



UNIVERSITÀ DI PARMA

ARCHIVIO DELLA RICERCA

University of Parma Research Repository

Completeness of string analysis for dynamic languages

This is the peer reviewed version of the following article:

Original

Completeness of string analysis for dynamic languages / Arceri, V.; Olliaro, M.; Cortesi, A.; Mastroeni, I.. - In: INFORMATION AND COMPUTATION. - ISSN 0890-5401. - (2021), p. 104791.104791. [10.1016/j.ic.2021.104791]

Availability:

This version is available at: 11381/2899227 since: 2024-11-18T20:11:03Z

Publisher:

Elsevier Inc.

Published

DOI:10.1016/j.ic.2021.104791

Terms of use:

Anyone can freely access the full text of works made available as "Open Access". Works made available

Publisher copyright

note finali coverpage

(Article begins on next page)

03 January 2025

Completeness of String Analysis for Dynamic Languages

Vincenzo Arceri^{b,*}, Martina Olliaro^{b,c}, Agostino Cortesi^b,
Isabella Mastroeni^a

^a*University of Verona. Department of Computer Science, Strada le Grazie 15,
37134 Verona, Italy*

^b*Ca' Foscari University of Venice. Scientific Campus, Via Torino 155,
30172 Mestre, Venice, Italy*

^c*Masaryk University. Faculty of Informatics, Botanická 68A,
60200 Brno, Czech Republic*

Abstract

In Abstract Interpretation, completeness ensures that the analysis does not lose information with respect to the property of interest. In particular, for dynamic languages like JavaScript, completeness of string analysis is a key security issue, as poorly managed string manipulation code may easily lead to significant security flaws. In this paper, we provide a systematic and constructive approach for generating the completion of string domains for dynamic languages, and we apply it to the refinement of existing string abstractions. We also provide an effective procedure to measure the precision improvement obtained when lifting the analysis to complete domains.

Keywords: Abstract Interpretation, String Analysis, Completeness

1. Introduction

Despite the growth of support for string manipulation in programming languages, string manipulation errors still lead to code vulnerabilities that can be exploited by malicious agents, causing potential catastrophic damages. This is even more true in the context of web applications, where common programming languages used for web-based software development (e.g.,

*Corresponding author

Email addresses: vincenzo.arceri@unive.it (Vincenzo Arceri),
martina.olliaro@unive.it (Martina Olliaro), cortesi@unive.it (Agostino Cortesi),
isabella.mastroeni@univr.it (Isabella Mastroeni)

JavaScript) offer a wide range of dynamic features that make string manipulation dangerous.

String analysis is a static program analysis technique that computes, for each execution trace of the program given as input, the set of the possible string values that may reach a certain program point. String analysis, like any other non-trivial program analysis, is an undecidable task. Thus, a certain degree of approximation is necessary to find evidence of bugs and vulnerabilities in string manipulating code. In the recent literature, different approximation techniques for string analysis have been developed, such as [8]: automata-based [6, 11, 49, 48], abstraction-based [16, 2, 4, 15, 3, 50], constraint-based [1, 42, 45, 32], and grammar-based [36, 47], and have been used, in particular, to detect web application vulnerabilities [47, 48, 49].

In this paper¹ we focus on string analysis by means of the Abstract Interpretation theory [17, 18]. Cousot and Cousot have proposed Abstract Interpretation in the 70s as a theory of sound abstraction (or approximation) of the semantics of computer programs, and nowadays, it is widely integrated in software verification tools, and it is used to prove approximations correctness by means of rigorous mathematical methods. Since the introduction of the Abstract Interpretation theory, many abstract domains representing properties of interest about numerical domains values have been designed [17, 19, 25, 37, 10, 12, 38, 26, 43]. On the other hand, just in the last few years, the scientific community has taken an interest in the development of abstract domains for string analysis [16, 15, 2, 4, 34, 29, 39], some of them language-specific, such as those defined as part of the JavaScript static analysers: TAJIS [27], SAFE [31], and JSAI [28].

Important features of Abstract Interpretation are *soundness* and *completeness* [18]. If soundness (or correctness), as a basic requirement, should always be guaranteed by static analysis tools to avoid the presence of false negatives, completeness is frequently not met. If completeness is satisfied, it means that the abstract computations do not lose information, during the abstraction process, with respect to a property of interest, and so the anal-

¹This paper is a revised and extended version of [7]. Precisely, we have added an auxiliary procedure to facilitate the reader in understanding the fundamentals theorems of completeness. Moreover, we have introduced the analysis relative precision, which quantifies the increment of the analysis precision gained by analysing a program with a complete abstract domain with respect to when it is analysed with its original version. Finally, we have presented an experimental evaluation.

ysis can be considered optimal. In [24], Giacobazzi et al. highlighted the fact that completeness is an abstract domain property, and they presented a methodology to obtain complete abstract domains with respect to operations by minimally extending or restricting the underlying domains.

1.1. Paper contribution

The goal of this paper is to provide a way-to-proceed in the context of imprecise string abstractions. In particular, we exploit the complete shells' theoretical framework, constructively showing how to improve the precision of incomplete abstractions, without designing new string abstract domains.

We consider two JavaScript-specific string abstract domains defined as part of TAJIS [27] and SAFE [31] static analysers, focusing on their completeness with respect to the main string-manipulating operations. We compute their complete versions, and we discuss the benefits of guaranteeing completeness in the context of Abstract Interpretation-based string analysis of dynamic languages. Finally, we present an effective procedure to measure the precision increment when analysing a program with a complete abstract domain.

1.2. Paper structure

Section 2 gives basics in mathematics and Abstract Interpretation. Section 3 recalls relevant concepts related to the completeness property in Abstract Interpretation that we will use throughout the whole paper. Moreover, a motivating example is given to show the importance of guaranteeing completeness in an Abstract Interpretation-based analysis with respect to strings. Section 4 defines our core language. Section 5 presents the completion of the string abstract domains integrated into TAJIS and SAFE static analysers with respect to two operations of interest. Section 6 highlights the strengths and usefulness of the completeness approach to static analysis of JavaScript string manipulating programs. Section 7 defines the measurement procedure of the analysis precision increment. Section 8 concludes.

2. Background: Foundations of Abstract Interpretation

Mathematical notation. Given a set S , we denote by S^* the set of all the finite sequences of elements of S . If $s = s_0 \dots s_n \in S^*$, we denote by s_i the i -th element of s , and by $|s| = n + 1$ its length. We denote by $s[x/i]$ the sequence obtained replacing s_i in s with x . Given two sets S and T , we denote with $\mathcal{P}(S)$

the powerset of S , with $S \setminus T$ the set difference, with $S \subset T$ the strict inclusion relation, with $S \subseteq T$ the inclusion relation and with $S \times T$ the Cartesian product between S and T . We denote by S^n the n -ary Cartesian product of a set S , with $n \geq 2$. We denote by $f : S \rightarrow T$ a function from elements of S (domain) to elements of T (image). Given $f : S \rightarrow T$ and $g : T \rightarrow Q$ we denote with $g \circ f : S \rightarrow Q$ their composition, i.e., $g \circ f = \lambda x. g(f(x))$. When f is a function of arity n , i.e., $f : S^n \rightarrow T$, with $s \in S^n$ and $i \in [0, n)$, $f_s^i = \lambda z. f(s[z/i]) : S \rightarrow T$ is the same function where all the parameters but the i -th are fixed by s , namely $f_s^i = f(s_0, \dots, s_{i-1}, x, s_{i+1}, \dots, s_n)$.

Ordered structures. A set L with a partial ordering relation \sqsubseteq is a poset and it is denoted by $\langle L, \sqsubseteq \rangle$. Given a poset $\langle L, \sqsubseteq \rangle$ and a $X \subseteq L$, an upper bound of X is an element $y \in L$ such that $x \sqsubseteq y, \forall x \in X$. An upper bound y for X is the least upper bound (lub) of X if and only if for every other upper bound y' of X it holds that $y \sqsubseteq y'$. The least upper bound, when it exists, it is unique. Dually, a lower bound of X is an element $y \in L$ such that $y \sqsubseteq x, \forall x \in X$. A lower bound y of X is the greatest lower bound (glb) of X if and only if for every other lower bound y' of X it holds that $y' \sqsubseteq y$. A poset $\langle L, \sqsubseteq, \sqcup, \sqcap \rangle$, where \sqcup and \sqcap are respectively the lub and glb operators of L , is a lattice if $\forall x, y \in L$ we have that $x \sqcap y$ and $x \sqcup y$ belong to L , and we say that it is also complete when for each $X \subseteq L$ we have that $\sqcap X, \sqcup X \in L$. Given a poset $\langle L, \sqsubseteq \rangle$ and $S \subseteq L$, we denote by $\max^{\sqsubseteq}(S) = \{x \in S \mid \forall y \in S. x \sqsubseteq y \Rightarrow x = y\}$ the set of the maximal elements of S in L w.r.t. \sqsubseteq . As usual, a complete lattice L , with ordering \sqsubseteq , lub \sqcap , glb \sqcup , greatest element (top) \top , and least element (bottom) \perp is denoted by $\langle L, \sqsubseteq, \sqcup, \sqcap, \top, \perp \rangle$. An *upper closure operator* on a poset $\langle L, \sqsubseteq \rangle$ is an operator $\rho : L \rightarrow L$ which is monotone (i.e., $x \sqsubseteq y \Rightarrow \rho(x) \sqsubseteq \rho(y)$), idempotent (i.e., $\rho(\rho(x)) = \rho(x)$), and extensive (i.e., $x \sqsubseteq \rho(x)$). The set of all closure operators on a poset L is denoted by $uco(L)$.

Abstract Interpretation. Abstract Interpretation [17, 18] is a theoretical framework for sound reasoning about program semantic properties of interest, and can be equivalently formalized either as Galois connections or closure operators on a given concrete domain [18]. In this paper, we shall assume that the concrete domain D is a complete lattice w.r.t. to a certain partial order \sqsubseteq_D . Observe that this condition is not restrictive, as if it is not matched it is sufficient to lift the concrete domain to its powerset, yielding a complete lattice w.r.t. the set inclusion partial order. This lifting operation in liter-

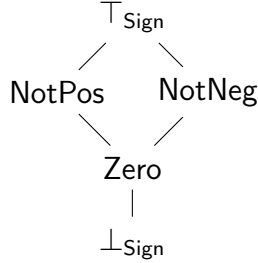


Figure 1: Sign abstract domain

ature is applied in the construction of a *collecting semantics* starting from any concrete semantics [17].

Let \mathbf{D} and $\bar{\mathbf{D}}$ be complete lattices, which respectively represent a concrete and an abstract domain, a pair of monotone functions $\alpha : \mathbf{D} \rightarrow \bar{\mathbf{D}}$ and $\gamma : \bar{\mathbf{D}} \rightarrow \mathbf{D}$ forms a *Galois Connection* (GC) between \mathbf{D} and $\bar{\mathbf{D}}$ if for every $d \in \mathbf{D}$ and for every $\bar{d} \in \bar{\mathbf{D}}$ we have $\alpha(d) \sqsubseteq_{\bar{\mathbf{D}}} \bar{d} \Leftrightarrow d \sqsubseteq_{\mathbf{D}} \gamma(\bar{d})$. The function α (resp. γ) is the *left-adjoint* (resp. *right-adjoint*) to γ (resp. α), and it is additive (resp. co-additive). We denote a GC as $\mathbf{D} \xleftrightarrow[\alpha]{\gamma} \bar{\mathbf{D}}$. Given $\mathbf{D} \xleftrightarrow[\alpha]{\gamma} \bar{\mathbf{D}}$, then $\gamma \circ \alpha \in uco(\mathbf{D})$. In particular, an upper closure operator (over \mathbf{D}) uniquely identifies a GC, and viceversa. For example, let us consider the abstract domain of **Sign**, abstracting integer sets (i.e., $\mathcal{P}(\mathbb{Z})$ is the concrete domain) reported in Figure 1. We can also see the GC between $\mathcal{P}(\mathbb{Z})$ and **Sign** by means of the corresponding upper closure operator $\rho_{\text{Sign}} : \mathcal{P}(\mathbb{Z}) \rightarrow \mathcal{P}(\mathbb{Z})$, where its image can be either \emptyset , $\{0\}$, $\{n \leq 0 \mid n \in \mathbb{Z}\}$, $\{n \geq 0 \mid n \in \mathbb{Z}\}$ or \mathbb{Z} , corresponding to the concretization of the abstract elements **Zero**, **NotPos**, **NotNeg**, \perp_{Sign} and \top_{Sign} , respectively, of the abstract domain **Sign**.

If \mathbf{D} is a complete lattice, then $\langle uco(\mathbf{D}), \sqsubseteq, \sqcup, \sqcap, \lambda d. \mathbf{d}, \text{id} \rangle$ forms a complete lattice [46], which is the set of all possible abstractions of \mathbf{D} , where the bottom element is $\text{id} = \lambda d. \mathbf{d}$, and for every $\rho, \eta \in uco(\mathbf{D})$, ρ is *more concrete than* η iff $\rho \sqsubseteq \eta$ iff $\forall d' \in \mathbf{D}. \rho(d') \sqsubseteq \eta(d')$, and $(\rho \sqcap \eta)(d) = \rho(d) \sqcap \eta(d)$ and $(\rho \sqcup \eta)(d) = d$ iff $\rho(d) = \eta(d) = d$. The operator $\rho \in uco(\mathbf{D})$ is disjunctive when $\rho(\mathbf{D})$ is a join-sublattice of \mathbf{D} which holds iff ρ is additive [18]. Let L be a complete lattice, then $X \subseteq L$ is a Moore family of L if $X = \mathfrak{M}(X) = \{\sqcap S \mid S \subseteq X\}$, where $\sqcap \emptyset = \top$. The condition that any concrete element of \mathbf{D} has the best abstraction in the abstract domain $\bar{\mathbf{D}}$, implies that $\bar{\mathbf{D}}$ is a Moore family of \mathbf{D} . We denote by $\mathfrak{M}(X)$ the Moore closure of $X \subseteq \mathbf{D}$, that is the least subset of \mathbf{D} , which is a Moore family of \mathbf{D} , and

contains X . Given $D \xleftrightarrow[\alpha]{\gamma} \bar{D}$, a concrete function $f : D \rightarrow D$ may be not computable. We use a function $f^\sharp : \bar{D} \rightarrow \bar{D}$ to *correctly* approximate f , namely f^\sharp must be *sound*.

