptyset \\ s_2 & \text{otherwise} \end{cases} \\
& =_\tau^{rta}(s_1)(s_2) = s_2 \quad (\text{scope}_\tau^{mr,v_1,\dots,v_n})^{rta}(s) = s \\
& (\text{unscope}_\tau^{mr})^{rta}(s_1)(s_2) = s_2 \quad (\text{restrict}_\tau^{vs})^{rta}(s) = s \\
& (\text{expand}_\tau^{v:t})^{rta}(s) = s \quad (\text{new}_\tau^{cr})^{rta}(s) = s \cup \{cr\} \quad \cup_\tau^{rta}(s_1)(s_2) = s_1 \cup s_2 \\
& (\text{lookup}_\tau^{id,mr})^{rta}(s) = \begin{cases} s & \text{if there exists } \kappa \in s \cap \downarrow\tau(\text{res}) \\ & \text{such that } d(\kappa).v(id) = mr \\ \emptyset & \text{otherwise.} \end{cases}
\end{aligned}$$

Fig. 4. The instantiation over  $State^{rta}$  of the operations of Figure 3 and of  $\cup$ .

## 4 Rapid Type Analysis (*rta*-Analysis)

Bacon and Sweeney defined [2] a simple and efficient *rapid type analysis* (*rta*) of little precision. They collect the set  $N$  of classes instantiated by a **new** command. The classes of an expression  $e$  are  $(\downarrow D(e)) \cap N$ , where  $D(e)$  is the declared type of  $e$ . In terms of abstract interpretation, this technique amounts to defining an abstraction which collects the classes of the objects in the memory of the states. A priori, abstract interpretation can distinguish this set of classes in different program points. Then it results in a more precise analysis than that in [2].

**Definition 4.** For  $\tau \in Typing$ , the abstract domain is  $State_\tau^{rta} = \wp(Class\_ref)$  with the concretisation map  $\gamma^{rta} : State_\tau^{rta} \mapsto \wp(State_\tau)$  such that (here,  $o \in \text{codom}(\mu)$  means that  $o$  is an object in the memory  $\mu$ )  $\gamma^{rta}(s) = \{\langle \phi, \mu \rangle \in State_\tau \mid \text{for every } \langle \kappa, \phi' \rangle \in \text{codom}(\mu) \text{ we have } \kappa \in s\}$ .

**Proposition 1.** For  $\tau \in Typing$ , the set  $\gamma^{rta}(State_\tau^{rta})$  is a Moore family of  $\wp(State_\tau)$ . Hence  $State_\tau^{rta}$  is an abstract domain whose induced abstraction map is  $\alpha^{rta} : \wp(State_\tau) \mapsto State_\tau^{rta}$  such that

$$\alpha^{rta}(S) = \{\kappa \in Class\_ref \mid \langle \phi, \mu \rangle \in S \text{ and } \langle \kappa, \phi' \rangle \in \text{codom}(\mu)\} .$$

*Proof (Sketch).* The non-empty set  $\gamma^{rta}(State_\tau^{rta})$  is a Moore family since it is  $\cap$ -closed. Indeed, it can be shown that for every  $\{s_i\}_{i \in \mathbb{N}} \subseteq State_\tau^{rta}$  we have  $\bigcap_{i \in \mathbb{N}} \gamma_\tau^{rta}(s_i) = \gamma_\tau^{rta}(\bigcap_{i \in \mathbb{N}} s_i)$ . The  $\alpha^{rta}$  map is derived [8,9] as  $\alpha^{rta}(S) = \bigcap \{s \in State_\tau^{rta} \mid S \subseteq \gamma^{rta}(s)\} = \bigcap \{s \in \wp(Class\_ref) \mid \langle \phi, \mu \rangle \in S \text{ and } \langle \kappa, \phi' \rangle \in \text{codom}(\mu) \text{ entails } \kappa \in s\} = \{\kappa \in Class\_ref \mid \langle \phi, \mu \rangle \in S \text{ and } \langle \kappa, \phi' \rangle \in \text{codom}(\mu)\}$ .



**Proposition 2.** *The optimal abstract counterparts of the powerset extension of the operations of Figure 3 and of  $\cup$  are those given in Figure 4.*

The only operation in Figure 4 which enlarges the set of classes created during the execution of the program is `new`. The operations `get.field` and `put.field` check if a class compatible with the declared type of `res` has been instantiated, since otherwise `res` must be bound to `nil`, and the concrete operation fails. Similarly, `lookup` checks if it has been instantiated some class such that a call to the method `id` results in the execution of the method referenced by `mr`. Note how these operations allow the set of classes instantiated to shrink. This is a consequence of the use of abstract interpretation, which deals with different program points. It is a distinguishing feature w.r.t. the original technique in [2].

