

# TARSIS: An effective automata-based abstract domain for string analysis

Luca Negrini<sup>1</sup>  | Vincenzo Arceri<sup>2</sup>  | Agostino Cortesi<sup>1</sup> | Pietro Ferrara<sup>1</sup>

<sup>1</sup>Department of Environmental Sciences, Informatics and Statistics, Ca' Foscari University of Venice, Venice, Italy

<sup>2</sup>Department of Mathematical, Physical and Computer Sciences, University of Parma, Parma, Italy

## Correspondence

Vincenzo Arceri, Department of Mathematical, Physical and Computer Sciences, University of Parma, Parco Area delle Scienze, 53/A, 43124 Parma, Italy.

Email: [vincenzo.arceri@unipr.it](mailto:vincenzo.arceri@unipr.it)

## Funding information

University of Parma, Grant/Award Numbers: MUR\_DM737\_2022\_FIL\_PROGETTI\_B\_ARCERI\_COFIN, CUP: D91B21005370003, PE00000014; EU—NGEU; PNRR, Grant/Award Number: ECS00000043; Ca' Foscari University

## Abstract

In this paper, we introduce TARSIS, a new abstract domain based on the abstract interpretation theory that approximates string values through finite state automata. The main novelty of TARSIS is that it works over an alphabet of strings instead of single characters. On the one hand, such an approach requires a more complex and refined definition of the lattice operators and of the abstract semantics of string operators. On the other hand, it is in position to obtain strictly more precise results than state-of-the-art approaches. We compare TARSIS both with simpler domains and with the standard automata model, targeting case studies containing standard yet challenging string manipulations. The performance gain w.r.t. the standard automata model is also assessed, measuring the speed-up gained by TARSIS. Experiments confirm that TARSIS can obtain precise results without incurring in excessive computational costs.

## KEYWORDS

abstract interpretation, static analysis, string analysis

## 1 | INTRODUCTION

Nowadays, string values play a key role in any modern programming language, because they are adopted for a variety of purposes and tasks. For instance, they allow to dynamically access object properties, to hide the program code by using string-to-code statements and reflection, or to manipulate data-interchange formats, such as JSON, just to name a few. In this context, the correctness of string manipulations is therefore crucial. Sound static analysis<sup>1,2</sup> has been widely applied to prove the correctness of programs (e.g., the absence of bugs). Recently, a relevant effort was spent towards the static approximation of string values in different contexts, such as SQL queries programmatically built by code,<sup>3</sup> reflection,<sup>4,5</sup> string-to-code statement analysis,<sup>6</sup> and injection vulnerabilities.<sup>7,8</sup>

Despite the great effort spent in reasoning about strings, static analysis often failed to manage programs that heavily manipulate strings, mainly due to the inaccuracy of the results and/or the prohibitive amount of resources (time and space) required to retrieve useful information on strings. On the one hand, finite height string abstractions<sup>9</sup> are computable in a reasonable time, but precision is suddenly lost when using advanced string manipulations. On the other hand, more sophisticated abstractions (e.g., the ones reported in Arceri et al<sup>10</sup> and Cortesi and Olliaro<sup>11</sup>) compute precise results, but they require a huge, and sometimes unrealistic, computational cost, making the analysis of real code intractable. A good representation of the latter abstractions is the finite state automata (FSA) domain.<sup>10</sup> Over-approximating strings into FSAs has shown to increase string analysis accuracy in many scenarios, but it does not scale up to real world programs dealing with statically unknown inputs and long text manipulations.

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2024 The Authors. Journal of Software: Evolution and Process published by John Wiley & Sons Ltd.

In this paper, we introduce TARSIS, a new abstract domain for string values based on FSAs. Standard FSA has been shown to provide precise abstractions of string values when all the components of such strings are known but with high computational cost. Instead of considering finite automata built over the classical alphabet of single characters, TARSIS considers automata built over an alphabet of strings. The alphabet comprises a special value to represent statically unknown strings. This avoids the creation of self-loops with any possible character as input, which otherwise would significantly degrade performance. We define the TARSIS's abstract semantics on all string operations reported in Arceri et al,<sup>10</sup> that we use as reference, either defined directly on the automaton or on its equivalent regular expression.

TARSIS has been implemented in GoLiSA,<sup>12</sup> a static analyzer for Go based on Library for Static Analysis (LiSA).<sup>13,14</sup> By comparing TARSIS with other cutting-edge domains for string analysis, results show that (i) when applied to simple code that causes a precision loss in simpler domains, TARSIS correctly approximates string values within a comparable execution time; (ii) on code that makes the standard automata domain unusable due to the complexity of the analysis, TARSIS is in position to perform in a limited amount of time, making it a viable domain for complex and real codebases; and (iii) TARSIS is able to precisely abstract complex string operations that have not been addressed by state-of-the-art domains.

This paper is a revised and extended version of Negrini et al.<sup>15</sup> Specifically, we completed the previous version of the paper by covering all of the string operations considered in Arceri et al,<sup>10</sup> with the addition of the string equality operator. For all supported operations (both newly added and already formalized in Negrini et al.<sup>15</sup>), we also reported proofs of soundness and completeness (or incompleteness). Moreover, we repeated the original experimental evaluation using GoLiSA<sup>12</sup> instead of the prototypical analyzer used in the original paper. Finally, to clearly assess the performance gain of TARSIS w.r.t. the FSA abstract domain introduced in Arceri et al,<sup>10</sup> we extended our evaluation to deeply compare execution times of the two domains.

The rest of the paper is structured as follows. Section 2 introduces a motivating example. Section 3 defines the mathematical notation used throughout the paper. Section 4 formalizes TARSIS and its abstract semantics. Section 5 compares TARSIS with other domains. Section 6 discusses most related works, while Section 7 concludes. Appendix A reports soundness and completeness (or incompleteness) proofs for TARSIS abstract semantics.

## 2 | MOTIVATING EXAMPLE

Consider the code of Figure 1 that counts the occurrences of string `substr` into string `s`. This code is (a simplification of) the Go API function `strings.Count` (see <https://cs.opensource.google/go/go/+refs/tags/go1.20.1/src/strings/strings.go>). Proving properties about the value of `n` at Line 9, is particularly challenging, because it requires to correctly model a set of string operations (viz., `len`, `Index`, and `substring`) and their interaction. State-of-the-art string analyses fail to precisely model most of such operations, because their abstraction of string values is not rigorous enough to deal with them. This loss of precision usually leads to failure in proving string-based properties (also on non-string values) in real-world software, such as the numerical bounds of the value returned by `Count` when applied to a string.

The goal of this paper is to provide an abstract interpretation-based static analysis, in order to deal with complex and nested string manipulations similar to the one reported in Figure 1. As we will discuss in Section 5, TARSIS models (among the others) all string operations used in `Count`, and it is precise enough to infer, given the abstractions of `s` and `substr`, the precise range of values that `n` might have when the function returns.

```

1  func Count(s, substr string) int {
2      if len(substr) == 0 {
3          return len(s) + 1
4      }
5      n := 0
6      for true {
7          i := strings.Index(s, substr)
8          if i == -1 {
9              return n
10         }
11         n++
12         s = s[i+len(substr):]
13     }
14 }

```

**FIGURE 1** The `strings.Count` function of the Go API.

### 3 | PRELIMINARIES

#### 3.1 | Mathematical notation

Given a set  $S$ ,  $S^*$  is the set of all finite sequences of elements of  $S$ . If  $s = s_0 \dots s_n \in S^*$ ,  $s_i$  is the  $i$ th element of  $s$ ,  $|s| = n + 1$  is its length, and  $s[x/y]$  is the sequence obtained by replacing all occurrences of  $x$  in  $s$  with  $y$ . When  $s'$  is a subsequence of  $s$ , we write  $s'_gs$ . Given  $s \in S^*$  and  $i, j \in \mathbb{N}$ ,  $0 \leq i \leq j \leq |s|$ , we denote the subsequence  $s_i s_{i+1} \dots s_{j-1}$  by  $s[i:j]$ , with  $s[i:]$  denoting the subsequence  $s_i s_{i+1} \dots s_n$ . We denote by  $s^n$ , with  $n \geq 0$  the  $n$ -times repetition of the string  $s$ . Given two sets  $S$  and  $T$ ,  $\wp(S)$  is the powerset of  $S$ ,  $ST$  is the set difference,  $S \subset T$  is the strict inclusion relation between  $S$  and  $T$ ,  $S \subseteq T$  is the inclusion relation between  $S$  and  $T$ ,  $S \times T$  is the Cartesian product between  $S$  and  $T$ , and  $S \cdot T$  is the concatenation of  $S$  and  $T$ , that is,  $S \cdot T = \{s \cdot t \mid s \in S, t \in T\}$ . Given a set  $S$  and  $n \in \mathbb{N}$ ,  $S^n$  is recursively defined as  $S^0 \triangleq \{\epsilon\}$ , and  $S^{n>0} \triangleq S \cdot S^{n-1}$ .

#### 3.2 | Ordered structures

A set  $L$  with a partial ordering relation  $\leq \subseteq L \times L$  is a poset, denoted by  $\langle L, \leq \rangle$ . A poset  $\langle L, \leq, \vee, \wedge \rangle$ , where  $\vee$  and  $\wedge$  are, respectively, the least upper bound (lub) and greatest lower bound (glb) operators of  $L$ , is a lattice if  $\forall x, y \in L. x \vee y$  and  $x \wedge y$  belong to  $L$ . It is also complete if  $\forall X \subseteq L$  we have that  $\bigvee X, \bigwedge X \in L$ . A complete lattice  $L$ , with ordering  $\leq$ , lub  $\vee$ , glb  $\wedge$ , top element  $\top$ , and bottom element  $\perp$  is denoted by  $\langle L, \leq, \vee, \wedge, \top, \perp \rangle$ .

#### 3.3 | Abstract interpretation

Abstract interpretation<sup>1,2</sup> is a theoretical framework for sound reasoning about semantic properties of a program, establishing a correspondence between the concrete semantics of a program and an approximation of it, called abstract semantics. Let  $C$  and  $A$  be complete lattices, a pair of monotone functions  $\alpha: C \rightarrow A$  and  $\gamma: A \rightarrow C$  forms a *Galois connection* (GC) between  $C$  and  $A$  if  $\forall x \in C, \forall y \in A: \alpha(x) \leq_A y \Leftrightarrow x \leq_C \gamma(y)$ . We denote a GC as  $C \xrightarrow[\alpha]{\gamma} A$ . Given  $C \xrightarrow[\alpha]{\gamma} A$ , a concrete function  $f: C \rightarrow C$  is, in general, not computable. Hence, a function  $f^\sharp: A \rightarrow A$  that must *correctly* approximate the function  $f$  is needed. If so, we say that the function  $f^\sharp$  is *sound*. Given  $C \xrightarrow[\alpha]{\gamma} A$  and a concrete function  $f: C \rightarrow C$ , an abstract function  $f^\sharp: A \rightarrow A$  is sound w.r.t.  $f$  if  $\forall c \in C. \alpha(f(c)) \leq_A f^\sharp(\alpha(c))$ , or equivalently  $\forall a \in A. f(\gamma(a)) \leq_C \gamma(f^\sharp(a))$ . Completeness<sup>16</sup> can be obtained by enforcing the equality of the soundness conditions. Doing so, we obtain two notion of completeness. Given  $C \xrightarrow[\alpha]{\gamma} A$ , a concrete function  $f: C \rightarrow C$  and an abstract function  $f^\sharp: A \rightarrow A$ ,  $f^\sharp$  is *backward complete* w.r.t.  $f$  if  $\forall c \in C. \alpha(f(c)) = f^\sharp(\alpha(c))$ , and it is *forward complete* w.r.t.  $f$  if  $\forall a \in A. f(\gamma(a)) = \gamma(f^\sharp(a))$ .

#### 3.4 | FSAs and regular expression notation

We follow the notation reported in Arceri et al<sup>10</sup> for introducing FSA. A FSA is a tuple  $\mathbb{A} = \langle Q, \Sigma, \delta, q_0, F \rangle$ , where  $Q$  is a finite set of states,  $q_0 \in Q$  is the initial state,  $\Sigma$  is a finite alphabet of symbols,  $\delta \subseteq Q \times \Sigma \times Q$  is the transition relation, and  $F \subseteq Q$  is the set of final states. If  $\delta: Q \times \Sigma \rightarrow Q$  is a function, then  $\mathbb{A}$  is called deterministic FSA. The set of all the FSAs is  $\mathcal{F}\mathbb{A}$ . If  $\mathcal{L} \subseteq \Sigma^*$  is recognized by a FSA, we say that  $\mathcal{L}$  is a regular language. Given  $\mathbb{A} \in \mathcal{F}\mathbb{A}$ ,  $\mathcal{L}(\mathbb{A})$  is the language accepted by  $\mathbb{A}$ . From the Myhill-Nerode theorem, for each regular language there uniquely exists a minimum FSA (w.r.t. the number of states) recognizing the language. Given a regular language  $\mathcal{L}$ ,  $\text{Min}(\mathbb{A})$  is the minimum FSA  $\mathbb{A}$  s.t.  $\mathcal{L} = \mathcal{L}(\mathbb{A})$ . Abusing notation, given a regular language  $\mathcal{L}$ ,  $\text{Min}(\mathcal{L})$  is the minimal FSA recognizing  $\mathcal{L}$ . Given  $\mathbb{A}$ , we denote by  $\text{Kleene}(\mathbb{A})$  the automaton recognizing the Kleene closure of  $\mathcal{L}(\mathbb{A})$ .

We denote as  $\text{Paths}(\mathbb{A}) \in \wp(\delta^*)$  the set of sequences of transitions corresponding to all the possible paths from the initial state  $q_0$  to a final state  $q_n \in F$ . When  $\mathbb{A}$  is cycle free, the set  $\text{Paths}(\mathbb{A})$  is finite and computable. Given  $\pi \in \text{Paths}(\mathbb{A})$ ,  $|\pi|$  is its length, meaning the sum of the lengths of the symbols that appear on the transitions composing the path. Furthermore,  $|\text{minPath}(\mathbb{A})| \in \mathbb{N}$  denotes the (unique) length of a minimum path. If  $\mathbb{A}$  is a cycle-free automaton,  $|\text{maxPath}(\mathbb{A})| \in \mathbb{N}$  denotes the (unique) length of a maximum path. Given  $\pi = t_0 \dots t_n \in \text{Paths}(\mathbb{A})$ ,  $\sigma_{\pi_i}$  is the symbol read by the transition  $t_i$ ,  $i \in [0, n]$ , and  $\sigma_\pi = \sigma_{\pi_0} \dots \sigma_{\pi_n}$  is the string recognized by such path. Predicate  $\text{cyclic}(\mathbb{A})$  holds if and only if the given automaton contains a cycle. Throughout the paper, it could be more convenient to refer to a FSA by its regular expression (regex for short), being equivalent. Given two regexes  $r_1$  and  $r_2$ ,  $r_1 || r_2$  is the disjunction between  $r_1$  and  $r_2$ ,  $r_1 r_2$  is the concatenation of  $r_1$  with  $r_2$ ,  $(r_1)^*$  is the Kleene closure of  $r_1$ .

### 3.5 | The FSA abstract domain

Here, we report the necessary notions about the FSA abstract domain presented in Arceri et al,<sup>10</sup> over-approximating string properties as the minimum deterministic FSA recognizing them. Given an alphabet  $\Sigma$ , the FSA domain is defined as  $\langle F_{A/\equiv}, \sqsubseteq_{F_A}, \sqcup_{F_A}, \sqcap_{F_A}, \text{Min}(\emptyset), \text{Min}(\Sigma^*) \rangle$ , where  $F_{A/\equiv}$  is the quotient set of  $F_A$  w.r.t. the equivalence relation induced by language equality,  $\sqsubseteq_{F_A}$  is the partial order induced by language inclusion, and  $\sqcup_{F_A}$  and  $\sqcap_{F_A}$  are the lub and the glb, respectively. The minimum is  $\text{Min}(\emptyset)$ , that is, the automaton recognizing the empty language, and the maximum is  $\text{Min}(\Sigma^*)$ , that is, the automaton recognizing any possible string over  $\Sigma$ . We abuse notation by representing equivalence classes in  $F_{A/\equiv}$  by one of its automaton (usually the minimum), that is, when we write  $A \in F_{A/\equiv}$  we mean  $[A]_{\equiv}$ . Because  $F_{A/\equiv}$  does not satisfy the ascending chain condition (ACC), that is, it contains infinite ascending chains, it is equipped with the parametric widening  $\nabla_{F_A}^n$ . The latter is defined in terms of a state equivalence relation merging states that recognize the same language, up to a fixed length  $n \in \mathbb{N}$ , a parameter used for tuning the widening precision.<sup>17,18</sup> For instance, let us consider the automata  $A, A' \in F_{A/\equiv}$  recognizing the languages  $\mathcal{L} = \{\epsilon, a\}$  and  $\mathcal{L}' = \{\epsilon, a, aa\}$ , respectively. The result of the application of the widening  $\nabla_{F_A}^n$ , with  $n = 1$ , is  $A \nabla_{F_A}^n A' = A'$  s.t.  $\mathcal{L}(A') = \{a^n \mid n \in \mathbb{N}\}$ .

### 3.6 | Core language and semantics

We introduce a core language  $\text{IMP}$ , whose syntax is reported in Figure 2. Such language, besides supporting arithmetic expressions ( $\text{AE}$ ) and Boolean expressions ( $\text{BE}$ ), also supports all of the string expressions ( $\text{SE}$ ) discussed in Arceri et al,<sup>10</sup> that we use as reference. Primitives values are  $\text{VAL} = \mathbb{Z} \cup \Sigma^* \cup \{\text{true}, \text{false}\}$ , namely, integers, strings, and Booleans. Programs states  $\mathbb{M}: \text{ID} \rightarrow \text{VAL}$  map identifiers to primitives values, ranged over the meta-variable  $m$ . The concrete semantics of  $\text{IMP}$  statements is captured by the function  $[[ \text{st} ]]: \mathbb{M} \rightarrow \mathbb{M}$ . The semantics is defined in a standard way and for this reason has been omitted. Such semantics relies on one of the expressions that we capture, abusing notation, as  $[[ e ]]: \mathbb{M} \rightarrow \text{VAL}$ . We define the part concerning strings in Figure 3.

## 4 | THE TARSIS ABSTRACT DOMAIN

In this section, we recast the original finite state abstract domain working over an alphabet of characters  $\Sigma$ , reported in Section 3, to an augmented abstract domain based on FSAs over an alphabet of strings.

### 4.1 | Abstract domain and widening

The key idea of TARSIS is to adopt the same FSA-based domain, changing the alphabet on which automata are defined to a set of strings, namely,  $\Sigma^*$ . The main concern is that  $\Sigma^*$  is infinite and it would not permit us to adopt the FSA model that requires the alphabet to be finite. Thus, in order to solve this problem, we make this abstract domain *parametric* to the program we aim to analyze and in particular to its strings. Given an  $\text{IMP}$  program  $P$ , we denote by  $\Sigma_P^*$  any substring of strings appearing in  $P$  (the set  $\Sigma_P^*$  can be easily computed collecting the constant strings in  $P$  by visiting its abstract syntax tree and then computing their substrings), *delimiting* the space of string properties we aim to check only on  $P$ .