Definition 1 (Soundness). Given $D \xleftrightarrow[\alpha]{\gamma} \bar{D}$ and a concrete function $f : D \rightarrow D$, an abstract function $f^\sharp : \bar{D} \rightarrow \bar{D}$ is sound w.r.t. f if

$$\forall d \in D. \alpha(f(d)) \sqsubseteq_{\bar{D}} f^\sharp(\alpha(d))$$

Among all the sound abstract functions f^\sharp , we aim at the best one, namely the one that loses less information when approximating the function f . This property is captured by the notion of best correct approximation.

Definition 2 (Best correct approximation). Given $D \xleftrightarrow[\alpha]{\gamma} \bar{D}$ and a concrete function $f : D \rightarrow D$, the function $\alpha \circ f \circ \gamma : \bar{D} \rightarrow \bar{D}$ is the best correct approximation of f .

In Abstract Interpretation, there exist two notions of completeness. *Backward completeness* property focuses on complete abstractions of the inputs, while *forward completeness* [23, 22, 21] focuses on complete abstractions of the outputs, both w.r.t. an operation of interest. In this paper, we focus on the more typical and best known notion of completeness, i.e., the backward completeness. In particular, the notion is obtained by enforcing the equality in the soundness condition reported in Definition 1, as reported in Definition 3.

Definition 3 (Backward completeness). Given $D \xleftrightarrow[\alpha]{\gamma} \bar{D}$, a concrete function $f : D \rightarrow D$ and an abstract function $f^\sharp : \bar{D} \rightarrow \bar{D}$, f^\sharp is backward complete w.r.t. f if

$$\forall d \in D. \alpha(f(d)) = f^\sharp(\alpha(d))$$

Having backward complete abstract functions is a desirable property since, when backward completeness is met, we have the guarantee that no loss of information arises during the input abstraction process, w.r.t. an operation of interest. For instance, the **Sign** numerical abstract domain, depicted in Figure 1, is backward complete w.r.t. the multiplication numerical operation. Indeed, in order to compute the sign of the expression $e_1 * e_2$ it is enough to know the sign of e_1 and e_2 , without any loss of information during the abstraction of the operands.

In the rest of the paper, when we will talk about completeness, we mean backward completeness. While the best correct approximation always induces a sound abstraction, this is not the case of completeness. Indeed, given a concrete function $f : D \rightarrow D'$, a pair of abstract domains $\langle \bar{D}, \bar{D}' \rangle \in uco(D) \times uco(D')$, with $\gamma_{\bar{D}, D} : \bar{D} \rightarrow D$ the concretization function from \bar{D} to D and $\alpha_{D', \bar{D}'} : D' \rightarrow \bar{D}'$ the abstraction function from D' to \bar{D}' , then there exists a complete function $f^\sharp : \bar{D} \rightarrow \bar{D}'$ iff the best correct approximation $\alpha_{D', \bar{D}'} \circ f \circ \gamma_{\bar{D}, D} : \bar{D} \rightarrow \bar{D}'$ is complete (and it is equal to f^\sharp) [24]. Given two posets D and D' and a function $f : D^n \rightarrow D'$, we denote by $\Gamma(D, D', f)$ the set of pairs of abstract domains $\langle \rho, \eta \rangle \in uco(D) \times uco(D')$ which are complete for f . Hence, given a function $f : D^n \rightarrow D'$, we can have more than one pair of complete domains $\langle \rho, \eta \rangle \in uco(D) \times uco(D')$ for f .

3. Background: Making Abstract Interpretations complete

In this Section, we recall notions of methodologies introduced in [24] that we will use through the whole paper, to constructively build, from an initial abstract domain, a new abstract domain that is complete w.r.t. an operation of interest. Finally, a motivating example showing the usefulness of completion of abstract domains for string analysis is given.

As reported in [24], it is worth noting that completeness is a property related to the underlying abstract domain. Starting from this fact, in [24], authors proposed a constructive method to manipulate the underlying incomplete abstract domain in order to get a complete abstract domain w.r.t. a certain operation. In particular, given two abstract domains \bar{D} and \bar{D}' and an operator $f : \bar{D}^n \rightarrow \bar{D}'$ with $n \in \mathbb{N}$, the authors gave two different notions of completion of abstract domains w.r.t. f : the one that *adds* the minimal number of abstract points to the input abstract domain \bar{D} and the other that *removes* the minimal number of abstract points from the output abstract domain \bar{D}' . The first approach captures the notion of *complete shell of \bar{D}* , while the latter defines the *complete core of \bar{D}'* , both w.r.t. an operator f .

Complete shell vs complete core. We will focus on the construction of complete shells of string abstract domains, rather than complete cores. This choice is guided by the fact that a complete core for an operation f removes abstract points from a starting abstract domain. So, even if it is complete for f , the complete core could worsen the precision of other operations.

Conversely, complete shells augment the starting abstract domains (adding abstract points), and consequently, they cannot compromise the precision of other operations.

Below, we recall two important theorems proved in [24] that provide a constructive method to compute abstract domain complete shells, defined in terms of an upper closure operator ρ . Precisely, the latter theorems present two notions of complete shells: *i. complete shells of ρ relative to η* (where η is an upper closure operator), meaning that they are complete shells of operations defined on ρ that return results in η , and *ii. absolute complete shells of ρ* , meaning that they are complete shells of operations that are defined on ρ and return results in ρ .

Definition 4 (Complete shell of ρ relative to η [24]). Let $\langle A, \sqsubseteq_A, \sqcup_A \rangle$ and $\langle B, \sqsubseteq_B, \sqcup_B \rangle$ be two posets and $f : A^n \rightarrow B$ be a continuous function. Given $\rho \in uco(A)$ and $\eta \in uco(B)$, then let $\mathcal{S}_f^\eta : uco(A) \rightarrow uco(A)$ be the domain transformer:

$$\mathcal{S}_f^\eta(\rho) \stackrel{\text{def}}{=} \bigsqcup_{uco(A)} \{ \delta \in uco(A) \mid \delta \sqsubseteq \rho, \langle \delta, \eta \rangle \in \Gamma(A, B, f) \}.$$

If $\langle \mathcal{S}_f^\eta(\rho), \eta \rangle \in \Gamma(A, B, f)$, then $\mathcal{S}_f^\eta(\rho)$ is called *complete shell of ρ relative to η with respect to an operation f* .

As discussed in [24], Definition 4 does not offer a constructive methodology to compute $\mathcal{S}_f^\eta(\rho)$. Theorem 1 reports a constructive characterization of the complete shell of ρ relative to η w.r.t. f , making use of the Moore closure operator defined in Section 2.

Theorem 1 ([24]). Let $\langle A, \sqsubseteq_A, \sqcup_A \rangle$ and $\langle B, \sqsubseteq_B, \sqcup_B \rangle$ be two posets and $f : A^n \rightarrow B$ be a continuous function. Given $\rho \in uco(A)$, $\eta \in uco(B)$ and $\mathcal{S}_f^\eta(\rho)$ as in Definition 4, the following equality holds:

$$\mathcal{S}_f^\eta(\rho) = \mathfrak{M}(\rho \cup (\bigsqcup_{\substack{A \\ i \in [0, n) \\ x \in A^n, y \in \eta}} \max^{\sqsubseteq_A}(\{z \in A \mid (f_x^i)(z) \sqsubseteq_B y\}))).$$

As already mentioned above, the idea behind the complete shell of ρ (input abstraction) relative to η (output abstraction) is to refine ρ adding the minimum number of abstract points to make ρ complete w.r.t. an operation f . By Theorem 1, this is obtained by adding to ρ the maximal elements in

A , whose f image is dominated by elements in η , at least in one dimension i . Clearly, the so-obtained abstraction may not be an upper closure operator for A . Hence, the Moore closure operator is applied. On the other hand, absolute complete shells are involved in the case in which the operator f of interest has same input and output abstract domain, i.e., $f : A^n \rightarrow A$. In this case, given $\rho \in uco(A)$, absolute complete shells of ρ can be obtained as the greatest fix-point (*gfp*) of the domain transformer presented in Definition 4, as stated below.

Definition 5 (Absolute complete shell of ρ [24]). Let $\langle A, \sqsubseteq_A, \sqcup_A \rangle$ be a poset and $f : A^n \rightarrow A$ be a continuous function. Given $\rho \in uco(A)$, then let $\bar{\mathcal{S}}_f : uco(A) \rightarrow uco(A)$ be the domain transformer:

$$\bar{\mathcal{S}}_f(\rho) \stackrel{\text{def}}{=} \bigsqcup_{uco(A)} \{ \delta \in uco(A) \mid \delta \sqsubseteq \rho, \langle \delta, \delta \rangle \in \Gamma(A, A, f) \}.$$

If $\langle \bar{\mathcal{S}}_f(\rho), \rho \rangle \in \Gamma(A, A, f)$, then $\bar{\mathcal{S}}_f(\rho)$ is called *absolute complete shell of ρ with respect to an operation f* .

Theorem 2 ([24]). Let A be a poset and $f : A^n \rightarrow A$ be a continuous function. For $\rho \in uco(A)$, if $\mathcal{S}_f^p(\rho)$ is the complete shell of ρ relative to ρ w.r.t. f , and if $\bar{\mathcal{S}}_f(\rho)$ is the absolute complete shell of ρ w.r.t. f , then the following equality holds:

$$\bar{\mathcal{S}}_f(\rho) = \text{gfp}(\lambda \rho. \mathcal{S}_f^p(\rho)).$$

In [24], Giacobazzi et al. have discussed the completeness and incompleteness of the **Sign** abstract domain, approximating numerical values, depicted in Figure 1. **Sign** is complete for the product operation. Let $*$: $\mathcal{P}(\mathbb{Z}) \times \mathcal{P}(\mathbb{Z}) \rightarrow \mathcal{P}(\mathbb{Z})$ be the concrete product operation and $*_{\text{Sign}} : \text{Sign} \times \text{Sign} \rightarrow \text{Sign}$ be the corresponding abstract product operation, following the well known sign rules (e.g., **NotPos** $*_{\text{Sign}}$ **NotPos** = **NotNeg**). Given the expression $e_1 * e_2$, the equality $\alpha_{\text{Sign}}(e_1 * e_2) = \alpha_{\text{Sign}}(e_1) *_{\text{Sign}} \alpha_{\text{Sign}}(e_2)$ holds, with α_{Sign} being the abstraction function of **Sign**. As an example, consider the concrete expression $\{2, 5\} * \{-1, -3\}$, then $\alpha_{\text{Sign}}(\{2, 5\} * \{-1, -3\}) = \alpha_{\text{Sign}}(\{-2, -5, -6, -15\}) = \text{NotPos}$ is equal to $\alpha_{\text{Sign}}(\{2, 5\}) *_{\text{Sign}} \alpha_{\text{Sign}}(\{-1, -3\}) = \text{NotNeg} *_{\text{Sign}} \text{NotPos} = \text{NotPos}$.

On the other hand, **Sign** is not complete for the sum operation. Let $+$: $\mathcal{P}(\mathbb{Z}) \times \mathcal{P}(\mathbb{Z}) \rightarrow \mathcal{P}(\mathbb{Z})$ be the concrete sum operation and $+_{\text{Sign}} : \text{Sign} \times \text{Sign} \rightarrow$

Algorithm 1: Relative complete shell procedure pseudo-code

Input: $\langle A, \sqsubseteq_A \rangle, \langle B, \sqsubseteq_B \rangle$ (posets),
 $\rho \in uco(A), \eta \in uco(B)$ (abstractions),
 $f : A^n \rightarrow B$

Output: $\mathcal{S}_f^\eta(\rho)$

```
1  $X \leftarrow \emptyset$ ;  
2 foreach  $i \in [0, n)$  and  $x \in A$  and  $y \in \eta$  do  
3   let  $z \in A$  be the maximum element such that  $f_x^i(z) \sqsubseteq_B y$ ;  
4   add  $z$  to  $X$ ;  
5 let  $\mathcal{S}_f^\eta(\rho)$  be the Moore closure of  $\rho \cup X$ ;  
6 return  $\mathcal{S}_f^\eta(\rho)$ ;
```

Sign be the corresponding abstract sum operation. Consider the concrete expression $\{2\} + \{-1, -2\}$, then $\alpha_{\text{Sign}}(\{2\} + \{-1, -2\}) = \alpha_{\text{Sign}}(\{0, 1\}) = \text{NotNeg}$ is not equal to $\alpha_{\text{Sign}}(\{2\}) +_{\text{Sign}} \alpha_{\text{Sign}}(\{-1, -2\}) = \text{NotNeg} +_{\text{Sign}} \text{NotPos} = \mathbb{Z}$.

In [24], the absolute complete shell of **Sign** w.r.t. the sum operation has been computed, which corresponds to the interval abstract domain [17].