*Example 3.* We execute over  $State_\tau^{rta}$  the same sequence of instructions of Ex. 2. The initial state is the abstraction (Prop. 1) of the initial concrete state. The final abstract state says that  $v_1$  and  $v_2$ , both of declared type  $a$ , can be bound to objects of class  $a$  or  $b$ , which is a correct but rather imprecise approximation of the concrete final state of Example 2.

Operation	State
	$\{\kappa_a\}$ (initial state)
$(\text{expand}_{[v_1 \mapsto \kappa_a]}^{v_2: \kappa_a})^{rta}$	$\{\kappa_a\}$
$(\text{new}_{[v_1, v_2 \mapsto \kappa_a]}^{\kappa_b})^{rta}$	$\{\kappa_a, \kappa_b\}$
$(\text{get\_field}_{[v_1, v_2 \mapsto \kappa_a, res \mapsto \kappa_b]}^n)^{rta}$	$\{\kappa_a, \kappa_b\}$
$(\text{put\_var}_{[v_1, v_2, res \mapsto \kappa_a]}^{v_2})^{rta}$	$\{\kappa_a, \kappa_b\}$

## 5 Dataflow Analysis (*df*-Analysis)

In the dataflow analysis of [10] only program variables are analysed, while fields are approximated with the downwards closure  $\downarrow(t)$  of their declared type  $t$ .

Compare this with Definition 2. An abstract value is a set of types which share a supertype. An abstract frame maps variables into abstract values consistent with the given typing. The special variable `this` cannot be unbound.

**Definition 5.** *Given  $\tau \in Typing$ , we define  $Value^{df} = \{S \in \wp(Type) \mid S \subseteq \downarrow t \text{ for some } t \in Type\}$  and*

$$Frame_\tau^{df} = \left\{ \phi \left| \begin{array}{l} \phi \in \text{dom}(\tau) \mapsto Value^{df} \text{ and for every } v \in \text{dom}(\tau) \\ \text{if } \tau(v) = int \text{ then } \phi(v) = \{int\} \\ \text{if } \tau(v) = bool \text{ then } \phi(v) = \{bool\} \\ \text{if } \tau(v) \in Class\_ref \text{ then } \phi(v) \subseteq \downarrow \tau(v) \\ \text{and if } this \in \text{dom}(\tau) \text{ then } \phi(this) \neq \emptyset \end{array} \right. \right\}.$$

In  $Frame_\tau^{df}$ , integer and boolean variables are always mapped to  $\{int\}$  and  $\{bool\}$ , respectively. We consider them just to simplify the presentation.

There is no abstract memory. Then abstract states are just abstract frames (compare with Definition 3). The special abstract state  $\emptyset$  represents the empty set of concrete states and improves the precision of the analysis.

$\text{nop}_\tau^{\text{df}}(\phi) = \phi$	$(\text{get\_int}_\tau^i)^{\text{df}}(\phi) = \phi[\{int\}/res]$
$(\text{get\_nil}_\tau^{\text{cr}})^{\text{df}}(\phi) = \phi[\emptyset/res]$	$(\text{get\_bool}_\tau^b)^{\text{df}}(\phi) = \phi[\{bool\}/res]$
$(\text{get\_var}_\tau^v)^{\text{df}}(\phi) = \phi[\phi(v)/res]$	$(\text{put\_var}_\tau^v)^{\text{df}}(\phi) = \phi[\phi(res)/v]_{-res}$
$(\text{get\_field}_\tau^f)^{\text{df}}(\phi) = \begin{cases} \emptyset & \text{if } \phi(res) = \emptyset \\ \phi[\downarrow(d(\tau(res)).\tau(f))/res] & \text{otherwise} \end{cases}$	
$(\text{put\_field}_{\tau,\tau'}^f)^{\text{df}}(\phi_1)(\phi_2) = \begin{cases} \emptyset & \text{if } \phi_1(res) = \emptyset \\ \phi_2 _{-res} & \text{otherwise} \end{cases}$	
$=_{\tau}^{\text{df}}(\phi_1)(\phi_2) = \phi_2[\{bool\}/res]$	$\cup_{\tau}^{\text{df}}(\phi_1)(\phi_2) = \lambda v \in \text{dom}(\tau). \phi_1(v) \cup \phi_2(v)$
$(\text{scope}_{\tau}^{mr,v_1,\dots,v_n})^{\text{df}}(\phi) = [\iota_1 \mapsto \phi(v_1), \dots, \iota_n \mapsto \phi(v_n), \text{this} \mapsto \phi(res)]$	$\text{where } \langle \iota_1, \dots, \iota_n \rangle = p(mr).s \setminus \{n(mr), \text{this}\}$
$(\text{unscope}_{\tau}^{mr})^{\text{df}}(\phi_1)(\phi_2) = \phi_1[\phi_2(n(mr))/res]$	
$(\text{restrict}_{\tau}^{vs})^{\text{df}}(\phi) = \phi _{-vs}$	$(\text{expand}_{\tau}^{v:t})^{\text{df}}(\phi) = \phi[\text{init}^{\text{df}}(t)/v]$
$(\text{lookup}_{\tau}^{id,mr})^{\text{df}}(\phi) = \begin{cases} \phi[S/res] & \text{if } S = \{cr \in \phi(res) \mid d(cr).\nu(id) = mr\} \neq \emptyset \\ \emptyset & \text{otherwise} \end{cases}$	
$(\text{new}_{\tau}^{\text{cr}})^{\text{df}}(\phi) = \phi[\{cr\}/res]$	$\text{is\_true}_{\tau}^{\text{df}}(\phi) = \text{is\_false}_{\tau}^{\text{df}}(\phi) = \phi,$