At this point, we can instantiate the automata-based framework proposed in Arceri et al<sup>10</sup> with the new alphabet as

```

a ∈ AE ::= x ∈ ID | n ∈ ℤ | a + a | a - a | a * a | a / a
          | length(s) | indexOf(s,s)
b ∈ BE ::= x ∈ ID | true | false | b && b | b || b | ! b | e < e
          | e == e | contains(s1,s2) | startsWith(s1,s2) | endsWith(s1,s2)
s ∈ SE ::= x ∈ ID | "σ" | substr(s,a,a) | charAt(s,a)
          | repeat(s,a) | concat(s,s) | replace(s,s,s)
          | trim(s) | trimLeft(s) | trimRight(s)    (σ ∈ Σ*)
e ∈ E ::= a | b | s
st ∈ STMT ::= st ; st | skip | x = e | if (b) { st } else { st } | while (b) { st }
P ∈ IMP ::= st ;

```

FIGURE 2 IMP syntax.

$$\begin{aligned}
\llbracket \text{substr}(\mathbf{s}, \mathbf{a}, \mathbf{a}') \rrbracket_{\mathfrak{m}} &= \sigma_i \dots \sigma_j \quad \text{if } i \leq j < |\sigma|, i = \llbracket \mathbf{a} \rrbracket_{\mathfrak{m}}, j = \llbracket \mathbf{a}' \rrbracket_{\mathfrak{m}} \\
\llbracket \text{charAt}(\mathbf{s}, \mathbf{a}) \rrbracket_{\mathfrak{m}} &= \llbracket \text{substr}(\mathbf{s}, \mathbf{a}, \mathbf{a} + 1) \rrbracket_{\mathfrak{m}} \\
\llbracket \text{indexOf}(\mathbf{s}, \mathbf{s}') \rrbracket_{\mathfrak{m}} &= \begin{cases} \min\{i \mid \sigma_i \dots \sigma_j = \sigma'\} & \text{if } \exists i, j \in \mathbb{N}. \sigma_i \dots \sigma_j = \sigma' \\ -1 & \text{otherwise} \end{cases} \\
\llbracket \text{replace}(\mathbf{s}, \mathbf{s}', \mathbf{s}'') \rrbracket_{\mathfrak{m}} &= \begin{cases} \sigma[\sigma'/\sigma''] & \text{if } \sigma' \curvearrowright_{\mathfrak{s}} \sigma \\ \sigma & \text{otherwise} \end{cases} \\
\llbracket \text{trimLeft}(\mathbf{s}) \rrbracket_{\mathfrak{m}} &= \sigma' \text{ where } \sigma = \psi\sigma', \psi = \max\{\psi' \in \{\_ \}^* \mid \sigma = \psi'\sigma'\} \\
\llbracket \text{trimRight}(\mathbf{s}) \rrbracket_{\mathfrak{m}} &= \sigma' \text{ where } \sigma = \sigma'\phi, \phi = \max\{\phi' \in \{\_ \}^* \mid \sigma = \sigma''\phi'\} \\
\llbracket \text{trim}(\mathbf{s}) \rrbracket_{\mathfrak{m}} &= \llbracket \text{trimLeft}(\text{trimRight}(\mathbf{s})) \rrbracket_{\mathfrak{m}} \\
\llbracket \text{concat}(\mathbf{s}, \mathbf{s}') \rrbracket_{\mathfrak{m}} &= \sigma \cdot \sigma' \quad \llbracket \text{length}(\mathbf{s}) \rrbracket_{\mathfrak{m}} = |\sigma| \quad \llbracket \mathbf{s} == \mathbf{s}' \rrbracket_{\mathfrak{m}} = \sigma == \sigma' \\
\llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket_{\mathfrak{m}} &= \begin{cases} \text{true} & \text{if } \exists i, j \in \mathbb{N}. \sigma_i \dots \sigma_j = \sigma' \\ \text{false} & \text{otherwise} \end{cases} \\
\llbracket \text{repeat}(\mathbf{s}, \mathbf{a}) \rrbracket_{\mathfrak{m}} &= \sigma^n \quad \text{if } n = \llbracket \mathbf{a} \rrbracket_{\mathfrak{m}}, n \geq 0 \\
\llbracket \text{startsWith}(\mathbf{s}, \mathbf{s}') \rrbracket_{\mathfrak{m}} &= \begin{cases} \text{true} & \text{if } \exists \phi \in \Sigma^*. \sigma = \sigma'\phi \\ \text{false} & \text{otherwise} \end{cases} \\
\llbracket \text{endsWith}(\mathbf{s}, \mathbf{s}') \rrbracket_{\mathfrak{m}} &= \begin{cases} \text{true} & \text{if } \exists \phi \in \Sigma^*. \sigma = \phi\sigma' \\ \text{false} & \text{otherwise} \end{cases}
\end{aligned}$$

**FIGURE 3** Concrete semantics of IMP string expressions, where  $\sigma = \llbracket \mathbf{s} \rrbracket_{\mathfrak{m}}, \sigma' = \llbracket \mathbf{s}' \rrbracket_{\mathfrak{m}}, \sigma'' = \llbracket \mathbf{s}'' \rrbracket_{\mathfrak{m}}$

$$\langle \mathcal{TF}_{A/\equiv}, \sqsubseteq_{\mathcal{T}}, \sqcup_{\mathcal{T}}, \sqcap_{\mathcal{T}}, \text{Min}(\emptyset), \text{Min}(\mathbb{A}_{\mathbb{P}}^*) \rangle.$$

The alphabet on which FSAs are defined is  $\mathbb{A}_{\mathbb{P}} \triangleq \Sigma_{\mathbb{P}}^* \cup \{\mathbf{T}\}$ , where  $\mathbf{T}$  is a special symbol that we intend as “any possible string.” Let  $\mathcal{TF}_{\mathbb{A}}$  be the set of any deterministic FSA over the alphabet  $\mathbb{A}_{\mathbb{P}}$ . Because we can have more automata recognizing a language,  $\mathcal{TF}_{A/\equiv}$  is the quotient set of  $\mathcal{TF}_{\mathbb{A}}$  w.r.t. the equivalence relation induced by language equality, that is, the elements of domain are equivalence classes. For simplicity, when we write  $A \in \mathcal{TF}_{A/\equiv}$ , we intend the equivalence class of  $A$ .  $\sqsubseteq_{\mathcal{T}}$  is the partial order induced by language inclusion, and  $\sqcup_{\mathcal{T}}$  and  $\sqcap_{\mathcal{T}}$  are the lub and the glb over elements of  $\mathcal{TF}_{A/\equiv}$ , computing the equivalence class of the union and the intersection of the two automata representing the corresponding classes, respectively. The bottom element is  $\text{Min}(\emptyset)$ , corresponding to the automaton recognizing the empty language, and the maximum is  $\text{Min}(\mathbb{A}_{\mathbb{P}}^*)$ , namely, the automaton recognizing any string over  $\mathbb{A}_{\mathbb{P}}$ .

Similarly to the standard FSA domain  $F_{A/\equiv}$ , also  $\mathcal{TF}_{A/\equiv}$  is not a complete lattice and, consequently, it does not form a GC with the string concrete domain  $\wp(\Sigma^*)$ . This comes from the non-existence, in general, of the best abstraction of a string set in  $\mathcal{TF}_{A/\equiv}$  (e.g., a context-free language has no best abstract element in  $\mathcal{TF}_{A/\equiv}$  approximating it). Nevertheless, this is not a concern because weaker forms of abstract interpretation are still possible<sup>19</sup> still guaranteeing soundness relations between concrete and abstract elements (e.g., polyhedra<sup>20,21</sup>). In particular, we can still ensure soundness comparing the concretizations of our abstract elements (cf. Section 8 of Cousot and Cousot<sup>19</sup>). Hence, we define the concretization function  $\gamma_{\mathcal{T}}: \mathcal{TF}_{A/\equiv} \rightarrow \wp(\Sigma^*)$  as  $\gamma_{\mathcal{T}}(A) \triangleq \bigcup_{\sigma \in \mathcal{L}(A)} \text{Flat}(\sigma)$ , where  $\text{Flat}$  converts a string in  $\mathbb{A}_{\mathbb{P}}^*$  into a set of strings in  $\Sigma^*$ . For instance,  $\text{Flat}(a \mathbf{T} \mathbf{T} b b c) = \{a \sigma b b c \mid \sigma \in \Sigma^*\}$ . Note that the language of strings recognized by  $A$  corresponds to the concretization function reported above, namely,  $\mathcal{L}(A) = \gamma_{\mathcal{T}}(A)$ .

#### 4.1.1 | Widening

Similarly to the standard automata domain  $F_{A/\equiv}$ , also  $\mathcal{TF}_{A/\equiv}$  does not satisfy ACC, meaning that fixpoint computations over  $\mathcal{TF}_{A/\equiv}$  may not converge in a finite time. Hence, we need to equip  $\mathcal{TF}_{A/\equiv}$  with a widening operator to ensure the convergence of the analysis. We define the widening  $\nabla_{\mathcal{T}}^n: \mathcal{TF}_{A/\equiv} \times \mathcal{TF}_{A/\equiv} \rightarrow \mathcal{TF}_{A/\equiv}$ , parametric in  $n \in \mathbb{N}$ , taking two automata as input and returning an over-approximation of the least upper bounds between them, as required by widening definition. We rely on the standard automata widening reported in Section 3 that, informally speaking, can be seen as a *subset construction* algorithm<sup>22</sup> up to languages of strings of length  $n$ .

To explain the widening  $\nabla_{\mathcal{T}}^n$ , consider the following function manipulating strings; for the sake of readability, in the program examples presented in this paper, the plus operation between strings corresponds to the string concatenation:

```

function f(v) {
  res = "";
  while (?)
    res = res + "id = " + v;
  return res;
}

```

Function  $f$  takes as input parameter  $v$  and returns variable  $res$ . Let us suppose that  $v$  is a statically unknown string, corresponding to the automaton recognizing  $T$  (i.e.,  $\text{Min}(\{T\})$ ). The result of the function  $f$  is a string of the form  $id = T$ , repeated zero or more times. Because the `while` guard is unknown, the number of iterations is statically unknown, and in turn, also the number of concatenations performed inside the loop body. The goal here is to over-approximate the value returned by the function  $f$ , that is, the value of  $res$  at the end of the function. Let  $A$ , reported in Figure 4A, be the automaton abstracting the value of  $res$  before starting the second iteration of the loop, and let  $A'$ , reported in Figure 4B, be the automaton abstracting the value of  $res$  at the end of the second iteration. At this point, we want to apply the widening operator  $\nabla_T^n$ , between  $A$  and  $A'$ , working as follows. We first compute  $A \sqcup_T A'$  (corresponding to the automaton reported in Figure 4C except that also  $q_0$  is also a final state). On this automaton, we merge any state that recognizes the same  $A_P$ -strings of length  $n$ , with  $n \in \mathbb{N}$ . In our example, let  $n$  be 2. The resulting automaton is reported in Figure 4C, where  $q_0$  and  $q_2$  are put together, and the other states are left as singletons. Figure 4D depicts the minimized version of Figure 4C.

The widening  $\nabla_T^n$  has been proved to meet the widening requirements (i.e., over-approximation of the least upper bounds and convergence on infinite ascending chains) in D'Silva.<sup>18</sup> The parameter  $n$ , tuning the widening precision, is arbitrary and can be chosen by the user. As highlighted in Arceri et al.,<sup>10</sup> the higher  $n$  is, the more the corresponding widening operator is precise in over-approximating lubs of infinite ascending chains (i.e., in fixpoint computations). A classical improvement on widening-based fixpoint computations is to integrate a threshold,<sup>23</sup> namely, widening is applied to over-approximate lubs when a certain threshold (usually over some property of abstract values) is overcome. In fixpoint computations, we decide to apply the previously defined widening  $\nabla_T^n$  only when the number of the states of the lubbed automata overcomes the threshold  $\tau \in \mathbb{N}$ . This permits us to postpone the widening application, getting more precise abstractions when the automata sizes do not overcome the threshold. At the moment, the threshold  $\tau$  is not automatically inferred, because it surely requires further investigations.

## 4.2 | String abstract semantics of IMP

In this section, we define the abstract semantics of the string operators defined in Section 3 over the new string domain  $\mathcal{TF}_{A/\equiv}$ . Soundness and completeness (or incompleteness) proofs of the TARSI's abstract semantics are reported in Appendix A. While TARSI also implements `startsWith` and `endsWith` operators, their abstract semantics is not discussed in this section because we adopt the same ones reported in Arceri et al.<sup>10</sup>

Because IMP supports strings, integers, and Booleans values, we need a way to merge the corresponding abstract domains. In particular, we abstract integers with the well-known interval abstract domain<sup>1</sup> defined as  $\text{Intv} \triangleq \{ [a, b] \mid a \in \mathbb{Z} \cup \{-\infty\}, b \in \mathbb{Z} \cup \{+\infty\}, a \leq b \} \cup \{ \perp_{\text{Intv}} \}$ , and Booleans with  $\text{Bool} \triangleq \wp(\{\text{true}, \text{false}\})$ . As usual, we denote by  $\sqcup_{\text{Intv}}$  and  $\sqcup_{\text{Bool}}$  the lubs between intervals and Booleans, respectively. In particular, we merge such abstract domains in  $\text{VAL}^\sharp$  by the smashed sum abstract domain<sup>24</sup>  $\text{VAL}^\sharp \triangleq \mathcal{TF}_{A/\equiv} \oplus \text{Intv} \oplus \text{Bool}$  that *smashes* the bottom elements of the involved domains into a single one, and adds a new top above the ones from the domains.

The program state is represented through abstract program memories  $\mathbb{M}^\sharp : \text{Id} \rightarrow \text{VAL}^\sharp$  from identifiers to abstract values. The abstract semantics is captured by the function  $[[st]]^\sharp : \mathbb{M}^\sharp \rightarrow \mathbb{M}^\sharp$ , relying on the abstract semantics of expressions defined by, abusing notation,  $[[e]]^\sharp : \mathbb{M}^\sharp \rightarrow \text{VAL}^\sharp$ . We focus on the abstract semantics of string operations, while the semantics of the other expressions is standard and does not involve strings.

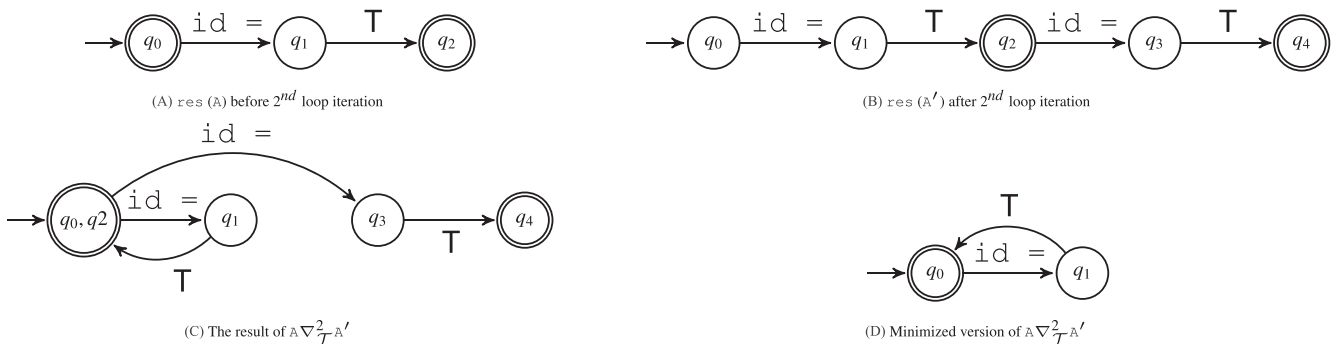


FIGURE 4 Example of widening application.

## 4.2.1 | Concat

Given  $\mathbb{A}, \mathbb{A}' \in \mathcal{TF}_{\mathbb{A}/\equiv}$ , the abstract semantics of `concat` returns a new automaton recognizing the language  $\{ \sigma \cdot \sigma' \mid \sigma \in \mathcal{L}(\mathbb{A}), \sigma' \in \mathcal{L}(\mathbb{A}') \}$ , that is, the concatenation between the strings of  $\mathcal{L}(\mathbb{A})$  with the strings of  $\mathcal{L}(\mathbb{A}')$ . This is achievable relying on the standard automata concatenation.<sup>22</sup> Let  $\mathbf{s}, \mathbf{s}' \in \mathbb{SE}$  and suppose that  $[[\mathbf{s}]]^{\sharp \text{mm}^{\sharp}} = \langle Q, \mathbb{A}_P, \delta, q_0, F \rangle \in \mathcal{TF}_{\mathbb{A}/\equiv}$ ,  $[[\mathbf{s}']]^{\sharp \text{mm}^{\sharp}} = \langle Q', \mathbb{A}'_P, \delta', q'_0, F' \rangle \in \mathcal{TF}_{\mathbb{A}'/\equiv}$ . The abstract semantics of `concat` is

$$[[\text{concat}(\mathbf{s}, \mathbf{s}')] ]^{\sharp \text{mm}^{\sharp}} \triangleq \text{Min}(\langle Q \cup Q', \mathbb{A}_P, \delta \cup \delta' \cup \{ (q_f, \epsilon, q'_0) \mid q_f \in F \}, q_0, F' \rangle).$$

Following the standard automata concatenation, the semantics merges the automata introducing an  $\epsilon$ -transition from each final state of  $\mathbb{A}$  to the initial state of  $\mathbb{A}'$ . The result's initial state is the initial state of  $\mathbb{A}$ , while its final states are the ones of  $\mathbb{A}'$ .

## 4.2.2 | Length

Given  $\mathbb{A} \in \mathcal{TF}_{\mathbb{A}/\equiv}$ , the abstract semantics of `length` returns an interval  $[c_1, c_2]$  such that  $\forall \sigma \in \mathcal{L}(\mathbb{A}). c_1 \leq |\sigma| \leq c_2$ . We recast the original idea of the abstract semantics of `length` over standard FSAs. Let  $\mathbf{s} \in \mathbb{SE}$ , supposing that  $[[\mathbf{s}]]^{\sharp \text{mm}^{\sharp}} = \mathbb{A} \in \mathcal{TF}_{\mathbb{A}/\equiv}$ . The `length` abstract semantics is

$$[[\text{length}(\mathbf{s})]]^{\sharp \text{mm}^{\sharp}} \triangleq \begin{cases} [[\text{minPath}(\mathbb{A})], +\infty] & \text{if } \text{cyclic}(\mathbb{A}) \vee \text{readsTop}(\mathbb{A}), \\ [[\text{minPath}(\mathbb{A})], |\text{maxPath}(\mathbb{A})|] & \text{otherwise,} \end{cases}$$

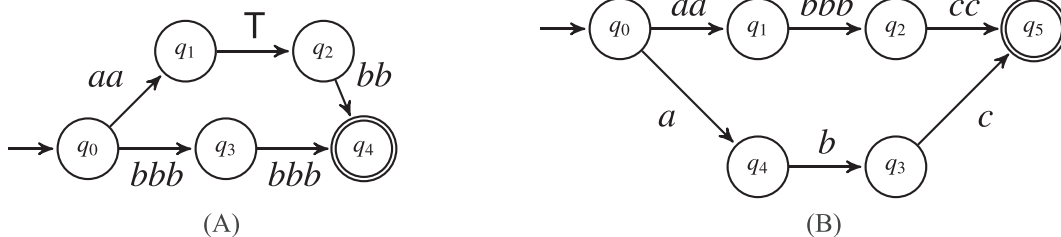
where  $\text{readsTop}(\mathbb{A}) \Leftrightarrow \exists q, q' \in Q. (q, T, q') \in \delta$ . Note that, when evaluating the length of the minimum path, T is considered to have a length of 0. For instance, consider the automaton  $\mathbb{A}$  reported in Figure 5A. The minimum path of  $\mathbb{A}$  is  $(q_0, aa, q_1), (q_1, T, q_2), (q_2, bb, q_4)$ , and its length is 4. Because a transition labeled with T is in  $\mathbb{A}$  (and its length cannot be statically determined), the abstract `length` of  $\mathbb{A}$  is  $[4, +\infty]$ . Consider the automaton  $\mathbb{A}'$  reported in Figure 5B. In this case,  $\mathbb{A}'$  has no cycles and has no transitions labeled with T and the length of every string recognized by  $\mathbb{A}'$  can be determined. The length of the minimum path of  $\mathbb{A}'$  is 3 (below path of  $\mathbb{A}'$ ), the length of the maximum path of  $\mathbb{A}'$  is 7 (above path of  $\mathbb{A}'$ ), and consequently, the abstract `length` of  $\mathbb{A}'$  is  $[3, 7]$ .

## 4.2.3 | Contains

Given  $\mathbb{A}, \mathbb{A}' \in \mathcal{TF}_{\mathbb{A}/\equiv}$ , the abstract semantics of `contains` should return `true` if every string of  $\mathbb{A}'$  is contained into every string of  $\mathbb{A}$ , `false` if no string of  $\mathbb{A}'$  is contained into any string of  $\mathbb{A}$ , and  $\{\text{true}, \text{false}\}$  in the other cases. For instance, consider the automaton  $\mathbb{A}$  depicted in Figure 7A and suppose we check if it contains the automaton  $\mathbb{A}'$  recognizing the language  $\{aa, a\}$ . The automaton  $\mathbb{A}'$  is a *single-path automaton*,<sup>25</sup> meaning that every string of  $\mathbb{A}'$  is a prefix of its longest string. In this case, the containment of the longest string (on each automaton path) implies the containment of the others, such as in our example, namely, it is enough to check that the longest string of  $\mathbb{A}'$  is contained into  $\mathbb{A}$ . Note that a single-path automaton cannot read the symbol T. We rely on the predicate `singlePath`( $\mathbb{A}$ ) when  $\mathbb{A}$  is a non-cyclic single-path automaton, and we denote by  $\sigma_{\text{sp}}$  its longest string.

Let  $\mathbf{s}, \mathbf{s}' \in \mathbb{SE}$ , supposing that  $[[\mathbf{s}]]^{\sharp \text{mm}^{\sharp}} = \mathbb{A} \in \mathcal{TF}_{\mathbb{A}/\equiv}$ ,  $[[\mathbf{s}']]^{\sharp \text{mm}^{\sharp}} = \mathbb{A}' \in \mathcal{TF}_{\mathbb{A}'/\equiv}$ . The `contains` abstract semantics is

$$[[\text{contains}(\mathbf{s}, \mathbf{s}')] ]^{\sharp \text{mm}^{\sharp}} \triangleq \begin{cases} \{\text{false}\} & \text{if } \mathbb{A}' \not\sqsubseteq_{\mathcal{T}} \text{FA}(\mathbb{A}) = \text{Min}(\emptyset), \\ \{\text{true}\} & \text{if } \text{singlePath}(\mathbb{A}') \wedge \forall \pi \in \text{Paths}(\mathbb{A}^{ac}). \sigma_{\text{sp}} \sigma_{\pi}, \\ \{\text{true}, \text{false}\} & \text{otherwise.} \end{cases}$$



**FIGURE 5** (A)  $\mathbb{A}$  s.t.  $\mathcal{L}(\mathbb{A}) = \{bbb\ bbb, aa\ T\ bb\}$  and (B)  $\mathbb{A}'$  s.t.  $\mathcal{L}(\mathbb{A}') = \{a\ b\ c, aa\ bbb\ cc\}$

In the first case, we denote by  $FA(A)$  the *factor* automaton of  $A$ , that is, the automaton recognizing any substring of  $A$ . In particular, if none of  $A$ 's substrings is part of  $A'$ , the abstract semantics safely returns `false` (checking the emptiness of the greatest lower bound between  $FA(A)$  and  $A'$ ). Then, if  $A'$  is a single-path automaton, the abstract semantics returns `true` if every path of  $A^{ac}$  reads the longest string of  $A'$ , with  $A^{ac}$  being a copy of  $A$  where all the cycles have been removed. Considering  $A^{ac}$  is necessary not only to make Paths computable but also to exclude *optional* strings recognized as part of loops. Here, we abuse notation denoting with  $\sigma_{sp} \sigma_\pi$  the fact that  $\sigma_{sp}$  is a substring of each string in  $Flat(\sigma_\pi)$ . Otherwise,  $\{\text{true}, \text{false}\}$  is returned.

#### 4.2.4 | String equality

Given  $A, A' \in \mathcal{TF}_{A/\equiv}$ , the abstract semantics of string equality returns `true` when  $A$  and  $A'$  recognize a singleton string and they are equal, `false` if no string recognized by  $A$  is equal to any string recognized by  $A'$ ,  $\{\text{true}, \text{false}\}$  in the other cases. Before defining the abstract semantics of string equality, we define equality between two strings over the alphabet  $\Sigma^* \cup \{T\}$ . For example, the strings  $aTb$  and  $abb$  may be equal, while  $aTb$  and  $abd$  definitely are not. Algorithm 1 defines the function  $\text{eq}: \{\Sigma \cup \{T\}\}^* \times \{\Sigma \cup \{T\}\}^* \rightarrow \text{Bool}$ , working on the expanded alphabet  $\{\Sigma \cup \{T\}\}^*$ .