Domain completion procedure. To improve the understanding about how to obtain a complete domain w.r.t. an operation of interest we provide a step-by-step reading of the formula in Theorem 1 (i.e., *relative complete shell*) by means of the procedure reported in Algorithm 1. The algorithm takes as input two posets $\langle A, \sqsubseteq_A \rangle, \langle B, \sqsubseteq_B \rangle$, two closures $\rho \in uco(A), \eta \in uco(B)$, which respectively correspond to the input and output abstraction, and a continuous function $f : A^n \rightarrow B$. The procedure returns the complete shell of ρ relative to η w.r.t. f . Algorithm 1 follows Theorem 1 and collects in X , for each dimension i and element of ρ , any element $z \in A$ whose f image is dominated by elements in η (lines 2-4). Then, the starting input abstraction ρ is joined with the new elements collected in X and since the so obtained result may not be a closure, the Moore closure operator is applied (line 5). Finally, the complete shell of ρ relative to η is returned at line 6.

We reported the procedure only to improve the understanding of the complete shell since, unfortunately, this is not an effective (decidable) procedure. Indeed, Algorithm 1 might diverge (at lines 2-4) when A is an infinite set or η is a not finite closure.

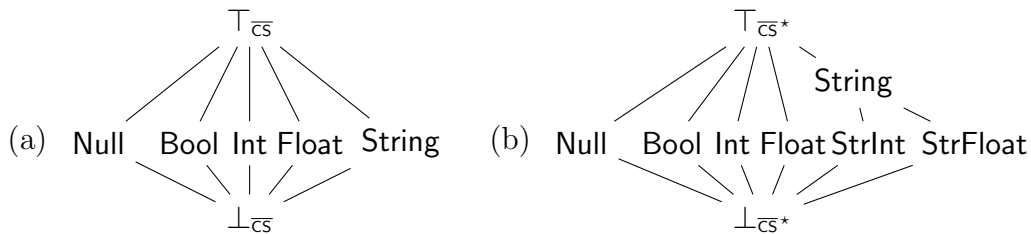


Figure 2: (a) *Coalesced Sum* abstract domain for PHP. (b) Complete shell of *Coalesced Sum* abstract domain w.r.t. the sum operation

3.1. Motivating example

A common feature of dynamic languages, such as PHP or JavaScript, is to be weakly typed. Hence, in those languages, it is allowed to change the variable type through the program execution. For example, in PHP, it is completely legal to write fragments such as `$\$x=1; \$x=true;$` , where the type of the variable `x` changes from integer to boolean. The first attempt for statically reasoning about variable types was adopting the so-called *Coalesced Sum* abstract domain (\overline{CS}) [5, 30], to detect whether a certain variable has a constant type through the whole program execution. In Figure 2a, we report the *Coalesced Sum* abstract domain for an intra-procedural version of PHP [5], that tracks null, boolean, integer, float and string types². Consider the formal semantics of the sum operation in PHP [20]. When one of the operands is a string, since the sum operation is feasible only between numbers, *implicit type conversion* occurs and converts the operand string to a number. In particular, if the prefix of the string is a number, it is converted to the maximum prefix of the string corresponding to a number, otherwise it is converted to 0. For example, the expression $e = "2.4hello" + "4"$ returns 6.4. Let $\alpha_{\overline{CS}}$ and $+\overline{CS}$ be the abstraction function and the abstract sum operation on the *Coalesced Sum* abstract domain respectively. The type of the expression e is given by: $\alpha_{\overline{CS}}(\{"2.4hello"\}) +_{\overline{CS}} \alpha_{\overline{CS}}(\{"4"\}) = \mathbf{String} +_{\overline{CS}} \mathbf{String} = \top_{\overline{CS}}$. The static type analysis based on the *Coalesced Sum* abstract domain returns $\top_{\overline{CS}}$ (i.e., any possible value), since the sum between two strings may

²By closing the *Coalesced Sum* abstract domain with the powerset operation we get a more precise domain called *union type* abstract domain [30], that tracks the set of types of a certain variable during program execution.

return either an integer or a float value. Precisely, the *Coalesced Sum* abstract domain is not complete w.r.t. the PHP sum operation, since for any string σ and σ' , it does not meet the completeness condition, i.e., $\alpha_{\overline{\text{CS}}}(\sigma + \sigma') = \alpha_{\overline{\text{CS}}}(\sigma) +_{\overline{\text{CS}}} \alpha_{\overline{\text{CS}}}(\sigma')$. Indeed, $\alpha_{\overline{\text{CS}}}(\{"2.4\text{hello}" + "4"}) = \text{Float}$ is different from $\alpha_{\overline{\text{CS}}}(\{"2.4\text{hello}"}) +_{\overline{\text{CS}}} \alpha_{\overline{\text{CS}}}(\{"4"}) = \top_{\overline{\text{CS}}}$.

Intuitively, the *Coalesced Sum* abstract domain is not complete w.r.t. the sum operation due to the loss of precision that occurs during the abstraction process of the inputs. Indeed, the domain is not precise enough to distinguish between strings that may be implicitly converted to integers or floats.

Figure 2b shows the complete shell of the *Coalesced Sum* abstract domain w.r.t. the sum ($\overline{\text{CS}}^*$). The latter adds two abstract values to the original domain, namely **StrFloat** and **StrInt**, that correspond to the abstractions of the strings that may be implicitly converted to floats and to integers, respectively. Notice that the type analysis on the new abstract domain is now complete w.r.t. the sum operation, e.g., $\alpha_{\overline{\text{CS}}^*}(\{"2.4\text{hello}" + "4"}) = \text{Float}$ that is equal to $\alpha_{\overline{\text{CS}}^*}(\{"2.4\text{hello}"}) +_{\overline{\text{CS}}^*} \alpha_{\overline{\text{CS}}^*}(\{"4"}) = \text{StrFloat} +_{\overline{\text{CS}}^*} \text{StrInt} = \text{Float}$.

As pointed out above, guaranteeing completeness in Abstract Interpretation is a desirable property that an abstract domain should aim to, since it ensures that no loss of precision occurs during the input abstraction process of the operation of interest. It is worth noting that *guessing* a complete abstract domain for a certain operation becomes particularly hard when the operation has tricky semantics, such as in our example or, more in general, in dynamic languages operations. For this reason, complete shells become important since they can mathematically guarantee completeness for a certain operation, starting from an abstract domain of interest.

4. The core dynamic language μDyn

We define μDyn , an imperative toy language expressive enough to handle some interesting behaviors related to strings in dynamic languages, e.g., implicit type conversion, inspired by the JavaScript programming language [35]. μDyn syntax is reported in Figure 3. The μDyn basic values are represented by the set

$$\text{VAL} = \text{INT} \cup \text{FLOAT} \cup \text{BOOL} \cup \text{STR}$$

```

a ::= n | x | a + a | a - a | a * a | a \ a
    | toNum(s) | length(s)
b ::= x | true | false | b && b | b || b | ! b
s ::= x | "s" | concat(s1,s2)
e ::= x | a | b | s
bl ::= { } | { S }
S ::= x = e; | if (b) bl1 else bl2 | while (b) bl | S1 S2

```

where $x \in \text{ID}$, $s \in \Sigma^*$, $n \in \text{INT} \cup \text{FLOAT}$

Figure 3: μDyn syntax

such that $\text{INT} = \mathbb{Z}$ denotes the set of signed integers, FLOAT denotes the set of signed decimal numbers³, $\text{BOOL} = \{\text{true}, \text{false}\}$ denotes the set of booleans, and $\text{STR} = \Sigma^*$ denotes the set of strings over an alphabet Σ . Then, we consider Σ^* composed of two sets, $\Sigma^* = \Sigma_{\text{Num}}^* \cup \Sigma_{\text{NotNum}}^*$, where Σ_{Num}^* is the set of numeric strings (e.g., "42", "-7.2") Σ_{NotNum}^* is the set of non numeric strings (e.g., "foo", "-2a"). Moreover, we consider Σ_{Num}^* additionally composed of four sets:

$$\Sigma_{\text{Num}}^* = \Sigma_{\text{UInt}}^* \cup \Sigma_{\text{UFloat}}^* \cup \Sigma_{\text{SInt}}^* \cup \Sigma_{\text{SFloat}}^*$$

From left to right, they correspond to the set of unsigned integer strings, unsigned float strings, signed integer strings, and signed float strings, respectively.

μDyn programs are elements generated by S syntax rules. Program states $\text{STATE} = \{\xi \mid \xi : \text{ID} \rightarrow \text{VAL}\}$ are maps from identifiers to values. The concrete semantics of μDyn statements follows⁴ [5], and it is given by the function $\llbracket \cdot \rrbracket : \text{STMT} \times \text{STATE} \rightarrow \text{STATE}$, inductively defined on the structure of the statements, as reported in Figure 4. We abuse notation in defining the

³In general, floats are represented in programming languages in the IEEE 754 double precision format. For the sake of simplicity, we use instead decimal numbers.

⁴Note that in this paper the semantics of the operations is detailed only for the ones that are relevant for string analysis. The complete concrete semantics of μDyn can be found in [5].

$$\begin{aligned}
\llbracket x = e; \rrbracket \xi &= \xi[x \leftarrow \llbracket e \rrbracket \xi] \\
\llbracket \text{if } (b) \text{ bl}_1 \text{ else bl}_2 \rrbracket \xi &= \begin{cases} \llbracket \text{bl}_1 \rrbracket \xi & \llbracket b \rrbracket \xi = \text{true} \\ \llbracket \text{bl}_2 \rrbracket \xi & \llbracket b \rrbracket \xi = \text{false} \end{cases} \\
\llbracket \text{while } (b) \text{ bl} \rrbracket \xi &= \llbracket \text{if } (b) \{ \text{bl while } (b) \text{ bl} \} \text{ else } \{ \} \rrbracket \xi \\
\llbracket \{ \} \rrbracket \xi &= \xi \quad \llbracket \{ S \} \rrbracket \xi = \llbracket S \rrbracket \xi \\
\llbracket S_1 S_2 \rrbracket \xi &= \llbracket S_2 \rrbracket (\llbracket S_1 \rrbracket \xi) \\
\llbracket \text{concat}(s, s') \rrbracket \xi &= \llbracket s \rrbracket \xi \cdot \llbracket s' \rrbracket \xi \quad \llbracket \text{length}(s) \rrbracket \xi = |\llbracket s \rrbracket \xi| \\
\llbracket \text{toNum}(s) \rrbracket \xi &= \begin{cases} \mathcal{N}(\llbracket s \rrbracket \xi) & \llbracket s \rrbracket \xi \in \Sigma_{\text{Num}}^* \\ 0 & \text{otherwise} \end{cases}
\end{aligned}$$

Figure 4: μDyn semantics

concrete semantics of expressions: $\llbracket \cdot \rrbracket : \text{EXP} \times \text{STATE} \rightarrow \text{VAL}$. Figure 4 shows the formal semantics of two relevant expressions involving strings we focus on: `concat`, that concatenates two strings, and string-to-number operation, namely `toNum`, that takes a string as input and returns the number that it represents if the input string corresponds to a numerical strings, 0 otherwise. Given $\sigma \in \text{STR}$, we denote by $\mathcal{N}(\sigma) \in \text{INT} \cup \text{FLOAT}$ the numeric value of a given string. For example, `toNum("4.2") = 4.2` and `toNum("asd") = 0`.

5. Making JavaScript string abstract domains complete

In this section, we study the completeness of two string abstract domains integrated into two state-of-the-art JavaScript static analysers based on Abstract Interpretation, namely TAJs [27] and SAFE [31]. Both the abstract domains track important information on JavaScript strings, e.g., TAJs can infer when a string corresponds to an unsigned integer, that may be used as array index, and SAFE tracks numeric strings, such as "2.5" or "+5".

For the sake of readability, we recast the original string abstract domains for μDyn , following the notation adopted in [4]. Figure 5 depicts them. Notice that the original abstract domain part of SAFE analyser treats the string "NaN" as a numeric string. Since our core language does not provide

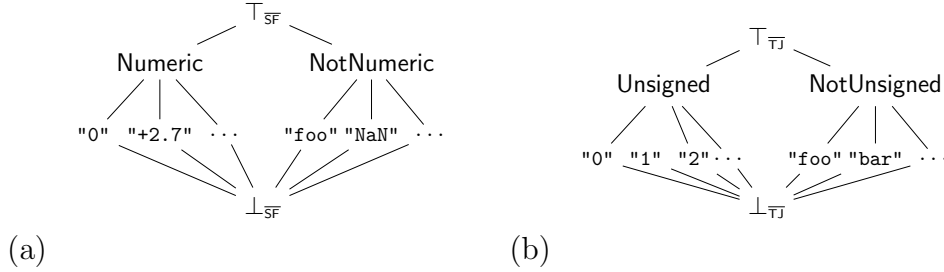


Figure 5: (a) SAFE, (b) TAJs string abstract domains recasted for μ Dyn

the primitive value `NaN`, the corresponding string, i.e., `"NaN"`, has no particular meaning here, and it is treated as a non-numerical string.

For each string abstract domain \overline{D} , we denote by $\alpha_{\overline{D}} : \mathcal{P}(\Sigma^*) \rightarrow \overline{D}$ its abstraction function, by $\gamma_{\overline{D}} : \overline{D} \rightarrow \mathcal{P}(\Sigma^*)$ its concretization function, and by $\rho_{\overline{D}} : \mathcal{P}(\Sigma^*) \rightarrow \mathcal{P}(\Sigma^*) \in uco(\overline{D})$ the associated upper closure operator.