where  $\text{init}^{\text{df}}(int) = \{int\}$ ,  $\text{init}^{\text{df}}(bool) = \{bool\}$ ,  $\text{init}^{\text{df}}(cr) = \emptyset$  for  $cr \in \text{Class\_ref}$ .

**Fig. 5.** The instantiation over  $\text{State}^{\text{df}}$  of the operations of Figure 3 and of  $\cup$ .

**Definition 6.** For  $\tau \in \text{Typing}$ , the abstract domain is  $\text{State}_{\tau}^{\text{df}} = \{\emptyset\} \cup \text{Frame}_{\tau}^{\text{df}}$ .

The relation  $\approx$  says when an abstract value (i.e., a set of types) approximates a concrete value. The set  $\{int\}$  approximates the integers, the set  $\{bool\}$  approximates the booleans and a set of classes  $S$  approximates  $nil$  and all locations containing an object of a class in  $S$ .

**Definition 7.** Let  $\mu \in \text{Memory}$ . We define  $\approx_{\mu}: \text{Value}^{\text{df}} \times \text{Value}$  as the minimal relation such that  $\{int\} \approx_{\mu} i$ , with  $i \in \text{Int}$ ,  $\{bool\} \approx_{\mu} b$ , with  $b \in \text{Bool}$ ,  $S \approx_{\mu} nil$ , with  $S \neq \{int\}$  and  $S \neq \{bool\}$ , and  $S \approx_{\mu} l$  with  $l \in \text{Loc}$  and  $\mu(l).\kappa \in S$ . This relation is pointwise extended to  $\approx_{\mu}: \text{Frame}_{\tau}^{\text{df}} \times \text{Frame}_{\tau}$ , for  $\tau \in \text{Typing}$ .

The concretisation of an abstract state  $\phi$  is the set of concrete states  $\langle \phi', \mu \rangle$  whose frame  $\phi'$  binds every variable  $v \in \text{dom}(\phi)$  to a concrete value  $\phi'(v)$  approximated ( $\approx_{\mu}$ ) by the abstract value  $\phi(v)$ , i.e.,  $\phi(v) \approx_{\mu} \phi'(v)$ .

**Definition 8.** For  $\tau \in \text{Typing}$ , we define the concretisation map  $\gamma^{\text{df}}: \text{State}_{\tau}^{\text{df}} \mapsto \wp(\text{State}_{\tau})$  such that  $\gamma^{\text{df}}(\emptyset) = \emptyset$  and  $\gamma^{\text{df}}(\phi) = \{\langle \phi', \mu \rangle \in \text{State}_{\tau} \mid \phi \approx_{\mu} \phi'\}$  for every  $\phi \in \text{State}_{\tau}^{\text{df}} \setminus \{\emptyset\}$ .

**Proposition 3.** For  $\tau \in \text{Typing}$ , the set  $\gamma^{\text{df}}(\text{State}_{\tau}^{\text{df}})$  is a Moore family of  $\wp(\text{State}_{\tau})$ . Hence  $\text{State}_{\tau}^{\text{df}}$  is an abstract domain whose induced abstraction map is  $\alpha^{\text{df}}: \wp(\text{State}_{\tau}) \mapsto \text{State}_{\tau}^{\text{df}}$  such that  $\alpha^{\text{df}}(\emptyset) = \emptyset$  and, for  $S \neq \emptyset$ ,  $\alpha^{\text{df}}(S) = a$  such that for  $v \in \text{dom}(\tau)$ ,  $a(v) = \{int\}$  if  $\tau(v) = int$ ,  $a(v) = \{bool\}$  if  $\tau(v) = bool$  and  $a(v) = \{\mu(\phi(v)).\kappa \mid \langle \phi, \mu \rangle \in S \text{ and } \phi(v) \in \text{Loc}\}$  if  $\tau(v) \in \text{Class\_ref}$ . In this last case,  $\alpha^{\text{df}}$  collects the classes of the objects bound to  $v$  in some concrete state.

