



An abstract domain for objects in dynamic programming languages

Vincenzo Arceri, Michele Pasqua and Isabella Mastroeni

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

September 12, 2019

An abstract domain for objects in dynamic programming languages

Vincenzo Arceri, Michele Pasqua, and Isabella Mastroeni

University of Verona, Department of Computer Science, Italy
{vincenzo.arceri | michele.pasqua | isabella.mastroeni}@univr.it

Abstract. Dynamic languages, such as JavaScript, PHP, Python or Ruby, provide a memory model for objects data structures allowing programmers to dynamically create, manipulate, and delete objects' properties. Moreover, in dynamic languages it is possible to access and update properties by using strings: this represents a hard challenge for static analysis. In this paper, we exploit the finite state automata abstract domain, approximating strings, in order to define a novel abstract domain for objects. We design an abstract interpreter useful to analyze objects in a toy language, inspired by real-world dynamic programming languages. We then show, by means of minimal yet expressive examples, the precision of the proposed abstract domain.

1 Introduction

In the last years, dynamic languages such as JavaScript or PHP have gained a huge success in a very wide range of applications. This mainly happened due to the several features that such languages provide to developers, making the writing of programs easier and faster. One of this features is the way strings may be used to interact with programs' objects. Indeed, it is popular, especially in dynamic languages, to create, manipulate, and delete objects' properties at run-time, interacting with them using strings. If, on the one hand, this may help developers to simplify coding and to build applications faster, on the other hand, this may lead to misunderstandings and bugs in the produced code. Furthermore, because of these dynamic features, reasoning about dynamic programs by means of static analysis is quite hard, producing very often imprecise results.

For instance, let us consider the simple yet expressive example reported in Fig. 1, supposing that the value of `par` is statically unknown. The value of `idx` is indeterminate after line 2 and it is updated at each iteration of the while loop (line 6). The loop guard is also statically unknown and at each iteration we access `obj` with `idx`, incrementally saving the results in `n`. The goal is to statically retrieve the value of `idx` and `n` at the end of the program. It is worth noting that a crucial role here is played by the string abstraction used to approximate the value of `idx`, that is used to access `obj`. Indeed, adopting finite abstract domains, such as [13–15], will lead to infer that `idx` could be any possible string. Consequently, when `idx` is used to access `obj`, in order to guarantee soundness,

```

1  if (par == ?) { idx = "a"; }
2  else { idx = "b"; };
3  n = 0; obj = new {a:1, aa:2, ab:3, ac:"world"};
4  while (?) {
5      n = n + obj[idx];
6      idx = concat(idx, "a");
7  }
8  obj[idx] = n; // value of idx and n ?

```

Fig. 1

we need to access all properties of `obj`. For instance, we also have to consider the property `ac`, which is never used to access `obj` during the execution of the program. This ends up in an imprecise approximation of `idx` and, in turn, of `n`.

In this paper, we employ a more precise abstraction for string values. In particular, we abstract strings with the finite state automata abstract domain [2], able to derive precise results also when strings are modified in iterative constructs. Then we define a novel abstract domain for objects, exploiting finite state automata. The idea is to abstract the objects' properties in the same domain used to abstract string values, namely the finite state automata abstract domain. We show that exploiting finite automata to abstract string values and objects properties produces precise result in abstract computations, in particular in objects' properties lookup and in objects' manipulation inside iterative constructs. We will formally present the objects abstract domain in Sec. 3.1.

Moreover, we use strings and objects abstract domains together with integers and booleans abstractions, presenting an abstract interpreter built upon the combination of these domains for a toy language, expressive enough to handle string operations, object declarations, objects' properties lookup and assignments.

2 Background

Notation. Given a finite set of symbols Σ , we denote by Σ^* the Kleene-closure of Σ , i.e., the set of all finite sequences of symbols in Σ . We denote an element of Σ^* , called *string*, by $s \in \Sigma^*$. If $s = s_0s_1 \dots s_n$, the length of s is $|s| = n + 1$ and the element in the i -th position is s_i . Given two strings s and s' , ss' is their concatenation. We use the following notations: $\Sigma^i \triangleq \{s \in \Sigma^* \mid |s| = i\}$ and $\Sigma^{<i} \triangleq \bigcup_{0 \leq j < i} \Sigma^j$, for $i \in \mathbb{N}$. We follow [12] for automata notation. A finite state automaton is a tuple $\mathbf{A} = \langle Q, q_0, \Sigma, \delta, F \rangle$ where Q is a finite set of states, $q_0 \in Q$ is the initial state, Σ is the (finite) alphabet, $\delta \subseteq Q \times \Sigma \times Q$ is the transition relation and $F \subseteq Q$ is a set of final states. In particular, if $\delta \in Q \times \Sigma \rightarrow Q$ is a function, then \mathbf{A} is called deterministic finite state automata (DFA). The class of languages recognized by finite state automata is the class of regular languages. Given an automaton \mathbf{A} , we denote the language accepted by \mathbf{A} as $\mathcal{L}(\mathbf{A})$. A language \mathcal{L} is regular iff there exists a finite state automaton \mathbf{A} such that $\mathcal{L} = \mathcal{L}(\mathbf{A})$. From the Myhill-Nerode theorem [9], we have that for each regular language there exists a unique minimum automaton, i.e., with the minimum number of states, recognizing the language. Given a regular language \mathcal{L} , we denote by $\text{Min}(\mathcal{L})$ the minimum DFA \mathbf{A} such that $\mathcal{L} = \mathcal{L}(\mathbf{A})$. For space

limitations, in the following we will refer to finite state automata by using the corresponding regular expressions, which are isomorphic to regular languages and, in turn, to finite state automata. Given two regular expressions \mathbf{r}_1 and \mathbf{r}_2 , we denote by $\mathbf{r}_1 \parallel \mathbf{r}_2$ the disjunction between \mathbf{r}_1 and \mathbf{r}_2 , by $(\mathbf{r}_1)^*$ the Kleene-closure of \mathbf{r}_1 , and by $(\mathbf{r}_1)^+$ the Kleene-closure of \mathbf{r}_1 with at least one repetition.

Given a partial function $f \in X \rightarrow Y$, we can define an equivalent total function $g \in X \rightarrow Y_\uparrow$, where $Y_\uparrow \triangleq Y \cup \{\uparrow\}$ and $\uparrow \notin Y$ denotes indefiniteness. The function g is defined as: $g(x) \triangleq f(x)$ when $f(x)$ is defined, and $g(x) \triangleq \uparrow$ otherwise. When we describe extensionally a function we omit the elements mapped to \uparrow , namely $g \in X \rightarrow Y_\uparrow$, described as $[x_1 \mapsto y_1 \ x_2 \mapsto y_2 \ \dots \ x_n \mapsto y_n]$, is such that $g(x_i) = y_i$ for every $i \in \{1, 2, \dots, n\}$ and $g(x_i) = \uparrow$ otherwise.

Abstract interpretation. The (concrete) semantic of a program is a representation of all its possible executions by means of a set of mathematical objects. This set is, in general, not computable. It is well known, due to Rice's theorem, that all non trivial properties of the concrete semantics of a program are undecidable. Abstract interpretation is born as a theory for soundly approximating the semantics of discrete dynamic systems. The approximation consists in the observation of the semantics at a specified level of abstraction, focusing only on some important aspects of computations. In this setting, abstract interpretation allows us to compute an abstract semantics of the program, depending on the properties of interest. The approximation is correct by design, in the sense that what holds in the abstract holds also in the concrete (no false negatives).

A theory of domains for abstract interpretation was defined in [7], based on the notion of *Galois insertion*. A Galois insertion (C, α, γ, A) consists of two partially ordered sets $\langle C, \leq_C \rangle$, $\langle A, \leq_A \rangle$ and two monotone functions $\alpha \in C \rightarrow A$, $\gamma \in A \rightarrow C$ such that for all c in C and a in A it holds: $\alpha(c) \leq_A a \Leftrightarrow c \leq_C \gamma(a)$ and $\alpha \circ \gamma = \text{id}$ (the identity function $\lambda x. x$). C is the concrete domain, A is the abstract domain, α is the abstraction function and γ is the concretization function. Sometimes, abstract interpretations are given by means of Galois connections (instead of Galois insertions), relaxing the constraints $\alpha \circ \gamma = \text{id}$. Let $f \in C \rightarrow C$ be a function on the concrete domain and $f^\# \in A \rightarrow A$ be a function on the abstract domain. $f^\#$ is a sound (or correct) approximation of f if $f \circ \gamma \leq_C \gamma \circ f^\#$ or, equivalently, if $\alpha \circ f \leq_A f^\# \circ \alpha$ [7].