---

#### Algorithm 1 eq algorithm

---

```

  let  $\sigma, \sigma' \in \{\Sigma \cup \{T\}\}^*$ 
1: if  $\sigma = \sigma' \wedge \sigma = \epsilon$  then
2:   return true
3: else if  $(\sigma = \epsilon \wedge \sigma' = T) \vee (\sigma' = \epsilon \wedge \sigma = T)$  then
4:   return  $\{\text{true}, \text{false}\}$ 
5: else if  $(\sigma = \epsilon \wedge \sigma' \neq T) \vee (\sigma' = \epsilon \wedge \sigma \neq T)$  then
6:   return false
7: else if  $\sigma = T \vee \sigma' = T$  then
8:   return  $\{\text{true}, \text{false}\}$ 
9: else if  $\sigma_0 \neq \sigma'_0 \wedge \sigma_0 \neq T \wedge \sigma'_0 \neq T$  then
10:  return false
11: else if  $\sigma_0 = \sigma'_0 \wedge \sigma_0 \neq T$  then
12:  return  $\text{eq}(\sigma[1:], \sigma'[1:])$ 
13: else if  $\sigma_0 = T \vee \sigma'_0 = T$  then
14:  return false  $\sqcup$   $\text{eq}(\sigma, \sigma'[1:])$   $\sqcup$   $\text{eq}(\sigma[1:], \sigma'[1:])$   $\sqcup$   $\text{eq}(\sigma[1:], \sigma')$ 
15: end if

```

---

Intuitively, `eq` checks string equality by recursively inspecting smaller suffixes of the given strings (Lines 11–15), returning definite answers only when `T` characters do not appear (Lines 2, 6, and 10). Note that, when one of the given strings begins with `T` (Lines 13 and 14), the algorithm can only prove inequality.

Let  $s, s' \in \mathcal{SE}$  and suppose  $[[s]]^{\#mm\#} = A$  and  $[[s']]^{\#mm\#} = A'$ . The abstract semantics of string equality is defined as

$$[[s == s']]^{\#mm\#} \triangleq \begin{cases} \{\text{false}\} & \text{if } A \sqcap_T A' = \text{Min}(\emptyset), \\ \{\text{true}, \text{false}\} & \text{if } \text{cyclic}(A) \vee \text{cyclic}(A'), \\ \bigsqcup_{\sigma \in \mathcal{L}(A), \sigma' \in \mathcal{L}(A')} \text{eq}(\sigma, \sigma') & \text{otherwise.} \end{cases}$$

In the first case, if the greatest lower bound between  $A$  and  $A'$  is  $\text{Min}(\emptyset)$ , then the automata do not share any common string; hence, the abstract semantics returns  $\{\text{false}\}$ . In the second case, if either  $A$  or  $A'$  are cyclic, the abstract semantics of string equality returns  $\{\text{true}, \text{false}\}$ . Otherwise, we rely on `eq` to compare the  $A_P$ -strings recognized by  $A$  and  $A'$  and we lub the results. To avoid cluttering the notation, the conversion from strings over  $\Sigma^* \cup \{T\}$  to strings over  $\{\Sigma \cup \{T\}\}^*$  when calling the function `eq` is implicit.

#### 4.2.5 | IndexOf

Given  $A, A' \in \mathcal{TF}_{A/\equiv}$ , the `indexOf` abstract semantics returns an interval of the first indexes of the strings of  $\mathcal{L}(A')$  inside strings of  $\mathcal{L}(A)$ , recalling that when at least one string of  $\mathcal{L}(A')$  is not a substring of any string of  $\mathcal{L}(A)$ , the resulting interval must take into account  $-1$  as well.

Let  $s, s' \in \mathcal{SE}$  and suppose  $[[s]]^{\#mm\#} = A$  and  $[[s']]^{\#mm\#} = A'$ . The abstract semantics of `indexOf` is defined as



$$\llbracket \text{indexOf}(s, s') \rrbracket^{\#}_{\text{mm}} \triangleq \begin{cases} [-1, +\infty] & \text{if } \text{cyclic}(\mathbb{A}) \vee \text{cyclic}(\mathbb{A}') \vee \text{readsTop}(\mathbb{A}'), \\ [-1, -1] & \text{if } \forall \sigma' \in \mathcal{L}(\mathbb{A}') \nexists \sigma \in \mathcal{L}(\mathbb{A}). \sigma'_s \sigma, \\ \bigsqcup_{\sigma \in \mathcal{L}(\mathbb{A}')}^{\text{Intv}} \text{IO}(\mathbb{A}, \sigma) & \text{otherwise.} \end{cases}$$

If one of the automata has cycles or the automaton abstracting strings we aim to search for (i.e.,  $\mathbb{A}'$  has a T-transition, we return  $[-1, +\infty]$ . Moreover, if none of the strings recognized by  $\mathbb{A}'$  is contained in a string recognized by  $\mathbb{A}$  (note that this is a decidable check because  $\mathbb{A}$  and  $\mathbb{A}'$  are cycle free; otherwise, the interval  $[-1, +\infty]$  would be returned in the first case), we can safely return the precise interval  $[-1, -1]$ . Otherwise, we rely on the auxiliary function  $\text{IO} : \mathcal{TFA}_{\equiv} \times \Sigma^* \rightarrow \text{Intv}$  that, given an automaton  $\mathbb{A}$  and a string  $\sigma \in \Sigma^*$ , returns an interval corresponding to the possible first positions of  $\sigma$  in strings recognized by  $\mathbb{A}$ . Because  $\mathbb{A}'$  recognizes a finite language, we apply  $\text{IO}(\mathbb{A}, \sigma)$  to each  $\sigma \in \mathcal{L}(\mathbb{A}')$  and to return the least upper bound of the resulting intervals.

In particular, the function  $\text{IO}(\mathbb{A}, \sigma)$  returns an interval  $[i, j] \in \text{Intv}$ , where  $i$  and  $j$  are computed as follows.

$$i = \begin{cases} -1 & \text{if } \exists \pi \in \text{Paths}(\mathbb{A}). \sigma_s \sigma_\pi, \\ \min_{\pi \in \text{Paths}(\mathbb{A})} \{i \mid \sigma_f \in \text{Flat}(\sigma_\pi) \wedge \sigma_{f_1} \dots \sigma_{f_{i+n}} = \sigma\} & \text{otherwise,} \end{cases}$$

$$j = \begin{cases} -1 & \text{if } \forall \pi \in \text{Paths}(\mathbb{A}). \sigma_s \sigma_\pi, \\ +\infty & \text{if } \exists \pi \in \text{Paths}(\mathbb{A}). \sigma_s \sigma_\pi \wedge \pi \text{ reads T before } \sigma, \\ \max_{\pi \in \text{Paths}(\mathbb{A})} \{j \mid \sigma_f \in \text{Flat}(\sigma_\pi) \wedge \sigma_{f_1} \dots \sigma_{f_{i+n}} = \sigma \wedge \sigma_s \sigma_{f_0} \dots \sigma_{f_{i+n-1}}\} & \text{otherwise.} \end{cases}$$

As for the abstract semantics of `contains`, we abuse notation denoting with  $\sigma_s \sigma_\pi$  the fact that  $\sigma$  is a substring of each string in  $\text{Flat}(\sigma_\pi)$ . Given  $\text{IO}(\mathbb{A}, \sigma) = [i, j] \in \text{Intv}$ ,  $i$  corresponds to the minimal position where the first occurrence of  $\sigma$  can be found in  $\mathbb{A}$ , while  $j$  to the maximal one. Let us first focus on the computation of the minimal position. If there exists a path  $\pi$  of  $\mathbb{A}$  s.t.  $\sigma$  is not recognized by  $\sigma_\pi$ , then the minimal position where  $\sigma$  can be found in  $\mathbb{A}$  does not exist and  $-1$  is returned. Otherwise, the minimal position where  $\sigma$  begins across  $\pi$  is returned. Let us consider now the computation of the maximal position. If all paths of the automaton do not recognize  $\sigma$ , then  $-1$  is returned. If there exists a path where  $\sigma$  is recognized but the character T appears earlier in the path, then  $+\infty$  is returned. Otherwise, the maximal index of the first occurrences of  $\sigma$  across the paths of  $\mathbb{A}$  is returned.

## 4.2.6 | Repeat

Given  $\mathbb{A} \in \mathcal{TFA}_{\equiv}$  and  $[i, j] \in \text{Intv}$ , the abstract semantics of `repeat` should return an automaton recognizing the language of every string recognized by  $\mathbb{A}$  repeated  $k$ -times, with  $i \leq k \leq j$ . We first define the auxiliary function `repeat`, reported in Algorithm 2, that inputs an automaton  $\mathbb{A}$  and  $n \in \mathbb{N}$  and returns an automaton recognizing the strings of  $\mathbb{A}$  repeated  $n$ -times.

---

### Algorithm 2 repeat algorithm

---

```

let  $\mathbb{A} \in \mathcal{TFA}_{\equiv}$ ,  $n \in \mathbb{N}$ 
1: if  $n = 0$  then
2:   return  $\text{Min}(\{\epsilon\})$ 
3: else if  $\text{cyclic}(\mathbb{A})$  then
4:    $\mathbb{A}' \leftarrow \text{Min}(\{\epsilon\})$ 
5:   for  $i \in [0, n-1]$  do
6:      $\mathbb{A}' \leftarrow \llbracket \text{concat}(\mathbb{A}', \mathbb{A}) \rrbracket^{\#}$ 
7:   end for
8:   return  $\mathbb{A}'$ 
9: else
10:   $\mathbb{A}_r \leftarrow \text{Min}(\emptyset)$ 
11:  for  $\pi \in \text{paths}(\mathbb{A})$  do
12:     $\mathbb{A}' \leftarrow \text{Min}(\{\epsilon\})$ 
13:    for  $i \in [0, n-1]$  do
14:       $\mathbb{A}' \leftarrow \llbracket \text{concat}(\mathbb{A}', \text{Min}(\{\sigma_\pi\})) \rrbracket^{\#}$ 
15:    end for
16:     $\mathbb{A}_r \leftarrow \mathbb{A}_r \sqcup_{\mathcal{T}} \mathbb{A}'$ 
17:  end for
18:  return  $\mathbb{A}_r$ 
19: end if

```

---

Lines 1 and 2 handle the case when  $n$  is zero and return the automaton recognizing the empty string. Lines 3–8 handle the case when  $A'$  is a cyclic automaton, and they build the automaton  $A'$  returned at Line 8, corresponding to the  $n$ -concatenation of the automaton  $A$ . Otherwise, for each path  $\pi$  of  $A$ , Lines 12–15 build the automaton  $A'$  that corresponds to the  $n$ -repetition of  $\sigma_\pi$ , that is, the string read by the path  $\pi$ . This operation is repeated for each path of the automaton (Lines 11–17), and the obtained automata are lubbed together in  $A_r$  at Line 16, which is finally returned.

Let  $s \in \text{SE}$  and  $a \in \text{AE}$ , supposing that  $[[s]]^{\#mm\#} = A \in \mathcal{TF}_{A/\equiv}$  and  $[[a]]^{\#mm\#} = [i,j] \in \text{Intv}$ . W.l.o.g., let us suppose that  $[i,j] \subseteq [0, +\infty]$ ; when negative values are met, the automaton recognizing the empty language is returned. The `repeat` abstract semantics is

$$[[\text{repeat}(s,a)]]^{\#mm\#} = \begin{cases} \text{Kleene}(A) & \text{if } [i,j] = [0, +\infty], \\ \text{repeat}(A,i) & \text{if } i = j \wedge i \in \mathbb{N}, \\ [[\text{concat}(\text{repeat}(A,i), \text{Kleene}(A))] ]^{\#} & \text{if } j = +\infty, \\ \bigsqcup_{k \in [i,j]} \text{repeat}(A,k) & \text{otherwise.} \end{cases}$$

When the input interval is  $[0, +\infty]$ , the first case is matched and the Kleene closure of  $A$  is returned. The second case is when the interval concretizes to a single value (e.g.,  $[2,2]$ ) and the abstract semantics returns  $\text{repeat}(A,i)$ . If the interval is  $[i, +\infty]$ , with  $i \in \mathbb{N} \setminus \{0\}$ , because it is excluded in the first case, the abstract semantics returns the concatenation between the  $i$ -repetition of  $A$  with its Kleene closure. An example of this case is reported in Figure 6. Otherwise, that is, the interval  $[i,j]$  is finite, the `repeat`'s abstract semantics returns the lub of the  $k$ -repetition of  $A$ , for each  $i \leq k \leq j$ .

#### 4.2.7 | TrimLeft, TrimRight, and Trim

The concrete semantics of `trimLeft` removes leading whitespaces (i.e., at the beginning) from a string. Similarly, its abstract semantics inputs an automaton  $A \in \mathcal{TF}_{A/\equiv}$  and removes leading whitespaces from the begin of each string recognized by  $A$ . Let  $A \in \mathcal{TF}_{A/\equiv}$  and  $r$  be the regex corresponding to the language recognized by  $A$ . The `trimLeft`'s abstract semantics is captured by the function `trimL` inductively defined on the structure of regexes as follows.

$$\text{trimL}(r) = \begin{cases} r & \text{if } r = \top \vee r = \emptyset, \\ [[\text{trimLeft}(\sigma)]] & \text{if } r = \sigma, \\ \text{trimL}(r_1) \parallel \text{trimL}(r_2) & \text{if } r = r_1 \parallel r_2, \\ \text{trimL}(r_2) & \text{if } r = r_1 r_2 \wedge \text{trimL}(r_1) = \epsilon, \\ \text{trimL}(r_1)(r_2 \parallel \text{trimL}(r_2)) & \text{if } r = r_1 r_2 \wedge \text{readWS}(r_1), \\ \text{trimL}(r_1)r_2 & \text{if } r = r_1 r_2 \wedge \neg \text{readWS}(r_1), \\ \epsilon & \text{if } r = (r_1)^* \wedge \text{trimL}(r_1) = \epsilon, \\ \epsilon \parallel \text{trimL}(r_1) r & \text{if } r = (r_1)^*. \end{cases}$$

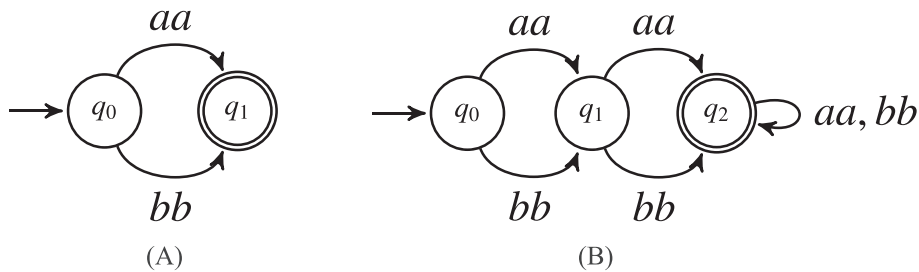


FIGURE 6 (A)  $A$  and (B)  $[[\text{repeat}(A, [2, +\infty])]]^{\#}$

The predicate  $\text{readWS}(r)$  holds if the language of  $r$  recognizes a whitespace string, that is,  $\text{readWS}(r) \Leftrightarrow \exists ws \in \{ \}^* . ws \in \mathcal{L}(r)$ . If the regex is empty or  $T$ ,  $\text{trimL}$  behaves as the identity function (first case). If the regex is an atom, we rely on the concrete semantics of  $\text{trimLeft}$  (second case). If the regex is a disjunction, the result is the disjunction of the application of  $\text{trimL}$  of the operands (third case). Then, three cases are needed for the regex concatenation  $r_1 r_2$ . If  $r_1$  recognizes only whitespace strings (i.e.,  $\text{trimL}(r_1) = \epsilon$ ), then we need to trim left also  $r_2$  (fourth case). If  $r_1$  recognizes at least one whitespace string, it *might* be necessary to also trim  $r_2$ , hence  $\text{trimL}(r_1)$  is concatenated with the disjunction of  $r_2$  and  $\text{trimL}(r_2)$  (fifth case). For instance, let us consider the regex  $( \| a ) b$ . The concrete semantics of  $\text{trimLeft}$  would return  $\{ a b, b \}$ . Our abstract semantics, in this case, correctly returns  $\text{trimL}( \| a ) ( b \| \text{trimL}( b ) ) = a ( b \| b ) = a b \| b$ . Lastly, if  $r_1$  does not recognize any whitespace string,  $r_1$  is left-trimmed and concatenated with  $r_2$  (sixth case). Finally, if  $r = (r_1)^*$ , two cases are identified. If  $r_1$  recognizes only whitespace strings, the empty string is returned (seventh case); otherwise, the result may be the empty string (in the case of 0-repetition) or the whole regex  $r$  appended to trimmed-left  $r_1$  (eighth case). Similarly, we can define  $\text{trimR}$  that removes trailing whitespace from the input regex. The definition of  $\text{trimR}$  is left implicit because it is analogous to one of the  $\text{trimL}$ .

Let  $s \in \text{SE}$ , supposing that  $[[ s ]]^{\#} \text{mm}^{\#} = A \in \mathcal{TF}_{A/\equiv}$  and let  $r$  be the regex equivalent to  $A$ . The abstract semantics of  $\text{trimLeft}$ ,  $\text{trimRight}$ , and  $\text{trim}$  are

$$[[ \text{trimLeft}(s) ]]^{\#} \text{mm}^{\#} = \text{trimL}(r) [[ \text{trimRight}(s) ]]^{\#} \text{mm}^{\#} = \text{trimR}(r) [[ \text{trim}(s) ]]^{\#} \text{mm}^{\#} = \text{trimL}(\text{trimR}(r)).$$

#### 4.2.8 | Replace

To give the intuition about how the abstract semantics of  $\text{replace}$  works, consider three automata  $A, A_s, A_r \in \mathcal{TF}_{A/\equiv}$ . Let us refer to  $A_s$  as the *search automaton* and to  $A_r$  as the *replace automaton*. Roughly speaking, the abstract semantics of  $\text{replace}$  substitutes strings of  $A_s$  with strings of  $A_r$  inside strings of  $A$ . We need to specify two types of possible replacements, by means of the following example. Consider  $A \in \mathcal{TF}_{A/\equiv}$  that is depicted in Figure 7A and suppose that the search automaton  $A_s$  is the one recognizing the string  $bbb$  and the replace automaton  $A_r$  is a random automaton. In this case, the  $\text{replace}$  abstract semantics performs a *must-replace* over  $A$ , namely, substituting the sub-automaton composed by  $q_1$  and  $q_2$  with the replace automaton  $A_r$ . Instead, let us suppose that the search automaton  $A_s$  is the one recognizing the language  $\{bbb, cc\}$ . Because it is unknown which string *must* be replaced (between  $bbb$  and  $cc$ ), the  $\text{replace}$  abstract semantics needs to perform a *may-replace*: When a string recognized by the search automaton is met inside a path of  $A$ , it is left unaltered in the automaton, and in the same position where the string is met, the abstract  $\text{replace}$  only extends  $A$  with the replace automaton. An example of may replacement is reported in Figure 7, where  $A$  is the one reported in Figure 7A, the search automaton  $A_s$  is the one recognizing the language  $\{bbb, cc\}$  and the replace automaton  $A_r$  is the one recognizing the string  $rr$ .

Before introducing the abstract semantics of  $\text{replace}$ , we define how to replace a string into an automaton's path. In particular, we define algorithm RP in Algorithm 3, that given a path  $\pi$  of an arbitrary automaton, a replace automaton  $A^r$ , and  $\sigma^s \in \Sigma^* \cup \{T\}$  returns a new automaton built starting from the path, but where portions of the path that recognize  $\sigma^s$  have been replaced with  $A^r$ .

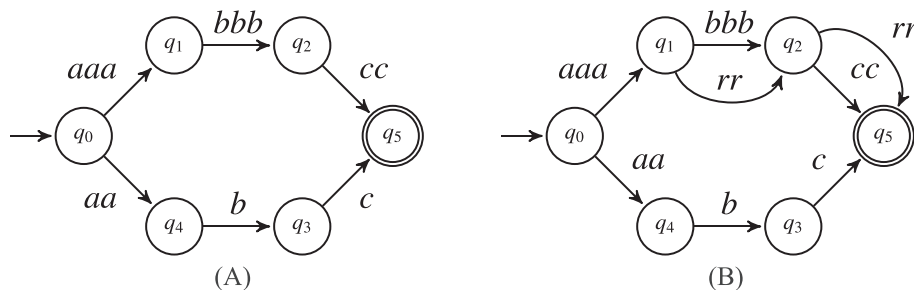


FIGURE 7 Example of may-replacement.

**Algorithm 3** RP algorithm

---

```

1: let  $\pi = (q_0, \sigma_0, q_1), \dots, (q_{n-1}, \sigma_{n-1}, q_n), A^r = \langle Q^r, \mathbb{A}, \delta^r, q_0^r, F^r \rangle \in \mathcal{TF}_{A/\equiv}, \sigma^s \in \Sigma^* \cup \{T\}$ 
2:  $Q^{result} \leftarrow \{ q \mid (q, \sigma, q') \in \pi \vee (q', \sigma, q) \in \pi \}$ 
3:  $\delta^{result} \leftarrow \pi$ 
4: for  $(q_i, \sigma_0^s, q_{i+1}), \dots, (q_{i+n-1}, \sigma_n^s, q_{i+n}) \in \pi$  do
5:    $\langle Q', \mathbb{A}, \delta', q_0', F' \rangle \leftarrow \text{clone}(A^r)$ 
6:    $\delta^{result} \leftarrow \delta^{result} \cup (q_i, \epsilon, q_0')$ 
7:    $\delta^{result} \leftarrow \delta^{result} \cup \{ (q_f, \epsilon, q_{i+n}) \mid q_f \in F' \}$ 
8:    $Q^{result} \leftarrow Q^{result} \setminus \{q_{i+1}, \dots, q_{i+n-1}\}$ 
9:    $\delta^{result} \leftarrow \delta^{result} \setminus \{(q_i, \sigma_0^s, q_{i+1}), \dots, (q_{i+n-1}, \sigma_n^s, q_{i+n})\}$ 
10: end for
11: return  $\langle Q^{result}, \mathbb{A}, \delta^{result}, q_0^o, F^o \rangle$ 

```

---

Algorithm 3 searches the given string  $\sigma^s$  across path  $\pi$ , collecting the sequences of transitions that recognize the search string  $\sigma^s$  and extracting them from  $\pi$  (Line 3). Whenever a matching sequence is found,  $A^r$  is cloned to  $A'$  to ensure that all additions target a different set of nodes (Line 4). Then, an  $\epsilon$ -transition is introduced going from the first state of the sequence to the initial state of  $A'$ , and one such transition is also introduced for each final state of  $A'$ , connecting that state with the ending state of the sequence (Lines 5 and 6). The list of states composing the sequence of transitions is then removed from the result (Line 7), together with the transitions connecting them (Line 8), because those were needed only to recognize the string that has been replaced. Note that RP corresponds to a must-replace. At this point, we are ready to define the `replace` abstract semantics. In particular, if either  $A$  or  $A_s$  have cycles or if one of them has a T-transition, we return  $\text{Min}(\{T\})$ , namely, the automaton recognizing T. Otherwise, the `replace` abstract semantics is

$$[[ \text{replace}(s, s_s, s_r) ]]^{\sharp \text{mm}^{\sharp}} \triangleq \begin{cases} A & \text{if } \forall \sigma_s \in \mathcal{L}(A_s). \nexists \sigma \in \mathcal{L}(A). \sigma_s s \sigma, \\ \bigsqcup_{\pi \in \text{Paths}(A)} \text{RP}(\pi, \sigma_s, A_r) & \text{if } \mathcal{L}(A_s) = \{\sigma_s\}, \\ \bigsqcup_{\substack{\sigma \in \mathcal{L}(A_s) \\ \pi \in \text{Paths}(A)}} \text{RP}(\pi, \sigma, A_r \sqcup_{\mathcal{T}} \text{Min}(\{\sigma\})) & \text{otherwise.} \end{cases}$$

In the first case, if none of the strings recognized by the search automaton  $A_s$  is contained in strings recognized by  $A$ , we can safely return the original automaton  $A$  without any replacement. In the special case where  $\mathcal{L}(A_s) = \{\sigma_s\}$ , we return the automaton obtained by replacing  $\sigma_s$  across all paths of  $A$  using function  $\text{RP}(\pi, \sigma_s, A_r)$ . In the last case, for each string  $\sigma \in \mathcal{L}(A_s)$  and for each path  $\pi \in \text{Paths}(A)$ , we perform a may-replace of  $\sigma$  with  $A_r$ : Note that this exactly corresponds to a call to RP where the replace automaton is  $A_r \sqcup_{\mathcal{T}} \text{Min}(\{\sigma\})$ . The so far obtained automata are finally lubbed together.