**Proposition 4.** *The optimal abstract counterparts of the powerset extension of the operations of Fig. 3 and of  $\cup$  are those given in Fig. 5. All of them, except  $\cup_\tau^{df}$ , are strict on all their arguments. For instance,  $(\text{unscope}_\tau^{nr})^{df}(\emptyset)(\phi_2) = (\text{unscope}_\tau^{nr})^{df}(\phi_1)(\emptyset) = (\text{unscope}_\tau^{nr})^{df}(\emptyset)(\emptyset) = \emptyset$ . For  $\cup_\tau^{df}$ , we define  $\cup_\tau^{df}(\emptyset)(\phi_2) = \phi_2$ ,  $\cup_\tau^{df}(\phi_1)(\emptyset) = \phi_1$  and  $\cup_\tau^{df}(\emptyset)(\emptyset) = \emptyset$ .*

*Example 4.* We execute over  $\text{State}_\tau^{df}$  the same sequence of instructions of Ex. 2. The initial state is the abstraction (Prop. 3) of the initial concrete state.

Operation	State
	$[v_1 \mapsto \{\kappa_a\}]$ (initial state)
$(\text{expand}_{[v_1 \mapsto \kappa_a]}^{v_2:\kappa_a})^{df}$	$[v_1 \mapsto \{\kappa_a\}, v_2 \mapsto \emptyset]$
$(\text{new}_{[v_1, v_2 \mapsto \kappa_a]}^{\kappa_b})^{df}$	$[v_1 \mapsto \{\kappa_a\}, v_2 \mapsto \emptyset, \text{res} \mapsto \{\kappa_b\}]$
$(\text{get\_field}_{[v_1, v_2 \mapsto \kappa_a, \text{res} \mapsto \kappa_b]}^n)^{df}$	$[v_1 \mapsto \{\kappa_a\}, v_2 \mapsto \emptyset, \text{res} \mapsto \{\kappa_a, \kappa_b\}]$
$(\text{put\_var}_{[v_1, v_2, \text{res} \mapsto \kappa_a]}^{v_2})^{df}$	$[v_1 \mapsto \{\kappa_a\}, v_2 \mapsto \{\kappa_a, \kappa_b\}]$

The final abstract state says that  $v_1$  can be bound to objects of class  $a$  and  $v_2$  to objects of class  $a$  or  $b$ , which is a correct approximation of the concrete final state of Example 2. It is strictly more precise than that obtained in Example 3.

## 6 0-CFA (*ps*-Analysis)

In [11], Palsberg and Schwartzbach defined a class analysis as a constraint problem. Every variable and every field is given a set of classes. These sets are then related by constraints which model the dataflow of the program. This means that the information about program variables is that used in the case of the dataflow analysis of Section 5, but this analysis deals with fields too. However, all objects of the same class are identified. The analysis of a method is the same for every call point and hence this analysis is a 0-CFA *analysis* [13].

That analysis leads to an abstract domain made of an abstract frame, identical to those used in Section 5, and of an abstract memory. The abstract memory is an abstract frame for the fields of the classes. We can assume without any loss of generality that fields in different classes have different names.

The difference with the original technique in [11] is that abstract interpretation does not *a priori* identify different occurrences of the same variable, while in [11] the same variable has assigned the same set of classes in every program point. This, together with the formally proved optimality of our abstract operations (Proposition 6), suggests that our analysis should be more precise than that in [11], without using any technique like multiple variables for different uses of the same variable or method splitting [12].

**Definition 9.** *We define the typing  $\bar{\tau} = \cup_{\kappa \in \text{Class\_refd}(\kappa).\tau}$  of the fields of all classes. It makes sense since fields have different names. For  $\tau \in \text{Typing}$  such that  $\text{dom}(\tau) \subseteq \text{dom}(\bar{\tau})$  and  $\phi \in \text{Frame}_\tau$ , we define  $\bar{\phi} \in \text{Frame}_{\bar{\tau}}$  as  $\bar{\phi}(v) = \phi(v)$  if  $v \in \text{dom}(\tau)$ , and  $\bar{\phi}(v) = \text{init}(\bar{\tau}(v))$  otherwise ( $\text{init}$  has been defined in Fig. 3).*