Nevertheless, Galois insertions/connections represent the optimal case: sometimes we have to settle for weaker forms of abstract interpretation, as in the case of the Polyhedra abstract domain [8], where we have only the concretization function γ . In this setting, the soundness is expressed as: $f \circ \gamma \leq_C \gamma \circ f^\#$.

Finite state automata abstract domain. We report here the finite state automata abstract domain presented in [2], that over-approximates strings as regular languages, represented by the minimum deterministic finite state automaton [9] recognizing them. The domain is $(\text{DFA}/\equiv, \sqsubseteq_{\text{DFA}}, \sqcup_{\text{DFA}}, \sqcap_{\text{DFA}}, \text{Min}(\emptyset), \text{Min}(\Sigma^*))$, where DFA/\equiv is the quotient set of DFA w.r.t. the equivalence relation induced by language equality, \sqsubseteq_{DFA} is the partial order induced by language inclusion, \sqcup_{DFA} and \sqcap_{DFA} are the least upper bound and the greatest lower bound, respectively. The minimum is $\text{Min}(\emptyset)$, corresponding to the automaton recognizing the empty

$ \begin{aligned} a \in \text{AE} &::= x \mid n \mid a + a \mid a - a \mid a * a \mid a / a \\ &\mid \text{length}(s) \mid \text{indexOf}(s_1, s_2) \\ b \in \text{BE} &::= x \mid \text{true} \mid \text{false} \mid b \ \&\& \ b \mid b \ \ \ \ b \mid ! \ b \mid a < a \\ &\mid a == a \mid s == s \\ s \in \text{SE} &::= x \mid "s" \mid \text{substr}(s, a_1, a_2) \mid \text{charAt}(s, a) \mid \text{concat}(s_1, s_2) \\ o \in \text{OE} &::= \{ \} \mid \{ s_0 : e_0, s_1 : e_1, \dots, s_n : e_n \} \\ e \in \text{E} &::= a \mid b \mid s \mid x[s] \\ \text{st} \in \text{STMT} &::= \text{st} ; \text{st} \mid \text{ski} \mid x = e \mid x = \text{new } o \mid x[s] = e \\ &\mid \text{if } b \{ \text{st} \} \text{ else } \{ \text{st} \} \mid \text{while } b \{ \text{st} \} \\ \text{where } x \in \text{ID} & \text{ (identifiers), } n \in \mathbb{Z} \text{ and } s, s_0, s_1, \dots, s_n \in \Sigma^* \end{aligned} $
--

Fig. 2: μJS syntax.

language and the maximum is $\text{Min}(\Sigma^*)$, corresponding to the automaton recognizing any possible string over Σ . We abuse notation by representing equivalence classes in the domain $\text{DFA}_{/\equiv}$ by one of its automaton (usually the minimum), i.e., when we write $A \in \text{DFA}_{/\equiv}$ we mean $[A]_{\equiv}$. Since the domain $\text{DFA}_{/\equiv}$ is infinite, and it is not ACC, i.e., it contains infinite ascending chains, it is equipped with the parametric widening ∇_{DFA}^n . The latter is defined in terms of a state equivalence relation merging states that recognize the same language, up to a fixed length $n \in \mathbb{N}$, a parameter used for tuning the widening precision [4, 10]. For instance, let us consider the automata $A, A' \in \text{DFA}_{/\equiv}$ recognizing the languages $\mathcal{L} = \{\epsilon, a\}$ and $\mathcal{L}' = \{\epsilon, a, aa\}$, respectively. The result of the application of the widening ∇_{DFA}^n , with $n = 1$, is $A \nabla_{\text{DFA}}^n A' = A''$ such that $\mathcal{L}(A'') = \{a^n \mid n \in \mathbb{N}\}$.

μJS language. In this paper, we adopt as core language μJS [2], whose syntax is reported in Fig. 2. This simple toy language is able to express arithmetic (AE), boolean (BE) and string expressions (SE). There is not implicit type conversion, since the problem of analyzing programs with implicit conversions had been already addressed in [1, 2]. Anyway, it is straightforward to merge our analysis with the ones proposed in [1, 2]. In addition, we have augmented μJS with objects (OE), where an object can be empty, denoted $\{ \}$, or a set of comma-separated property-expression associations, denoted $\{a:1, b:2, c:3\}$.

Concerning the language's semantics, the execution of a μJS program relies on the notion of state, which is composed by environments and heaps, namely states $\sigma \in \text{STATE}$ are pairs $\langle \xi, \rho \rangle \in \text{ENV} \times \text{HEAP}$. An environment is a map from identifiers to values, namely $\text{ENV} \triangleq \text{ID} \rightarrow \text{VAL}$, while a heap is a map from addresses to objects, namely $\text{HEAP} \triangleq \text{ADDR} \rightarrow \text{OBJ}$. Values v have domain $\text{VAL} \triangleq \text{INT} \cup \text{BOOL} \cup \text{STR} \cup \text{ADDR} \cup \{\uparrow\}$, where $\text{INT} \triangleq \mathbb{Z}$, $\text{BOOL} \triangleq \{\text{true}, \text{false}\}$, $\text{STR} \triangleq \Sigma^*$, $\text{ADDR} \triangleq \{\underline{n} \mid n \in \mathbb{N}\}$ and \uparrow denotes indefiniteness. An object $o \in \text{OBJ}$ is represented as a map that associates strings to values, namely $\text{OBJ} \triangleq \text{STR} \rightarrow \text{VAL}$. It is worth noting that there is no order relation between objects' properties, as it happens in standard programming languages. Environments update is defined as usual: $\xi[x \leftarrow v](y) \triangleq v$ when $x = y$, and $\xi[x \leftarrow v](y) \triangleq \xi(y)$ otherwise. The update for heaps and objects is analogous. The big-step semantics

of a μJS program (i.e., a statement) is standard, following [1, 2], and it is captured by the function $\llbracket \text{st} \rrbracket \in \text{STATE} \rightarrow \text{STATE}$. After showing the concrete semantics of object-related expressions, we will focus on the semantics of assignments, that slightly changes w.r.t. the standard one. As far as expression semantics is concerned, it is also standard [2]. We abuse notation denoting the semantics of an expression as $\llbracket e \rrbracket \in \text{STATE} \rightarrow \text{VAL}$. The evaluation of an object takes each association string-expression and it recursively evaluates the expressions. The result is a map containing the string-value associations.

$$\llbracket \{s_0 : e_0, s_1 : e_1, \dots, s_n : e_n\} \rrbracket \sigma \triangleq [s_n \mapsto \llbracket e_n \rrbracket \sigma] \bullet \dots \bullet [s_1 \mapsto \llbracket e_1 \rrbracket \sigma] \bullet [s_0 \mapsto \llbracket e_0 \rrbracket \sigma]$$

$$\text{where } f \bullet g(s) \triangleq g(s) \text{ if } g(s) \neq \uparrow \wedge f(s) = \uparrow \text{ and } f \bullet g(s) \triangleq f(s) \text{ otherwise}$$

For example, the expression $\{\mathbf{a}:1, \mathbf{b}:\text{length}(\text{"foo"}), \mathbf{c}:5+3\}$ evaluates to the object $[\mathbf{a} \mapsto 1, \mathbf{b} \mapsto 3, \mathbf{c} \mapsto 8]$. Following the JavaScript semantics, it is worth noting that, for instance, $\{\mathbf{a}:1, \mathbf{a}:2\}$ evaluates to $[\mathbf{a} \mapsto 2]$, saving only the last association with the same property \mathbf{a} . The semantics of objects' properties lookup checks whether the input object contains an identifier-value association where the identifier corresponds to the input string. Hence, its definition is the following, supposing that $\llbracket s \rrbracket \langle \xi, \rho \rangle = s \in \text{STR}$:

$$\llbracket x[s] \rrbracket \langle \xi, \rho \rangle \triangleq \rho(\llbracket x \rrbracket \langle \xi, \rho \rangle)(s) \text{ if } \llbracket x \rrbracket \langle \xi, \rho \rangle \in \text{ADDR} \text{ and } \llbracket x[s] \rrbracket \langle \xi, \rho \rangle \triangleq \uparrow \text{ otherwise}$$

In our core language, we allow only to access already stored objects (condition $\llbracket o \rrbracket \sigma \in \text{ADDR}$). Moreover, it is worth noting that when we try to access a property s not present in o , then $\rho(\llbracket o \rrbracket \langle \xi, \rho \rangle)$ returns \uparrow .