#### 4.2.9 | Substr and CharAt

Given  $A \in \mathcal{TF}_{A/\equiv}$  and two intervals  $i, j \in \text{Intv}$ , the abstract semantics of `substr` returns a new automaton  $A'$  soundly approximating any substring from  $i$  to  $j$  of strings recognized by  $A$ , for any  $i \in i, j \in j$  s.t.  $i \leq j$ .

Given  $A \in \mathcal{TF}_{A/\equiv}$ , in the definition of the `substr` semantics, we rely on the corresponding regex  $r$  because the two representations are equivalent and regexes allow us to define a more intuitive formalization of the semantics of `substr`. Let us suppose that  $[[ [s] ]]^{\sharp \text{mm}^{\sharp}} = A \in \mathcal{TF}_{A/\equiv}$ , and let us denote by  $r$  the regex corresponding to the language recognized by  $A$ . At the moment, let us consider exact intervals representing one integer value, namely,  $[[ a_1 ]]^{\sharp \text{mm}^{\sharp}} = [i, i]$  and  $[[ a_2 ]]^{\sharp \text{mm}^{\sharp}} = [j, j]$ , with  $i, j \in \mathbb{Z}$ . In this case, the abstract semantics is defined as

$$[[ \text{substr}(s, a_1, a_2) ]]^{\sharp \text{mm}^{\sharp}} \triangleq \bigsqcup \text{Min}(\{ \sigma \mid (\sigma, 0, 0) \in \text{Sb}(r, i, j - i) \}),$$

where  $\text{Sb}$  takes as input a regex  $r$ , two indexes  $i, j \in \mathbb{N}$ , and computes the set of substrings from  $i$  to  $j$  of all the strings recognized by  $r$ . In particular,  $\text{Sb}$  is defined by Algorithm 4, and given a regex  $r$  and  $i, j \in \mathbb{N}$ , it returns a set of triples of the form  $(\sigma, n_1, n_2)$ , such that  $\sigma$  is the *partial substring* that Algorithm 4 has computed up to now,  $n_1 \in \mathbb{N}$  tracks how many characters have still to be skipped before the substring can be computed and  $n_2 \in \mathbb{N}$  is the number of characters Algorithm 4 still needs to look for to successfully compute a substring.

**Algorithm 4** Sb algorithm

```

let r regex over  $\mathbb{A}$ ,  $i, j \in \mathbb{N}$ 
1: if  $j=0 \vee r = \emptyset$  then
2:   return  $\emptyset$ 
3: else if  $r = \sigma \in \Sigma^*$  then
4:   if  $i > |\sigma|$  then
5:     return  $\{(\epsilon, i - |\sigma|, j)\}$ 
6:   else if  $i + j > |\sigma|$  then
7:     return  $\{(\sigma_i \dots \sigma_{|\sigma|}, 0, j - |\sigma| + i)\}$ 
8:   else
9:     return  $\{(\sigma_i \dots \sigma_{i+j}, 0, 0)\}$ 
10:  end if
11: else if  $r = \top$  then
12:  result  $\leftarrow \{(\epsilon, i - k, j) : 0 \leq k \leq i, k \in \mathbb{N}\}$ 
13:  result  $\leftarrow \text{result} \cup \{(\bullet^k, 0, j - k) \mid 0 \leq k \leq j, k \in \mathbb{N}\}$ 
14:  return result
15: else if  $r = r_1 r_2$  then
16:  result  $\leftarrow \emptyset$ 
17:  subs1  $\leftarrow \text{Sb}(r_1, i, j)$ 
18:  for  $(\sigma_1, i_1, j_1) \in \text{subs}_1$  do
19:    if  $j_1 = 0$  then
20:      result  $\leftarrow \text{result} \cup \{(\sigma_1, i_1, j_1)\}$ 
21:    else
22:      result  $\leftarrow \text{result} \cup \{(\sigma_1 \cdot \sigma_2, i_2, j_2) \mid (\sigma_2, i_2, j_2) \in \text{Sb}(r_2, i_1, j_1)\}$ 
23:    end if
24:  end for
25:  return result
26: else if  $r = r_1 \| r_2$  then
27:  return  $\text{Sb}(r_1, i, j) \cup \text{Sb}(r_2, i, j)$ 
28: else if  $r = (r_1)^*$  then
29:  result  $\leftarrow \{(\epsilon, i, j)\}$ 
30:  partial  $\leftarrow \emptyset$ 
31:  repeat
32:    result  $\leftarrow \text{result} \cup \text{partial}$ ;
33:    partial  $\leftarrow \emptyset$ 
34:    for  $(\sigma_n, i_n, j_n) \in \text{result}$  do
35:      for  $(\text{suff}, i_s, j_s) \in \text{Sb}(r_1, i_n, j_n)$  do
36:        if  $(\sigma_n \cdot \text{suff}, i_s, j_s) \notin \text{result}$  then
37:          partial  $\leftarrow \text{partial} \cup \{(\sigma_n \cdot \text{suff}, i_s, j_s)\}$ 
38:        end if
39:      end for
40:    end for
41:  until partial  $\neq \emptyset$ 
42:  return result
43: end if

```

Hence, given  $\text{Sb}(r, i, j)$ , the result is a set of such triples; note that given an element of the resulting set  $(\sigma, n_1, n_2)$ ,  $n_2 = 0$  means that no more characters are needed and  $\sigma$  corresponds to a proper substring of  $r$  from  $i$  to  $j$ . Thus, from the resulting set, we can filter out the partial substrings, and retrieve only proper substrings of  $r$  from  $i$  to  $j$ , by only considering the value of  $n_2$ . Algorithm 4 is defined by case on the structure of the input regex  $r$ :

1.  $j=0$  or  $r = \emptyset$  (Lines 1 and 2):  $\emptyset$  is returned because we either completed the substring or we have no more characters to add.
2.  $r = \sigma \in \Sigma^*$  (Lines 3–10): if  $i > |\sigma|$ , the requested substring happens after this atom, and we return a singleton set  $\{(\epsilon, i - |\sigma|, j)\}$ , thus tracking the consumed characters before the start of the requested substring; if  $i + j > |\sigma|$ , the substring begins in  $\sigma$  but ends in subsequent regexes, and we return a singleton set containing the substring of  $\sigma$  from  $i$  to its end, with  $n_1 = 0$  because we began collecting characters, and  $n_2 = j - |\sigma| + i$  because we collected  $|\sigma| - i$  characters; otherwise, the substring is fully inside  $\sigma$ , and we return the substring of  $\sigma$  from  $i$  to  $i + j$ , setting both  $n_1$  and  $n_2$  to 0.
3.  $r = \top$  (Lines 11–14): because  $r$  might have any length, we generate substrings that (a) gradually consume all the missing characters before the substring can begin (Line 12) and (b) gradually consume all the characters that make up the substring, adding the unknown character  $\bullet$  (Line 13).

4.  $r = r_1 r_2$  (Lines 15–25): the desired substring can either be fully found in  $r_1$  or  $r_2$  or could overlap them; thus, we compute all the partial substrings of  $r_1$ , recursively calling  $Sb$  (Line 17); for all  $\{\sigma_1, i_1, j_1\}$  returned, substrings that are fully contained in  $r_1$  (i.e., when  $j_1 = 0$ ) are added to the result (Line 20) while the remaining ones are joined with ones computed by recursively calling  $Sb$  on  $r_2$  with  $n_1 = j_1$  and  $n_2 = j_2$ .
5.  $r = r_1 || r_2$  (Lines 26 and 27): we return the partial substring of  $r_1$  and the ones of  $r_2$ , recursively calling  $Sb$  on both of them;
6.  $r = (r_1)^*$  (Lines 28–42): we construct the set of substrings through fixpoint iteration, starting by generating  $\{\epsilon, i, j\}$  (corresponding to  $r_1$  repeated 0 times—Line 29) and then, at each iteration, by joining all the partial results obtained until now with the ones generated by a further recursive call to  $Sb$ , keeping only the joined results that are new (Lines 31–42).

Above, we have defined the abstract semantics of `substr` when intervals are constant. When  $[[a_1]]^{\#mm\#} = [i, j]$  and  $[[a_2]]^{\#mm\#} = [l, k]$ , with  $i, j, l, k \in \mathbb{Z}$ , the abstract semantics of `substr` is

$$[[\text{substr}(s, a_1, a_2)]]^{\#mm\#} \triangleq \bigsqcup_{a \in [i, j], b \in [l, k], a \leq b} \text{Min}(\{\sigma \mid (\sigma, 0, 0) \in \text{Sb}(r, a, b - a)\}).$$

We do not report the cases when input intervals are unbounded (e.g.,  $[1, +\infty]$ ). Nevertheless, these cases have been already considered in Arceri et al.<sup>10</sup> and treated analogously in our implementation.

We exploit the abstract semantics of `substr` to instantiate the one of `charAt` as a special case of the former:

$$[[\text{charAt}(s, a_1)]]^{\#mm\#} \triangleq [[\text{substr}(s, a_1, a_1 + 1)]]^{\#mm\#}.$$

## 5 | EXPERIMENTAL RESULTS

TARSIS has been compared with five string abstract domains, namely, the prefix (Pr), suffix (Su), char inclusion (Ci), bricks (Br) domains (all defined in Costantini et al.<sup>9</sup>), and  $FA_{\equiv}$  (defined in Arceri et al.<sup>10</sup>). All domains have been implemented in GoLiSA, which we will briefly introduce before presenting our experimental results. TARSIS and  $FA_{\equiv}$  share a common implementation for the automata structure that is parametric to the alphabet they use. This ensures that performance differences can be accounted on the different size of the automata, eliminating biases that could be introduced with separate implementations having different degrees of optimization.

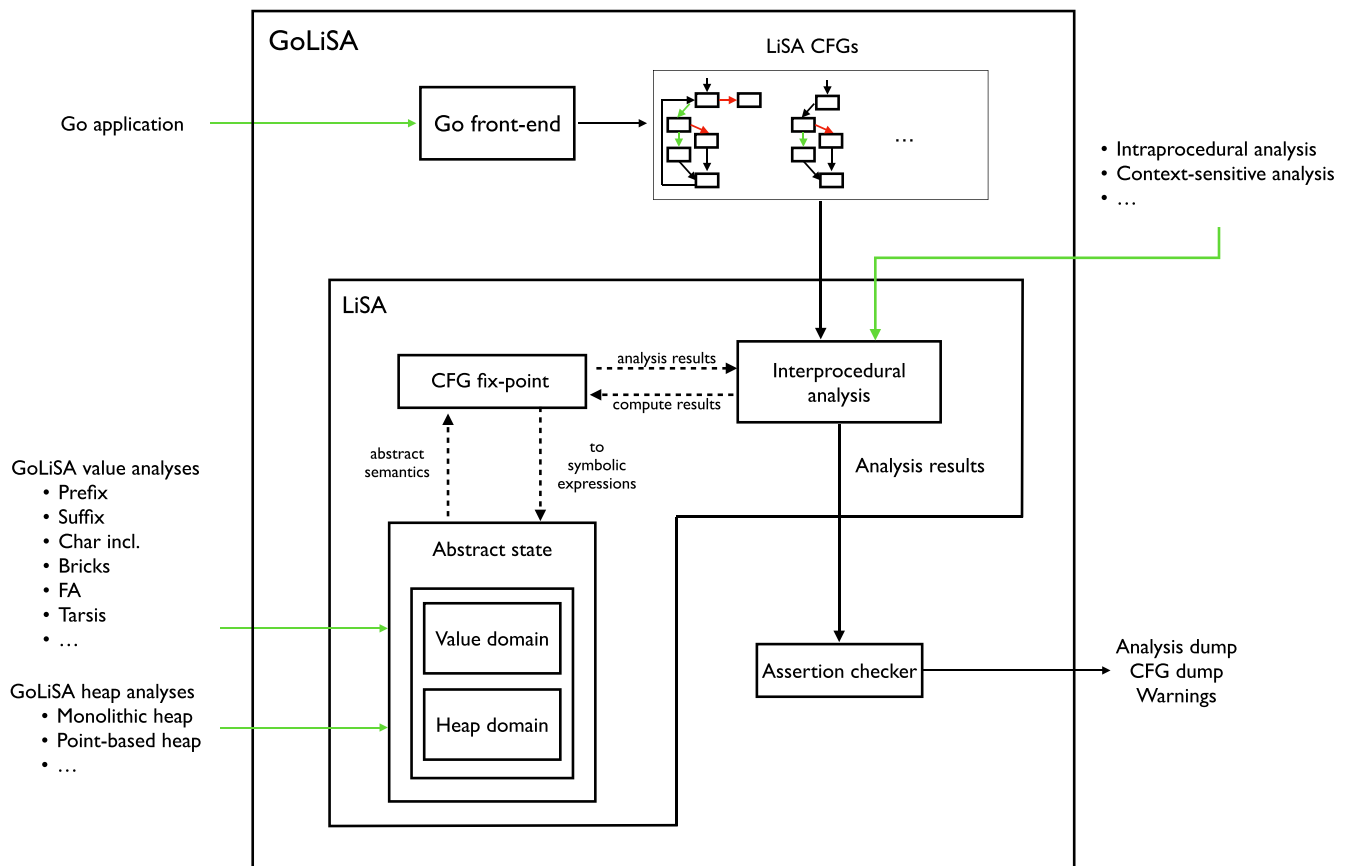
Comparisons have been performed by analyzing the code through the coalesced sum domain specified in Section 4.2 with trace partitioning<sup>26</sup> (note that all traces are merged when evaluating an assertion), plugging in the various string domains. All experiments have been performed on a HP EliteBook G6 machine, with an Intel Core i7-8565U @ 1.8 GHz processor and 16 GB of RAM memory.

To achieve a fair comparison with the other string domains, the subjects of our evaluation are small hand-crafted code fragments that represent standard string manipulations that occur regularly in software. Pr, Su, Ci, and Br have been built to model simple properties and to work with integers instead of intervals and have been evaluated on small programs: Section 5.2 compares them with TARSIS and  $FA_{\equiv}$  without expanding the scope of such evaluations. Section 5.3 instead focuses on slightly more advanced and complex string manipulations that are not modeled by the aforementioned domains, but that  $FA_{\equiv}$  and TARSIS can indeed tackle, highlighting differences between them. Finally, Section 5.5 focuses on the performance difference between  $FA_{\equiv}$  and TARSIS, benchmarking their lattice operations and abstract transformers.

### 5.1 | LiSA and GoLiSA

Experiments presented in this paper have been performed using GoLiSA (<https://github.com/lisa-analyzer/go-lisa>), a static analyzer for Go based on LiSA,<sup>13,14</sup> whose high level infrastructure is visible in Figure 8.

LiSA is a modular framework for developing static analyzers based on the abstract interpretation theory. LiSA analyzes CFGs whose statements do not have a predefined semantics: Instead, users of the framework define custom statement instances implementing language-specific semantic functions, enabling the analysis of a wide range of programming languages and the development of multilanguage analyses. The analysis infrastructure is partitioned into three main areas: call evaluation, memory modeling, and value analysis. Each area corresponds to a separate configurable analysis component that operates agnostically w.r.t. how the others are implemented. The analysis begins in the *Interprocedural Analysis* that executes a program-wide fixpoint by computing each individual CFG's fixpoint. Whenever a call is encountered, the computation of its result is delegated back to the *Interprocedural Analysis*. Instead, non-calling statements are decomposed into a sequence of atomic operations, called *symbolic expressions*, each with a precise semantics that the abstract domains can interpret. Memory-dealing expressions are handled by the *Heap Domain*, tracking their effect and rewriting them as abstract identifiers representing possible memory locations. Finally, the *Value Domain* tracks properties about variables (either program variables or abstract identifiers) and computes invariants for each program point.



**FIGURE 8** Schema of GoLiSA's architecture (taken from Olivieri et al.<sup>12</sup>)

Code parsing logic and the definition of language-specific statements are provided by *Frontends* such as GoLiSA,<sup>12</sup> which can also provide implementations for LiSA's components. Several frontends exist for different languages,<sup>12,27,28</sup> with each constituting effective static analyzers for individual languages that can be combined to obtain multilanguage analyses. In particular, GoLiSA provides a Go-specific *Heap Domain* to accurately track memory operations, while it exploits a context-based *Interprocedural Analysis* provided out-of-the-box by LiSA. All string abstractions considered in the evaluation are implemented as *Value Domains*. Finally, after an analysis is completed, GoLiSA executes the *Assertion Checker*, that is, a program visitor that can access the results of string analyses to raise *definite* alarms (DA for short) when a failing assert (i.e., whose condition is definitely false) is met, or *possible* alarms (PA for short) when the assertion *might* fail (i.e., the assertion's condition evaluates to  $T_{\text{Bool}}$ ). Note that Go does not have built-in `assert` instructions, and we simulate them by invoking a function that `panics` when the given condition is false.

## 5.2 | Precision of the various domains on test cases

We start by considering programs `SUBS` (Figure 9A) and `LOOP` (Figure 9B).

`SUBS` calls `substr` on the concatenation between two strings, where the first is constant and the second one is chosen in a non-deterministic way (i.e., `nondet` condition is statically unknown, Lines 3–7). `LOOP` builds a string by repeatedly appending a suffix, which contains a user input (i.e., an unknown string), to a constant value. Table 1 reports the value approximation for `res` for each abstract domain and analyzed program when the first assertion of each program is met, as well as if the abstract domain precisely dealt with the program assertions. For the sake of readability, `TARSIS` and `FA≡` approximations are expressed as regexes.

When analyzing `SUBS`, both `PR` and `SU` lose precision because the string to append to `res` is statically unknown.

This leads, at Line 8, to a partial substring of the concrete one with `PR` and to an empty string with `SU`. Instead, the `substring` semantics of `CI` moves every character of the receiver in the set of possibly contained ones, thus the abstract value at Line 8 is composed by an empty set of included characters, and a set of possibly included characters containing the ones of both strings. Finally, `BR`, `FA≡`, and `TARSIS` are expressive enough to track any string produced by any concrete execution of `SUBS`.

```

1 func Subs(nondet bool) {
2   res := "substring test"
3   if (nondet) {
4     res = res + " passed"
5   } else {
6     res = res + " failed"
7   }
8   res = res[5:18]
9   assert (strings.Contains(res, "g"));
10  assert (strings.Contains(res, "p"));
11  assert (strings.Contains(res, "f"));
12  assert (strings.Contains(res, "d"));
13 }

```

(A) Program SUBS

```

1 func Loop(value string, nondet bool) {
2   res := "Repeat: "
3   for nondet {
4     res = res + value + "!"
5   }
6   assert (strings.Contains(res, "t"));
7   assert (strings.Contains(res, "!"));
8   assert (strings.Contains(res, "f"));
9 }

```

(B) Program LOOP

**FIGURE 9** Program samples used for domain comparison.**TABLE 1** Values of `res` at the first assert of each program.

Domain	Program Subs		Program Loop	
Pr	ring test	×	Repeat :	×
Su	ε	×	ε	×
Ci	[[abdefgilnrstu]	×	[ : aepRt ] [ ! : aepRt T ]	✓
Br	[[ring test fai, ring test pas]](1,1)	✓	[[T]](0, +∞)	×
FA <sub>≡</sub>	ring test (pas    fai)	✓	Repeat : (T)*	✓
TARSIS	(ring test pas    ring test fai)	✓	Repeat : (T!)*	✓

When evaluating the assertions of `Subs`, a PA should be raised on Lines 10 and 11, because  $\sim p \sim$  or  $\sim f \sim$  might be in `res`, together with a DA alarm on Line 12, because  $\sim d \sim$  is surely not contained in `res`. No alarm should be raised on Line 9 instead, because  $\sim g \sim$  is part of the common prefix of both branches and thus will be included in the substring. Such behavior is achieved when using `Br`, `FA≡`, or `TARSIS`. Because the `substring` semantics of `Ci` moves all characters to the set of possibly contained ones, PAs are raised on all four assertions.

Moreover, `Su` loses all information about `res`, PAs are raised on Lines 9–12 when using such domain. `Pr` instead tracks the definite prefix of `res`; thus, the PA at Line 9 is avoided.

When analyzing `Loop`, we expect to obtain no alarm at Line 6 (because character  $\sim t \sim$  is always contained in the resulting string value) and PAs at Lines 7 and 8. `Pr` infers as prefix of `res` the string  $\sim \text{Repeat} : \sim$ , keeping such value for the analysis of the whole program. This allows the analyzer to prove the assertion at Line 6, but it raises PAs when it checks the ones at Lines 7 and 8.