5.1. Completing TAJs string abstract domain

Figure 5b depicts the string abstract domain \overline{TJ} , the recasted version of the domain integrated into the TAJs static analyser [27]. \overline{TJ} splits the strings into **Unsigned**, that denotes the strings corresponding to unsigned numbers, and **NotUnsigned**, any other string. Hence, for example, $\alpha_{\overline{TJ}}(\{"9", "+9"\}) = \top_{\overline{TJ}}$ and $\alpha_{\overline{TJ}}(\{"9.2", "foo"\}) = \text{NotUnsigned}$. Before reaching these abstract values, \overline{TJ} precisely tracks each string.

Here we focus on the `toNum` (i.e., string-to-number) operation. Since this operation clearly involves numbers, in Figure 6 we report the TAJs numerical abstract domain, denoted by \overline{TJ}_N . The latter domain behaves similarly to \overline{TJ} , distinguishing between unsigned and not unsigned integers.

Below we define the abstract semantics of the string-to-number operation for \overline{TJ} . In particular, we define the function:

$$\llbracket \text{toNum}(\bullet) \rrbracket^{\overline{TJ}} : \overline{TJ} \rightarrow \overline{TJ}_N$$

that takes as input a string abstract value in \overline{TJ} , and returns an integer abstract value in \overline{TJ}_N .

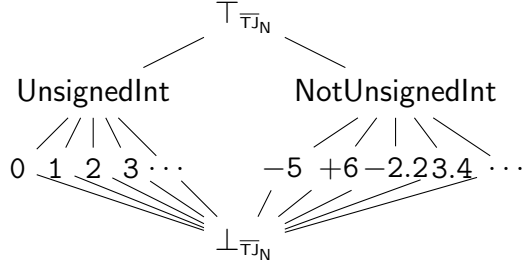


Figure 6: TAJN numerical abstract domain

$$\llbracket \text{toNum}(s) \rrbracket^{\overline{TJ}} = \begin{cases} \perp_{\overline{TJ}_N} & \llbracket s \rrbracket^{\overline{TJ}} = \perp_{\overline{TJ}} \\ \llbracket \text{toNum}(\sigma) \rrbracket & \llbracket s \rrbracket^{\overline{TJ}} = \sigma \in \Sigma^* \\ \text{UnsignedInt} & \llbracket s \rrbracket^{\overline{TJ}} = \text{Unsigned} \\ \top_{\overline{TJ}_N} & \llbracket s \rrbracket^{\overline{TJ}} = \text{NotUnsigned} \vee \llbracket s \rrbracket^{\overline{TJ}} = \top_{\overline{TJ}} \end{cases}$$

When the input evaluates to $\perp_{\overline{TJ}}$, the bottom is propagated, and $\perp_{\overline{TJ}_N}$ is returned (first row). If the input evaluates to a single string value, i.e., a concrete string, the abstract semantics relies on its concrete one (second row), as single strings are precisely captured by \overline{TJ} . When the input evaluates to the string abstract value **Unsigned**, the integer abstract value **UnsignedInt** is returned (third row). Indeed, in the concrete scenario, an unsigned and not float numeric string, which exactly represents the strings approximated by **Unsigned** in \overline{TJ} , is converted into the correspondent numeric value by the string-to-number operation. Therefore, the abstraction in \overline{TJ}_N of the numeric value of all the strings approximated by **Unsigned** is **UnsignedInt**. Finally, when the input evaluates to **NotUnsigned** or $\top_{\overline{TJ}}$, the top abstract value $\top_{\overline{TJ}_N}$ is returned (fourth row). In the second case, i.e., when the input is evaluated to $\top_{\overline{TJ}}$, it is trivial to notice that the result of the abstract string-to-number operation represents the best correct approximation. However, in the first case, it is not straightforward. Note that **NotUnsigned** approximates not numeric strings and signed and/or float numeric strings. Thus the only safe abstraction of the numeric value of strings approximated by **NotUnsigned** is $\top_{\overline{TJ}_N}$.

Lemma 1. The \overline{TJ} abstract domain is not complete with respect to the

`toNum` operation. In particular⁵, $\exists S \in \mathcal{P}(\Sigma^*)$ such that:

$$\alpha_{\overline{TJ}_N}(\llbracket \text{toNum}(S) \rrbracket) \not\sqsubseteq \llbracket \text{toNum}(\alpha_{\overline{TJ}}(S)) \rrbracket^{\overline{TJ}}$$

Proof of Lemma 1. As a counterexample to completeness, consider the set $S = \{ "2.3", "3.4" \}$. We can show that $\alpha_{\overline{TJ}_N}(\llbracket \text{toNum}(S) \rrbracket) \neq \llbracket \text{toNum}(\alpha_{\overline{TJ}}(S)) \rrbracket^{\overline{TJ}}$. Indeed,

- $\alpha_{\overline{TJ}_N}(\llbracket \text{toNum}(S) \rrbracket) = \text{NotUnsignedInt}$
- $\llbracket \text{toNum}(\alpha_{\overline{TJ}}(S)) \rrbracket^{\overline{TJ}} = \llbracket \text{toNum}(\text{NotUnsigned}) \rrbracket^{\overline{TJ}} = \top_{\overline{TJ}_N}$.

□

The completeness condition does not hold because the \overline{TJ} string abstract domain loses too much information during the abstraction process, and the latter information cannot be retrieved during the abstract `toNum` operation. In particular, when non-numeric strings and unsigned integer strings are converted to numbers by `toNum`, they are mapped to the same value, namely 0. Indeed, \overline{TJ} does not differentiate between non-numeric and unsigned integer string values, and this is the principal cause of the \overline{TJ} incompleteness w.r.t. `toNum`. Additionally, more precision can be obtained if we could differentiate numeric strings holding float numbers from those holding integer numbers. Thus, to make \overline{TJ} complete w.r.t. `toNum`, we have to derive the complete shell of the \overline{TJ} string abstract domain relative to the \overline{TJ}_N numerical abstract domain, by applying Theorem 1.

Definition 6 (Complete shell of \overline{TJ}). Let $\rho_{\overline{TJ}} \in \text{uco}(\mathcal{P}(\Sigma^*))$ and $\rho_{\overline{TJ}_N} \in \text{uco}(\mathcal{P}(\text{INT} \cup \text{FLOAT}))$ be the upper closure operators related to \overline{TJ} and \overline{TJ}_N abstract domains respectively. Given $\llbracket \text{toNum}(\bullet) \rrbracket^{\overline{TJ}} : \overline{TJ} \rightarrow \overline{TJ}_N$, we define by \overline{TJ}^* the transformer $\mathcal{S}_{\text{toNum}}^{\rho_{\overline{TJ}_N}}(\rho_{\overline{TJ}})$.

By Definition 4, \overline{TJ}^* is the complete shell of $\rho_{\overline{TJ}}$ relative to $\rho_{\overline{TJ}_N}$ with respect to the `toNum` operation. As already argued in Section 5, to compute \overline{TJ}^* we use the constructive characterization given by Theorem 1.

Figure 7 depicts \overline{TJ}^* . In particular, the abstract points inside dashed boxes are the abstract values added during the computation of \overline{TJ}^* , the

⁵We abuse notation denoting with $\llbracket \cdot \rrbracket$ the additive lift to set of basic values of the concrete semantics, i.e., the collecting semantics.

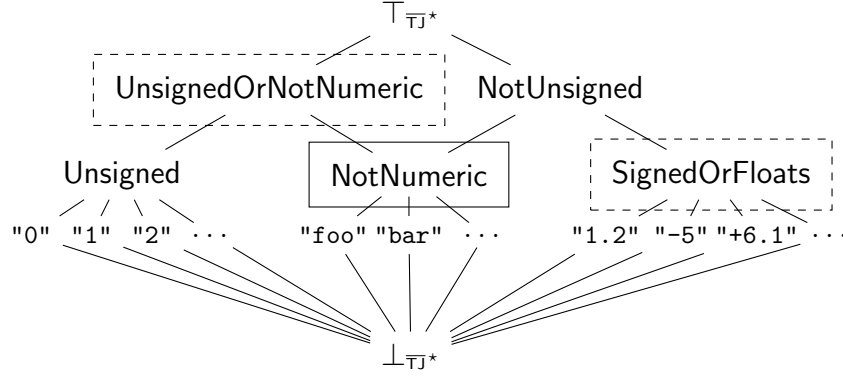


Figure 7: Complete shell of $\rho_{\overline{TJ}}$ relative to $\rho_{\overline{TJ}_N}$ w.r.t. toNum

point inside the solid box is instead obtained by the Moore closure of the other points of the domain, while the remaining abstract values were already in \overline{TJ} . A non-intuitive point added in \overline{TJ}^* is **SignedOrFloats**, namely the abstract value s.t. its concretization contains any float string and the signed integers. This abstract point is added during the computation of \overline{TJ}^* . In particular, following Algorithm 1, instantiated with $A = \mathcal{P}(\Sigma^*)$ and $\eta = \rho_{\overline{TJ}_N}$, this abstract point is computed at lines 2-4 at the iteration when y is $\gamma_{\overline{TJ}_N}(\text{NotUnsignedInt})$, namely

$$\text{SignedOrFloats} \in \max^{\subseteq}(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{toNum}(Z) \rrbracket \subseteq \gamma_{\overline{TJ}_N}(\text{NotUnsignedInt})\})$$

Informally: which is the maximal set of strings Z s.t. $\text{toNum}(Z)$ is dominated by **NotUnsignedInt**? To obtain from $\text{toNum}(Z)$ only values dominated by **NotUnsignedInt**, the maximal set doing so is exactly the set of the float strings and the signed strings. Other strings, such that: unsigned integer strings or not numerical strings are excluded, since they are both converted to unsigned integers, and they would violate the dominance relation. Similarly, the abstract point **UnsignedOrNotNumeric** is added to the complete shell \overline{TJ}^* . Following again Algorithm 1, instantiated with $A = \mathcal{P}(\Sigma^*)$ and $\eta = \rho_{\overline{TJ}_N}$, the above abstract element is computed at lines 2-4 at the iteration when y is $\gamma_{\overline{TJ}_N}(\text{UnsignedInt})$, namely

$$\text{UnsignedOrNotNumeric} \in \max^{\subseteq}(\{Z \in \mathcal{P}(\Sigma^*) \mid \text{toNum}(Z) \subseteq \gamma_{\overline{TJ}_N}(\text{UnsignedInt})\})$$

To obtain from $\text{toNum}(Z)$ only values dominated by **UnsignedInt**, the maximal set doing so is exactly the set of the unsigned integer strings and the non-

numerical strings, since the latter are converted to 0. Any other string, such as a signed or float string, would violate the dominance relation.

The last point added in the complete shell is **NotNumeric**. In particular, this abstract point is added by the Moore closure between **Unsigned** and **UnsignedOrNotNumeric**. The remaining abstract points, namely **Unsigned** and **NotUnsigned**, were already present in the original TAJIS string abstract domain, and they are leaved unaltered in its complete shell w.r.t. **toNum**.

Example 1. Consider again the string set $S = \{"2.3", "3.4"\}$ of Example 5.1. We can show that in \overline{TJ}^* , $\alpha_{\overline{TJ}_N}(\llbracket \text{toNum}(S) \rrbracket) = \llbracket \text{toNum}(\alpha_{\overline{TJ}^*}(S)) \rrbracket^{\overline{TJ}^*}$. Indeed, $\alpha_{\overline{TJ}_N}(\llbracket \text{toNum}(S) \rrbracket) = \text{NotUnsignedInt}$ and,

$$\begin{aligned} \llbracket \text{toNum}(\alpha_{\overline{TJ}^*}(S)) \rrbracket^{\overline{TJ}^*} &= \llbracket \text{toNum}(\text{SignedOrFloats}) \rrbracket^{\overline{TJ}^*} \\ &= \text{NotUnsignedInt} \end{aligned}$$

Note that, the new abstract semantics $\llbracket \text{toNum}(\bullet) \rrbracket^{\overline{TJ}^*}$, handling the abstract points added by the complete shell, corresponds to the best correction approximation, namely $\alpha_{\overline{TJ}^*} \circ \llbracket \text{toNum}(\bullet) \rrbracket \circ \gamma_{\overline{TJ}^*} : \overline{TJ}^* \rightarrow \overline{TJ}_N$ (see Definition 2).

5.2. Completing SAFE string abstract domain

Figure 5a depicts the string abstract domain \overline{SF} , i.e., the recasted version of the domain involved into SAFE [31] static analyser. It splits strings into the abstract values: **Numeric** (i.e., numerical strings) and **NotNumeric** (i.e., all the other strings). As for \overline{TJ} , before reaching these abstract values, \overline{SF} precisely tracks single string values. For instance, $\alpha_{\overline{SF}}(\{"+9.6", "7"\}) = \text{Numeric}$, and $\alpha_{\overline{SF}}(\{"+9.6", "\text{bar}"\}) = \perp_{\overline{SF}}$.