The semantics of generic statements is standard, here we explain only the semantics for assignments, which is also used for objects allocation and update. We have three cases: $x = e$, where e evaluates to a value; $x = \text{new } o$, where o evaluates to an object; $x[s] = e$, where s evaluates to a string and e evaluates to a value. In the first case, we only update the environment, following the typical concrete semantics of assignments. In the second case, we need to allocate the object into a new address which x will point to. Then, both environment and heap are properly updated. In the third case, we update the object pointed by x in the heap. Formally, let $\underline{n} \in \text{ADDR}$ be a fresh, i.e., not-used, address:

$$\begin{aligned} \llbracket x = e \rrbracket \langle \xi, \rho \rangle &\triangleq \langle \xi[x \leftarrow \llbracket e \rrbracket \langle \xi, \rho \rangle], \rho \rangle & \llbracket x = \text{new } o \rrbracket \langle \xi, \rho \rangle &\triangleq \langle \xi[x \leftarrow \underline{n}], \rho[\underline{n} \leftarrow \llbracket o \rrbracket \langle \xi, \rho \rangle] \rangle \\ \llbracket x[s] = e \rrbracket \langle \xi, \rho \rangle &\triangleq \langle \xi, \rho[\xi(x) \leftarrow \rho(\xi(x))[\llbracket s \rrbracket \langle \xi, \rho \rangle \leftarrow \llbracket e \rrbracket \langle \xi, \rho \rangle]] \rangle \end{aligned}$$

As a final remark, we point out that in our extension of μJS we do not model features such as pointer arithmetic, objects comparisons and implicit type conversion (e.g., $\mathbf{x} = 1$; $\mathbf{y} = \text{true}$; $\mathbf{z} = \mathbf{x} == \mathbf{y}$ leads to an error).

3 Static Analysis of μJS

In order to reason about a μJS program we need to take into account all its possible executions, by means of the so called collecting semantics. Our concrete collecting semantics is a classic post-condition semantics, computing states invariants at every statement. It is defined as the direct-image lift of the big-step semantics of μJS , hence it is a function from sets of states to sets of states.

We denote by $\llbracket \text{st} \rrbracket \in \wp(\text{STATE}) \rightarrow \wp(\text{STATE})$ the concrete collecting semantics. For instance, the collecting semantics for assignments involving expressions, is defined as $\llbracket x = e \rrbracket X \triangleq \{\llbracket x = e \rrbracket \sigma \mid \sigma \in X\}$. The semantics is similarly defined for the other constructs and for assignments involving objects. In particular, the collecting semantics for conditionals and loops is defined, as usual, as:

$$\begin{aligned} \llbracket \text{if } b \{ \text{st}_1 \} \text{ else } \{ \text{st}_2 \} \rrbracket X &\triangleq \llbracket \text{st}_1 \rrbracket \text{filter}_b(X) \cup \llbracket \text{st}_2 \rrbracket \text{filter}_{!b}(X) \\ \llbracket \text{while } b \{ \text{st} \} \rrbracket X &\triangleq \text{filter}_{!b}(\text{lfp } \lambda T. X \cup \llbracket \text{st} \rrbracket \text{filter}_b(T)) \end{aligned}$$

Here $\text{filter}_b \in \wp(\text{STATE}) \rightarrow \wp(\text{STATE})$ is a filtering function, namely it filters out the states that do not fulfill the boolean condition b . Unfortunately, we are not able to compute the concrete collecting semantics since it is an infinite mathematical object. Hence, in order to perform static analysis, we approximate the collecting semantics, following the abstract interpretation framework. In order to make the computation, and in turn the analysis, feasible we need an abstract semantics $\llbracket \text{st} \rrbracket^\#$ computer-representable and ensuring termination of the analysis. Ideally, the abstract semantics computes on abstract states in $\text{STATE}^\#$, approximations of the concrete ones. Precisely, $\text{STATE}^\#$ is an approximation of $\wp(\text{STATE})$, with a concretization $\gamma \in \text{STATE}^\# \rightarrow \wp(\text{STATE})$. The abstract semantics must be sound, meaning that what we prove in the abstract also holds for the concrete semantics. Put it in abstract interpretation terms, this means that for every $\sigma^\# \in \text{STATE}^\#$ we have that $\llbracket \text{st} \rrbracket \gamma(\sigma^\#) \subseteq \gamma(\llbracket \text{st} \rrbracket^\# \sigma^\#)$. Before defining the abstract semantics, we focus on the *objects abstract domain*, which is the core of our paper and it is used to represent, possibly infinite, sets of concrete objects.

3.1 Abstract Objects

As previously introduced, in order to make the analysis feasible, we need to finitely represent an infinite set of states. We start here with our representation of infinite sets of objects, namely we define an abstract domain approximating $\wp(\text{OBJ})$. First, we have a non-relational abstraction between objects-properties and values, i.e., we abstract $\wp(\text{OBJ})$ in $\wp(\text{STR}) \rightarrow \wp(\text{VAL})$.

Then we abstract $\wp(\text{STR})$ with the automata domain, while for $\wp(\text{VAL})$ we abstract separately each type of values in its abstract domain, obtaining the product domain $\text{VAL}^\# \triangleq \text{INT}^\# \times \text{BOOL}^\# \times \text{STR}^\# \times \wp(\text{ADDR}^\#) \times \{\text{def}, ?\}$. For numeric values we can use any non-relational domain, such as integer intervals. $\text{BOOL}^\# \triangleq \{\perp, \text{tt}, \text{ff}, \top\}$ is isomorphic to $\wp(\text{BOOL})$ and for sets of strings we use the automata domain, namely $\text{STR}^\# \triangleq \text{DFA}_{/\equiv}$. As we will see in the next subsection, we approximate heaps with an allocation-site abstraction of ADDR . So, possibly infinite sets of addresses are abstracted into finite sets of allocation sites, namely $\text{ADDR}^\# \triangleq \text{LINES}$, where LINES is the finite set of lines of code of a given program. Here we abstract $\wp(\text{ADDR})$ in $\wp(\text{ADDR}^\#)$, since an abstract object could have more than one allocation site. The domain $\{\text{def}, ?\}$ is isomorphic to $\wp(\{\uparrow\})$ and def represents the absence of indefiniteness while $?$ represents potential indefiniteness. An abstract value $v^\# = \langle i^\#, b^\#, s^\#, A, u^\# \rangle \in \text{VAL}^\#$ represents the union of the elements taken from every single-type abstraction:

$$\gamma_V(v^\#) = \gamma_I(i^\#) \cup \gamma_B(b^\#) \cup \gamma_S(s^\#) \cup \bigcup_{l \in A} \gamma_A(l) \cup \gamma_U(u^\#)$$

```

1  o = new {x:1, y:2, z:3};
2  idx = "x";
3  while (?) {
4      if (?) { idx = concat(idx, "x") }
5      else { idx = concat(idx, "y") }
6  };
7  o[idx] = 7;

```

Fig. 3: μJS program example.

where γ_I is the concretization defined in the numerical non-relational domain, $\gamma_B(\perp) \triangleq \emptyset$, $\gamma_B(\text{tt}) \triangleq \{\text{true}\}$, $\gamma_B(\text{ff}) \triangleq \{\text{false}\}$, $\gamma_B(\top) \triangleq \{\text{true}, \text{false}\}$, γ_S is the concretization for the automata domain (i.e., the language recognized by the given automaton) and $\gamma_U(\text{def}) = \emptyset$, $\gamma_U(?) = \{\uparrow\}$. The concretization for addresses will be introduced in Sec. 3.2, when we deal with abstract heaps. Briefly, the concretization of a given allocation site is the set of all possible addresses that can be allocated at that line of code. The abstract join $\sqcup_V^\#$ and the partial order $\sqsubseteq_V^\#$ for $\text{VAL}^\#$ are defined pointwise.