Again, `Su` loses all information about `res` because the lub operation occurring at Line 3 cannot find a common suffix between  $\sim \text{Repeat} : \sim$  and  $\sim ! \sim$ ; hence, PAs are raised on Lines 6–8. Because the set of possible characters contains `T`, `Ci` can correctly state that any character might appear in the string. For this reason, two PAs are reported on Lines 7 and 8, while no alarm is raised on Line 6 (again, this is possible because the string used in the `contains` call has length 1). The alternation of `T` and  $\sim ! \sim$  prevents `Br` normalization algorithm from merging similar bricks. This will eventually lead to overcoming the length threshold  $k_L$ , hence resulting in the  $[[T]](0, +\infty)$  abstract value. In such a situation, `Br` returns  $T_{\text{Bool}}$  on all `contains` calls, resulting in PAs on Lines 6–8. The parametric widening of `FA≡` collapses the colon into `T`. In `TARSIS`, because the automaton representing `res` grows by two states each iteration, the parametric widening defined in Section 4.1 can collapse the whole content of the loop into a two-state loop recognizing `T!`. The precise approximation of `res` of both domains enable the analyzer to detect that the assertion at Line 6 always holds, while PAs are raised on Lines 7 and 8.

In summary, `Pr` and `Su` failed to produce the expected results on both `Subs` and `Loop`, while `Ci` and `Br` produced exact results in one case (`Loop` and `Subs`, respectively) but not in the other. Hence, `FA≡` and `TARSIS` were the two only domains that produced the desired behavior in these rather simple test cases.

### 5.3 | Evaluation on realistic code samples

In this section, we explore two real world code samples. Method `ToSTRING` (Figure 10A) transforms an array of names that come as string values into a single string. While it resembles the code of `Loop` in Figure 9B (thus, results of all the analyses show the same strengths and weaknesses), now assertions check `contains` predicates with a multi-character string.



Method `COUNTMATCHES` (Figure 10B) makes use of `strings.Count` (reported in Section 2) to prove properties about its return value.

Table 2 reports the results of both methods (stored in `res` and `count`, respectively) evaluated by each analysis at the first assertion, as well as if the abstract domain precisely dealt with the program assertions.

As expected, when analyzing `TOSTRING`, each domain showed results similar to those of `LOOP`. In particular, we expect to obtain no alarm at Line 12 (because  $\sim\text{People}\sim$  is surely contained in the resulting string) and two PAs at Lines 13 and 14. `PR`, `SU`, `CI`, and `BR` raise PAs on all the three assert statements. `FA≡` and `TARSIS` detect that the assertion at Line 12 always holds. Thus, when using them, the analyzer raises PAs on Lines 13 and 14 because (i) comma character is part of `res` if the loop is iterated at least once and (ii) `T` might match  $\sim\text{not}\sim$ .

If `COUNTMATCHES` was to be executed, `count` would be either 2 or 3 when the first assertion is reached, depending on the choice of `str`. Thus, no alarm should be raised at Line 9, while a DA should be raised on Line 10 and a PA on Line 11. Because `PR`, `SU`, `CI`, and `BR` do not define most of the operations used in the code, the analyzer does not have information about the string on which `strings.Count` is executed and thus abstract `count` with the interval  $[0, +\infty]$ . Thus, PAs are raised on Lines 9–11. Instead, `FA≡` and `TARSIS` are instead able to detect that `sub` is present in all the possible strings represented by `str`. Thus, thanks to trace partitioning, the trace where the loop is skipped and `count` remains 0 gets discarded. Then, when the first `indexOf` call happens,  $[0,0]$  is stored into `idx`, because all possible values of `str` start with `sub`. Because the call to `length` yields  $[10,17]$ , all possible substrings from  $[2,2]$  (`idx` plus the length of `sub`) to  $[10,17]$  are computed (viz.,  $\sim\text{e throat}\sim$ ,  $\sim\text{is is th}\sim$ ,  $\sim\text{is is the}\sim$ , ...,  $\sim\text{is is the thing}\sim$ ), and the resulting automaton is the one that recognizes all of them. Because the value of `sub` is still contained in every path of such automaton, the loop guard still holds and the second iteration is analyzed, repeating the same operations. When the loop guard is reached for the third time, the remaining substring of the shorter starting string (viz.,  $\sim\text{roat}\sim$ ) recognized by the automaton representing `str` will no longer contain `sub`: a trace where `count` equals  $[2,2]$  will leave the loop. A further iteration is then analyzed, after which `sub` is no longer contained in any of the strings that `str` might hold. Thus, a second and final trace where `count` equals  $[3,3]$  will reach the assertions and will be merged by interval `lub`, obtaining  $[2,3]$  as final value for `count`. This allows `TARSIS` and `FA≡` to identify that the assertion at Line 10 never holds, raising a DA, while the one at Line 11 might not hold, raising a PA.

## 5.4 | Efficiency w.r.t. simpler string domains

The detailed analysis of two test cases and two examples taken from real-world code underlined that `TARSIS` and `FA≡` are the only ones able to obtain precise results on them. We now discuss the efficiency of the analyses. Table 3 reports the execution times for all the domains on the case

<pre> 1 func ToString(names []string) { 2   res := "People : {" 3   i := 0; 4   for i &lt; len(names) { 5     res = res + names[i]; 6     if i != len(names) - 1 { 7       res = res + ","; 8     } 9     i++; 10  } 11  res = res + "}"; 12  assert (strings.Contains(res, "People")); 13  assert (strings.Contains(res, ",")); 14  assert (strings.Contains(res, "not")); 15 } </pre> <p style="text-align: center;">(A) Program <code>TOSTRING</code></p>	<pre> 1 func CountMatches(nondet boolean) { 2   str := "" 3   if nondet { 4     str = "this is the thing"; 5   } else { 6     str = "the throat"; 7   } 8   count := strings.Count(str, "th"); 9   assert (count &gt; 0); 10  assert (count == 0); 11  assert (count == 3); 12 } </pre> <p style="text-align: center;">(B) Program <code>COUNTMATCHES</code></p>
--	--

**FIGURE 10** Programs used for assessing domain precision.

**TABLE 2** Values of `res` and `CountMatches` at the first assert of the respective program.

Domain	Program <code>TOSTRING</code>		Program <code>COUNTMATCHES</code>	
<code>PR</code>	<code>People : {</code>	×	$[0, +\infty]$	×
<code>SU</code>	$\epsilon$	×	$[0, +\infty]$	×
<code>CI</code>	$\{\{\}: \text{Peopl} \} \{\{\}: \text{Peopl } T\}$	×	$[0, +\infty]$	×
<code>BR</code>	$\{\{T\}\}(0, +\infty)$	×	$[0, +\infty]$	×
<code>FA<sub>≡</sub></code>	<code>People : {(T)*T}</code>	✓	$[2,3]$	✓
<code>TARSIS</code>	<code>People : {}    People : {(T)*T}</code>	✓	$[2,3]$	✓

**TABLE 3** Execution times of the domains on each program.

Domain	Subs	Loop	ToSTRING	COUNTMATCHES
Pr	5 ms	28 ms	1 s 509 ms	129 ms
Su	5 ms	32 ms	1 s 599 ms	110 ms
Ci	4 ms	45 ms	1 s 594 ms	138 ms
Br	8 ms	78 ms	3 s 911 ms	345 ms
FA <sub>≡</sub>	11 ms	20 s 895 ms	22 m 37 s 621 ms	9 s 175 ms
TARSIS	7 ms	443 ms	12 s 437 ms	123 ms

studies analyzed in this section, by taking the time that GoLiSA reports as the one needed to compute a program-wide fixpoint (thus excluding the time for booting up the analysis, parsing the code and dumping the results). Note that the times reported here are higher than the ones of the original paper due to the usage of a complete static analyzer: Memory abstractions and call resolution were not performed by the prototypical analyzer used in Negrini et al.<sup>15</sup> Overall, Pr, Su, Ci, and Br are the fastest domains with execution times usually in the order of milliseconds, with the exception of ToSTRING that proved more challenging for all domains. Thus, if on the one hand these domains failed to prove some of the properties of interest, they are quite efficient and they might be helpful to prove simple properties. TARSIS execution times are sometimes higher but still comparable with them. Instead, FA<sub>≡</sub> blows up on three out of the four test cases (and in particular on ToSTRING). Hence, TARSIS is the only domain that executes the analysis in a limited time while being able to prove all the properties of interest on these four case studies.

The reason behind the performance gap between TARSIS and FA<sub>≡</sub> can be accounted on the alphabets underlying the automata. In FA<sub>≡</sub>, automata are built over an alphabet of single characters. While this simplifies the semantic operations, it also causes state and transition blow up w.r.t. the size of the string that needs to be represented. This does not happen in TARSIS, because atomic strings (not built through concatenation or other string manipulations) are part of the alphabet and can be used as transition symbol. Having less states and transitions to operate upon drastically lowers the time and memory requirements of automata operations, making TARSIS faster than FA<sub>≡</sub>.

TARSIS's alphabet has another peculiarity w.r.t. FA<sub>≡</sub>'s: It has a special symbol for representing the unknown string. Having such a symbol requires some fine-tuning of the algorithms to have them behave differently when the symbol is encountered, but without additional tolls on their performances. FA<sub>≡</sub>'s alphabet does not have such a symbol, thus representing the unknown string is achieved through a state having one self-loop for each character in the alphabet (including the empty string). This requires significantly more resources for automata algorithms, leading to higher execution times.

It is important to notice that performances of programs relying on automata are heavily dependent on their implementation. Both FA<sub>≡</sub> and TARSIS come as non-optimized implementations whose performances can be greatly improved, thus further reducing the gap between them and the simpler string abstractions. The source code of FA<sub>≡</sub> is available at <https://github.com/SPY-Lab/fsa>, while TARSIS's implementation is published at <https://github.com/UniVE-SSV/tarsis>. Instead, implementations used in the experiments are part of LiSA available at <https://github.com/lisa-analyzer/lisa>.

## 5.5 | Performance benchmark of TARSIS and FA<sub>≡</sub>

As automata-based abstractions have already been proved to be effective in string program analysis,<sup>10</sup> we focus the final part of our evaluation on the difference in resource consumption between TARSIS and FA<sub>≡</sub>. In fact, our main goal is proving that the adoption of TARSIS can make automata-based abstractions viable for the analysis of non-trivial code.

To measure the performance of TARSIS and FA<sub>≡</sub>, one could track the resource consumption of full program analyses employing the two domains and reason about their difference. However, this could produce misleading results: the measured performance would be affected by the ones of other analysis components (e.g., memory abstractions and interprocedural analyses) and would hence be affected by their running time and memory consumption. The different semantics of TARSIS and FA<sub>≡</sub> would thus only account for a small portion of each measurement, hindering the purpose of our experiments. Hence, to ensure that we can precisely measure only the performance of interest, that is, the ones concerning TARSIS and FA<sub>≡</sub>, we directly compare each lattice operation and abstract transformer in isolation, benchmarking their execution times. This enables accurate speedup measurements, because such times are not affected by external factors like memory saturation caused by the remaining analysis components.

We compare lattice operators  $\sqsubseteq$ ,  $\sqcup$ ,  $\sqcap$ , and  $\nabla$ , together with the semantics of string operators (`concat`, `contains`, `length`, `substr`, `trim`, `repeat`, and `indexOf`); because `trim`'s semantics relies on `trimLeft` and `trimRight`, the comparison of the latter operations have been omitted. Moreover, concerning the `trimLeft`, `trimRight`, and `trim` abstract semantics of FA<sub>≡</sub>, we adopted the same one reported in Section 4.2 for TARSIS, recasted to standard automata, because a bug was found in the original FA<sub>≡</sub> abstract semantics when automata with cycles

were involved, leading to unsound results. For instance, let us consider the automaton recognizing the language  $\{ (a)^n \mid n \in \mathbb{N} \}$ , the original `trim` abstract semantics of  $FA_{/\equiv}$  returns the automaton recognizing the language  $\{ a^n \mid n \in \mathbb{N} \}$ , which is an unsound result (e.g., the string  $a$  is not recognized by the resulting automaton).

During an actual program analysis, the target operators would be invoked by a fixpoint engine on operands built through several string manipulations. Each such operand would thus be an automaton with an arbitrarily complex structure, depending on the combination of operators used. To ensure that our experiment measures the performance of each domain's operands on all possible input structures, we identify and define here seven automata classes, each containing automata sharing a common structure, on which lattice operators and string operations can be applied to. Then, we compose a benchmark by using automata from all classes, thus ensuring that the measurements will take into account performance on different automata structures. In the following, we describe each class reporting an exemplification of a Go fragment that may generate an automaton belonging to each class (where  $b$  is a statically unknown Boolean value), together with the automaton computed by TARSIS and  $FA_{/\equiv}$  for the string variable  $x$  at the end of the fragment.

1. *Statically unknown strings*. Such class contains the abstractions of statically unknown strings, such as a user input. Because both domains have a unique minimum automaton for such strings, this class is made of a single automaton.

Go fragment	TARSIS	$FA_{/\equiv}$
<code>_, err := fmt.Sprintf("%s", &amp;x);</code>		

2. *Constant strings*. This class contains automata generated as abstractions of string literals appearing in the program.

Go fragment	TARSIS	$FA_{/\equiv}$
<code>x := "foo";</code>		

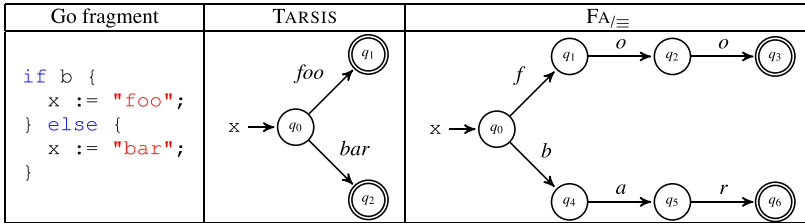
3. *Concatenated strings*. This class contains automata modeling the concatenation of simple strings, with each being either statically unknown or constant. Automata in this class are effective concatenations of ones from the previous classes.

Go fragment	TARSIS	$FA_{/\equiv}$
<code>z := "in:"; _, err := fmt.Sprintf("%s", &amp;y); x := z + y;</code>		

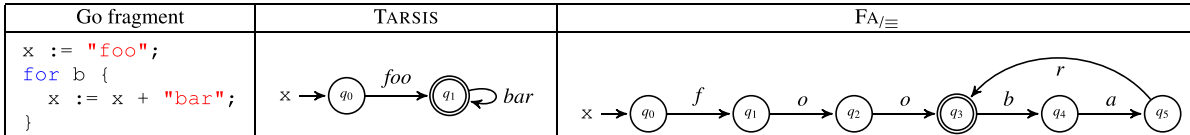
4. *Increasing strings*. The automata in this class model sets of strings built by optionally appending a finite number of strings at the end of an existing one, thus leading to single path automata. Each of such strings can be either statically unknown or constant (i.e., each automaton is part of one of the first two classes).

Go fragment	TARSIS	$FA_{/\equiv}$
<code>x := "foo"; _, err := fmt.Sprintf("%s", &amp;y); if b {   x := x + y; }</code>		

5. *Disjoint strings*. This class contains automata built as the union of up to four different automata coming from the third class, thus modeling their least upper bound. Automata in this class represent the result of merging branches where string variables were assigned to different values. Each string used in the lub can thus be either statically unknown, constant, or resulting from a concatenation.



6. *Looping strings.* The automata in this class are built by inserting up to two loops in the automata of the third class, thus modeling strings that are built inside a loop.



7. *Random strings.* This class contains automata that are built using an arbitrary set of manipulations and thus have no predefined structure. They represent worst-case scenarios for both domains, and thus, automata of this class can be significantly complex. As this class does not follow any particular structure, examples automata with their generating code are omitted.

### 5.5.1 | Benchmark composition

We benchmark each lattice and string operators mentioned above by executing them 100 times (with each execution referred to as a *round*), each one using randomly generated automata from one of the above classes for all the required inputs (e.g., a round executing `concat` on automata of Class 2 will generate two automata  $A_1$  and  $A_2$  from that class and would then run `concat(A_1, A_2)`). Specifically, one round is executed using the single automaton from Class 1, five rounds with automata from Class 2, 10 rounds are executed for Classes 3 to 6, while the remaining 54 rounds use automata from Class 7. Each round consists of (i) the generation of the TARSIS automata, (ii) their conversion to the equivalent  $FA_{\equiv}$  ones, and (iii) the execution of each operation individually, measuring the required time. Each operation's execution has a 30 s timeout, as we deem it a reasonable time bound for an operation to complete when running as part of an analysis. Before executing the benchmark, a warm-up iteration is executed to ensure that setup operations of the JVM would not impact the actual measurement. Finally, integer indexes for the substring operation were randomly chosen between 0 and 20.

As the size of  $FA_{\equiv}$  automata grows fast, we impose some limits on the structure of each generated TARSIS automaton, such that, when we generate the equivalent  $FA_{\equiv}$  automaton, its size is limited. This prevents an excessive amount of timeouts that would invalidate our experiments. Specifically:

- all atomic strings (i.e., the ones used as symbols in TARSIS's transitions) have a random length of up to 10 characters;
- each atomic string has a 10% chance of being statically unknown;
- intermediate states of single path automata (Class 4) have a 50% chance of being an additional accepting state; and
- automata of Class 7 can contain up to five states with up to three transitions per state; each state has a 25% chance of being an additional final state.

### 5.5.2 | Benchmark results

Table 4 reports, for each operation and domain, the number of successful rounds and the one of timed-out ones, together with the total, average, minimum, and maximum execution times of the successful rounds. Finally, the speedup of TARSIS w.r.t.  $FA_{\equiv}$  is determined by only comparing execution times of the rounds where both domains ran successfully, reporting the ratio  $t_{FA_{\equiv}}/t_{TARSIS}$ .

The benchmark confirms that TARSIS is overall reliably faster than  $FA_{\equiv}$ . TARSIS operations time out significantly less than the ones of  $FA_{\equiv}$  (two timeouts instead of 137), reporting lower total times despite the higher number of automata tackled. Worst-case performances are also improved, with most operations requiring far less time with TARSIS even on complex automata, as highlighted by column *Max. time* of Table 4. Moreover,

**TABLE 4** Lattice operators and abstract transformers benchmark results.

Operation	Domain	# Successes	# Timeouts	Total time	Avg. time	Min. time	Max. time	Time ratio
$\sqsubseteq$	TARSIS	100	0	2 s 37 ms	20 ms	<1 ms	403 ms	920.74×
	$FA_{\perp}$	88	12	2 m 27 s 815 ms	1 s 679 ms	<1 ms	18 s 685 ms	
$\sqcup$	TARSIS	100	0	29 ms	<1 ms	<1 ms	10 ms	114.98×
	$FA_{\perp}$	100	0	3 s 430 ms	34 ms	<1 ms	945 ms 31 ms	
$\cap$	TARSIS	100	0	4 s 905 ms	<1 ms	<1 ms	8 ms	29.37×
	$FA_{\perp}$	100	0	760 ms	7 ms	<1 ms	426 ms	
$\nabla$	TARSIS	100	0	148 ms	1 ms	<1 ms	27 ms	97.70×
	$FA_{\perp}$	96	4	13 s 485 ms	140 ms	<1 ms	3 s 522 ms	
concat	TARSIS	100	0	40 ms	<1 ms	<1 ms	5 ms	528.42×
	$FA_{\perp}$	100	0	21 s 537 ms	215 ms	<1 ms	14 s 177 ms	
contains	TARSIS	100	0	2 s 187 ms	21 ms	<1 ms	668 ms	7.16×
	$FA_{\perp}$	99	1	10 s 888 ms	109 ms	<1 ms	1 s 922 ms	
length	TARSIS	100	0	143 ms	1 ms	<1 ms	36 ms	2206.42×
	$FA_{\perp}$	74	26	1 m 19 s 427 ms	1 s 73 ms	<1 ms	15 s 454 ms	
indexOf	TARSIS	100	0	6 s 575 ms	65 ms	<1 ms	2 s 222 ms	1.72×
	$FA_{\perp}$	100	0	11 s 352 ms	113 ms	<1 ms	10 s 693 ms	
substr	TARSIS	100	0	22 s 757 ms	227 ms	<1 ms	13 s 807 ms	465.37×
	$FA_{\perp}$	53	47	3 m 49 s 127 ms	4 s 323 ms	<1 ms	32 s 998 ms	
replace	TARSIS	99	1	326 ms	3 ms	<1 ms	47 ms	7.98×
	$FA_{\perp}$	99	1	2 s 440 ms	24 ms	<1 ms	1 s 470 ms	
trim	TARSIS	99	1	20 s 759 ms	209 ms	<1 ms	18 s 270 ms	486.42×
	$FA_{\perp}$	57	43	1 m 38 s 269 ms	1 s 724 ms	<1 ms	22 s 760 ms	
repeat	TARSIS	100	0	37 ms	209 ms	<1 ms	895 ms	3.93×
	$FA_{\perp}$	97	3	12 s 488 ms	128 ms	<1 ms	2 s 290 ms	

when focusing only on rounds where both domains succeeded within the required time bound, the speedup (column *Time ratio* of Table 4) is still noticeable: Even in the worst case of `indexOf`, TARSIS still performs almost twice as fast as  $FA_{\perp}$ .

Let us comment the results concerning `length` and `trim`. For these operations, TARSIS and  $FA_{\perp}$  use the same algorithms to implement the corresponding abstract semantics. Still, TARSIS incurs in a single timeout when executing `trim` (against 43 timeouts of  $FA_{\perp}$ ) and in none when executing `length` (while  $FA_{\perp}$  incurs in 26 timeouts), also enabling significant speedup (more than 2000× for `length` and almost 500× for `trim`). This highlights that even when using the same abstract semantics for both domains, the alphabet chosen for TARSIS is still able to drastically reduce the required resources that a static analysis needs to analyze strings.

## 6 | RELATED WORK

The problem of statically analyzing strings has been already tackled in different contexts in the literature during the last two decades.<sup>9–11,29–32</sup> As already discussed, the TARSIS abstract domain built upon the FSA abstract domain defined in Arceri et al<sup>10</sup> in the context of dynamic languages, providing automata-based abstract semantics for common ECMAScript string operations. The same abstract domain has been integrated also for defining a sound-by-construction analysis for string-to-code statements.<sup>33</sup> As reported by the experimental results in Section 5 (Table 3 in particular), TARSIS is quite more efficient than  $FA_{\perp}$  while keeping (almost) the same level of precision.

Generally speaking, the practical comparison of different string static analyses is particularly challenging because of (i) the lack of standard benchmarks, (ii) the variety of the information that can be tracked over string values (e.g., included characters, their order, regular expressions, and substring relations), and (iii) the availability of the implementations of existing analyses (often formalized several years ago and not actively maintained). For these reasons, in the rest of this section, we qualitatively discuss the differences between TARSIS and other related works, but we do not experimentally compare them.