$\text{nop}_\tau^{ps}(\langle \phi, \mu \rangle) = \langle \phi, \mu \rangle$	$(\text{get\_int}_\tau^i)^{ps}(\langle \phi, \mu \rangle) = \langle \phi[\{int\}/res], \mu \rangle$
$(\text{get\_nil}_\tau^{cr})^{ps}(\langle \phi, \mu \rangle) = \langle \phi[\emptyset/res], \mu \rangle$	$(\text{get\_bool}_\tau^b)^{ps}(\langle \phi, \mu \rangle) = \langle \phi[\{bool\}/res], \mu \rangle$
$(\text{get\_var}_\tau^v)^{ps}(\langle \phi, \mu \rangle) = \langle \phi[\phi(v)/res], \mu \rangle$	$(\text{put\_var}_\tau^v)^{ps}(\langle \phi, \mu \rangle) = \langle \phi[\phi(res)/v]_{-res}, \mu \rangle$
$(\text{get\_field}_\tau^f)^{ps}(\langle \phi, \mu \rangle) = \begin{cases} \langle \emptyset, \mu \rangle & \text{if } \phi(res) = \emptyset \\ \langle \phi[\mu(f)/res], \mu \rangle & \text{otherwise} \end{cases}$	
$(\text{put\_field}_\tau^f)^{ps}(\langle \phi_1, \mu_1 \rangle)(\langle \phi_2, \mu_2 \rangle) = \begin{cases} \langle \emptyset, \mu_2 \rangle & \text{if } \phi_1(res) = \emptyset \\ \langle \phi_2 _{-res}, \mu_2[\mu_2(f) \cup \phi_2(res)/f] \rangle & \text{otherwise} \end{cases}$	
$\quad =_{\tau}^{ps}(\langle \phi_1, \mu_1 \rangle)(\langle \phi_2, \mu_2 \rangle) = \langle \phi_2[\{bool\}/res], \mu_2 \rangle$	
$(\text{scope}_\tau^{mr, v_1, \dots, v_n})^{ps}(\langle \phi, \mu \rangle) = \langle [t_1 \mapsto \phi(v_1), \dots, t_n \mapsto \phi(v_n), \text{this} \mapsto \phi(res)], \mu \rangle$	where $\langle t_1, \dots, t_n \rangle = p(mr).s \setminus \{n(mr), \text{this}\}$
$(\text{unscope}_\tau^{mr})^{ps}(\langle \phi_1, \mu_1 \rangle)(\langle \phi_2, \mu_2 \rangle) = \langle \phi_1[\phi_2(n(mr))/res], \mu_2 \rangle$	
$(\text{restrict}_\tau^{vs})^{ps}(\langle \phi, \mu \rangle) = \langle \phi _{-vs}, \mu \rangle$	$(\text{expand}_\tau^{v:t})^{ps}(\langle \phi, \mu \rangle) = \langle \phi[\text{init}^{df}(t)/v], \mu \rangle$
$(\text{lookup}_\tau^{id, mr})^{ps}(\langle \phi, \mu \rangle) = \begin{cases} \langle \phi[S/res], \mu \rangle & \text{if } S = \{cr \in \phi(res) \mid d(cr).v(id) = mr\} \neq \emptyset \\ \langle \emptyset, \mu \rangle & \text{otherwise} \end{cases}$	
$(\text{new}_\tau^{cr})^{ps}(\langle \phi, \mu \rangle) = \langle \phi[\{cr\}/res], \mu \rangle$	$\text{is\_true}_\tau^{ps}(\langle \phi, \mu \rangle) = \text{is\_false}_\tau^{ps}(\langle \phi, \mu \rangle) = \langle \phi, \mu \rangle$
$\cup_\tau^{ps}(\langle \phi_1, \mu_1 \rangle)(\langle \phi_2, \mu_2 \rangle) = \langle \cup_\tau^{df}(\phi_1)(\phi_2), \cup_\tau^{df}(\mu_1)(\mu_2) \rangle,$	

where  $\text{init}^{df}$  and  $\cup^{df}$  are defined in Figure 5.

**Fig. 6.** The instantiation over  $State^{ps}$  of the operations of Fig. 3 and of  $\cup$ .

**Definition 10.** For  $\tau \in Typing$ , the abstract domain is  $State_\tau^{ps} = (Frame_\tau^{ps} \times Memory^{ps})$ , with  $Frame_\tau^{ps} = State_\tau^{df}$  and  $Memory^{ps} = Frame_\tau^{df}$  (Defs. 5 and 6).

Compare the following definition with Definition 8. The relation  $\approx$  has been defined in Definition 7. An element  $\langle \phi, \mu \rangle \in State_\tau^{ps}$  represents those concrete states whose frame is compatible with  $\phi$ , i.e., their class variables are bound to objects of a class allowed by  $\phi$ , and whose memory contains objects with an internal frame compatible with  $\mu$ , i.e., their fields are bound to objects of a class allowed by  $\mu$ . Since the frame of an object is for its instance variables only, we extend it to the whole set of fields before its comparison with  $\mu$ .