The partial order $\sqsubseteq_O^\#$ for $\text{OBJ}^\#$ is the pointwise ordering between functions, i.e., $o_1^\# \sqsubseteq_O^\# o_2^\# \triangleq (\forall A \in \text{DFA}_{/\equiv} . o_1^\#(A) \sqsubseteq_V^\# o_2^\#(A))$. This order is not optimal but it does not harm the analysis since, as we can see in Sect. 3.1, the order can be strengthened. Analogously, the join for $\text{OBJ}^\#$ is defined as $\sqcup_O^\# X \triangleq \lambda A . \sqcup_V^\# \{o^\#(A) \mid o^\# \in X\}$. It is straightforward to see that $(\text{OBJ}^\#, \sqsubseteq_O^\#)$ is a complete lattice, with minimum mapping every automaton to the tuple composed by the minimum of each value-type domain, and maximum mapping every automaton to the tuple composed by the maximum of each value-type domain. The concretization $\gamma_O \in \text{OBJ}^\# \rightarrow \wp(\text{OBJ})$ is defined as:

$$\gamma_O(o^\#) \triangleq \{o \in \text{OBJ} \mid \forall s \in \text{STR} \exists A \in \text{DFA}_{/\equiv} . (s \in \gamma_S(A) \wedge o(s) \in \gamma_V(o^\#(A)))\}$$

In order to optimize the implementation of the abstract domain, we represent singleton sets of string as they are, instead of converting them into an automaton. Indeed, it is worth noting that we can partition the finite state automata abstract domain as $\text{DFA}_{/\equiv} = \text{DFA}_{/\equiv}^1 \cup \text{DFA}_{/\equiv}^\omega$, where $\text{DFA}_{/\equiv}^1 \triangleq \{A \in \text{DFA}_{/\equiv} \mid |\mathcal{L}(A)| = 1\}$, namely the set of finite state automata that recognize singleton languages, and $\text{DFA}_{/\equiv}^\omega \triangleq \text{DFA}_{/\equiv} \setminus \text{DFA}_{/\equiv}^1$, namely the set of finite state automata that recognizes languages of size 0 or size greater than 1 (possibly infinite). Clearly $\text{DFA}_{/\equiv}^1$ is isomorphic to STR , hence we can equivalently define abstract objects as maps in $\text{OBJ}^\# \triangleq (\text{STR} \cup \text{DFA}_{/\equiv}^\omega) \rightarrow \text{VAL}^\#$.

In order to show how our objects abstract domain works, we consider a simple yet expressive μJS example (Fig. 3, where we suppose that the boolean guards of the `while` and `if` statements are statically unknown). The fragment declares the object `o` at line 1, and its abstract value at lines 1-7 is reported in Fig. 4a. Then, it indefinitely iterates over the string variable `idx` at lines 3-6 appending either the strings "x" or "y". Finally, `idx` is used to access the object `o` at line 7. Let us suppose to statically analyze the above program with the abstract domain previously presented. Since the number of iterations of the `while`-loop is statically unknown, the string value of `idx`, abstracted as a finite state automaton, may diverge. In order to enforce termination, the automata widening ∇_{DFA}^n is applied. Tuning ∇_{DFA}^n with $n = 3$, the abstract value of `idx` at

line 7, after the `while` computation, corresponds to the automaton expressed by the regular expression $\mathbf{x}(\mathbf{x} \parallel \mathbf{y})^*$. Since `idx` does not represent just a single string, when we analyze $\mathbf{o}[\mathbf{idx}]$ we may have to overwrite an object property (e.g., \mathbf{x}) and add new properties to \mathbf{o} (e.g., \mathbf{xy}). Since the abstract value of `idx` expresses an infinite number of object properties, we call this property *summary property*. The abstract value after the execution of the statement at line 7 (Fig. 4b, where the summary property is added to the object reported in Fig. 4a). Note that in the abstract object updated after line 7, the abstract properties \mathbf{x} and $\mathbf{x}(\mathbf{x} \parallel \mathbf{y})^*$ share the common concrete property \mathbf{x} . In particular, the value of $\mathbf{o}["\mathbf{x}"]$ may be either 1 or 7. We aim at an objects' representation where every property does not share any property with the others, namely when objects are in normal form.

Normalization. We now formally define the notion of abstract object normal form. Given an abstract object $o^\# \in \text{OBJ}^\#$, we denote by $\text{props}(o^\#) \subseteq \text{STR}^\#$ the set of its abstract properties, namely the properties which are not undefined. We remind that $\text{STR}^\#$ is the optimized version of the automata domain, i.e., $\text{STR}^\# = \text{STR} \cup \text{DFA}_{\equiv}^\omega$. Formally, $\text{props}(o^\#) \triangleq \{p \in \text{STR}^\# \mid o^\#(p) = \langle i^\#, b^\#, s^\#, A, u^\# \rangle \wedge u^\# = \text{def}\}$. Abstract properties represent sets of concrete properties. Hence, given $p \in \text{props}(o^\#)$, we abuse notation denoting by $\mathcal{L}(p)$ the language of the concrete properties captured by p . $\mathcal{L}(p)$ is the language recognized by the corresponding automaton, when $p \in \text{DFA}_{\equiv}^\omega$ and it is the language $\{p\}$ when $p \in \text{STR}$.

Definition 1 (Abstract object normal form). *An abstract object $o^\# \in \text{OBJ}^\#$ is in normal form when:*

$$\forall p \in \text{props}(o^\#). |\mathcal{L}(p)| \in \{1, \omega\} \wedge \forall p_1, p_2 \in \text{props}(o^\#). \mathcal{L}(p_1) \cap \mathcal{L}(p_2) = \emptyset$$

Informally, we say that an abstract object is in normal form when each property p represents only a single string (i.e., $|\mathcal{L}(p)| = 1$) or an infinite language (i.e., $|\mathcal{L}(p)| = \omega$) and it does not share any concrete property with other abstract properties. Hence, a normal form abstract object has two kind of properties: p is a *non-summary property*, if $|\mathcal{L}(p)| = 1$, and p is a *summary property*, if $|\mathcal{L}(p)| = \omega$. For instance, the abstract object in Fig. 4a is in normal form, since any abstract property expresses concrete properties that are not expressed by other abstract properties and it only contains summary and non-summary properties. Instead, the abstract object in Fig. 4b is not in formal form, despite it has only summary and non-summary properties, since the string \mathbf{x} is expressed by the non-summary property \mathbf{x} and by the summary property $\mathbf{x}(\mathbf{x} \parallel \mathbf{y})^*$. During abstract

$$(a) \left[\begin{array}{l} \mathbf{x} \mapsto [1,1] \\ \mathbf{y} \mapsto [2,2] \\ \mathbf{z} \mapsto [3,3] \\ \hline - \end{array} \right] \quad (b) \left[\begin{array}{l} \mathbf{x} \mapsto [1,1] \\ \mathbf{y} \mapsto [2,2] \\ \mathbf{z} \mapsto [3,3] \\ \hline \mathbf{x}(\mathbf{x} \parallel \mathbf{y})^* \mapsto [7,7] \end{array} \right] \quad (c) \left[\begin{array}{l} \mathbf{x} \mapsto [1,7] \\ \mathbf{y} \mapsto [2,2] \\ \mathbf{z} \mapsto [3,3] \\ \hline \mathbf{x}(\mathbf{x} \parallel \mathbf{y})^+ \mapsto [7,7] \end{array} \right]$$

Fig. 4: (a) Abstract value of \mathbf{o} after line 1 of the fragment reported in Fig. 3 (b) Abstract value of \mathbf{o} after line 7. (c) Its normal form.