## 6.1 | Static analysis of string values

Our experimental results compared TARSIS with several simple abstract domains introduced in Costantini et al.<sup>9,34</sup> Those domains track relatively simple information about string values, such as characters included or not, prefixes and suffixes, and concatenation of constant string values. As reported in Tables 1–3, these abstractions are quite more efficient but less precise than TARSIS. Overall, TARSIS exposed execution times comparable with the ones of Costantini et al.<sup>9,34</sup> and precision similar to Arceri et al.<sup>10,33</sup>

Static analysis was applied to string values in many different contexts, such as SQL queries programmatically built by code,<sup>3</sup> reflection,<sup>4,5</sup> and injection vulnerabilities.<sup>7,8</sup> Instead, the main goal of our work is to introduce a novel abstract domain (i.e., approximation of string values) that outperforms state-of-the-art approaches in terms of precision or efficiency, and this can be applied to different contexts.

## 6.2 | Automata-based approaches

Several approaches to statically analyze string values through automata (like TARSIS) have been proposed in the literature. For instance, the authors of Almashfi and Lu<sup>35</sup> provided an automata abstraction merged with interval abstractions for analyzing JavaScript arrays and objects. Another interesting automata-based model is symbolic automata,<sup>36</sup> which differs from the standard one having an alphabet of predicates (that can potentially be infinite) instead of single characters. Examples of applications of symbolic automata in the context of static analysis are regex processing, sanitizer analysis,<sup>37</sup> and their usage as program models for mixing syntactic and semantic abstractions over the program.<sup>38</sup> Generally speaking, automata-based approaches are usually not efficient, because they need to perform algorithmically complex operations on the automata. Our approach is aimed at (partially) overcoming such limits by working on an alphabet of strings instead of single characters.

## 6.3 | Regular expressions

A major stream of research focused on the abstraction of string values through regular expressions. In particular, Christensen et al.<sup>30</sup> proposed a static analysis of Java strings based on the abstraction of the control-flow graph as a context-free grammar. By relying on regular expressions, this approach is in position to track information not only on constants inside the string values like TARSIS but also to approximate ranges of possible characters such as  $0|(-?[1-9][0-9]^*)$  (i.e., the string representation of integer values). However, such approximation is strictly more complex and led to less efficient analyses.

Similarly, regular strings<sup>39</sup> is an abstraction of the FSA domain and approximates strings as a strict subset of regular expressions. The authors introduced an aggressive widening operator that improves the efficiency but worsens the precision of the analysis.

## 6.4 | String constraints verification

Another major research effort was spent on the context of string constraints verification by the investigation and development of various techniques and tools focused on the study of decidable fragments of string constraint formulas<sup>40</sup> or proposing new efficient decidable procedures or string constraints representations<sup>41,42</sup> also based on automata, such as Wang et al.<sup>43</sup> and Yu et al.,<sup>44</sup> or involving type conversion string constraints.<sup>45</sup> For instance, Z3-str<sup>46</sup> extended the SMT-solver Z3<sup>47</sup> by treating strings as primitive types supporting the most common operations over strings (e.g., concatenation and substring). All those approaches allow to track very precise information over string values but usually require manually annotating some portions of the code (e.g., through loop invariants and pre- and post-conditions) and require solving NP-complete problems (e.g., SAT solving) causing slowdowns (or timeouts) of the analyses in some (hopefully corner) cases.

## 7 | CONCLUSION

In this paper, we introduced TARSIS, an abstract domain for sound abstraction of string values. TARSIS is based on FSAs paired with their equivalent regular expression: a representation that allows precise modeling of complex string values. Experiments show that TARSIS achieves great precision also on code that heavily manipulate string values, while the time needed for the analysis is comparable with the one of other simpler domains.

In order to enforce loop convergence, our analysis has been equipped with a widening with threshold operator over TARSIS automata. As it usually happens in abstract interpretation, in order to retrieve some information lost by the widening application, the analysis can be equipped also with a narrowing.<sup>23</sup> Hence, a narrowing operator for TARSIS will be studied, in order to get more precise results on loop analyses. Another improvement with respect to the efficiency of the analysis can be reached defining a split operator<sup>48</sup> for TARSIS, when analyzing conditional guards;

this would allow TARSIS to achieve significant efficiency improvements with respect to the classical approach based on the filter operator. Moreover, TARSIS' semantics can also be employed as summaries for library functions that manipulate strings: When one such function is used, it can be automatically identified by the analysis engine,<sup>49</sup> and our proposed abstract semantics can be applied instead of analyzing the function's code.

Finally, TARSIS provides a non-relational string abstraction: As such, information relating different variables is not modeled inside the domain. We thus intend on working on a relational extension of TARSIS that is also able to relate different variables, also extending existing solutions based on combining TARSIS with other domains.<sup>50</sup>

## ACKNOWLEDGMENTS

This work is partially supported by Bando di Ateneo per la Ricerca 2022, funded by University of Parma (MUR\_DM737\_2022\_FIL\_PROGETTI\_B\_ARCERI\_COFIN and CUP: D91B21005370003); "Formal verification of GPLs blockchain smart contracts," SERICS (PE00000014) under the NRRP MUR program funded by the EU–NGEU; iNEST-Interconnected NordEst Innovation Ecosystem funded by PNRR (Mission 4.2, Investment 49 1.5) NextGeneration EU—Project ID: ECS00000043; and SPIN-2021 "Static Analysis for Data Scientists" funded by Ca' Foscari University.

## DATA AVAILABILITY STATEMENT

The data that support the findings of this study are openly available in LiSA at <https://github.com/lisa-analyzer/lisa>.

## ORCID

Luca Negrini  <https://orcid.org/0000-0001-9930-8854>

Vincenzo Arceri  <https://orcid.org/0000-0002-5150-0393>

## REFERENCES

- Cousot P, Cousot R. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In: Graham RM, Harrison MA, Sethi R, eds. *Conference Record of the Fourth ACM Symposium on Principles of Programming Languages, Los Angeles, California, USA*: ACM; 1977:238-252. <https://doi.org/10.1145/512950.512973>
- Cousot P, Cousot R. Systematic design of program analysis frameworks. In: Aho AV, Zilles SN, Rosen BK, eds. *Conference Record of the Sixth Annual ACM Symposium on Principles of Programming Languages, San Antonio, Texas, USA*: ACM Press; 1979:269-282. <https://doi.org/10.1145/567752.567778>
- Gould C, Su Z, Devanbu PT. JDBC checker: a static analysis tool for SQL/JDBC applications. In: Finkelstein A, Estublier J, Rosenblum DS, eds. *26th International Conference on Software Engineering (ICSE 2004), 23-28 May 2004, Edinburgh, United Kingdom*: IEEE Computer Society; 2004:697-698. <https://doi.org/10.1109/ICSE.2004.1317494>
- Landman D, Serebrenik A, Vinju JJ. Challenges for static analysis of Java reflection: literature review and empirical study. In: Uchitel S, Orso A, Robillard MP, eds. *Proceedings of the 39th International Conference on Software Engineering, ICSE 2017, Buenos Aires, Argentina*: IEEE/ACM; 2017:507-518. <https://doi.org/10.1109/ICSE.2017.53>
- Bodden E, Sewe A, Sinschek J, Oueslati H, Mezini M. Taming reflection: aiding static analysis in the presence of reflection and custom class loaders. In: Taylor RN, Gall HC, Medvidovic N, eds. *Proceedings of the 33rd International Conference on Software Engineering, ICSE 2011, Waikiki, Honolulu, HI, USA*: ACM; 2011:241-250. <https://doi.org/10.1145/1985793.1985827>
- Arceri V, Mastroeni I. A sound abstract interpreter for dynamic code. In: Hung C-C, Cerný T, Shin D, Bechini A, eds. *SAC '20: The 35th ACM/SIGAPP Symposium on Applied Computing, Online Event, [Brno, Czech Republic]*: ACM; 2020:1979-1988. <https://doi.org/10.1145/3341105.3373964>
- Fu X, Lu X, Peltsverger B, Chen S, Qian K, Tao L. A static analysis framework for detecting SQL injection vulnerabilities. In: 31st Annual International Computer Software and Applications Conference, COMPSAC 2007, Beijing, China, Vol. 1. IEEE Computer Society; 2007:87-96. <https://doi.org/10.1109/COMPSAC.2007.43>
- Livshits VB, Lam MS. Finding security vulnerabilities in Java applications with static analysis. In: McDaniel PD, ed. *Proceedings of the 14th USENIX Security Symposium, Baltimore, MD, USA*: USENIX Association; 2005. <https://www.usenix.org/conference/14th-usenix-security-symposium/finding-security-vulnerabilities-java-applications-static>
- Costantini G, Ferrara P, Cortesi A. A suite of abstract domains for static analysis of string values. *Softw Pract Exp*. 2015;45(2):245-287. <https://doi.org/10.1002/spe.2218>
- Arceri V, Mastroeni I, Xu S. Static analysis for ECMAScript string manipulation programs. *Appl Sci*. 2020;10:3525. <https://doi.org/10.3390/app10103525>
- Cortesi A, Olliaro M. M-string segmentation: a refined abstract domain for string analysis in C programs. In: Pang J, Zhang C, He J, Weng J, eds. *2018 International Symposium on Theoretical Aspects of Software Engineering, TASE 2018, Guangzhou, China*: IEEE Computer Society; 2018:1-8.
- Olivieri L, Negrini L, Arceri V, Tagliaferro F, Ferrara P, Cortesi A, Spoto F. Information flow analysis for detecting non-determinism in blockchain. In: Ali K, Salvaneschi G, eds. *37th European Conference on Object-Oriented Programming, ECOOP 2023, Seattle, Washington, United States, LIPIcs, vol. 263: Schloss Dagstuhl—Leibniz-Zentrum für Informatik*; 2023:23:1-23:25. <https://doi.org/10.4230/LIPIcs.ECOOP.2023.23>
- Ferrara P, Negrini L, Arceri V, Cortesi A. Static analysis for dummies: experiencing LiSA. In: Do LNQ, Urban C, eds. *SOAP@PLDI 2021: Proceedings of the 10th ACM SIGPLAN International Workshop on the State of the Art in Program Analysis, Virtual Event, Canada*: ACM; 2021:1-6. <https://doi.org/10.1145/3460946.3464316>
- Negrini L, Ferrara P, Arceri V, Cortesi A. LiSA: a generic framework for multilanguage static analysis. In: Arceri V, Cortesi A, Ferrara P, Olliaro M, eds. *Challenges of Software Verification*: Springer Nature Singapore; 2023:19-42. [https://doi.org/10.1007/978-981-19-9601-6\\_2](https://doi.org/10.1007/978-981-19-9601-6_2)

15. Negrini L, Arceri V, Ferrara P, Cortesi A. Twinning automata and regular expressions for string static analysis. In: Verification, Model Checking, and Abstract Interpretation—22nd International Conference, VMCAI 2021, Proceedings Henglein F, Shoham S, Vizek Y, eds., Lecture Notes in Computer Science, vol. 12597. Springer; 2021:267-290. [https://doi.org/10.1007/978-3-030-67067-2\\_13](https://doi.org/10.1007/978-3-030-67067-2_13)
16. Giacobazzi R, Ranzato F, Scozzari F. Making abstract interpretations complete. *J ACM*. 2000;47(2):361-416. <https://doi.org/10.1145/333979.333989>
17. Bartzis C, Bultan T. Widening arithmetic automata. In: Computer Aided Verification, 16th International Conference, CAV 2004, Proceedings Alur R, Peled DA, eds., Lecture Notes in Computer Science, vol. 3114. Springer; 2004:321-333. [https://doi.org/10.1007/978-3-540-27813-9\\_25](https://doi.org/10.1007/978-3-540-27813-9_25)
18. D'Silva V. Widening for Automata. *MSc Thesis: Inst. Fur Inform.—UZH*; 2006.
19. Cousot P, Cousot R. Abstract interpretation frameworks. *J Log Comput*. 1992;2(4):511-547. <https://doi.org/10.1093/logcom/2.4.511>
20. Bagnara R, Hill PM, Zaffanella E. The Parma Polyhedra Library: toward a complete set of numerical abstractions for the analysis and verification of hardware and software systems. *Sci Comput Program*. 2008;72(1-2):3-21. <https://doi.org/10.1016/j.scico.2007.08.001>
21. Becchi A, Zaffanella E. PPLite: zero-overhead encoding of NNC polyhedra. *Inf Comput*. 2020;275:104620. <https://doi.org/10.1016/j.ic.2020.104620>
22. Davis MD, Sigal R, Weyuker EJ. *Computability, Complexity, and Languages: Fundamentals of Theoretical Computer Science*: Academic Press Professional, Inc.; 1994.
23. Cortesi A, Zanioli M. Widening and narrowing operators for abstract interpretation. *Comput Lang Syst Struct*. 2011;37(1):24-42. <https://doi.org/10.1016/j.cl.2010.09.001>
24. Arceri V, Maffei S. Abstract domains for type juggling. *Electron Notes Theor Comput Sci*. 2017;331:41-55. <https://doi.org/10.1016/j.entcs.2017.02.003>
25. Arceri V, Olliaro M, Cortesi A, Mastroeni I. Completeness of abstract domains for string analysis of JavaScript programs. In: Theoretical Aspects of Computing—ICTAC 2019—16th International Colloquium, Proceedings Hierons RM, Mosbah M, eds., Lecture Notes in Computer Science, vol. 11884. Springer; 2019:255-272. [https://doi.org/10.1007/978-3-030-32505-3\\_15](https://doi.org/10.1007/978-3-030-32505-3_15)
26. Rival X, Mauborgne L. The trace partitioning abstract domain. *ACM Trans Program Lang Syst*. 2007;29(5):26-es. <https://doi.org/10.1145/1275497.1275501>
27. Olivieri L, Jensen TP, Negrini L, Spoto F. MichelsonLiSA: a static analyzer for tezos. In: IEEE International Conference on Pervasive Computing and Communications Workshops and Other Affiliated Events, PerCom Workshops 2023. IEEE; 2023:80-85. <https://doi.org/10.1109/PerComWorkshops56833.2023.10150247>
28. Negrini L, Shabadi G, Urban C. Static analysis of data transformations in Jupyter notebooks. In: Proceedings of the 12th ACM SIGPLAN International Workshop on the State of the Art in Program Analysis, SOAP 2023 Ferrara P, Hadarean L, eds. ACM; 2023:8-13. <https://doi.org/10.1145/3589250.3596145>
29. Park C, Im H, Ryu S. Precise and scalable static analysis of jQuery using a regular expression domain. In: Proceedings of the 12th Symposium on Dynamic Languages, DLS 2016 Ierusalimschy R, ed. ACM; 2016:25-36.
30. Christensen AS, Møller A, Schwartzbach MI. Precise analysis of string expressions. In: Static Analysis, 10th International Symposium, SAS 2003, Proceedings Cousot R, ed., Lecture Notes in Computer Science, vol. 2694. Springer; 2003:1-18. [https://doi.org/10.1007/3-540-44898-5\\_1](https://doi.org/10.1007/3-540-44898-5_1)
31. Madsen M, Andreassen E. 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. Proceedings Cohen A, ed., Lecture Notes in Computer Science, vol. 8409. Springer; 2014:197-217. [https://doi.org/10.1007/978-3-642-54807-9\\_12](https://doi.org/10.1007/978-3-642-54807-9_12)
32. Abdulla PA, Atig MF, Chen Y-F, Holík L, Rezine A, Rümmer P, Stenman J. String constraints for verification. In: Computer Aided Verification—26th International Conference, CAV 2014, Held as Part of the Vienna Summer of Logic, VSL 2014. Proceedings Biere A, Bloem R, eds., Lecture Notes in Computer Science, vol. 8559. Springer; 2014:150-166. [https://doi.org/10.1007/978-3-319-08867-9\\_10](https://doi.org/10.1007/978-3-319-08867-9_10)
33. Arceri V, Mastroeni I. Analyzing dynamic code: a sound abstract interpreter for *evil eval*. *ACM Trans Priv Secur*. 2021;24(2):10:1-10:38. <https://doi.org/10.1145/3426470>
34. Costantini G, Ferrara P, Cortesi A. Static analysis of string values. In: Formal Methods and Software Engineering—13th International Conference on Formal Engineering Methods, ICFEM 2011. Proceedings Qin S, Qiu Z, eds., Lecture Notes in Computer Science, vol. 6991. Springer; 2011:505-521. [https://doi.org/10.1007/978-3-642-24559-6\\_34](https://doi.org/10.1007/978-3-642-24559-6_34)
35. Almahfi N, Lu L. Precise string domain for analyzing JavaScript arrays and objects. In: 3rd International Conference on Information and Computer Technologies, ICICT 2020 May Huang SZ, ed. IEEE; 2020:17-23. <https://doi.org/10.1109/ICICT50521.2020.00011>
36. D'Antoni L, Veanes M. Minimization of symbolic automata. In: The 41st Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '14 Jagannathan S, Sewell P, eds. ACM; 2014:541-554. <https://doi.org/10.1145/2535838.2535849>
37. Veanes M. Applications of symbolic finite automata. In: Implementation and Application of Automata—18th International Conference, CIAA 2013. Proceedings Konstantinidis S, ed., Lecture Notes in Computer Science, vol. 7982. Springer; 2013:16-23. [https://doi.org/10.1007/978-3-642-39274-0\\_3](https://doi.org/10.1007/978-3-642-39274-0_3)
38. Preda MD, Giacobazzi R, Lakhota A, Mastroeni I. Abstract symbolic automata: mixed syntactic/semantic similarity analysis of executables. In: Proceedings of the 42nd annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2015 Rajamani SK, Walker D, eds. ACM; 2015:329-341. <https://doi.org/10.1145/2676726.2676986>
39. Choi T-H, Lee O, Kim H, Doh K-G. A practical string analyzer by the widening approach. In: Programming Languages and Systems, 4th Asian Symposium, APLAS 2006, Proceedings Kobayashi N, ed., Lecture Notes in Computer Science, vol. 4279. Springer; 2006:374-388. [https://doi.org/10.1007/11924661\\_23](https://doi.org/10.1007/11924661_23)
40. Abdulla PA, Atig MF, Diep BP, Holík L, Janku P. Chain-free string constraints. In: Automated Technology for Verification and Analysis—17th International Symposium, ATVA 2019, Proceedings Chen Y-F, Cheng C-H, Esparza J, eds., Lecture Notes in Computer Science, vol. 11781. Springer; 2019: 277-293. [https://doi.org/10.1007/978-3-030-31784-3\\_16](https://doi.org/10.1007/978-3-030-31784-3_16)
41. Chen T, Hague M, Lin AW, Rümmer P, Wu Z. Decision procedures for path feasibility of string-manipulating programs with complex operations. *Proc ACM Program Lang*. 2019;3:49:1-49:30. <https://doi.org/10.1145/3290362>
42. Amadini R, Gange G, Stuckey PJ. Dashed strings for string constraint solving. *Artif Intell*. 2020;289:103368. <https://doi.org/10.1016/j.artint.2020.103368>
43. Wang H-E, Chen S-Y, Yu F, Jiang J-HR. A symbolic model checking approach to the analysis of string and length constraints. In: Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering, ASE 2018 Huchard M, Kästner C, Fraser G, eds. ACM; 2018:623-633.



44. Yu F, Alkhalaf M, Bultan T, Ibarra OH. Automata-based symbolic string analysis for vulnerability detection. *Formal Methods Syst Des.* 2014;44(1): 44-70.
45. Abdulla PA, Atig MF, Chen Y-F, et al. Efficient handling of string-number conversion. In: Proceedings of the 41st ACM SIGPLAN International Conference on Programming Language Design and Implementation, PLDI 2020 Donaldson AF, Torlak E, eds. ACM; 2020:943-957. <https://doi.org/10.1145/3385412.3386034>
46. Zheng Y, Zhang X, Ganesh V. Z3-str: a z3-based string solver for web application analysis. In: Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering, ESEC/FSE'13, Saint Petersburg, Russian Federation Meyer B, Baresi L, Mezini M, eds. ACM; 2013:114-124.
47. de Moura LM, Bjørner NS. Z3: an efficient SMT solver. In: Tools and Algorithms for the Construction and Analysis of Systems, 14th International Conference, TACAS 2008, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2008. Proceedings Ramakrishnan CR, Rehof J, eds., Lecture Notes in Computer Science, vol. 4963. Springer; 2008:337-340.
48. Arceri V, Dolcetti G, Zaffanella E. Speeding up static analysis with the split operator. In: Proceedings of the 12th ACM SIGPLAN International Workshop on the State of the Art in Program Analysis, SOAP 2023 Ferrara P, Hadarean L, eds. ACM; 2023:14-19. <https://doi.org/10.1145/3589250.3596141>
49. Ferrara P, Negrini L. SARL: OO framework specification for static analysis. In: Software Verification—12th International Conference, VSTTE 2020, and 13th International Workshop, NSV 2020, Revised Selected Papers Christakis M, Polikarpova N, Duggirala PS, Schrammel P, eds., Lecture Notes in Computer Science, vol. 12549. Springer; 2020:3-20. [https://doi.org/10.1007/978-3-030-63618-0\\_1](https://doi.org/10.1007/978-3-030-63618-0_1)
50. Arceri V, Olliaro M, Cortesi A, Ferrara P. Relational string abstract domains. In: Verification, Model Checking, and Abstract Interpretation—23rd International Conference, VMCAI 2022, Proceedings Finkbeiner B, Wies T, eds., Lecture Notes in Computer Science, vol. 13182. Springer; 2022:20-42. [https://doi.org/10.1007/978-3-030-94583-1\\_2](https://doi.org/10.1007/978-3-030-94583-1_2)

**How to cite this article:** Negrini L, Arceri V, Cortesi A, Ferrara P. TARSIS: An effective automata-based abstract domain for string analysis. *J Softw Evol Proc.* 2024;e2647. <https://doi.org/10.1002/smr.2647>

## APPENDIX A: SOUNDNESS AND COMPLETENESS PROOFS OF TARSIS'S SEMANTICS

We prove the soundness and completeness of TARSIS's abstract semantics by showing that their concretization is an over-approximation of the concrete one. As we formalized our transfer functions w.r.t. the smashed sum  $\mathbb{V}_{AL}^{\#} \triangleq \mathcal{T}F_{A/\equiv} \oplus \text{Intv} \oplus \text{Bool}$ , we compare concretizations of its elements with a concrete smashed sum  $\overline{\mathbb{V}}_{AL} \triangleq \wp(\Sigma^*) \cup \wp(\mathbb{Z}) \cup \wp(\{\text{true}, \text{false}\})$ , which is defined as a collecting semantics. We abuse notation denoting with  $\mathbb{M} : \text{Id} \rightarrow \overline{\mathbb{V}}_{AL}$  the set of collecting memories, ranging over meta-variable  $\text{mm}$ , that associate each identifier to a collecting value. The concrete expression semantics of such domain is defined as  $[[e]] : \mathbb{M} \rightarrow \overline{\mathbb{V}}_{AL}$ , evaluating  $e$  and returning its possible values. Such semantics is defined as the additive lift of the one in Figure 3. Function  $\gamma_{\mathbb{V}_{AL}^{\#}} : \mathbb{V}_{AL}^{\#} \rightarrow \overline{\mathbb{V}}_{AL}$  is the smashed sum concretization, and it is defined as

$$\gamma_{\mathbb{V}_{AL}^{\#}}(a) \triangleq \begin{cases} \emptyset & \text{if } a = \perp, \\ \gamma_{\text{Intv}}(a) & \text{if } a \in \text{Intv}, \\ \gamma_{\text{Bool}}(a) & \text{if } a \in \text{Bool}, \\ \gamma_{\mathcal{T}}(a) & \text{if } a \in \mathcal{T}F_{A/\equiv}, \\ \overline{\mathbb{V}}_{AL} & \text{otherwise,} \end{cases}$$

where  $\gamma_{\text{Intv}} : \text{Intv} \rightarrow \wp(\mathbb{Z})$  and  $\gamma_{\text{Bool}} : \text{Bool} \rightarrow \wp(\{\text{true}, \text{false}\})$  correspond to the concretization functions of intervals and Booleans, respectively. We can now define the abstract memories concretization  $\gamma : \mathbb{M}^{\#} \rightarrow \mathbb{M}$  as  $\gamma(\text{mm}^{\#}) \triangleq \{ (x, \gamma_{\mathbb{V}_{AL}^{\#}}(\text{mm}^{\#}(x))) \mid x \in \text{dom}(\text{mm}^{\#}) \}$ . With this setup, we prove the abstract semantics to be sound by proving that  $\forall \text{mm}^{\#} \in \mathbb{M}^{\#}. [[e]]\gamma(\text{mm}^{\#}) \subseteq \gamma_{\mathbb{V}_{AL}^{\#}}([[e]]^{\#}(\text{mm}^{\#}))$ . We also prove completeness by enforcing the equality on such relation, and incompleteness by providing a counterexample.