We study the completeness of \overline{SF} w.r.t. **concat** operation. Figure 8 presents the abstract semantics of the concatenation operation for \overline{SF} , that is:

$$\llbracket \text{concat}(\bullet, \bullet) \rrbracket^{\overline{SF}} : \overline{SF} \times \overline{SF} \rightarrow \overline{SF}$$

In particular, when both abstract values correspond to single strings, the standard string concatenation is applied (second row, second column). In the case in which one abstract value, involved in the concatenation, is a string and the other is **Numeric** (third row, second column and second row, third column) we distinguish two cases: if the string is empty or corresponds to an unsigned integer we can safely return **Numeric**, otherwise **NotNumeric** is returned. This

		$\llbracket \text{concat}(s_1, s_2) \rrbracket^{\overline{\text{SF}}}$				
$\llbracket s_1 \rrbracket^{\overline{\text{SF}}}$ \diagdown	$\llbracket s_2 \rrbracket^{\overline{\text{SF}}}$	$\perp_{\overline{\text{SF}}}$	$\sigma_2 \in \Sigma^*$	Numeric	NotNumeric	$\top_{\overline{\text{SF}}}$
$\perp_{\overline{\text{SF}}}$	$\perp_{\overline{\text{SF}}}$	$\perp_{\overline{\text{SF}}}$		$\perp_{\overline{\text{SF}}}$	$\perp_{\overline{\text{SF}}}$	$\perp_{\overline{\text{SF}}}$
$\sigma_1 \in \Sigma^*$	$\perp_{\overline{\text{SF}}}$	$\sigma_1 \cdot \sigma_2$	$\begin{cases} \text{Numeric} & \sigma_1 = "" \text{ or} \\ & \sigma_1 \in \Sigma_{\text{UInt}}^* \\ \text{NotNumeric} & \text{otherwise} \end{cases}$		$\top_{\overline{\text{SF}}}$	$\top_{\overline{\text{SF}}}$
Numeric	$\perp_{\overline{\text{SF}}}$	$\begin{cases} \text{Numeric} & \sigma_2 = "" \text{ or} \\ & \sigma_2 \in \Sigma_{\text{UInt}}^* \\ \text{NotNumeric} & \text{otherwise} \end{cases}$			$\top_{\overline{\text{SF}}}$	$\top_{\overline{\text{SF}}}$
NotNumeric	$\perp_{\overline{\text{SF}}}$	$\top_{\overline{\text{SF}}}$			$\top_{\overline{\text{SF}}}$	NotNumeric
$\top_{\overline{\text{SF}}}$	$\perp_{\overline{\text{SF}}}$	$\top_{\overline{\text{SF}}}$			$\top_{\overline{\text{SF}}}$	$\top_{\overline{\text{SF}}}$

Figure 8: SAFE concat abstract semantics

happens because, when two float strings (hence numerical strings) are concatenated, a non-numerical string is returned (e.g., $\text{concat}("1.1", "2.2") = "1.12.2"$). For the same reason, when both input abstract values are **Numeric**, the result is not guaranteed to be numerical. Indeed, $\llbracket \text{concat}(\text{Numeric}, \text{Numeric}) \rrbracket^{\overline{\text{SF}}} = \top_{\overline{\text{SF}}}$.

Lemma 2. The $\overline{\text{SF}}$ abstract domain is not complete with respect to the `concat` operation. In particular, $\exists S_1, S_2 \in \mathcal{P}(\Sigma^*)$ such that:

$$\alpha_{\overline{\text{SF}}}(\llbracket \text{concat}(S_1, S_2) \rrbracket) \not\sqsubseteq \llbracket \text{concat}(\alpha_{\overline{\text{SF}}}(S_1), \alpha_{\overline{\text{SF}}}(S_2)) \rrbracket^{\overline{\text{SF}}}$$

Proof of Lemma 2. As a counterexample to completeness, consider the sets $S_1 = \{"2.2", "2.3"\}$ and $S_2 = \{"2", "3"\}$. Then, in $\overline{\text{SF}}$, $\alpha_{\overline{\text{SF}}}(\llbracket \text{concat}(S_1, S_2) \rrbracket) \neq \llbracket \text{concat}(\alpha_{\overline{\text{SF}}}(S_1), \alpha_{\overline{\text{SF}}}(S_2)) \rrbracket^{\overline{\text{SF}}}$. Indeed, $\alpha_{\overline{\text{SF}}}(\llbracket \text{concat}(S_1, S_2) \rrbracket) = \text{Numeric}$ and,

$$\llbracket \text{concat}(\alpha_{\overline{\text{SF}}}(S_1), \alpha_{\overline{\text{SF}}}(S_2)) \rrbracket^{\overline{\text{SF}}} = \llbracket \text{concat}(\text{Numeric}, \text{Numeric}) \rrbracket^{\overline{\text{SF}}} = \top_{\overline{\text{SF}}}$$

□

The $\overline{\text{SF}}$ abstract domain loses too much information during the abstraction process, which can not be retrieved during the abstract concatenation. Intuitively, to gain completeness w.r.t. `concat` operation, $\overline{\text{SF}}$ should improve the precision of the numerical strings abstraction, e.g., discriminating between float and integer strings.

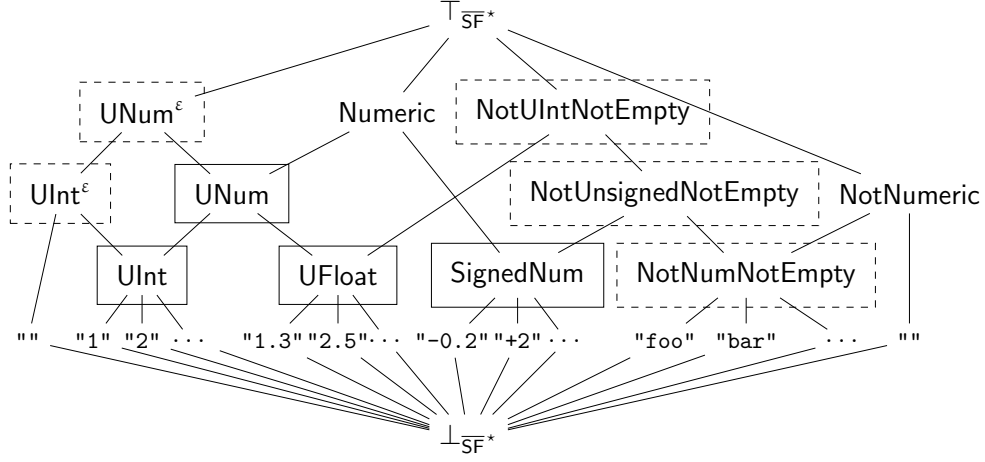


Figure 9: Absolute complete shell of $\rho_{\overline{SF}}$ w.r.t. `concat`

Definition 7 (Complete shell of \overline{SF}). Let $\rho_{\overline{SF}} \in uco(\mathcal{P}(\Sigma^*))$ be the upper closure operator related to \overline{SF} abstract domain. Given $\llbracket \text{concat}(\bullet, \bullet) \rrbracket^{\overline{SF}} : \overline{SF} \times \overline{SF} \rightarrow \overline{SF}$, we define by \overline{SF}^* the transformer $\overline{S}_{\text{concat}}(\rho_{\overline{SF}})$.

By Definition 5, \overline{SF}^* is the absolute complete shell of $\rho_{\overline{SF}}$ with respect to the `concat` operation.

By Theorem 2, the transformer $\overline{S}_{\text{concat}}(\rho_{\overline{SF}})$ is equal to the Moore closure of the union between \overline{SF} and the binary operator defined in Table 1. Table 1 depicts the first iteration of the fix-point computation of Theorem 2, where $\perp_{\overline{SF}}$ and $\top_{\overline{SF}}$ rows and columns are omitted for space limitations; note that Table 1 also corresponds to the whole fix-point result of Theorem 2, since the fix-point is reached at the next step. In particular, $\forall X \in \mathcal{P}(\Sigma^*). \max^{\subseteq}(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(Z, X) \rrbracket \subseteq \gamma_{\overline{SF}}(\perp_{\overline{SF}})\}) = \perp_{\overline{SF}}$, $\max^{\subseteq}(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(X, Z) \rrbracket \subseteq \gamma_{\overline{SF}}(\perp_{\overline{SF}})\}) = \perp_{\overline{SF}}$ and $\max^{\subseteq}(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(X, Z) \rrbracket \subseteq \gamma_{\overline{SF}}(\perp_{\overline{SF}})\}) = \top_{\overline{SF}}$, when $X \neq \perp_{\overline{SF}}$, $\max^{\subseteq}(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(Z, X) \rrbracket \subseteq \gamma_{\overline{SF}}(\perp_{\overline{SF}})\}) = \perp_{\overline{SF}}$, when $X \neq \top_{\overline{SF}}$. The complete shell \overline{SF}^* is depicted in Figure 9.⁶ In particular, the points inside dashed boxes are the abstract values added during the iterative computations of \overline{SF}^* , the points inside solid boxes are instead obtained by the Moore closure of the other points of the domain, while the remaining abstract values were

⁶The empty strings ("") at the most-left and most-right sides of Figure 9 are not distinct elements but they are only duplicated in order to declutter the figure.

X \ Y	Numeric	NotNumeric	"n" ∈ Z	"f" ∈ F	"s" ∈ NotNum
Numeric	{""} ∪ UInt	[NotNum \ {""}] ∪ NotUInt ∪ NotUFloat	⊥ _{SF}	⊥ _{SF}	⊥ _{SF}
NotNumeric	⊥ _{SF}	NotNumeric	⊥ _{SF}	⊥ _{SF}	⊥ _{SF}
"n" ∈ Z	{""} ∪ UInt ∪ UFloat	[NotNum \ {""}] ∪ NotUInt ∪ NotUFloat	⊥ _{SF} ∨ str	⊥ _{SF} ∨ str	⊥ _{SF} ∨ str
"f" ∈ F	{""} ∪ UInt	[NotNum \ {""}] ∪ Float ∪ NotUInt	⊥ _{SF}	⊥ _{SF} ∨ str	⊥ _{SF} ∨ str
"s" ∈ NotNum	⊥ _{SF}	NotNumeric	⊥ _{SF}	⊥ _{SF}	⊥ _{SF} ∨ str

Table 1: First iteration of $\max^{\subseteq}(\{Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(X, Z) \rrbracket \subseteq \gamma_{\overline{\text{SF}}}(\text{Y})\})$

already in $\overline{\text{SF}}$. The meaning of abstract values in $\overline{\text{SF}}^*$ is intuitive. In order to satisfy the completeness property, $\overline{\text{SF}}^*$ splits the **Numeric** abstract value, already taken into account in $\overline{\text{SF}}$, into all the strings corresponding to unsigned integer (**UInt**), unsigned floats (**UFloat**), and signed numbers (**SignedNum**). Moreover, particular importance is given to the empty string since the new abstract domain specifies whether each abstract value contains "". Indeed, the UInt^ε abstract value represents the strings corresponding to unsigned integer or the empty string. The UNum^ε abstract value represents the strings corresponding to unsigned numbers or the empty string. An unexpected abstract value considered in $\overline{\text{SF}}^*$ is **NotUnsignedNotEmpty**, such that:

$$\gamma_{\overline{\text{SF}}^*}(\text{NotUnsignedNotEmpty}) = \{\sigma \mid \sigma \in \Sigma_{\text{SInt}}^* \cup \Sigma_{\text{SFloat}}^* \cup (\Sigma_{\text{NotNum}}^* \setminus \{\text{""}\})\}$$

Namely, the abstract point whose concretization corresponds to the set of any non-numerical string, except the empty string, and any string corresponding to a signed number. This abstract point has been added to $\overline{\text{SF}}^*$, following Theorem 2. Since the absolute complete shell is the greatest fix-point of the relative one (that in our case is reached after one iteration), the corresponding procedure is the greatest fix-point of Algorithm 1, instantiated with $A = \mathcal{P}(\Sigma^*)$ and $\eta = \rho_{\overline{\text{SF}}}$. In particular, the abstract element is computed

at lines 2-4 at the iteration when y is $\gamma_{\overline{\text{SF}}}(\text{NotNumeric})$, x is $\gamma_{\overline{\text{SF}}}(\text{Numeric})$ and $i = 1$, namely

$$\text{NotUnsignedNotEmpty} \in \max^{\subseteq}(Z \in \mathcal{P}(\Sigma^*) \mid \llbracket \text{concat}(\gamma_{\overline{\text{SF}}}(\text{Numeric}), Z) \rrbracket \subseteq \gamma_{\overline{\text{SF}}}(\text{NotNumeric}))$$

Informally: which is the maximal set of strings s.t. concatenated to any possible numerical string will produce any possible non-numerical string? Indeed, to be sure to obtain non-numerical strings, the maximal set doing so is exactly the set of any non-numerical non-empty string, and any string corresponding to a signed number, that is **NotUnsignedNotEmpty**.

Example 2. Consider again the string sets $S_1 = \{"2.2", "2.3"\}$ and $S_2 = \{"2", "3"\}$ of Example 5.2. We can show that in $\overline{\text{SF}}^*$, $\alpha_{\overline{\text{SF}}}^*(\llbracket \text{concat}(S_1, S_2) \rrbracket) = \llbracket \text{concat}(\alpha_{\overline{\text{SF}}}^*(S_1), \alpha_{\overline{\text{SF}}}^*(S_2)) \rrbracket^{\overline{\text{SF}}^*}$. Indeed, $\alpha_{\overline{\text{SF}}}^*(\llbracket \text{concat}(S_1, S_2) \rrbracket) = \text{UFloat}$ and $\llbracket \text{concat}(\alpha_{\overline{\text{SF}}}^*(S_1), \alpha_{\overline{\text{SF}}}^*(S_2)) \rrbracket^{\overline{\text{SF}}^*} = \llbracket \text{concat}(\text{UFloat}, \text{UInt}) \rrbracket^{\overline{\text{SF}}^*} = \text{UFloat}$.