Algorithm 1: $\text{Norm} \in \text{OBJ}^\# \rightarrow \text{OBJ}^\#$ algorithm

Data: $o^\# \in \text{OBJ}^\#$
Result: $\text{Norm}(o^\#)$

- 1 **foreach** $p \in \text{props}(o^\#)$ **do**
- 2 $v^\# \leftarrow o^\#(p)$;
- 3 **if** $|\mathcal{L}(p)| \notin \{1, \omega\}$ **then**
- 4 $o^\# \leftarrow o^\#[p \leftarrow \langle \perp, \perp, \text{Min}(\emptyset), \emptyset, ? \rangle]$;
- 5 **foreach** $s \in \mathcal{L}(p)$ **do**
- 6 $o^\# \leftarrow o^\# \bullet^\# [s \mapsto v^\#]$;
- 7 **foreach** $p_1 \in \text{props}(o^\#)$ **do**
- 8 $v_1^\# \leftarrow o^\#(p_1)$; remove p_1 from $o^\#$; **normalized** \leftarrow **false**;
- 9 **foreach** $p_2 \in \text{props}(o^\#)$ **do**
- 10 $v_2^\# \leftarrow o^\#(p_2)$;
- 11 **if** $p_1 \sqcap_s^\# p_2 \neq \emptyset \wedge p_1 \neq p_2$ **then**
- 12 **normalized** \leftarrow **true**;
- 13 $o^\# \leftarrow o^\# \bullet^\# [p_1 \sqcap_s^\# p_2 \mapsto o^\#(p_1 \sqcap_s^\# p_2) \sqcup_V^\# v_1^\# \sqcup_V^\# v_2^\#]$;
- 14 $o^\# \leftarrow o^\# \bullet^\# [p_1 \searrow_s^\# p_2 \mapsto o^\#(p_1 \searrow_s^\# p_2) \sqcup_V^\# v_1^\#]$;
- 15 $o^\# \leftarrow o^\# \bullet^\# [p_2 \searrow_s^\# p_1 \mapsto o^\#(p_2 \searrow_s^\# p_1) \sqcup_V^\# v_2^\#]$;
- 16 remove p_2 from $o^\#$;
- 17 **if** **!normalized** **then** $o^\# \leftarrow o^\# \bullet^\# [p_1 \mapsto v_1^\#]$;
- 18 **return** $o^\#$;

computations, it may happen that abstract objects are not in normal form, so we need to normalize them. We rely on the function $\text{Norm} \in \text{OBJ}^\# \rightarrow \text{OBJ}^\#$ that normalizes an abstract object and its behaviour is captured by the algorithm reported by Alg. 1, where the $o_1^\# \bullet^\# o_2^\#$ is defined as, assuming $\langle i_1^\#, b_1^\#, s_1^\#, A_1, u_1^\# \rangle = o_1^\#(p)$, $\langle i_2^\#, b_2^\#, s_2^\#, A_2, u_2^\# \rangle = o_2^\#(p)$:

$$o_1^\# \bullet^\# o_2^\#(p) \triangleq o_2^\#(p) \text{ if } u_2^\# \neq ? \wedge u_1^\# = ? \text{ and } o_1^\# \bullet^\# o_2^\#(p) \triangleq o_1^\#(p) \text{ otherwise}$$

In the algorithm, the operators $\sqcap_s^\#$ and $\searrow_s^\#$ are the operators \sqcap_{DFA} and \searrow_{DFA} , respectively, of the automata domain adapted to its optimized versions $\text{STR} \cup \text{DFA}_{\neq}^\omega$. The first part of Alg. 1, namely lines 1-6, checks if any property of $o^\#$ is summary or non-summary. If it finds a property p such that $|\mathcal{L}(p)| \notin \{1, \omega\}$ then the algorithm first remove that property from the object, and then looks at its language (that is finite) and adds any single property captured by p with its old corresponding value. All the automata operations reported above and the check $|\mathcal{L}(p)| \notin \{1, \omega\}$ can be performed with linear complexity w.r.t. the number of state of the automata. For example, let consider the object $[x \mid y \mapsto [5, 5]]$, the algorithm returns as result the normal form abstract object $[x \mapsto [5, 5], y \mapsto [5, 5]]$. The idea of the second part of Alg. 1 (lines 7-17) is to check, for any $p_1 \in \text{props}(o^\#)$, if it shares at least a concrete property with any other $p_2 \in \text{props}(o^\#)$ (lines 11-16). The idea is check if p_1 is contained in p_2 and, if so, the previous abstract value of p_1 and the abstract value of p_2 are merged (with lub) and collected in p_1 . In particular, the old association with p_1 is removed (line 8) and

three new abstract properties are created in o . After that, p_1 and p_2 do not share any concrete property. In particular:

- the property $p_1 \sqcap_s^\# p_2$ points to the join of the previous values of p_1 and p_2 and the previous value (if present) of $p_1 \sqcap_s^\# p_2$ in $o^\#$ (line 13);
- the property $p_1 \searrow_s^\# p_2$ points to the previous value of p_1 and the previous value (if present) of $p_1 \searrow_s^\# p_2$ in $o^\#$ (line 14);
- the property $p_2 \searrow_s^\# p_1$ points to the previous value of p_2 and the previous value (if present) of $p_2 \searrow_s^\# p_1$ in $o^\#$ (line 15);

Then, also the property p_2 is removed from $o^\#$ (line 16). If the property p_1 does not share any property with other abstract properties of $o^\#$, the association $\langle p_1, o^\#(p_1) \rangle$ is simply added to o (line 17). For example, let us consider again the abstract object reported in Fig. 4b. The result obtained by applying Alg. 1 is the abstract object reported in Fig. 4c.

Proposition 1. *Given $o^\# \in \text{OBJ}^\#$, the abstract object $\text{Norm}(o^\#)$, computed by Alg. 1, is in normal form (Def. 1). Moreover, we have that $\gamma_O(o^\#) = \gamma_O(\text{Norm}(o^\#))$.*

As we have mentioned in Sect. 3, normalization strengthens the abstract order between objects. For example, the objects $[a \mapsto [1, 1], b \mapsto [1, 1]]$ and $[a \parallel b \mapsto [1, 2]]$ are not comparable, but, if we normalize the second object (i.e., in $[a \mapsto [1, 2], b \mapsto [1, 2]]$), then we have $[a \mapsto [1, 1], b \mapsto [1, 1]] \sqsubseteq_o^\# \text{Norm}([a \parallel b \mapsto [1, 2]])$.

3.2 Abstract Semantics

Abstract states in $\text{STATE}^\#$ are composed by abstract environments and abstract heaps, so we have an abstraction from $\wp(\text{ENV} \times \text{HEAP})$ to $\wp(\text{ENV}) \times \wp(\text{HEAP})$. As an abstract representation of the heap, we use a classic allocation-site abstraction of ADDR [16]. Possibly infinite sets of addresses are abstracted into finite sets of allocation sites, namely $\text{ADDR}^\# \triangleq \text{LINES}$, where LINES is the finite set of lines of code of a given program. Given a μJS program, we suppose to have a labeling assigning to each statement of the program a unique line of code (a natural number). Then, we define two functions, $\text{line} \in \text{STMT} \rightarrow \text{LINES}$ and $\text{code} \in \text{LINES} \rightarrow \text{STMT}$, returning the line of code of a given statement and the statement assigned to a given line of code, respectively. The concretization is

$$\gamma_A(l) \triangleq \left\{ \underline{n} \in \text{ADDR} \mid \exists \langle \xi, \rho \rangle \in \text{STATE}. \begin{array}{l} \llbracket \text{code}(l) \rrbracket \langle \xi, \rho \rangle = \langle \xi', \rho' \rangle \wedge \\ \rho(\underline{n}) = \lambda s. \uparrow \wedge \rho'(\underline{n}) \neq \lambda s. \uparrow \end{array} \right\}$$

meaning that the concretization of a given allocation site l is the set of all possible addresses that can be allocated at that line of code. An abstract heap is a map associating abstract addresses, i.e., lines of code, to abstract objects, namely $\text{HEAP}^\# \triangleq \text{ADDR}^\# \rightarrow \text{OBJ}^\#$. An abstract object is a map associating an automaton with an abstract value.