In the following, we remove the subscript from  $\gamma$  to avoid cluttering the notation, because it is clear from the context which concretization function applies. Moreover, we mark proof steps as *automata lift* if they represent the transition from a condition over languages (i.e., sets of strings) to its equivalent condition over automata. Finally, given the non-existence of a GC between the string concrete domain and TARSIS, from here on, when we refer to completeness, we intend forward completeness, defined in Section 3.

### A.1 | Concat

**Theorem A1.**  $[[\text{concat}(s, s')]]^{\#}$  is a sound and complete abstraction of  $[[\text{concat}(s, s')]]$ . Formally,

$$\forall \text{mm}^\sharp \in \mathbb{M}^\sharp, \forall \mathbf{s}, \mathbf{s}' \in \text{SE}. [[\text{concat}(\mathbf{s}, \mathbf{s}')]]\gamma(\text{mm}^\sharp) = \gamma([[ \text{concat}(\mathbf{s}, \mathbf{s}') ]])^\sharp \text{mm}^\sharp.$$

*Proof.* Soundness and completeness follow from the fact that finite state automata and regular languages are closed under finite concatenation.<sup>22</sup>  $\square$

## A.2 | Length

**Theorem A2.**  $[[\text{length}(\mathbf{s})]]^\sharp$  is a sound but not complete abstraction of  $[[\text{length}(\mathbf{s})]]$ . Formally,

$$\forall \text{mm}^\sharp \in \mathbb{M}^\sharp, \forall \mathbf{s} \in \text{SE}. [[\text{length}(\mathbf{s})]]\gamma(\text{mm}^\sharp) \subsetneq \gamma([[ \text{length}(\mathbf{s}) ]])^\sharp \text{mm}^\sharp.$$

*Proof.* The collecting semantics of `length` is defined as the additive lift of the concrete one reported in Figure 3, namely,  $[[\text{length}(\mathbf{s})]]_{\text{mm}} = \{ |\sigma| \mid \sigma \in \mathcal{L} \}$ , where  $[[\mathbf{s}]]_{\text{mm}} = \mathcal{L} \in \wp(\Sigma^*)$ . Let us suppose that  $[[\mathbf{s}]]^\sharp_{\text{mm}^\sharp} = \mathbb{A} \in \mathcal{TF}_{\mathbb{A}/\equiv}$  s.t.  $\gamma(\mathbb{A}) = \mathcal{L} \in \wp(\Sigma^*)$ , and let  $l = \{ |\sigma| \mid \sigma \in \mathcal{L} \}$ . Following the semantics definition, if  $\text{cyclic}(\mathbb{A}) \vee \text{readsTop}(\mathbb{A})$ , we prove the soundness as

$$\begin{aligned} & [[\text{length}(\mathbf{s})]]\gamma(\text{mm}^\sharp) \\ &= l && \text{? def. } \mathbb{E} \text{ } \S \\ &\subseteq \gamma([\min l, +\infty]) && \text{? def. } \min, \gamma \text{ } \S \\ &= \gamma([\min \text{Path}(\mathbb{A}), +\infty]) && \text{? def. } \min \text{Path} \text{ } \S \\ &= \gamma([[ \text{length}(\mathbf{s}) ]])^\sharp \text{mm}^\sharp && \text{? def. } \mathbb{E}, \text{ first case } \text{ } \S \end{aligned}$$

Otherwise, because  $\mathcal{L}$  is a finite language, soundness is proven as

$$\begin{aligned} & [[\text{length}(\mathbf{s})]]\gamma(\text{mm}^\sharp) \\ &= l && \text{? def. } \mathbb{E} \text{ } \S \\ &\subseteq \gamma_a([\min l, \max l]) && \text{? def. } \max, \gamma \text{ } \S \\ &= \gamma([\min \text{Path}(\mathbb{A}), \max \text{Path}(\mathbb{A})]) && \text{? def. } \min \text{Path}, \max \text{Path} \text{ } \S \\ &= \gamma([[ \text{length}(\mathbf{s}) ]])^\sharp \text{mm}^\sharp && \text{? def. } \mathbb{E}, \text{ second case } \text{ } \S \end{aligned}$$

As a counterexample for completeness, let us consider  $[[\mathbf{s}]]^\sharp_{\text{mm}^\sharp} = \mathbb{A} \in \mathcal{TF}_{\mathbb{A}/\equiv}$  s.t.  $\gamma(\mathbb{A}) = \{a, aaa\}$ .

$$[[\text{length}(\mathbf{s})]]\gamma(\text{mm}^\sharp) = \{1, 3\} \subsetneq \gamma([1, 3]) = \gamma([[ \text{length}(\mathbf{s}) ]])^\sharp \text{mm}^\sharp.$$

$\square$

## A.3 | Contains

**Theorem A3.**  $[[\text{contains}(\mathbf{s}, \mathbf{s}')]]^\sharp$  is a sound but not complete abstraction of  $[[\text{contains}(\mathbf{s}, \mathbf{s}')]]$ . Formally,

$$\forall \text{mm}^\sharp \in \mathbb{M}^\sharp, \forall \mathbf{s}, \mathbf{s}' \in \text{SE}. [[\text{contains}(\mathbf{s}, \mathbf{s}')]]\gamma(\text{mm}^\sharp) \subseteq \gamma([[ \text{contains}(\mathbf{s}, \mathbf{s}') ]])^\sharp \text{mm}^\sharp.$$

*Proof.* The collecting semantics of `contains` is defined as the additive lift of the concrete one, that is,  $[[\text{contains}(\mathbf{s}, \mathbf{s}')]]_{\text{mm}} = \{ \mathbf{b} \mid \mathbf{b} = \text{contains}(\sigma, \sigma'), \sigma \in \mathcal{L}, \sigma' \in \mathcal{L}' \}$ , where  $[[\mathbf{s}]]_{\text{mm}} = \mathcal{L} \in \wp(\Sigma^*)$ ,  $[[\mathbf{s}']]_{\text{mm}} = \mathcal{L}' \in \wp(\Sigma^*)$  and

contains:  $\Sigma^* \times \Sigma^* \rightarrow \{\text{true}, \text{false}\}$  corresponds to the concrete semantics of Figure 3. Let  $[[s]]^{\#} \text{mm}^{\#} = A \in \mathcal{TFA}_{/\equiv}$  s.t.  $\gamma(A) = \mathcal{L} \in \wp(\Sigma^*)$  and  $[[s']]^{\#} \text{mm}^{\#} = A' \in \mathcal{TFA}_{/\equiv}$  s.t.  $\gamma(A') = \mathcal{L}' \in \wp(\Sigma^*)$ . We split the proof following possible values produced by the concrete semantics.

$[[\text{contains}(s, s')]]^{\#} \gamma(\text{mm}^{\#}) = \{\text{false}\}$ . Thus, no substring of the strings in  $\mathcal{L}$  is in  $\mathcal{L}'$ .

$$\begin{aligned} [[\text{contains}(s, s')]]^{\#} \gamma(\text{mm}^{\#}) = \{\text{false}\} &\stackrel{\text{def. } \mathbb{E}}{\Leftrightarrow} \forall \sigma \in \mathcal{L} \forall \sigma' \in \mathcal{L}' . \sigma'_s \sigma \stackrel{\text{def. FA}}{\Leftrightarrow} \mathcal{L}(\text{FA}(A)) \cap \mathcal{L}' = \emptyset \\ &\stackrel{\text{automata lift}}{\Leftrightarrow} \text{FA}(A) \cap_{\mathcal{T}} A' = \text{Min}(\emptyset) \stackrel{\text{def. } \mathbb{E}, 1^{\text{st}} \text{ case}}{\Leftrightarrow} [[\text{contains}(s, s')]]^{\#} \text{mm}^{\#} = \{\text{false}\}. \end{aligned}$$

$[[\text{contains}(s, s')]]^{\#} \gamma(\text{mm}^{\#}) = \{\text{true}\}$ . Thus, all strings in  $\mathcal{L}$  contain all the strings of  $\mathcal{L}'$ . This invalidates the first case of our abstract semantics, as  $\exists \sigma \in \mathcal{L}(\text{FA}(A)) . \sigma \in \mathcal{L}'$ . If  $\text{singlePath}(A')$  holds, our semantics matches the concrete one:

$$\begin{aligned} [[\text{contains}(s, s')]]^{\#} \gamma(\text{mm}^{\#}) = \{\text{true}\} \wedge \text{singlePath}(A') &\stackrel{\text{def. } \mathbb{E}, \text{singlePath}}{\Leftrightarrow} \forall \sigma \in \mathcal{L} . \sigma_{\text{sps}} \sigma \stackrel{\mathcal{L}(A^{ac}) \subseteq \mathcal{L}}{\Leftrightarrow} \forall \sigma \in \mathcal{L}(A^{ac}) . \sigma_{\text{sps}} \sigma \\ &\stackrel{\text{automata lift}}{\Leftrightarrow} \forall \pi \in \text{Paths}(A^{ac}) . \sigma_{\text{sps}} \sigma_{\pi} \stackrel{\text{def. } \mathbb{E}, 2^{\text{nd}} \text{ case}}{\Leftrightarrow} [[\text{contains}(s, s')]]^{\#} \text{mm}^{\#} = \{\text{true}\}. \end{aligned}$$

Otherwise, if  $A'$  is not a single-path automaton, the semantics returns  $\{\text{true}, \text{false}\}$ , and soundness is met.

$[[\text{contains}(s, s')]]^{\#} \gamma(\text{mm}^{\#}) = \{\text{true}, \text{false}\}$ . In this case, soundness is met as none of the conditions appearing in the definition of  $[[\text{contains}(s, s')]]^{\#}$  are satisfied, and the latter returns  $\{\text{true}, \text{false}\}$  as well (third case). As a counterexample for the completeness of contains, let us consider  $[[s]]^{\#} \text{mm}^{\#} = A \in \mathcal{TFA}_{/\equiv}$  s.t.  $\gamma(A) = \{ab, bba\}$  and  $[[s']]^{\#} \text{mm}^{\#} = A' \in \mathcal{TFA}_{/\equiv}$  s.t.  $\gamma(A') = \{a, b\}$ .

$$[[\text{contains}(s, s')]]^{\#} \gamma(\text{mm}^{\#}) = \{\text{true}\} \subsetneq \{\text{true}, \text{false}\} = \gamma([[ \text{contains}(s, s') ]])^{\#} \text{mm}^{\#}.$$

□

## A.4 | IndexOf

**Theorem A4.**  $[[\text{indexOf}(s, s')]]^{\#}$  is a sound but not complete abstraction of  $[[\text{indexOf}(s, s')]]$ . Formally,

$$\forall \text{mm}^{\#} \in \mathbb{M}^{\#}, \forall s, s' \in \text{SE}. [[\text{indexOf}(s, s')]]^{\#} \gamma(\text{mm}^{\#}) \subsetneq \gamma([[ \text{indexOf}(s, s') ]])^{\#} \text{mm}^{\#}.$$

*Proof.* The collecting semantics of `indexOf` is defined as the additive lift of the concrete one:  $[[\text{indexOf}(s, s')]]^{\#} \text{mm}^{\#} = \{i \mid i = \text{indexOf}(\sigma, \sigma'), \sigma \in \mathcal{L}, \sigma' \in \mathcal{L}'\}$ , where  $[[s]]^{\#} \text{mm}^{\#} = \mathcal{L} \in \wp(\Sigma^*)$ ,  $[[s']]^{\#} \text{mm}^{\#} = \mathcal{L}' \in \wp(\Sigma^*)$  and  $\text{indexOf}: \Sigma^* \times \Sigma^* \rightarrow \mathbb{N}$  corresponds to the concrete semantics of Figure 3. Let  $[[s]]^{\#} \text{mm}^{\#} = A \in \mathcal{TFA}_{/\equiv}$  s.t.  $\gamma(A) = \mathcal{L} \in \wp(\Sigma^*)$  and  $[[s']]^{\#} \text{mm}^{\#} = A' \in \mathcal{TFA}_{/\equiv}$  s.t.  $\gamma(A') = \mathcal{L}' \in \wp(\Sigma^*)$ . By definition of the concrete semantics,  $[[\text{indexOf}(s, s')]]^{\#} \gamma(\text{mm}^{\#}) \subseteq \gamma([-1, +\infty])$ . When  $A$  or  $A'$  are cyclic or  $A'$  has a T transition, the abstract semantics returns the interval  $[-1, +\infty]$ , guaranteeing soundness. We thus continue by assuming that  $A$  and  $A'$  are not cyclic and  $A'$  has no T transitions (i.e.,  $\mathcal{L}'$  is finite). We split the proof following the possible concrete values.

$[[\text{indexOf}(s, s')]]^{\#} \gamma(\text{mm}^{\#}) = \{-1\}$ . No string of  $\mathcal{L}'$  is contained in any string of  $\mathcal{L}$ . Formally,

$$\begin{aligned} [[\text{indexOf}(s, s')]]^{\#} \gamma(\text{mm}^{\#}) = \{-1\} &\stackrel{\text{necessary condition}}{\Leftrightarrow} \forall \sigma' \in \mathcal{L}' \nexists \sigma \in \mathcal{L} . \sigma'_s \sigma \\ &\stackrel{\text{def. } \mathbb{E}, 2^{\text{nd}} \text{ case}}{\Leftrightarrow} [[\text{indexOf}(s, s')]]^{\#} \text{mm}^{\#} = [-1, -1] = . \end{aligned}$$

As  $\gamma_{\text{Intv}}([-1, -1]) = \{-1\}$ , soundness is met.

$[[\text{indexOf}(s, s')]]\gamma(\text{mm}^\#) = I \subseteq \{n \in \mathbb{Z} \mid n \geq -1\} \cdot \exists i \in I. i \geq 0$ . This implies  $\exists \sigma \in \mathcal{L}(\mathbb{A}), \sigma' \in \mathcal{L}(\mathbb{A}'). \sigma'_s \sigma$ , as the collecting semantics returns at least one value that is not  $-1$ . Here, the abstract semantics relies on function IO that computes an interval for each string  $\sigma' \in \mathcal{L}(\mathbb{A}')$ , lubbing the results together. Hence, it is enough to prove the correctness of IO. Given  $\sigma' \in \mathcal{L}'$ , let us define the set  $I_{\sigma'} \subseteq I = \{i \mid i = \text{indexOf}(\sigma, \sigma'), \sigma \in \mathcal{L}\}$  of positions where  $\sigma'$  can be found in  $\mathcal{L}$  and let  $m, M \in I_{\sigma'}$  be the minimal and the maximal elements of  $I_{\sigma'}$ . Therefore, it is sufficient to prove that  $\gamma_{\text{Intv}}([m, M]) \subseteq \gamma_{\text{Intv}}([i, j])$ , where  $[i, j] = \text{IO}(\mathbb{A}, \sigma')$ . For soundness to hold,  $i \leq m$  and  $M \leq j$  must be true, according to  $\gamma_{\text{Intv}}$ . We first prove  $i \leq m$ , identifying two cases. If  $m = -1$ :

$$-1 \in I_{\sigma'} \stackrel{\text{necessary condition}}{\Leftrightarrow} \exists \sigma \in \mathcal{L}. \sigma'_s \sigma \stackrel{\text{automata lift}}{\Leftrightarrow} \exists \pi \in \text{Paths}(\mathbb{A}). \sigma'_s \sigma_\pi \stackrel{\text{def. } i, 1^{\text{st}} \text{ case}}{\Leftrightarrow} i = -1.$$

Instead, if  $m > -1$ :

$$\begin{aligned} m = \min I_{\sigma'}, m \neq -1 &\stackrel{\text{necessary cond.}}{\Leftrightarrow} \exists \sigma \in \mathcal{L}. \sigma_m \dots \sigma_{m+|\sigma'|} = \sigma' \wedge \forall \sigma \in \mathcal{L} \nexists n < m. \sigma_n \dots \sigma_{n+|\sigma'|} = \sigma' \\ \text{automata lift} &\Leftrightarrow \exists \pi \in \text{Paths}(\mathbb{A}). \exists \sigma_f \in \text{Flat}(\sigma_\pi). \sigma_{f_m} \dots \sigma_{f_{m+|\sigma'|}} = \sigma' \wedge \forall \pi \in \text{Paths}(\mathbb{A}) \forall \sigma_f \in \text{Flat}(\sigma_\pi). \sigma_{f_k} \dots \sigma_{f_{k+|\sigma'|}} = \sigma' \Rightarrow k \geq m \\ \text{def. } i, 2^{\text{nd}} \text{ case} &\Leftrightarrow i = \min k = m. \end{aligned}$$

We now prove that  $M \leq j$ , identifying three cases. If  $M = -1$ :

$$I_{\sigma'} = \{-1\} \stackrel{\text{necessary condition}}{\Leftrightarrow} \forall \sigma \in \mathcal{L}. \sigma'_s \sigma \stackrel{\text{automata lift}}{\Leftrightarrow} \forall \pi \in \text{Paths}(\mathbb{A}). \sigma'_s \sigma_\pi \stackrel{\text{def. } j, 1^{\text{st}} \text{ case}}{\Leftrightarrow} j = 1.$$

Instead, if  $M > -1$  and  $\forall \pi \in \text{Paths}(\mathbb{A}). \pi \text{ reads } \sigma \Rightarrow \pi \text{ does not read T before } \sigma$ :

$$\begin{aligned} M = \max I_{\sigma'} &\stackrel{\text{necessary cond.}}{\Leftrightarrow} \exists \sigma \in \mathcal{L}. \sigma_M \dots \sigma_{M+|\sigma'|} = \sigma' \wedge \forall \sigma \in \mathcal{L} \nexists n > M. \sigma_n \dots \sigma_{n+|\sigma'|} = \sigma' \\ \text{automata lift} &\Leftrightarrow \exists \pi \in \text{Paths}(\mathbb{A}). \exists \sigma_f \in \text{Flat}(\sigma_\pi). \sigma_{f_M} \dots \sigma_{f_{M+|\sigma'|}} = \sigma' \wedge \forall \pi \in \text{Paths}(\mathbb{A}) \forall \sigma_f \in \text{Flat}(\sigma_\pi). \sigma_{f_k} \dots \sigma_{f_{k+|\sigma'|}} = \sigma' \Rightarrow k \leq M \\ \text{def. } j, 3^{\text{rd}} \text{ case} &\Leftrightarrow j = \max k = M. \end{aligned}$$

Finally, if  $M > -1$  and  $\exists \pi \in \text{Paths}(\mathbb{A}). \pi \text{ reads T before } \sigma$ ,  $j = +\infty$  by the second case of the definition of  $j$ , that is, thus greater than  $M$ . As both inequalities are always satisfied, we can conclude that soundness is met in all cases.  $\square$

As a counterexample for completeness, let us consider  $[[s]]^\# \text{mm}^\# = \mathbb{A} \in \mathcal{TF}_{\mathbb{A}/\equiv}$  s.t.  $\gamma(\mathbb{A}) = \{a, bba\}$  and  $[[s']]^\# \text{mm}^\# = \mathbb{A}' \in \mathcal{TF}_{\mathbb{A}'/\equiv}$  s.t.  $\gamma(\mathbb{A}') = \{a\}$ .

$$[[\text{indexOf}(s, s')]]\gamma(\text{mm}^\#) = \{0, 2\} \subsetneq \gamma([0, 2]) = \gamma([[ \text{indexOf}(s, s') ]])^\# \text{mm}^\#.$$

## A.5 | Repeat

The abstract semantics of `repeat` relies of the auxiliary function `repeat`, described by Algorithm 2. First, we prove the soundness of `repeat`.

**Theorem A5.** Given  $\mathbb{A} \in \mathcal{TF}_{\mathbb{A}/\equiv}$ ,  $i \in \mathbb{N}$ , `repeat` is sound. Namely,

$$\text{repeat}(\mathcal{L}(\mathbb{A}), i) \subseteq \gamma(\text{repeat}(\mathbb{A}, i)),$$

where we abuse notation defining the collecting semantics  $\text{repeat}(\mathcal{L}, i) \triangleq \{\sigma' \mid \sigma \in \mathcal{L}\}$ .