**Definition 11.** Given  $\tau \in Typing$ , we define the concretisation map  $\gamma^{ps} : State_\tau^{ps} \mapsto \wp(State_\tau)$  as

$$\gamma^{ps}(\langle \phi, \mu \rangle) = \gamma^{df}(\phi) \cap \left\{ \langle \phi', \mu' \rangle \in State_\tau \mid \begin{array}{l} \text{for every } \langle \kappa, \phi'' \rangle \in \text{codom}(\mu') \\ \text{we have } \mu \approx_{\mu'} \overline{\phi''} \end{array} \right\}.$$

The following abstraction map formalises the idea of the analysis. A set of concrete states  $S$  is abstracted through  $\alpha^{df}$ , as in Section 5, but a second component provides more information about the fields. Namely, it is the abstraction through  $\alpha^{df}$  of the states of the objects in memory.

**Proposition 5.** For  $\tau \in \text{Typing}$ , the set  $\gamma^{ps}(\text{State}_\tau^{ps})$  is a Moore family of  $\wp(\text{State}_\tau)$ . Hence  $\text{State}_\tau^{ps}$  is an abstract domain whose induced abstraction map is  $\alpha^{ps} : \wp(\text{State}_\tau) \mapsto \text{State}_\tau^{ps}$  such that  $(\text{init}(\bar{\tau}))$  avoids empty abstract memories

$$\alpha^{ps}(S) = \langle \alpha^{df}(S), \alpha^{df}(\{\text{init}(\bar{\tau})\} \cup \{\langle \bar{\phi}', \mu \rangle \mid \langle \phi, \mu \rangle \in S \text{ and } \langle \kappa, \phi' \rangle \in \text{codom}(\mu)\}) \rangle.$$

**Proposition 6.** The optimal abstract counterparts of the powerset extension of the operations of Fig. 3 and of  $\cup$  are given in Fig. 6 for the case when the frame of their arguments is not  $\emptyset$ . Otherwise they yield  $\langle \emptyset, \mu \rangle$  for an arbitrary  $\mu \in \text{Memory}^{ps}$ , except  $\cup^{ps}$  which is such that  $\cup_{\tau}^{ps}(\langle \emptyset, \mu \rangle)(e) = \cup_{\tau}^{ps}(e)(\langle \emptyset, \mu \rangle) = e$ .

In Figure 6, the frame component of the abstract states behaves like in Section 5, except in `get.field`. The memory component does not change if the memory of the concrete operations (Figure 3) does not change. The operations `get.field` and `put.field` allow a flow of information between abstract frames and abstract memories. Namely, `get.field` loads in the abstract frame the set of classes  $\mu(f)$  of the objects stored in the field  $f$ . Conversely, `put.field` adds the classes contained in  $res$ , i.e.,  $\phi_2(res)$ , to the set of classes already stored in the field, i.e.,  $\mu_2(f)$ . In this way, we accumulate all classes stored in the field during the execution.

*Example 5.* We execute over  $\text{State}_\tau^{ps}$  the same sequence of instructions of Ex. 2. The initial state is the abstraction (Prop. 5) of the initial concrete state.

Operation	State
	$\langle [v_1 \mapsto \{\kappa_a\}], [\mathbf{n} \mapsto \emptyset] \rangle$ (initial state)
$(\text{expand}_{[v_1 \mapsto \kappa_a]}^{v_2: \kappa_a})^{ps}$	$\langle [v_1 \mapsto \{\kappa_a\}, v_2 \mapsto \emptyset], [\mathbf{n} \mapsto \emptyset] \rangle$
$(\text{new}_{[v_1, v_2 \mapsto \kappa_a]}^{\kappa_b})^{ps}$	$\langle [v_1 \mapsto \{\kappa_a\}, v_2 \mapsto \emptyset, res \mapsto \{\kappa_b\}], [\mathbf{n} \mapsto \emptyset] \rangle$
$(\text{get.field}_{[v_1, v_2 \mapsto \kappa_a, res \mapsto \kappa_b]}^{\mathbf{n}})^{ps}$	$\langle [v_1 \mapsto \{\kappa_a\}, v_2 \mapsto \emptyset, res \mapsto \emptyset], [\mathbf{n} \mapsto \emptyset] \rangle$
$(\text{put.var}_{[v_1, v_2, res \mapsto \kappa_a]}^{v_2})^{ps}$	$\langle [v_1 \mapsto \{\kappa_a\}, v_2 \mapsto \emptyset], [\mathbf{n} \mapsto \emptyset] \rangle$

The final state says that  $v_1$  can be bound to objects of class  $a$ , and that  $v_2$  and the field  $\mathbf{n}$  are bound to *nil*, the exact approximation of the final state of Ex. 2.

## 7 Comparison

We compare here the precision of the domains defined in the previous sections.

**Proposition 7.** For  $\tau \in \text{Typing}$ , we have  $\gamma^{df}(\text{State}_\tau^{df}) \subseteq \gamma^{ps}(\text{State}_\tau^{ps})$  and the *ps-analysis* computes a more precise information than the *df-analysis*.