As in the TAJIS case, the new abstract semantics $\llbracket \text{concat}(\bullet, \bullet) \rrbracket^{\overline{\text{SF}}^*}$, handling the abstract points added by the complete shell, corresponds the best correct approximation, namely $\alpha_{\overline{\text{SF}}}^* \circ \llbracket \text{concat}(\bullet, \bullet) \rrbracket \circ \gamma_{\overline{\text{SF}}}^* : \overline{\text{SF}}^* \times \overline{\text{SF}}^* \rightarrow \overline{\text{SF}}^*$ (see Definition 2).

6. What do we gain from using a complete abstract domain?

Now, we discuss and evaluate the benefits of adopting the complete shells reported in Section 5 and, more in general, complete domains, w.r.t. a certain operation. In particular, we compare the μDyn versions of the string abstract domains adopted by SAFE and TAJIS with their corresponding complete shells. We discuss the complexity of the complete shells, and finally we argue how adopting complete abstract domains can be useful into static analysers.

Precision. In the previous section, we focused on the completeness of the string abstract domains integrated into SAFE and TAJIS, for μDyn , w.r.t. two string operations, namely `concat` and `toNum`, respectively. While string concatenation is common in any programming language, `toNum` assumes critical importance for dynamic languages, where implicit type conversion is provided. It is worth noting that, in this paper we have only considered `toNum` operation but the completion construction of other implicit type conversion operations (e.g., `toBool` or `toString`) is similar to the one reported in Section 5.1. Since type conversion is often hidden from the developer, aim


```

1 var obj = {
2     "foo" : 1,
3     "bar" : 2,
4     "1.2" : 3,
5     "2.2" : "hello"
6 }
7
8 y = obj[idx];

```

Figure 10: JavaScript object access example

to completeness of the analysis increases the precision of such operations. For instance, let x be a variable, at a certain program execution point. x may have concrete value in the set $S = \{\text{"foo"}, \text{"bar"}\}$. If S is abstracted into the original TAJs string abstract domain, its abstraction will correspond to **NotUnsigned**, losing the information about the fact that the concrete value of x surely does not contain numerical values. Hence, when the abstract value of S is used as input of `toNum`, the result will return $\top_{\overline{TJ}_N}$, i.e., any possible concrete integer value. Conversely, abstracting S in \overline{TJ}^* (the absolute complete shell of \overline{TJ} relative to `toNum` discussed in Section 5.1) leads to a more precise abstraction, since \overline{TJ}^* is able to differentiate between non-numerical and numerical strings. In particular, the abstract value of S in \overline{TJ}^* is **NotNumeric**, and $\llbracket \text{toNum}(\text{NotNumeric}) \rrbracket^{\overline{TJ}^*}$ will precisely return 0.

Adopting a complete shell w.r.t. a certain operation does not compromise the precision of the others. For example, the JavaScript fragment reported in Figure 10 and let us analyze it with the domain of TAJs. Suppose that the value of `idx` is the abstraction, in the starting TAJs string abstract domain, of the string set $S = \{\text{"foo"}, \text{"bar"}\}$, namely the abstract value **NotUnsigned**. The variable `idx` is used to access the property of the object `obj` at line 8 and, to guarantee soundness, it accesses *all* the properties of `obj`, including the fields "1.2" and "2.2", introducing noise in the abstract computation, since "1.2" and "2.2" are false positives values introduced by the abstraction of the values of `idx`. If we analyse the same JavaScript fragment with the absolute complete shell (w.r.t. `toNum` operation) of the TAJs string abstract domain defined in Section 5.1, we obtain more precise results. Indeed, in this case, the value of `idx` corresponds to the abstract value **NotNumeric**, and when it is used to access the object `obj` at line 8, only "foo" and "bar" are accessed, since they are the only non-numerical string properties of `obj`.

Qualitative evaluation of complete shells. We evaluate the complete shells we have provided in the previous section from a qualitative point of view. As usual in static analysis by Abstract Interpretation, there is a trade-off between precision and efficiency: choosing a more precise abstract domain may compromise the efficiency of the abstract computations. A representative example is reported in [24]: the complete shell of the sign abstract domain w.r.t. addition is the interval abstract domain. Hence, starting from a finite height abstract domain (signs) we obtain an infinite height abstract domain (intervals). In particular, fix-point computations on signs converge, while on intervals may diverge. Indeed, after the completion, the interval abstract domain should also be equipped with a widening operator [17] to still guarantee termination. A worst-case scenario is when the complete shells w.r.t. a certain operation exactly corresponds to the collecting abstract domain, i.e., the concrete domain. Clearly, we cannot use the concrete domain due to undecidability reasons, but this suggests us to change the starting abstract domain, since it is not able to track any information related to the operation of interest. An example is the suffix abstract domain [16] with `substring` operation: since this abstract domain tracks only the common suffix of a strings set, it can not track the information about the indexes of the common suffix, and the complete shell of the suffix abstract domain w.r.t. `substring` would lead to the concrete domain. Hence, if the focus of the abstract interpreter is to improve the precision of the `substring` operation, we should change the abstract domain with a more precise one for `substring`, such as the finite state automata [6] abstract domain.

Consider now the complete shells reported in Section 5. The obtained complete shells still have finite height, hence termination is still guaranteed without the need to equip the complete shells with widening operators. Moreover, the complexity of the string operations of interest is preserved after completion. Indeed, in both TAJs and SAFE starting abstract domains, `concat` and `toNum` operations have constant complexity, respectively, and the same complexity is preserved in the corresponding complete shells. It is worth noting that also the complexity of the abstract domain-related operations, such as least upper bound, greatest lower bound and the ordering operator, is preserved in the complete shells. Hence, to conclude, as far as the complete shells we have reported for TAJs and SAFE are concerned, there is no worsening when we substitute the original string abstract domains with the corresponding complete shells, and this leads, as we have already mentioned before, to completeness during the input abstraction process w.r.t. the

relative operations, namely `concat` for SAFE and `toNum` for TAJIS.

False positives reduction. In static analysis, a certain degree of abstraction must be added in order to obtain decidable procedures to infer invariants on a generic program. Clearly, using less precise abstract domains lead to an increase of *false positive* values of the computed invariants. In particular, after a program is analysed, this burdens the phase of false positive detection: when a program is analysed, the following phase consists in detecting which values of the invariants derived by the static analyser are spurious, namely those values that certainly are not computed by the concrete execution of the program of interest.

In particular, using imprecise (i.e., not complete) abstract domains clearly increases the number of false positives in the abstract computation of the static analyser, burdening the next phase of detection of the spurious values. Conversely, adopting (backward) complete abstract domains w.r.t. a certain operation reduce the numbers of false positives introduced during the abstract computations, at least in the input abstraction process. Clearly, in this way, the next phase of detecting false positives will be lighter since less noise has been introduced during the abstract computation of the invariants. For example, let us consider the following simplified example, inspired by [41].

```
if (valid(in))
    safeOnlyIfValid(in);
else
    safeOnlyIfInvalid(in);
```

It is common, in web developing, to check and sanitize user inputs before performing some action in a secure way, e.g., for preventing XSS or injections attacks. For instance, we can suppose that the function `valid` in the above example checks whether the input string `in` does not contain the string `" < script"`: if so, it returns `true`, `false` otherwise. Hence, if `in` is valid, some safe action with valid inputs is performed (`true` branch), otherwise another safe action with non-valid input is performed (`false` branch). Let us suppose that the value of `in` may be one of the strings set $S = \{"+1", "+2"\}$. Note that, in any case, the program execution always lead to the execution of the `true` branch, since all the possible strings of `in` are valid, i.e., do not contain the string `" < script"`. The abstract value of `in`, using the starting TAJIS abstract domain, is `NotUnsigned`. Hence, the `if`-guard cannot be statically determined, since a `NotUnsigned` string may be contain or not the

string `" < script"`. Consequently, both `if` branches must be taken into account, leading to false alarms both in the `true` and `false` branches. Using the complete shell of TAJIS w.r.t. `toNum` instead leads to a more precise abstract value for `in`, namely `SignedOrFloat`. In this case, the `if` boolean guard is satisfied (since any string approximated by `SignedOrFloat` does not contain the string `" < script"`) and no false alarms are raised during the execution of the program.

7. Relative Precision

In this section, we recall the notion of pseudo-distance on abstract domains, firstly defined in [33], which takes into account the order relation between abstract elements (together with the fact that different abstract elements might approximate the same set of concrete values) and their possible incomparability. Then we formalise the increment of the precision obtained when analysing a program with a complete abstract domain with respect to the original version of the domain.

7.1. Abstract domains precision: an overview

It is well known that abstract domains precision can be qualitatively compared by exploiting the information they are able to capture [13, 14]. However, quantitatively evaluating the precision of abstract domains has been proven to be quite challenging and frequently resulted in ad hoc measures.

Di Pierro and Wiklicky [40], introduced the notion of probabilistic Abstract Interpretation, used to numerically estimate the incompleteness of numerical abstract domains. Sotin [44] presented the notion of precision of a numerical abstract value, measuring the volume it describes, with the purpose of quantitatively comparing the precision of numerical abstract domains. Finally, Casso et al. [9] wanted to compare the precision of different analyses on logic programs. Thus, they proposed distances in two abstract domains used in constraint logic programming and they extended them to distances between the results of different analyses of a given program.

Logozzo et al. [33] gave a more general definition of pseudo-distance between abstract domain elements that allows quantifying the error of approximating a concrete element in an abstract domain.

Definition 8 (Pseudo-distance [33]). Let $\langle \overline{\mathbb{D}}, \sqsubseteq_{\overline{\mathbb{D}}}, \perp_{\overline{\mathbb{D}}}, \top_{\overline{\mathbb{D}}}, \sqcap_{\overline{\mathbb{D}}}, \sqcup_{\overline{\mathbb{D}}} \rangle$ be an abstract domain. A function $\delta : \overline{\mathbb{D}} \times \overline{\mathbb{D}} \rightarrow \mathbb{R} \cup \{+\infty\}$ is a *pseudo-distance* over

the abstract domain \bar{D} if and only if, for any $\bar{d}, \bar{d}', \bar{d}'' \in \bar{D}$, it satisfies the following axioms:

- *non-negativity*: $\delta(\bar{d}, \bar{d}') \geq 0$
- *if-identity*: $\bar{d} =_{\bar{D}} \bar{d}' \Rightarrow \delta(\bar{d}, \bar{d}') = 0$
- *symmetry*: $\delta(\bar{d}, \bar{d}') = \delta(\bar{d}', \bar{d})$
- *weak triangle inequality*: $\bar{d} \sqsubseteq_{\bar{D}} \bar{d}'' \sqsubseteq_{\bar{D}} \bar{d}' \Rightarrow \delta(\bar{d}, \bar{d}'') \leq \delta(\bar{d}, \bar{d}') + \delta(\bar{d}', \bar{d}'')$

We recall the path length distance definition.

Definition 9 (Path length distance [33]). Let $\langle \bar{D}, \sqsubseteq_{\bar{D}}, \perp_{\bar{D}}, \top_{\bar{D}}, \sqcap_{\bar{D}}, \sqcup_{\bar{D}} \rangle$ be an abstract domain. The path length distance $\delta_{plen} : \bar{D} \times \bar{D} \rightarrow \mathbb{R} \cup \{+\infty\}$ is defined as

$$\delta_{plen}(\bar{d}, \bar{d}') = \begin{cases} plen(\bar{d}, \bar{d}') & \bar{d} \sqsubseteq_{\bar{D}} \bar{d}' \\ plen(\bar{d}', \bar{d}) & \bar{d}' \sqsubseteq_{\bar{D}} \bar{d} \\ +\infty & \text{otherwise} \end{cases}$$

where $plen(\bar{d}, \bar{d}')$ computes the distance between two abstract values $\bar{d}, \bar{d}' \in \bar{D}$, when they are comparable, and it is defined as follows.⁷

$$plen(\bar{d}, \bar{d}') = \min\{n \in \mathbb{N} \mid \{\bar{d}_0, \bar{d}_1, \dots, \bar{d}_n\} \in \mathcal{P}(\bar{D}), \bar{d}_0 = \bar{d}, \bar{d}_n = \bar{d}', \\ \forall i \in [0, n). \bar{d}_i \sqsubseteq_{\bar{D}} \bar{d}_{i+1} \wedge \nexists \bar{d}'' \in \bar{D} \setminus \{\bar{d}_i, \bar{d}_{i+1}\}. \bar{d}_i \sqsubseteq_{\bar{D}} \bar{d}'' \sqsubseteq_{\bar{D}} \bar{d}_{i+1}\}$$

Observe that, if $\bar{d} \sqsubseteq_{\bar{D}} \bar{d}'$, $plen(\bar{d}, \bar{d}')$ corresponds to the length of the minimal path from \bar{d} to \bar{d}' in the Hasse diagram representation of \bar{D} w.r.t. $\sqsubseteq_{\bar{D}}$. For instance, let us consider the **Sign** abstract domain reported in Figure 1 and suppose to compute the path length distance between \perp_{Sign} and \top_{Sign} . In the abstract domain, two paths exist between \perp_{Sign} and \top_{Sign} , namely $\{\perp_{\text{Sign}}, \text{NotPos}, \top_{\text{Sign}}\}$ and $\{\perp_{\text{Sign}}, \text{NotNeg}, \top_{\text{Sign}}\}$. Following the definition of $plen$ reported above, we have to take the minimum path length of the paths set, hence $plen(\perp_{\text{Sign}}, \top_{\text{Sign}}) = 2$ (path length distance starts counting from 0).