For what concerns environments, we consider a non-relational abstraction, approximating every identifier separately. This means that we abstract from $\wp(\text{ID} \rightarrow \text{VAL})$ to $\text{ID} \rightarrow \wp(\text{VAL})$. Abstract environments are maps from identifiers to abstract values, namely $\text{ENV}^\# \triangleq \text{ID} \rightarrow \text{VAL}^\#$, exploiting the abstraction between $\wp(\text{VAL})$ and $\text{VAL}^\#$ we have introduced in the previous subsection. Finally,

$$\begin{aligned}
\langle n \rangle_A^\# &\triangleq \alpha_I(\{n\}) & \langle x \rangle_A^\# \langle \xi^\#, \rho^\# \rangle &\triangleq \xi^\#(x) & \langle \mathbf{a}_1 + \mathbf{a}_2 \rangle_A^\# \sigma^\# &\triangleq \langle \mathbf{a}_1 \rangle_A^\# \sigma^\# +^I \langle \mathbf{a}_2 \rangle_A^\# \sigma^\# \\
\langle \mathbf{true} \rangle_B^\# \sigma^\# &\triangleq \mathbf{tt} & \langle x \rangle_B^\# \langle \xi^\#, \rho^\# \rangle &\triangleq \xi^\#(x) & \langle \mathbf{b}_1 \parallel \mathbf{b}_2 \rangle_B^\# \sigma^\# &\triangleq \langle \mathbf{b}_1 \rangle_B^\# \sigma^\# \sqcup_B^\# \langle \mathbf{b}_2 \rangle_B^\# \sigma^\# \\
\langle \{\} \rangle_O^\# \sigma^\# &\triangleq \lambda p. \langle \perp, \perp, \text{Min}(\emptyset), \emptyset, ? \rangle \\
\langle \{s_0 : \mathbf{e}_0, s_1 : \mathbf{e}_1, \dots, s_n : \mathbf{e}_n\} \rangle_O^\# \sigma^\# &\triangleq \\
& [s_n \mapsto \langle \mathbf{e}_n \rangle_E^\# \sigma^\#] \bullet^\# \dots [s_1 \mapsto \langle \mathbf{e}_1 \rangle_E^\# \sigma^\#] \bullet^\# [s_0 \mapsto \langle \mathbf{e}_0 \rangle_E^\# \sigma^\#] \\
\langle x[s] \rangle_E^\# \langle \xi^\#, \rho^\# \rangle &\triangleq \\
& \begin{cases} \sqcup_V^\# \{ \rho^\#(l)(p) \mid l \in \xi^\#(x) \wedge \mathcal{L}(p) \cap \mathcal{L}(\langle \mathbf{s} \rangle_S^\# \langle \xi^\#, \rho^\# \rangle) \neq \emptyset \} & \text{if } \xi^\#(x) \in \text{ADDR}^\# \\ \langle \perp, \perp, \text{Min}(\emptyset), \emptyset, ? \rangle & \text{otherwise} \end{cases} \\
\text{filter}_{\mathbf{true}}^\#(\sigma^\#) &\triangleq \sigma^\# & \text{filter}_x^\#(\langle \xi^\#, \rho^\# \rangle) &\triangleq \begin{cases} \langle \xi^\#, \rho^\# \rangle & \text{if } \xi^\#(x) \in \{\mathbf{tt}, \top\} \\ \sigma_\perp^\# & \text{otherwise} \end{cases} \\
\text{filter}_{\mathbf{b}_1 \parallel \mathbf{b}_2}^\#(\sigma^\#) &\triangleq \text{filter}_{\mathbf{b}_1}^\#(\sigma^\#) \sqcup^\# \text{filter}_{\mathbf{b}_2}^\#(\sigma^\#) & \text{filter}_{\mathbf{false}}^\#(\sigma^\#) &\triangleq \sigma_\perp^\#
\end{aligned}$$

Fig. 5: Abstract semantics for expressions and objects and the abstract filter

abstract states are, as in the concrete, pairs of abstract environments and abstract heaps, namely $\text{STATE}^\# \triangleq \text{ENV}^\# \times \text{HEAP}^\#$. The definition of the abstract join $\sqcup^\#$ and the partial order $\sqsubseteq^\#$ for $\text{STATE}^\#$ is straightforward.

The abstract semantics is then a function $(\text{st})^\# \in \text{STATE}^\# \rightarrow \text{STATE}^\#$, computing on abstract states. It relies on the abstract semantics for expressions $\langle \mathbf{e} \rangle_E^\# \in \text{STATE}^\# \rightarrow \text{VAL}^\#$, on the abstract semantics for objects $\langle \mathbf{o} \rangle_O^\# \in \text{STATE}^\# \rightarrow \text{OBJ}^\#$ and on the abstract filtering function $\text{filter}_b^\# \in \text{STATE}^\# \rightarrow \text{STATE}^{\#1}$. All of them must be sound w.r.t. their concrete counterparts, namely $\langle \mathbf{e} \rangle \gamma(\sigma^\#) \subseteq \gamma_V(\langle \mathbf{e} \rangle_E^\# \sigma^\#)$, $\langle \mathbf{o} \rangle \gamma(\sigma^\#) \subseteq \gamma_O(\langle \mathbf{o} \rangle_O^\# \sigma^\#)$ and $\text{filter}_b(\gamma(\sigma^\#)) \subseteq \gamma(\text{filter}_b^\#(\sigma^\#))$, for every $\sigma^\# \in \text{STATE}^\#$. In Fig. 5 we have a part of the definition of the abstract semantics for expressions and objects and the abstract filter, where $\sigma_\perp^\#$ is the minimum of the lattice $\langle \text{STATE}^\#, \sqsubseteq^\# \rangle$. The abstract semantics for statements is quite standard:

$$\begin{aligned}
\langle \text{st}_1 ; \text{st}_2 \rangle^\# \sigma^\# &\triangleq \langle \text{st}_2 \rangle^\# \langle \text{st}_1 \rangle^\# \sigma^\# & \langle \text{skip} \rangle^\# \sigma^\# &\triangleq \sigma^\# \\
\langle \text{if } \mathbf{b} \{ \text{st} \} \text{ else } \{ \text{st} \} \rangle^\# \sigma^\# &\triangleq \langle \text{st}_1 \rangle^\# \text{filter}_b^\#(\sigma^\#) \sqcup^\# \langle \text{st}_2 \rangle^\# \text{filter}_{!b}^\#(\sigma^\#) \\
\langle \text{while } \mathbf{b} \{ \text{st} \} \rangle^\# \sigma^\# &\triangleq \text{filter}_{!b}^\#(\text{lfp } \lambda \sigma_w^\# . \sigma^\# \sqcup^\# \langle \text{st} \rangle^\# \text{filter}_b^\#(\sigma_w^\#))
\end{aligned}$$

Concerning generic assignments, the abstract semantics follows the definition of the concrete one, so we have three cases: $x = \mathbf{e}$, where \mathbf{e} evaluates to a value; $x = \mathbf{o}$, where \mathbf{o} evaluates to an object; $x[s] = \mathbf{e}$, where \mathbf{s} evaluates to a string and \mathbf{e} evaluates to a value. In the first, we have to modify the abstract environment, setting x to the (abstract) evaluation of \mathbf{e} . In the second, we need to update the abstract address pointed by the identifier x , with the line of code of the assignment. Then we have to update the abstract object pointed, in the abstract heap, by the new line of code with the (abstract) evaluation of \mathbf{o} . Formally:

$$\begin{aligned}
\langle x = \mathbf{e} \rangle^\# \langle \xi^\#, \rho^\# \rangle &\triangleq \langle \xi^\#[x \leftarrow \langle \mathbf{e} \rangle_E^\# \langle \xi^\#, \rho^\# \rangle], \rho^\# \rangle \\
\langle x = \mathbf{new } \mathbf{o} \rangle^\# \langle \xi^\#, \rho^\# \rangle &\triangleq \langle \xi^\#[x \leftarrow \{\text{line}(x = \mathbf{new } \mathbf{o})\}], \rho^\#[\text{line}(x = \mathbf{o}) \leftarrow \langle \mathbf{o} \rangle_O^\# \langle \xi^\#, \rho^\# \rangle] \rangle
\end{aligned}$$

¹ We assume that all negations ! have been removed using DeMorgan's laws and usual arithmetic laws: $!(\mathbf{b}_1 \parallel \mathbf{b}_2) \equiv !\mathbf{b}_1 \&\& !\mathbf{b}_2$, $!(\mathbf{a}_1 < \mathbf{a}_2) \equiv (\mathbf{a}_2 < \mathbf{a}_1 \parallel \mathbf{a}_2 == \mathbf{a}_1)$, etc.

$$(a) \left[\begin{array}{l} \mathbf{b} \mapsto [2,2] \\ \mathbf{c} \mapsto [3,3] \\ \hline \mathbf{a}(\mathbf{a})^* \mapsto [4,4] \end{array} \right] \qquad (b) \left[\begin{array}{l} \mathbf{a} \mapsto [1,1] \\ \mathbf{b} \mapsto [2,2] \\ \mathbf{c} \mapsto [3,3] \\ \hline (\mathbf{a})^* \mapsto [4,4] \end{array} \right]$$

Fig. 6: Example of materialization.