*Proof.* Since  $\gamma^{df}(\emptyset) = \emptyset = \gamma^{ps}(\langle \emptyset, \mu \rangle) \in \gamma^{ps}(\text{State}_\tau^{ps})$  for  $\mu \in \text{Memory}^{ps}$ , it is enough to show that  $\gamma^{df}(\text{Frame}_\tau^{df}) \subseteq \gamma^{ps}(\text{State}_\tau^{ps})$ . This will entail the desired result, since the abstract operations (Figs. 5 and 6) are optimal (§8 of [9]). The idea is that every  $\phi \in \text{Frame}_\tau^{df}$  is the frame of a pair  $\langle \phi, \mu \rangle \in \text{State}_\tau^{ps}$  where  $\mu$  does not introduce any restriction. Let  $\mu(v) = \{int\}$  if  $\bar{\tau}(v) = int$ ,  $\mu(v) = \{bool\}$  if  $\bar{\tau}(v) = bool$  and  $\mu(v) = \downarrow(\bar{\tau}(v))$  if  $\bar{\tau}(v) \in \text{Class.ref}$  for  $v \in \text{dom}(\bar{\tau})$ . Then  $\langle \phi', \mu' \rangle \in \gamma^{df}(\phi)$  iff  $\phi \approx_{\mu'} \phi'$  iff  $(\phi \approx_{\mu'} \phi'$  and for every  $\langle \kappa, \phi'' \rangle \in \text{codom}(\mu')$  we have  $\mu \approx_{\mu'} \phi''$ ), iff  $\langle \phi', \mu' \rangle \in \gamma^{ps}(\langle \phi, \mu \rangle)$ , i.e.,  $\gamma^{df}(\phi) = \gamma^{ps}(\langle \phi, \mu \rangle)$ .

The *rta*- and the *df*-analysis are incomparable. In the final states of Exs. 3 and 4 the *df*- is more precise than the *rta*-analysis. But consider the initial states of those examples, i.e., the respective abstractions of the initial concrete state of Ex. 2. In Ex. 3 the possible classes for the field  $\mathbf{n}$  are  $\{\kappa_a\}$ , while in Ex. 4 they are  $\{\kappa_a, \kappa_b\}$ , because that abstraction does not provide any information about the fields. Hence the *df*-analysis provides a coarser approximation of the set of classes for the field  $\mathbf{n}$  than the *rta*-analysis.

We strongly believe that the *ps*-analysis is more precise than the *rta*-analysis. Indeed, the *ps*-analysis distinguishes between the classes stored in different variables, while the *rta*-analysis merges all classes in just one set. For instance, in the context of Ex. 1, the concrete states  $\sigma_1 = \langle [v_1 \mapsto l, v_2 \mapsto \text{nil}], [l \mapsto \langle \kappa_a, [\mathbf{n} \mapsto \text{nil}]] \rangle$  and  $\sigma_2 = \langle [v_1 \mapsto \text{nil}, v_2 \mapsto l], [l \mapsto \langle \kappa_a, [\mathbf{n} \mapsto \text{nil}]] \rangle$  are such that  $\alpha^{rta}(\{\sigma_1\}) = \alpha^{rta}(\{\sigma_2\}) = \{\kappa_a\}$  while  $\alpha^{ps}(\{\sigma_1\}) = \langle [v_1 \mapsto \{\kappa_a\}, v_2 \mapsto \emptyset], [\mathbf{n} \mapsto \emptyset] \rangle$  and  $\alpha^{ps}(\{\sigma_2\}) = \langle [v_1 \mapsto \emptyset, v_2 \mapsto \{\kappa_a\}], [\mathbf{n} \mapsto \emptyset] \rangle$ . In this case, the *ps*-analysis distinguishes  $\sigma_1$  and  $\sigma_2$  which do bind the variables to objects of different classes, while the *rta*-analysis does not. When trying to generalise this result into a formal proof, we are faced with the problem that the *rta*-analysis considers all objects in memory, while the *ps*-analysis only those reachable from the current frame or that of another object (Def. 11). For instance, in the context of Ex. 1, the concrete states  $\sigma_3 = \langle [v \mapsto \text{nil}], [l \mapsto \langle \kappa_a, [\mathbf{n} \mapsto \text{nil}]] \rangle$  and  $\sigma_4 = \langle [v \mapsto \text{nil}], [] \rangle$  are such that  $\alpha^{ps}(\{\sigma_3\}) = \alpha^{ps}(\{\sigma_4\}) = \langle [v \mapsto \emptyset], [\mathbf{n} \mapsto \emptyset] \rangle$ , while  $\alpha^{rta}(\{\sigma_3\}) = \{\kappa_a\}$  and  $\alpha^{rta}(\{\sigma_4\}) = \emptyset$ . However, this ability to distinguish such states, which contain *the same* class information, is of no use for the class analysis. Indeed, objects in memory can affect an expression of the program only if they are reachable from some variable. A formal proof of this statement can be obtained, e.g., by quotienting [5] both domains w.r.t. class information.