⁷Notice that the original definition of $plen$ in [33] contains a typo, as it refers to the partial order in a chain by disregarding its transitive closure.

```

1 nums = "";
2 notnums = "";
3 i = 0;
4 while (i < length(s)) {
5   if (toNum(charAt(str, i)) == 0) {
6     notnums = concat(notnums, charAt(s, i));7
7   } else {
8     nums = concat(nums, charAt(s, i));9
9   }
10  i = i + 1;11
11 }12

```

Figure 11: Example of μ Dyn annotated program

7.2. Measuring precision gained by complete shells

We aim at computing the distance and, in turn, the increment of precision, between abstract points of an abstract domain \bar{D} and its complete shell, w.r.t. a certain operation of interest f , noting that the abstract values of \bar{D} are all contained in its complete shell. In order not to clutter the notation, given an abstract domain \bar{D} , we denote by $\text{Shell}_f(\bar{D})$ its complete shell w.r.t the function f .

Let P be a μ Dyn program and let denote by $\text{Lab}(P)$ the program points of P . We denote by $l_i \in \text{Lab}(P)$ its i -th program point and by $\text{Vars}(P)$ the program variables. An example of program annotated with its program points is reported in Figure 11. Given such a program, as usual in static program analysis, the goal is to compute the abstract values associated to each variable $x \in \text{Vars}(P)$ at each program point $l \in \text{Lab}(P)$. We denote by $\xi^{\bar{D}} : \text{Vars}(P) \rightarrow \bar{D}$ the abstract state associating each variable to an abstract value of \bar{D} . When it is clear from the context, we denote the abstract state by ξ . Hence, given an abstract domain \bar{D} and a program P , the result of the analysis of P using the abstract domain \bar{D} (and the corresponding abstract semantics $\llbracket \cdot \rrbracket^{\bar{D}}$) is defined by

$$\text{Analysis}(P, \bar{D}, \xi_{\emptyset}) = \{(l_i, \xi_i) \mid l_i \in \text{Lab}(P), \xi_i = \llbracket P \rrbracket^{\bar{D}} \xi_{\emptyset} \text{ at program point } l_i\}$$

where ξ_{\emptyset} is the (initial) empty abstract state. For example, let us analyse the μ Dyn program with the original TAJIS abstract domain reported in Figure 5b. The abstract state holding at program point 12 (namely the exit program point) is $\xi_{12} = \{i \mapsto \text{UnsignedInt}, \text{nums} \mapsto \top_{\overline{\text{TJ}}}, \text{notnums} \mapsto \top_{\overline{\text{TJ}}}\}$, indeed the pair $(12, \xi_{12}) \in \text{Analysis}(P, \overline{\text{TJ}}, \xi_{\emptyset})$.

We can exploit the definition of path length distance reported in Definition 9 to define, for each program point $l \in \text{Lab}(\mathbf{P})$, the distance from $\perp_{\bar{D}}$ and each abstract value associated with each $\mathbf{x} \in \text{Vars}(\mathbf{P})$.

Definition 10 (Set of \perp -distances). Let \mathbf{P} be a μDyn program, \bar{D} an abstract domain and δ_{plen} the path length distance defined in Definition 9. The set of $\perp_{\bar{D}}$ -distances for a given $l \in \text{Lab}(\mathbf{P})$ is defined as follows.

$$\Phi_l(\mathbf{P}, \bar{D}) = \{(\mathbf{x}, \delta_{plen}(\perp_{\bar{D}}, \xi(\mathbf{x}))) \mid (l, \xi) \in \text{Analysis}(\mathbf{P}, \bar{D}, \xi_{\emptyset}), \mathbf{x} \in \text{Vars}(\mathbf{P})\}$$

For example, let us consider again ξ_{12} , previously defined, that is the abstract state holding at the program point 12 of the program reported in Figure 11 analyzing it with \bar{TJ} abstract domain. Hence, the set of the $\perp_{\bar{TJ}}$ -distances at program point 12 is $\Phi_{12}(\mathbf{P}, \bar{TJ}) = \{(i, 2), (\text{nums}, 3), (\text{notnums}, 3)\}$.

Note that, by the weak triangle inequality, given two abstract elements $\bar{d}, \bar{d}' \in \bar{D}$, if $\bar{d} \sqsubseteq_{\bar{D}} \bar{d}' \Rightarrow \delta_{plen}(\perp_{\bar{D}}, \bar{d}) \leq \delta_{plen}(\perp_{\bar{D}}, \bar{d}')$. Moreover, since $\perp_{\bar{D}}$ is comparable with any abstract value $a \in \bar{D}$, $\delta(\perp_{\bar{D}}, \bar{d}) \neq \infty$.

Definition 11. [Precision entropy at a program point] Let \mathbf{P} be a μDyn program and \bar{D} an abstract domain. Given a program point $l \in \text{Lab}(\mathbf{P})$, the precision entropy of \bar{D} for the program \mathbf{P} at l is defined as

$$\mathbb{P}_l(\mathbf{P}, \bar{D}) = \sum_{(\mathbf{x}_i, d_i) \in \Phi_l(\mathbf{P}, \bar{D})} d_i$$

where d_i is the path length distance from $\perp_{\bar{D}}$ to the abstract value of \mathbf{x}_i at the program point l .

We can use Definition 11 to define the precision entropy of an analysis for a given program.

Definition 12. [Precision entropy] Let \mathbf{P} be a μDyn program and \bar{D} be an abstract domain. The precision entropy $\mathbb{P}(\mathbf{P}, \bar{D})$ of the abstract domain A for the program \mathbf{P} is defined as

$$\mathbb{P}(\mathbf{P}, \bar{D}) = \sum_{l \in \text{Lab}(\mathbf{P})} \mathbb{P}_l(\mathbf{P}, \bar{D})$$

The precision entropy $\mathbb{P}(\mathbf{P}, \bar{D})$ says how much information the analysis based on the abstract domain \bar{D} is able to express: the more $\mathbb{P}(\mathbf{P}, \bar{D})$ is low,

the more the analysis based on \bar{D} is precise for the program P . Using this metric, we are able to compare the analysis results based on a certain abstract domain and its corresponding complete shell. As we have already mentioned at the beginning of Section 7, any abstract point of \bar{D} is contained in its complete shell. Moreover, it is worth noting that, given a variable $x \in \text{Vars}(P)$ and a certain program point $l \in \text{Lab}(P)$, the abstract value computed by the analysis on \bar{D} associated to x is always comparable with the abstract value computed by the analysis on $\text{Shell}_f(\bar{D})$ associated to x . Formally, let consider $(l, \xi) \in \text{Analysis}(P, \bar{D}, \xi_\emptyset)$ and $(l, \xi') \in \text{Analysis}(P, \text{Shell}_f(\bar{D}), \xi_\emptyset)$, for some program point $l \in \text{Lab}(P)$, we have that $\forall x \in \text{Vars}(P). \xi'(x) \sqsubseteq_{\text{Shell}_f(\bar{D})} \xi(x)$. This can be expressed by the following predicate

$$\mathbb{P}(P, \text{Shell}_f(\bar{D})) \leq \mathbb{P}(P, \bar{D}).$$

Informally speaking, for a program point l , the analysis result of x on $\text{Shell}_f(\bar{D})$ is always dominated (i.e., less than or equal to) by the analysis result of x on \bar{D} (contained in $\text{Shell}_f(\bar{D})$). This fact directly comes from the dominance relation involved in the definition of complete shells given in Theorems 1 and 2. For this reason, we can always compare the analysis results produced by \bar{D} and $\text{Shell}_f(\bar{D})$ for any variable and any program point.

This leads us to an automatic procedure to compute how much the analysis performed by $\text{Shell}_f(\bar{D})$ is better than the one performed by \bar{D} .

The CLAM static analyzer. The CLAM static analyzer for μDyn programs⁸ implements the TAJIS and SAFE string abstract domains and their corresponding complete shells discussed in Section 5. The abstract interpreter is parametric, as it can be tuned by selecting the string abstract domain to analyze a given program. Other string abstract domains can be easily plugged into our static analyzer, without re-implementing the underlying abstract interpreter. Moreover, it is possible to check the precision entropy, discussed in Section 7, of an abstract domain and its complete shell. In this way, it is possible to check at which program point and for which variables the complete shell-based analysis gains precision w.r.t. the analysis on the original abstract domain. In the following, we will consider TAJIS and its complete shell: note that, in our analyzer, both the domains specify the path length distances, as defined in Definition 9, from the bottom element to

⁸Available at <https://github.com/VincenzoArceri/clam>

Variable	\overline{TJ}	\overline{TJ}^*	$\mathbb{P}_{12}(P, \overline{TJ}) - \mathbb{P}_{12}(P, \overline{TJ}^*)$
<code>s</code>	UnsignedStr	UnsignedStr	0
<code>i</code>	UnsignedInt	UnsignedInt	0
<code>nums</code>	$\top_{\overline{TJ}}$	UnsignedOrNotNumeric	1
<code>notnums</code>	$\top_{\overline{TJ}}$	UnsignedOrNotNumeric	1

Table 2: Analysis output results with \overline{TJ} and \overline{TJ}^*

each of their abstract points (that are computed off-line and *hard-coded*) by computing all the lengths of the complete chains linking the bottom element with that points and picking the minimal one.

Let us consider the μ Dyn program reported in Figure 11, where the value of `s` is statically unknown. The program takes the strings `s` and puts its non-zero numerical characters into `nums` and the others into `notnums`. If the variable `s` is initialized as "24kobe8", at the program point 12 the value of `nums` is "248" and the value of `notnums` is "kobe". Let us consider TAJs and its complete shell and let analyze the program with both the abstract domains. When the variable `s` is initialized with the abstract value `UnsignedStr`, the analysis output, for the exit point is reported in Table 2, where the second and third columns correspond to the result analysis of \overline{TJ} and \overline{TJ}^* for the corresponding variable, respectively, and the last column is the increment precision gained by performing the analysis on \overline{TJ}^* . Observe that, for the variables `s` and `i`, we have no precision improvement, as stated by the last column of Table 2. Concerning variables `nums` and `notnums`, the analysis on \overline{TJ} returns the top values, whereas the analysis on \overline{TJ}^* leads to a precision increment since it returns the abstract value `UnsignedOrNotNumeric`, for both variables. The precision increment is stated by the last column of Table 2 for the variables `nums` and `notnums`, since their abstract values (`UnsignedOrNotNumeric`) on the analysis on \overline{TJ}^* is distant 1 from the abstract values of the analysis on \overline{TJ} ($\top_{\overline{TJ}}$). The total precision entropy of the analysis on \overline{TJ} and the analysis on \overline{TJ}^* is 54 and 44, respectively, so the overall precision increment when considering the complete shell, in this case, is equal to 10.

8. Conclusion

This paper focused on backward completeness in JavaScript-specific string abstract domains, and provided in particular the complete shells of TAJIS and SAFE string abstract domains w.r.t. `toNum` and `concat` operations, together with an effective procedure to measure the precision improvement of the analysis when moving to the complete shell. Our results can be easily applied also to JSAI string abstract domain [28], as it can be seen as an extension of the SAFE domain. The next challenge is to investigate forward completeness [24], the property that guarantees that no loss of precision occurs during the output abstraction process of a given operation. We aim to integrate the two completeness methodologies within an industrial JavaScript static analyzer to deeply evaluate their actual cost and overall impact.

Acknowledgment. Special thanks to the very professional reviewers, that contributed to improve the technical quality of the submitted manuscript.

References

- [1] P. A. Abdulla, M. F. Atig, Y. Chen, L. Holík, A. Rezine, P. Rümmer, J. Stenman, Norn: An SMT Solver for String Constraints, in: Computer Aided Verification - 27th International Conference, CAV 2015, San Francisco, CA, USA, July 18-24, 2015, Proceedings, Part I, 462–469, URL https://doi.org/10.1007/978-3-319-21690-4_29, 2015.
- [2] R. Amadini, G. Gange, F. Gauthier, A. Jordan, P. Schachte, H. Søndergaard, P. J. Stuckey, C. Zhang, Reference Abstract Domains and Applications to String Analysis, *Fundam. Inform.* 158 (4) (2018) 297–326, URL <https://doi.org/10.3233/FI-2018-1650>.
- [3] R. Amadini, G. Gange, P. J. Stuckey, G. Tack, A Novel Approach to String Constraint Solving, in: Principles and Practice of Constraint Programming - 23rd International Conference, CP 2017, Melbourne, VIC, Australia, August 28 - September 1, 2017, Proceedings, 3–20, URL https://doi.org/10.1007/978-3-319-66158-2_1, 2017.
- [4] R. Amadini, A. Jordan, G. Gange, F. Gauthier, P. Schachte, H. Søndergaard, P. J. Stuckey, C. Zhang, Combining String Abstract

- Domains for JavaScript Analysis: An Evaluation, in: Tools and Algorithms for the Construction and Analysis of Systems - 23rd International Conference, TACAS 2017, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2017, Uppsala, Sweden, April 22-29, 2017, Proceedings, Part I, 41–57, URL https://doi.org/10.1007/978-3-662-54577-5_3, 2017.
- [5] V. Arceri, S. Maffei, Abstract Domains for Type Juggling, *Electron. Notes Theor. Comput. Sci.* 331 (2017) 41–55, URL <https://doi.org/10.1016/j.entcs.2017.02.003>.
- [6] V. Arceri, I. Mastroeni, Static Program Analysis for String Manipulation Languages, in: Proceedings Seventh International Workshop on Verification and Program Transformation, VPT@Programming 2019, Genova, Italy, 2nd April 2019, 19–33, URL <https://doi.org/10.4204/EPTCS.299.5>, 2019.
- [7] V. Arceri, M. Olliaro, A. Cortesi, I. Mastroeni, Completeness of Abstract Domains for String Analysis of JavaScript Programs, in: Theoretical Aspects of Computing - ICTAC 2019 - 16th International Colloquium, Hammamet, Tunisia, October 31 - November 4, 2019, Proceedings, 255–272, URL https://doi.org/10.1007/978-3-030-32505-3_15, 2019.
- [8] T. Bultan, F. Yu, M. Alkhalaf, A. Aydin, String Analysis for Software Verification and Security, Springer, ISBN 978-3-319-68668-4, URL <https://doi.org/10.1007/978-3-319-68670-7>, 2017.
- [9] I. Casso, J. F. Morales, P. López-García, R. Giacobazzi, M. V. Hermenegildo, Computing Abstract Distances in Logic Programs, in: M. Gabbrielli (Ed.), Logic-Based Program Synthesis and Transformation - 29th International Symposium, LOPSTR 2019, Porto, Portugal, October 8-10, 2019, Revised Selected Papers, vol. 12042 of *Lecture Notes in Computer Science*, Springer, 57–72, URL https://doi.org/10.1007/978-3-030-45260-5_4, 2019.
- [10] L. Chen, A. Miné, P. Cousot, A Sound Floating-Point Polyhedra Abstract Domain, in: Programming Languages and Systems, 6th Asian Symposium, APLAS 2008, Bangalore, India, December 9-11, 2008. Proceedings, 3–18, URL https://doi.org/10.1007/978-3-540-89330-1_2, 2008.