As a third case, the so called *materialization*, we have the abstract semantics of object-property update, namely $x[s] = \mathbf{e}$. As we have already mentioned before, we allow to update only the objects that have been already stored into the heap. Suppose that $v^\sharp = \langle \mathbf{e} \rangle_{\mathbf{e}}^\sharp \langle \xi^\sharp, \rho^\sharp \rangle$, $p = \langle \mathbf{s} \rangle_{\mathbf{s}}^\sharp \langle \xi^\sharp, \rho^\sharp \rangle$ and $\{l_1, \dots, l_n\} = \xi^\sharp(x)$:

$$\begin{aligned} \text{let } o_i^\sharp &= \text{Norm}(\rho^\sharp(l_i)[p \leftarrow v^\sharp \sqcup_{\text{VAL}^\sharp} \rho^\sharp(l_i)(p)]), \text{ with } i \in \{1, \dots, n\} \text{ in} \\ \langle x[s] = \mathbf{e} \rangle^\sharp \langle \xi^\sharp, \rho^\sharp \rangle &\triangleq \langle \xi^\sharp, \rho^\sharp[l_1 \leftarrow o_1^\sharp, \dots, l_n \leftarrow o_n^\sharp] \rangle \end{aligned}$$

The abstract semantics of $x[s] = \mathbf{e}$ does not update the environment, since it only needs to update properties of abstract objects stored into the heap. For each location $l \in \text{ADDR}^\sharp$, associated to the identifier x , i.e., the one contained in $\xi^\sharp(x)$, the abstract semantics updates $\rho^\sharp(l)$, at the property p , with the value v^\sharp lubbed with the previous value of $\rho^\sharp(l)(p)$. This corresponds to a *weak update* of the object contained in x [3]. Before storing the updated abstract object in $\rho^\sharp(l)$, the latter is normalized. In this paper, we only perform weak updates. We could improve the precision of the analysis performing a *must-may analysis* in order to differentiate between properties that certainly point to some value and properties that may point to others. This can be done improving the proposed analysis using standard techniques, such as the ones reported in [3, 16, 17].

For example, let us suppose that $\rho^\sharp(l)$ is the object reported in Fig. 6(a) and we want to update the property expressed by p such that $\mathcal{L}(p) = \{\mathbf{a}\}$, with the interval $[1, 1]$. Applying these values to the previously defined abstract semantics, we obtain, at the allocation site l , the abstract object reported in Fig. 6(b). We say that the property \mathbf{a} has been *materialized*, since, before the update, it was part of a summary property, and after the update it is a non-summary property. We say that a (concrete) property is materialized when a string of an abstract object passes, during the update, from a summary property to a non-summary property. It is worth noting that normalization take care of materialization. The abstract semantics is sound w.r.t. the concrete collecting semantics, i.e., it computes an over-approximation of state invariants at every statement.

Theorem 1 (Soundness). *For every μJS program $\text{st} \in \text{STMT}$ we have that:*

$$\forall \sigma^\sharp \in \text{STATE}^\sharp. (\text{st})\gamma(\sigma^\sharp) \subseteq \gamma((\text{st})^\sharp \sigma^\sharp)$$

3.3 Widening

The domain $\langle \text{STATE}^\sharp, \sqsubseteq^\sharp \rangle$ is not ACC, i.e., it contains infinite ascending chains, because of the intervals abstract domain, the automata abstract domain and the novel objects abstract domain. Hence, fix-point computations in our abstract interpreter may diverge without an extrapolation operator. Hence, in order to enforce termination, the abstract domain VAL^\sharp is equipped with the widening

(a)	<pre> 1 o = new {a:1}; 2 key = "a"; 3 while (?) { 4 key = concat(key, "a"); 5 o[key] = 1; 6 }; </pre>	(b) $\left[\frac{\begin{array}{l} \mathbf{a} \mapsto [1,1] \\ \mathbf{aa} \mapsto [1,1] \\ \mathbf{aaa(a)^*} \mapsto [1,1] \end{array}}{\quad} \right]$
-----	---	--

 Fig. 7: (a) μ JS fragment, (b) Value of \mathbf{o} after **while**-loop.

(a)	<pre> 1 o = new {a:1}; 2 while (?) { 3 o["a"] = o["a"] + 1; 4 }; </pre>	(b) $\left[\frac{\mathbf{a} \mapsto [1, +\infty]}{-} \right]$
-----	---	--

 Fig. 8: (a) μ JS fragment, (b) Value of \mathbf{o} after **while**-loop.

operator $\nabla_v \in \text{VAL}^\sharp \times \text{VAL}^\sharp \rightarrow \text{VAL}^\sharp$ that is defined point-wise. In particular, intervals domain is equipped with its well-known widening defined in [7], finite state automata abstract domain is equipped with the widening ∇_{DFA}^n , reported in Sec. 2, while addresses and booleans are equipped with their least upper bound (they are finite). We can define the widening operator $\nabla_\varepsilon \in \text{ENV}^\sharp \times \text{ENV}^\sharp \rightarrow \text{ENV}^\sharp$ between environments upon ∇_v , applied point-wisely. For instance, suppose to use the widening ∇_{DFA}^n , with $n = 3$, for the finite state automata. We have that $[x \mapsto \langle [1, 1], \perp, \text{Min}(\mathbf{aaa}), \emptyset, \mathbf{def} \rangle] \nabla_\varepsilon [x \mapsto \langle [2, 2], \perp, \text{Min}(\mathbf{aaaa}), \emptyset, \mathbf{def} \rangle]$ is equal to the abstract environment $[x \mapsto \langle [1, +\infty], \perp, \text{Min}(\mathbf{a}^*), \emptyset, \mathbf{def} \rangle]$. Fix-point computations may also diverge on heaps, since also HEAP^\sharp is not ACC, due to the objects abstract domain. In particular, this happens because we model objects' properties with the finite state automata domain, which is not ACC. Anyway, a slight extension of the join $\sqcup_{\mathbf{o}}^\sharp$ is enough to guarantee termination of heap computations, exploiting the widening of the finite state automata domain. Informally speaking, abstract string values, in **while**-loop computations, always converge since finite state automata domain is equipped by a widening.

Let us consider the μ JS fragment reported in Fig. 7a and suppose that the boolean guard value is statically unknown. At each iteration on the **while**-loop, the string "a" is concatenated to the string value of **key** and then it is used to add a new property of the object \mathbf{o} . If the DFA_{\equiv} were not equipped with a widening, the value of **key** would diverge, as the properties of \mathbf{o} . Since convergence of strings is enforced by the widening ∇_{DFA}^n (with $n = 3$), also the objects' properties of \mathbf{o} converge. Indeed, the **while**-loop converges and the abstract interpreter produces, for the variable \mathbf{o} , the (normalized) object reported in Fig. 7b. Clearly, the simple object join is enough for objects' properties convergence but it is not for the associated value. For example, let consider the μ JS fragment reported in Fig. 8a. In this case, the number of properties of the object \mathbf{o} does not increase in the **while**-loop but it only increase the value of the property **a**. The idea behind the widening for objects is to apply the widening of values point-wisely between the properties of the two objects. Hence, we define the widening on OBJ^\sharp as: $o_1^\sharp \nabla_{\mathbf{o}} o_2^\sharp \triangleq \lambda p. o_1^\sharp(p) \nabla_v o_2^\sharp(p)$. Coming back to the example, applying the widening defined above, the abstract value of \mathbf{o} after the **while**-loop execution is reported in Fig. 8b. We then use this widening in order to define the widening for abstract heaps and, in turn, for abstract states.