## 8 Conclusions

Class analyses for object-oriented languages are many and varied. This paper shows that abstract interpretation leads to a formal framework for the development and comparison of such analyses. To demonstrate this, we have put three traditional techniques for class analysis inside that framework. This provides a systematic construction of the abstract operations of a particular analysis. Furthermore, it allows a formal comparison of the relative precision of the analyses. Due to space limitations, it has not been possible to show the details of how the set of operations can be used to give semantics to *e.g.*, Java or Java bytecode.

Several extensions to the present work are to be considered: There are more class analyses than the ones considered in this paper and the classification that we have initiated here should be carried further. Analyses can be combined so as to use the combination of two abstract domains in an analysis. Formally, this is characterised by the *reduced product* of two abstractions that defines a semi-lattice structure on the set of abstractions. An interesting problem is to construct a concrete representation of the reduced product of, say, the *rta*- and the *df*-analysis.

Since we have algorithms for the abstraction of a finite set of states (Props. 1, 3 and 5), and for the abstract operations (Props. 2, 4 and 6) and since our domains are finite (for a given  $\tau \in \textit{Typing}$ ), the framework provides a way of implementing the static class analyses. The practical aspects of such an implementation remain to be investigated.

## References

1. O. Agesen. Constraint-Based Type Inference and Parametric Polymorphism. In B. Le Charlier, editor, *Proc. of the 1st Int. Static Analysis Symp.*, volume 864 of *Lecture Notes in Computer Science*, pages 78–100. Springer-Verlag, 1994.
2. D. F. Bacon and P. F. Sweeney. Fast Static Analysis of C++ Virtual Function Calls. In *Proc. of OOPSLA '96*, volume 31(10) of *ACM SIGPLAN Notices*, pages 324–341, New York, 1996. ACM Press.
3. P. Bertelsen. Semantics of Java Byte Code. Technical report, Department of Information Technology, Technical University of Denmark, March 1997.
4. E. Börger and W. Schulte. Defining the Java Virtual Machine as Platform for Provably Correct Java Compilation. In L. Brim, J. Grunski, and J. Zlatusla, editors, *23rd Int. Symp. on Mathematical Foundations of Computer Science*. Springer LNCS vol. 1450, 1998.
5. A. Cortesi, G. Filé, and W. Winsborough. The Quotient of an Abstract Interpretation. *Theoretical Computer Science*, 202(1-2):163–192, 1998.
6. P. Cousot. Constructive Design of a Hierarchy of Semantics of a Transition System by Abstract Interpretation. In S. Brookes and M. Mislove, editors, *13th Conf. on Math. Found. of Programming Semantics*, volume 6 of *Electronic Notes on Theoretical Computer Science*, Pittsburgh, PA, USA, March 1997. Elsevier Science Publishers. Available at <http://www.elsevier.nl/locate/entcs/volume6.html>.
7. P. Cousot. Types as Abstract Interpretations. In *24th ACM Symposium on Principles of Programming Languages (POPL '97)*, pages 316–331. ACM Press, 1997.
8. P. Cousot and R. Cousot. Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints. In *4th ACM Symp. on Principles of Programming Languages*, pages 238–252, 1977.
9. P. Cousot and R. Cousot. Systematic Design of Program Analysis Frameworks. In *6th ACM Symp. on Principles of Programming Languages*, pages 269–282, 1979.
10. A. Diwan, J. E. B. Moss, and K. S. McKinley. Simple and Effective Analysis of Statically Typed Object-Oriented Programs. In *Proc. of OOPSLA '96*, volume 31(10) of *ACM SIGPLAN Notices*, pages 292–305, New York, 1996. ACM Press.
11. J. Palsberg and M. I. Schwartzbach. Object-Oriented Type Inference. In *Proc. of OOPSLA '91*, volume 26(11) of *ACM SIGPLAN Notices*, pages 146–161. ACM Press, November 1991.
12. J. Plevyak and A. A. Chien. Precise Concrete Type Inference for Object-Oriented Languages. In *Proc. of OOPSLA '94*, volume 29(10) of *ACM SIGPLAN Notices*, pages 324–340. ACM Press, October 1994.
13. O. Shivers. Control-Flow Analysis in Scheme. In *Proc. of the 1988 Conf. on Programming Languages Design and Implementation*, volume 23(7) of *ACM SIGPLAN Notices*, pages 164–174. ACM Press, July 1988.