- [11] A. S. Christensen, A. Møller, M. I. Schwartzbach, Precise Analysis of String Expressions, in: Static Analysis, 10th International Symposium, SAS 2003, San Diego, CA, USA, June 11-13, 2003, Proceedings, 1–18, URL https://doi.org/10.1007/3-540-44898-5_1, 2003.
- [12] R. Clarisó, J. Cortadella, The octahedron abstract domain, *Sci. Comput. Program.* 64 (1) (2007) 115–139, URL <https://doi.org/10.1016/j.scico.2006.03.009>.
- [13] A. Cortesi, G. Filé, W. H. Winsborough, Comparison of Abstract Interpretations, in: Automata, Languages and Programming, 19th International Colloquium, ICALP92, Vienna, Austria, July 13-17, 1992, Proceedings, 521–532, URL https://doi.org/10.1007/3-540-55719-9_101, 1992.
- [14] A. Cortesi, G. Filé, W. H. Winsborough, The Quotient of an Abstract Interpretation, *Theor. Comput. Sci.* 202 (1-2) (1998) 163–192, URL [https://doi.org/10.1016/S0304-3975\(97\)00137-0](https://doi.org/10.1016/S0304-3975(97)00137-0).
- [15] A. Cortesi, M. Oliaro, M-String Segmentation: A Refined Abstract Domain for String Analysis in C Programs, in: 2018 International Symposium on Theoretical Aspects of Software Engineering, TASE 2018, Guangzhou, China, August 29-31, 2018, 1–8, URL <https://doi.org/10.1109/TASE.2018.00009>, 2018.
- [16] G. Costantini, P. Ferrara, A. Cortesi, A suite of abstract domains for static analysis of string values, *Softw., Pract. Exper.* 45 (2) (2015) 245–287, URL <https://doi.org/10.1002/spe.2218>.
- [17] P. Cousot, R. Cousot, Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints, in: Conference Record of the Fourth ACM Symposium on Principles of Programming Languages, Los Angeles, California, USA, January 1977, 238–252, URL <https://doi.org/10.1145/512950.512973>, 1977.
- [18] P. Cousot, R. Cousot, Systematic Design of Program Analysis Frameworks, in: Conference Record of the Sixth Annual ACM Symposium on Principles of Programming Languages, San Antonio, Texas, USA, January 1979, 269–282, URL <https://doi.org/10.1145/567752.567778>, 1979.

- [19] P. Cousot, N. Halbwachs, Automatic Discovery of Linear Restraints Among Variables of a Program, in: Conference Record of the Fifth Annual ACM Symposium on Principles of Programming Languages, Tucson, Arizona, USA, January 1978, 84–96, URL <https://doi.org/10.1145/512760.512770>, 1978.
- [20] D. Filaretti, S. Maffei, An Executable Formal Semantics of PHP, in: ECOOP 2014 - Object-Oriented Programming - 28th European Conference, Uppsala, Sweden, July 28 - August 1, 2014. Proceedings, 567–592, URL https://doi.org/10.1007/978-3-662-44202-9_23, 2014.
- [21] R. Giacobazzi, I. Mastroeni, Transforming Abstract Interpretations by Abstract Interpretation, in: M. Alpuente, G. Vidal (Eds.), Static Analysis, 15th International Symposium, SAS 2008, Valencia, Spain, July 16-18, 2008. Proceedings, vol. 5079 of *Lecture Notes in Computer Science*, Springer, 1–17, URL https://doi.org/10.1007/978-3-540-69166-2_1, 2008.
- [22] R. Giacobazzi, I. Mastroeni, Making abstract models complete, *Math. Struct. Comput. Sci.* 26 (4) (2016) 658–701, URL <https://doi.org/10.1017/S0960129514000358>.
- [23] R. Giacobazzi, E. Quintarelli, Incompleteness, Counterexamples, and Refinements in Abstract Model-Checking, in: P. Cousot (Ed.), Static Analysis, 8th International Symposium, SAS 2001, Paris, France, July 16-18, 2001, Proceedings, vol. 2126 of *Lecture Notes in Computer Science*, Springer, 356–373, URL https://doi.org/10.1007/3-540-47764-0_20, 2001.
- [24] R. Giacobazzi, F. Ranzato, F. Scozzari, Making abstract interpretations complete, *J. ACM* 47 (2) (2000) 361–416, URL <https://doi.org/10.1145/333979.333989>.
- [25] P. Granger, Static Analysis of Arithmetical Congruences, *International Journal of Computer Mathematics - IJCM* 30 (1989) 165–190.
- [26] P. Granger, Static Analysis of Linear Congruence Equalities among Variables of a Program, in: TAPSOFT'91: Proceedings of the International Joint Conference on Theory and Practice of Software Development, Brighton, UK, April 8-12, 1991, Volume 1: Colloquium

- on Trees in Algebra and Programming (CAAP'91), 169–192, URL https://doi.org/10.1007/3-540-53982-4_10, 1991.
- [27] S. H. Jensen, A. Møller, P. Thiemann, Type Analysis for JavaScript, in: Static Analysis, 16th International Symposium, SAS 2009, Los Angeles, CA, USA, August 9-11, 2009. Proceedings, 238–255, URL https://doi.org/10.1007/978-3-642-03237-0_17, 2009.
- [28] V. Kashyap, K. Dewey, E. A. Kuefner, J. Wagner, K. Gibbons, J. Saracino, B. Wiedermann, B. Hardekopf, JSAI: a static analysis platform for JavaScript, in: Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering, (FSE-22), Hong Kong, China, November 16 - 22, 2014, 121–132, URL <https://doi.org/10.1145/2635868.2635904>, 2014.
- [29] S. Kim, W. Chin, J. Park, J. Kim, S. Ryu, Inferring Grammatical Summaries of String Values, in: Programming Languages and Systems - 12th Asian Symposium, APLAS 2014, Singapore, November 17-19, 2014, Proceedings, 372–391, URL https://doi.org/10.1007/978-3-319-12736-1_20, 2014.
- [30] E. Kneuss, P. Suter, V. Kuncak, Phantm: PHP analyzer for type mismatch, in: Proceedings of the 18th ACM SIGSOFT International Symposium on Foundations of Software Engineering, 2010, Santa Fe, NM, USA, November 7-11, 2010, 373–374, URL <https://doi.org/10.1145/1882291.1882355>, 2010.
- [31] H. Lee, S. Won, J. Jin, J. Cho, S. Ryu, SAFE: Formal Specification and Implementation of a Scalable Analysis Framework for ECMAScript, in: FOOL'12, URL <https://doi.org/10.1.1.387.1030>, 2012.
- [32] T. Liang, A. Reynolds, N. Tsiskaridze, C. Tinelli, C. W. Barrett, M. Deters, An efficient SMT solver for string constraints, in: Formal Methods in System Design 48 (3) (2016) 206–234, URL <https://doi.org/10.1007/s10703-016-0247-6>.
- [33] F. Logozzo, Towards a Quantitative Estimation of Abstract Interpretations, in: Workshop on Quantitative Analysis of Software, Microsoft, URL <https://www.microsoft.com/en-us/research/publication/>

towards-a-quantitative-estimation-of-abstract-interpretations/,
2009.

- [34] M. Madsen, E. Andreassen, String Analysis for Dynamic Field Access, in: Compiler Construction - 23rd International Conference, CC 2014, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2014, Grenoble, France, April 5-13, 2014. Proceedings, 197–217, URL https://doi.org/10.1007/978-3-642-54807-9_12, 2014.
- [35] S. Maffei, J. C. Mitchell, A. Taly, An Operational Semantics for JavaScript, in: Programming Languages and Systems, 6th Asian Symposium, APLAS 2008, Bangalore, India, December 9-11, 2008. Proceedings, 307–325, URL https://doi.org/10.1007/978-3-540-89330-1_22, 2008.
- [36] Y. Minamide, Static approximation of dynamically generated Web pages, in: Proceedings of the 14th international conference on World Wide Web, WWW 2005, Chiba, Japan, May 10-14, 2005, 432–441, URL <https://doi.org/10.1145/1060745.1060809>, 2005.
- [37] A. Miné, The octagon abstract domain, in: High. Order Symb. Comput. 19 (1) (2006) 31–100, URL <https://doi.org/10.1007/s10990-006-8609-1>.
- [38] R. Oucheikh, I. Berrada, O. E. Hichami, The 4-Octahedron Abstract Domain, in: Networked Systems - 4th International Conference, NETYS 2016, Marrakech, Morocco, May 18-20, 2016, Revised Selected Papers, 311–317, URL https://doi.org/10.1007/978-3-319-46140-3_25, 2016.
- [39] C. Park, H. Im, S. Ryu, Precise and scalable static analysis of jQuery using a regular expression domain, in: Proceedings of the 12th Symposium on Dynamic Languages, DLS 2016, Amsterdam, The Netherlands, November 1, 2016, 25–36, URL <https://doi.org/10.1145/2989225.2989228>, 2016.
- [40] A. D. Pierro, H. Wiklicky, Measuring the Precision of Abstract Interpretations, in: Logic Based Program Synthesis and Transformation,

- 10th International Workshop, LOPSTR 2000 London, UK, July 24-28, 2000, Selected Papers, 147–164, URL https://doi.org/10.1007/3-540-45142-0_9, 2000.
- [41] M. D. Preda, R. Giacobazzi, I. Mastroeni, Completeness in Approximate Transduction, in: X. Rival (Ed.), Static Analysis - 23rd International Symposium, SAS 2016, Edinburgh, UK, September 8-10, 2016, Proceedings, vol. 9837 of *Lecture Notes in Computer Science*, Springer, 126–146, URL https://doi.org/10.1007/978-3-662-53413-7_7, 2016.
- [42] P. Saxena, D. Akhawe, S. Hanna, F. Mao, S. McCamant, D. Song, A Symbolic Execution Framework for JavaScript, in: 31st IEEE Symposium on Security and Privacy, S&P 2010, 16-19 May 2010, Berkeley/Oakland, California, USA, 513–528, URL <https://doi.org/10.1109/SP.2010.38>, 2010.
- [43] A. Simon, A. King, J. M. Howe, Two Variables per Linear Inequality as an Abstract Domain, in: Logic Based Program Synthesis and Transformation, 12th International Workshop, LOPSTR 2002, Madrid, Spain, September 17-20,2002, Revised Selected Papers, 71–89, URL https://doi.org/10.1007/3-540-45013-0_7, 2002.
- [44] P. Sotin, Quantifying the Precision of Numerical Abstract Domains, Research Report, URL <https://hal.inria.fr/inria-00457324>, 2010.
- [45] M. Veanes, P. de Halleux, N. Tillmann, Rex: Symbolic Regular Expression Explorer, in: Third International Conference on Software Testing, Verification and Validation, ICST 2010, Paris, France, April 7-9, 2010, 498–507, URL <https://doi.org/10.1109/ICST.2010.15>, 2010.
- [46] M. Ward, The Closure Operators of a Lattice, *Annals of Mathematics* 43 (2) (1942) 191–196.
- [47] G. Wassermann, Z. Su, Sound and precise analysis of web applications for injection vulnerabilities, in: Proceedings of the ACM SIGPLAN 2007 Conference on Programming Language Design and Implementation, San Diego, California, USA, June 10-13, 2007, 32–41, URL <https://doi.org/10.1145/1250734.1250739>, 2007.
- [48] F. Yu, M. Alkhalaf, T. Bultan, O. H. Ibarra, Automata-based symbolic string analysis for vulnerability detection, *Formal Methods in*

System Design 44 (1) (2014) 44–70, URL <https://doi.org/10.1007/s10703-013-0189-1>.

- [49] F. Yu, T. Bultan, M. Cova, O. H. Ibarra, Symbolic String Verification: An Automata-Based Approach, in: Model Checking Software, 15th International SPIN Workshop, Los Angeles, CA, USA, August 10-12, 2008, Proceedings, 306–324, URL https://doi.org/10.1007/978-3-540-85114-1_21, 2008.
- [50] F. Yu, T. Bultan, B. Hardekopf, String Abstractions for String Verification, in: Model Checking Software - 18th International SPIN Workshop, Snowbird, UT, USA, July 14-15, 2011. Proceedings, 20–37, URL https://doi.org/10.1007/978-3-642-22306-8_3, 2011.