Motivating example. We now show the so far defined analysis on the example reported in the introduction (Fig. 1). It is worth noting that, in this example, object widening does not occur. We have already commented it with the fragments reported in Fig. 7 and Fig. 8. The goal is to reason about the value of \mathbf{n} at the end of the execution. At the beginning of the first iteration of the while loop, the value of \mathbf{idx} is $\langle \perp, \perp, (\mathbf{a} \parallel \mathbf{b}), \emptyset, \mathbf{def} \rangle$. The latter is used to access \mathbf{obj} and then the result is stored in \mathbf{n} (line 5). Since the property \mathbf{b} is not present in \mathbf{obj} , only the property \mathbf{a} is accessed by \mathbf{idx} , and the value of \mathbf{n} is $\langle [1, 1], \perp, \text{Min}(\emptyset), \emptyset, \mathbf{def} \rangle$. Before starting the next iteration, \mathbf{idx} is updated at line 6 and its value becomes $\langle \perp, \perp, (\mathbf{aa} \parallel \mathbf{ba}), \emptyset, \mathbf{def} \rangle$. Moreover, widening is applied before starting a new iteration. Supposing to apply the widening ∇_{DFA}^n , $n = 1$, the values of the variables before the second (and last) iteration are: $\mathbf{n} = \langle [0, +\infty], \perp, \text{Min}(\emptyset), \emptyset, \mathbf{def} \rangle$, $\mathbf{idx} = \langle \perp, \perp, (\mathbf{a} \parallel \mathbf{b}) \mathbf{a}^+, \emptyset, \mathbf{def} \rangle$ (other variables do not change). At the second iteration, \mathbf{obj} is accessed by the new abstract value of \mathbf{idx} , hence it will access the properties \mathbf{a} , \mathbf{aa} and \mathbf{ab} , updating the variable \mathbf{n} at line 5 with $\langle [1, +\infty], \perp, \text{Min}(\emptyset), \emptyset, \mathbf{def} \rangle$. The previous value of \mathbf{n} (before starting the second iteration) is widened with the abstract value of \mathbf{n} at the end of the second iteration, producing the value $\langle [0, +\infty], \perp, \text{Min}(\emptyset), \emptyset, \mathbf{def} \rangle$, that it is also the value holding after the while loop, at line 7, since with this value, the fix-point is reached. Finally, at line 8, the abstract value of \mathbf{n} is assigned to $\mathbf{obj}[\mathbf{idx}]$, updating the abstract object of \mathbf{obj} as follows (we omit bottom values): $[\mathbf{a} \mapsto [1, 1], \mathbf{aa} \mapsto [0, +\infty], \mathbf{ab} \mapsto [3, 3], \mathbf{ac} \mapsto \text{Min}(\{\text{"world"}\}), (\mathbf{a} \parallel \mathbf{b}) \mathbf{a}^+ \setminus \mathbf{aa} \mapsto [0, +\infty]]$. The summary property $(\mathbf{a} \parallel \mathbf{b}) \mathbf{a}^+ \setminus \mathbf{aa}$ is added and only the properties \mathbf{aa} and \mathbf{ab} are modified. Properties already present in \mathbf{obj} remain unaltered (e.g. \mathbf{a} and \mathbf{ac}).

4 Discussion and conclusion

We have proposed an abstract domain suitable for the analysis of objects' properties in dynamic programming languages. The novelty consists in exploiting *finite state automata*, in order to approximate objects' properties. This leads to a better precision (less false positives), compared to state-of-the-art domains approximating strings (for instance, [5, 6]). A key aspect of our abstract domain is the *normal form for objects* and, in the paper, we have presented a normalization algorithm: it transforms objects in their normal form. An object is in normal form if and only if it has only two kind of properties: *summary* and *non-summary*. The idea behind summarization, and hence materialization, is not new in static analysis, and comes from the well-known shape analysis [16]. For example, this idea has been adopted in [11], where the authors present a static analyzer for PHP that also involve heap analysis, where the heap, in their abstraction, is made of summary heap identifiers and non-summary heap identifiers. In particular, in [11], a summary heap identifier summarizes all the elements of the heap that could be updated by statically unknown assignments. We have adopted the same idea with the difference that we may have more summary properties, expressed by automata recognizing infinite languages, rather than a single summary property that merges together heap elements updated by statically unknown assignments. The idea of summarization has been also taken into

account in [3], where the authors propose the recency abstraction, which consists in representing each abstract allocation site with two memory regions, namely the *most recently allocated block* and the *not most recently allocated blocks*. The latter is basically a summary memory region, since more than one block may be allocated. Recency abstraction has been implemented also in TAJIS [13], an abstract interpretation-based static analyzer for JavaScript, showing that such abstraction outperforms other abstract allocation-based techniques. As future work, we aim to implement our objects' abstract domain upon TAJIS. We believe that combining our abstract domain and recency abstraction can produce good results and it would be interesting to make a comparison with TAJIS and other JavaScript static analyzers, such as SAFE [15] and JSAI [14].

References

1. Arceri, V., Maffei, S.: Abstract domains for type juggling. *Electr. Notes Theor. Comput. Sci.* **331** (2017)
2. Arceri, V., Mastroeni, I.: Static program analysis for string manipulation languages. In: VPT'19 (2019). <https://doi.org/10.4204/EPTCS.299.5>
3. Balakrishnan, G., Reps, T.W.: Recency-abstraction for heap-allocated storage. In: SAS'16 (2006). https://doi.org/10.1007/11823230_15
4. Bartzis, C., Bultan, T.: Widening Arithmetic Automata. In: CAV'04 (2004)
5. Cortesi, A., Olliaro, M.: M-String Segmentation: A Refined Abstract Domain for String Analysis in C Programs. In: TASE'18 (2018)
6. Costantini, G., Ferrara, P., Cortesi, A.: A Suite of Abstract Domains for Static Analysis of String Values. *Softw., Pract. Exper.* **45**(2) (2015)
7. Cousot, P., Cousot, R.: Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints. In: POPL'77
8. Cousot, P., Halbwachs, N.: Automatic discovery of linear restraints among variables of a program. In: POPL (1978)
9. Davis, M.D., Sigal, R., Weyuker, E.J.: *Computability, Complexity, and Languages: Fund. of Theor. CS.* Academic Press Professional, Inc. (1994)
10. D'Silva, V.: Widening for Automata. MsC Thesis, Inst. Fur Inform. - UZH (2006)
11. Haurar, D., Kofron, J.: Framework for static analysis of PHP applications. In: ECOOP'15 (2015). <https://doi.org/10.4230/LIPIcs.ECOOP.2015.689>
12. Hopcroft, J.E., Ullman, J.D.: *Introduction to Automata Theory, Languages and Computation.* Addison-Wesley (1979)
13. Jensen, S.H., Møller, A., Thiemann, P.: Type analysis for javascript. In: SAS '09 (2009). https://doi.org/10.1007/978-3-642-03237-0_17
14. Kashyap, V., Dewey, K., Kuefner, E.A., Wagner, J., Gibbons, K., Sarracino, J., Wiedermann, B., Hardekopf, B.: JSAI: a Static Analysis Platform for JavaScript. In: FSE '14 (2014)
15. Lee, H., Won, S., Jin, J., Cho, J., Ryu, S.: SAFE: Formal Specification and Implementation of a Scalable Analysis Framework for ECMAScript. In: FOOL (2012)
16. Nielson, F., Nielson, H.R., Hankin, C.: *Principles of program analysis.* Springer (1999). <https://doi.org/10.1007/978-3-662-03811-6>
17. Wilhelm, R., Sagiv, S., Reps, T.W.: Shape analysis. In: CC'00 (2000). https://doi.org/10.1007/3-540-46423-9_1