*Proof.* We split the proof in the following cases.

$i = 0$ . In this case, Algorithm 2 returns the automaton recognizing the empty string ( $\text{Min}(\{\epsilon\})$ ), and repeat is sound:

$$\text{repeat}(\mathcal{L}(\mathbb{A}), i) = \{ \sigma^0 \mid \sigma \in \mathcal{L}(\mathbb{A}) \} = \{ \epsilon \} = \gamma(\text{Min}(\{\epsilon\})) \quad \text{? Lines 1 and 2 of Algorithm 2}$$

$\mathbb{A}$  is cyclic. In the following we rely on the fact that, for some  $i \in \mathbb{N}$ :

$$\{ \sigma^i \mid \sigma \in \mathcal{L} \} \subseteq \mathcal{L}^i \quad (\text{A1})$$

$$\begin{aligned} \text{repeat}(\mathcal{L}(\mathbb{A}), i) &= \\ &= \{ \sigma^i \mid \sigma \in \mathcal{L}(\mathbb{A}) \} && \text{? def. repeat}(\mathcal{L}, i) \\ &= \{ \epsilon \} \cdot \{ \sigma^i \mid \sigma \in \mathcal{L}(\mathbb{A}) \} && \text{? neutrality of } \epsilon \text{ for concatenation} \\ &\subseteq \{ \epsilon \} \cdot \mathcal{L}(\mathbb{A})^i && \text{? Equation A1} \end{aligned}$$

The last formula is the language obtained by concatenating the empty string with the  $i$ -concatenation of  $\mathcal{L}(\mathbb{A})$ , corresponding the language of concatenating  $i$ -times the automaton  $\mathbb{A}$  with itself and the automaton recognizing  $\epsilon$ .

These operations are the ones performed at Lines 3–8 of Algorithm 2.

$i \neq 0$  and  $\mathbb{A}$  is not cyclic. In this case, we have

$$\begin{aligned} \text{repeat}(\mathcal{L}(\mathbb{A}), i) &= \\ &= \{ \sigma^i \mid \sigma \in \mathcal{L}(\mathbb{A}) \} && \text{? def. repeat}(\mathcal{L}, i) \\ &= \{ \sigma^i \mid \exists \pi \in \text{Paths}(\mathbb{A}). \sigma = \sigma_\pi \} && \text{? } \sigma \in \mathcal{L}(\mathbb{A}) \\ &= \bigcup_{\pi \in \text{Paths}(\mathbb{A})} \sigma_\pi^i && \text{? } \sigma = \sigma_\pi \\ &= \gamma\left(\bigsqcup_{\pi \in \text{Paths}(\mathbb{A})} \text{Min}(\{\sigma_\pi^i\})\right) && \text{? def. } \gamma, \sqcup \end{aligned}$$

The input of function  $\gamma$  in the last formula corresponds to Lines 10–18 of Algorithm 2, where  $\text{Min}(\{\sigma_\pi^i\})$  is computed by Lines 13–15.  $\square$

**Theorem A6.**  $[[ \text{repeat}(s, a) ]]$  is a sound but not complete abstraction of  $[[ \text{repeat}(s, a) ]]$ . Formally,

$$\forall \text{mm}^\# \in \mathbb{M}^\#, \forall s \in \text{SE}, a \in \text{AE}. [[ \text{repeat}(s, a) ]]\gamma(\text{mm}^\#) \subsetneq \gamma([[ \text{repeat}(s, a) ]])^\# \text{mm}^\#.$$

*Proof.* The collecting semantics of `replace` is defined as the additive lift of the concrete one, that is,  $[[ \text{repeat}(s, a) ]]\text{mm} = \{ \sigma^k \mid k \in [i, j], \sigma \in \mathcal{L} \}$ , where  $[[ s ]]\text{mm} = \mathcal{L} \in \wp(\Sigma^*)$  and  $[[ a ]]\text{mm} = [i, j] \in \wp(\mathbb{Z})$ . Let  $[[ s ]]^\# \text{mm}^\# = \mathbb{A} \in \mathcal{TF}_{\mathbb{A}/\equiv}$  s.t.  $\gamma(\mathbb{A}) = \mathcal{L} \in \wp(\Sigma^*)$ , and  $[[ a ]]^\# \text{mm}^\# = [i, j] \in \text{Intv}$ . By definition of the concrete semantics, we suppose that  $[i, j] \subseteq [0, +\infty]$ . We split the proof in the following cases.

$$[i, j] = [0, +\infty].$$

$$\begin{aligned} [[ \text{repeat}(s, a) ]]\gamma(\text{mm}^\#) &= \\ &= \{ \sigma^i \mid \sigma \in \mathcal{L}(\mathbb{A}), i \in \mathbb{N} \} && \text{? def. repeat}(s, a) \\ &\subseteq \bigcup_{i \in \mathbb{N}} \mathcal{L}(\mathbb{A})^i && \text{? Equation(A1)} \\ &= \gamma(\text{Kleene}(\mathbb{A})) && \text{? def. Kleene} \end{aligned}$$

$i = j \wedge i \in \mathbb{N}$ . Soundness follows from Theorem (A5).

$$j = +\infty.$$

$$\begin{aligned}
[[ \text{repeat}(s,a) ]]\gamma(\text{mm}^\#) &= \\
&= \{ \sigma^k \mid \sigma \in \mathcal{L}(A), k \geq i \} && \text{? def. repeat}(s,a) \} \\
&= \{ \sigma^i \sigma' \mid \sigma \in \mathcal{L}(A), i \in \mathbb{N} \} && \text{? def. } \geq, \text{ string concat. } \} \\
&= \{ \sigma^i \mid \sigma \in \mathcal{L}(A) \} \cdot \{ \sigma' \mid \sigma \in \mathcal{L}(A), i \in \mathbb{N} \} && \text{? def. lang. concat. } \} \\
&\subseteq \gamma([[ \text{concat}(\text{repeat}(A,i), \text{Kleene}(A)) ]]\#) && \text{? def. Kleene, repeat, } \gamma \}.
\end{aligned}$$

$i, j \in \mathbb{N}$ . Soundness follows from Theorem A5 and the definition of  $\sqcup_{\mathcal{T}}$ . As a counterexample for completeness, let  $[[s]]^\# \text{mm}^\# = A \in \mathcal{TFA}_{\equiv}$  s.t.  $\gamma(A) = \{ a^n \mid n \in \mathbb{N} \} \cup \{ b \}$  and  $[[a]]^\# \text{mm}^\# = [2, 2]$ .

$$[[ \text{repeat}(s,a) ]]\gamma(\text{mm}^\#) = \{ a^n \mid n \in \mathbb{N} \} \cup \{ bb \} \subsetneq \{ a^n b a^m b \mid n, m \in \mathbb{N} \} = \gamma([[ \text{repeat}(s,a) ]]\# \text{mm}^\#).$$

□

### A.6 | TrimLeft, TrimRight, and Trim

The abstract semantics of `trimLeft` relies of the auxiliary function `trimL`, working on regexes. In the following we prove the soundness of `trimL`. The soundness proof for `trimR` is analogous, while soundness of `trim` follows from the soundness of `trimL` and `trimR`.

**Theorem A7.** Given  $A \in \mathcal{TFA}_{\equiv}$ , let  $r$  be the regex equivalent to  $A$ . `trimL` is sound, namely,

$$\text{trimL}(\gamma(r)) \subseteq \gamma(\text{trimL}(r)),$$

where we abuse notation of `trimL` defining the collecting semantics  $\text{trimL}(\mathcal{L}) \triangleq \left\{ \sigma' \mid \begin{array}{l} \sigma \in \mathcal{L}, \sigma = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma = \psi' \sigma'' \} \end{array} \right\}$ .

*Proof.* The proof is done by induction on the structure of the regex  $r$ .

#### Base cases

$r = \emptyset$ . Soundness holds because  $\text{trimL}(\mathcal{L}(\emptyset)) = \emptyset = \gamma(\text{trimL}(\emptyset))$  (first case).

$r = T$ . Soundness holds because  $\text{trimL}(\mathcal{L}(T)) = \Sigma^* = \gamma(\text{trimL}(T))$  (first case).

$r = \sigma$ . If the regex is an atom, the abstract semantics relies on its concrete semantics; hence, soundness holds.

#### Inductive cases

$r = r_1 \parallel r_2$ .

$$\begin{aligned}
\text{trimL}(\gamma(r_1 \parallel r_2)) &= \\
&= \text{trimL}(\gamma(r_1) \cup \gamma(r_2)) && \text{? def. } \gamma \} \\
&= \text{trimL}(\gamma(r_1)) \cup \text{trimL}(\gamma(r_2)) && \text{? distrib. of trimL } \} \\
&= \gamma(\text{trimL}(r_1)) \cup \gamma(\text{trimL}(r_2)) && \text{? ind. hp. } \} \\
&= \gamma(\text{trimL}(r_1) \cup \text{trimL}(r_2)) && \text{? def. } \cup, \gamma \} \\
&= \gamma(\text{trimL}(r_1 \parallel r_2)) && \text{? distrib. of trimL, third case } \}.
\end{aligned}$$

$r = r_1 r_2$ . For regex concatenation, we split the proof in three sub-cases. •  $\gamma(r_1) \subseteq \{ \}^* \Rightarrow \text{trimL}(\gamma(r_1)) = \epsilon$

$$\begin{aligned}
\text{trimL}(\gamma(r_1 r_2)) &= \\
&= \left\{ \sigma' \mid \begin{array}{l} \sigma \in \gamma(r_1 r_2), \sigma = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma = \psi' \sigma' \} \end{array} \right\} && \text{? def. trimL } \S \\
&= \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \gamma(r_1), \sigma_2 \in \gamma(r_2), \sigma_1 \sigma_2 = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma_1 \sigma_2 = \psi' \sigma' \} \end{array} \right\} && \text{? def. } \gamma(r_1 r_2) \S \\
&= \left\{ \sigma' \mid \begin{array}{l} \sigma_2 \in \gamma(r_2), \sigma_2 = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma_2 = \psi' \sigma' \} \end{array} \right\} && \text{? trimL}(\gamma(r_1)) = \epsilon \Rightarrow \sigma_1 \in \psi \S \\
&= \text{trimL}(\gamma(r_2)) && \text{? def. trimL } \S \\
&\subseteq \gamma(\text{trimL}(r_2)) && \text{? ind. hp. } \S \\
&= \gamma(\text{trimL}(r_1 r_2)) && \text{? fourth case } \S.
\end{aligned}$$

- $\text{readWS}(r_1) \Rightarrow \gamma(r_1) = \mathcal{L}^{WS} \cup \mathcal{L}^{-WS}$  s.t.  $\mathcal{L}^{WS} \subseteq \{ \}^*$  and  $\mathcal{L}^{-WS} \cap \{ \}^* = \emptyset$

$$\begin{aligned}
\text{trimL}(\gamma(r_1 r_2)) &= \\
&= \left\{ \sigma' \mid \begin{array}{l} \sigma \in \gamma(r_1 r_2), \sigma = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma = \psi' \sigma' \} \end{array} \right\} && \text{? def. trimL } \S \\
&= \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \gamma(r_1), \sigma_2 \in \gamma(r_2), \sigma_1 \sigma_2 = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma_1 \sigma_2 = \psi' \sigma' \} \end{array} \right\} && \text{? def. } \gamma(r_1 r_2) \S \\
&= \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \mathcal{L}^{WS}, \sigma_2 \in \gamma(r_2), \sigma_1 \sigma_2 = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma_1 \sigma_2 = \psi' \sigma' \} \end{array} \right\} \\
&\cup \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \mathcal{L}^{-WS}, \sigma_2 \in \gamma(r_2), \sigma_1 \sigma_2 = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma_1 \sigma_2 = \psi' \sigma' \} \end{array} \right\} && \text{? def. } \mathcal{L}^{WS}, \mathcal{L}^{-WS} \S \\
&= \left\{ \sigma' \mid \begin{array}{l} \sigma_2 \in \gamma(r_2), \sigma_2 = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma_2 = \psi' \sigma' \} \end{array} \right\} \\
&\cup \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \mathcal{L}^{-WS}, \sigma_2 \in \gamma(r_2), \sigma_1 \sigma_2 = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma_1 \sigma_2 = \psi' \sigma' \} \end{array} \right\} && \text{? } \sigma_1 \in \{ \}^* \S \\
&= \left\{ \sigma' \mid \begin{array}{l} \sigma_2 \in \gamma(r_2), \sigma_2 = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma_2 = \psi' \sigma' \} \end{array} \right\} \\
&\cup \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \mathcal{L}^{-WS}, \sigma_2 \in \gamma(r_2), \sigma_1 \sigma_2 = \psi \sigma', \\ \psi = \max\{ \psi' \in \{ \}^* \mid \sigma_1 \sigma_2 = \psi' \sigma' \} \end{array} \right\} && \text{? } \sigma_1 \notin \{ \}^*, \psi \in \sigma_1 \S \\
&= \text{trimL}(\gamma(r_2)) \cup \text{trimL}(\gamma(r_1)) \gamma(r_2) && \text{? def. trimL } \S \\
&\subseteq \gamma(\text{trimL}(r_2)) \cup \gamma(\text{trimL}(r_1) r_2) && \text{? ind. hp., def. } \gamma \S \\
&= \gamma(\text{trimL}(r_2) \parallel \text{trimL}(r_1) r_2) && \text{? def. } \gamma(r_1 \parallel r_2) \S \\
&= \gamma(\text{trimL}(r_1 r_2)) && \text{? fifth case } \S.
\end{aligned}$$

- $\neg \text{readWS}(r_1)$ . The proof is analogous to the previous case.

$r = (r_1)^*$ . If  $\text{trimL}(r_1) = \epsilon$ , then  $\text{trimL}(r) = \epsilon = \gamma(\text{trimL}(r))$ ; hence, soundness trivially holds. Otherwise, the proof is analogous to the case  $r = r_1 r_2$ .  $\square$

As a counterexample for completeness, let  $[[s]]^{\#mm\#} = A = \text{Min}(\{\text{Tab}\})$ , thus  $\gamma(A) = \{ \sigma ab \mid \sigma \in \Sigma^* \}$ .

$$[[\text{trimLeft}(s)]] \gamma^{\#mm\#} = \{ \sigma ab \mid \sigma \in \Sigma^* \}^* \subsetneq \gamma(\text{Min}(\{\text{Tab}\})) = \gamma([[ \text{trimLeft}(s) ]]^{\#mm\#}).$$

Consequently, also the abstract semantics of  $\text{trimRight}$  and  $\text{trim}$  are not complete.

## A.7 | Replace

**Theorem A8.**  $[[\text{replace}(s, s_s, s_r)]]^{\#}$  is a sound but not complete abstraction of  $[[\text{replace}(s, s_s, s_r)]]$ . Formally,

$$\forall \text{mm}^\# \in \mathbb{M}^\#, \forall \mathbf{s}, \mathbf{s}_s, \mathbf{s}_r \in \text{SE}. [[ \text{replace}(\mathbf{s}, \mathbf{s}_s, \mathbf{s}_r) ]] \gamma(\text{mm}^\#) \subseteq \gamma([[ \text{replace}(\mathbf{s}, \mathbf{s}_s, \mathbf{s}_r) ]]^\# \text{mm}^\#).$$

*Proof.* The collecting semantics of `replace` is defined as the additive lift of the concrete one, that is,  $[[ \text{replace}(\mathbf{s}, \mathbf{s}_s, \mathbf{s}_r) ]] \text{mm} = \{ \sigma' \mid \sigma' = \text{replace}(\sigma, \sigma_s, \sigma_r), \sigma \in \mathcal{L}, \sigma_s \in \mathcal{L}_s, \sigma_r \in \mathcal{L}_r \}$ , where  $[[ \mathbf{s} ]] \text{mm} = \mathcal{L} \in \wp(\Sigma^*)$ ,  $[[ \mathbf{s}_s ]] \text{mm} = \mathcal{L}_s \in \wp(\Sigma^*)$ ,  $[[ \mathbf{s}_r ]] \text{mm} = \mathcal{L}_r \in \wp(\Sigma^*)$  and  $\text{replace}: \Sigma^* \times \Sigma^* \times \Sigma^* \rightarrow \Sigma^*$  corresponds to the concrete semantics of Figure 3. Let  $[[ \mathbf{s} ]]^\# \text{mm}^\# = \mathbf{A} \in \mathcal{TFA}_{\equiv}$  s.t.  $\gamma(\mathbf{A}) = \mathcal{L} \in \wp(\Sigma^*)$ ,  $[[ \mathbf{s}_s ]]^\# \text{mm}^\# = \mathbf{A}_s \in \mathcal{TFA}_{\equiv}$  s.t.  $\gamma(\mathbf{A}_s) = \mathcal{L}_s \in \wp(\Sigma^*)$ , and  $[[ \mathbf{s}_r ]]^\# \text{mm}^\# = \mathbf{A}_r \in \mathcal{TFA}_{\equiv}$  s.t.  $\gamma(\mathbf{A}_r) = \mathcal{L}_r \in \wp(\Sigma^*)$ . When  $\mathbf{A}$  or  $\mathbf{A}_s$  have a cycle or have a T-transition, our semantics returns  $\text{Min}(\{\mathbf{T}\})$  and is thus trivially sound. Otherwise, when no replacement happens (i.e.,  $[[ \text{replace}(\mathcal{L}, \mathcal{L}_s, \mathcal{L}_r) ]] \text{mm} = \mathcal{L}$ ):

$$[[ \text{replace}(\mathcal{L}, \mathcal{L}_s, \mathcal{L}_r) ]] = \mathcal{L} \stackrel{\text{necessary condition}}{\Leftrightarrow} \forall \sigma_s \in \mathcal{L}_s \nexists \sigma \in \mathcal{L}. \sigma_s \sigma \stackrel{\text{def. } \mathbb{B}_1^{\text{st}} \text{ case}}{\Leftrightarrow} [[ \text{replace}(\mathbf{A}, \mathbf{A}_s, \mathbf{A}_r) ]] = \mathbf{A}.$$

Otherwise, when at least one replacement happens, the semantics returns the lub of several applications of RP ranging over all possible combinations of strings in  $\mathcal{L}$  and  $\mathcal{L}_s$ , which can be thoroughly explored because  $\mathcal{L}$  and  $\mathcal{L}_s$  are finite sets. Once RP has been proven correct, soundness naturally follows according to the properties of lub. We thus prove that  $\forall \sigma \in \mathcal{L}, \forall \sigma_s \in \mathcal{L}_s, \forall \pi \in \text{Paths}(\mathbf{A}). \sigma_\pi = \sigma$ :

$$[[ \text{replace}(\{\sigma\}, \{\sigma_s\}, \mathcal{L}_r) ]] \subseteq \gamma(\text{RP}(\pi, \sigma_s, \mathbf{A}_r)).$$

Specifically, RP removes every occurrence of  $\sigma_s$  in  $\pi$  (Lines 7 and 8 of Algorithm 3, where states and transitions composing  $\sigma_s$  are removed from the resulting automaton), substituting them with a copy of the replace automaton (Line 4) that is connected to the path with  $\epsilon$ -transitions. This means that all  $\sigma'_s \sigma_\pi \cdot \sigma' = \sigma_s$  are substituted with *all* the strings recognized by  $\mathbf{A}_r$ . We can then characterize the language of the automaton returned by RP as  $\{ \sigma_\pi [\sigma_s / \sigma_r] \mid \sigma_r \in \mathcal{L}(\mathbf{A}_r) \}$ . Soundness is thus ensured:

$$\begin{aligned} [[ \text{replace}(\{\sigma\}, \{\sigma_s\}, \mathcal{L}_r) ]] &= \{ \sigma [\sigma_s / \sigma_r] \mid \sigma_r \in \mathcal{L}_r \} \\ &\subseteq \{ \sigma_\pi [\sigma_s / \sigma_r] \mid \sigma_r \in \mathcal{L}(\mathbf{A}_r) \} && \text{? automata lift ?} \\ &= \gamma(\text{RP}(\pi, \sigma_s, \mathbf{A}_r)) && \text{? def. RP ?} \end{aligned}$$

Soundness is thus proven as the result on individual strings can be lifted to languages, and because the  $\mathbf{A}_r$  passed to RP is an over-approximation of the concrete strings it represents (as the semantics performs a may-replacement whenever  $|\mathcal{L}_s| > 1$ ). The abstract semantics of `replace` is not complete, due to it returns  $\text{Min}(\{\mathbf{T}\})$  when either the input automaton or the search automaton contain cycles or read T.

As a counterexample for completeness, let  $[[ \mathbf{s} ]]^\# \text{mm}^\# = \mathbf{A} \in \mathcal{TFA}_{\equiv}$ ,  $[[ \mathbf{s}^s ]]^\# \text{mm}^\# = \mathbf{A}^s \in \mathcal{TFA}_{\equiv}$ ,  $[[ \mathbf{s}^r ]]^\# \text{mm}^\# = \mathbf{A}^r \in \mathcal{TFA}_{\equiv}$  s.t.  $\gamma(\mathbf{A}) = \{abc\}$ ,  $\gamma(\mathbf{A}^s) = \{a,z\}$  and  $\gamma(\mathbf{A}^r) = \{r\}$ .

$$[[ \text{replace}(\mathbf{s}, \mathbf{s}^s, \mathbf{s}^r) ]] \gamma(\text{mm}^\#) = \{rbc\} \not\subseteq \{abc, rbc\} = \gamma([[ \text{replace}(\mathbf{s}, \mathbf{s}^s, \mathbf{s}^r) ]]^\# \text{mm}^\#). \quad \square$$

## A.8 | Substring and CharAt

**Theorem A9.**  $[[ \text{substr}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_2) ]]^\#$  is a sound and complete abstraction of  $[[ \text{substr}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_2) ]]$ . Formally,

$$\forall \text{mm}^\# \in \mathbb{M}^\#, \forall \mathbf{s} \in \text{SE}, \forall \mathbf{a}_1, \mathbf{a}_2 \in \text{AE}. [[ \text{substr}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_2) ]] \gamma(\text{mm}^\#) = \gamma([[ \text{substr}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_2) ]]^\# \text{mm}^\#).$$

*Proof.* The collecting semantics of `substr` is defined as the additive lift of the concrete one, that is,  $[[ \text{substr}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_2) ]] \text{mm} = \{ \sigma \mid \sigma = \text{substr}(\sigma, i, j), \sigma \in \mathcal{L}, i \in I, j \in J \}$ , where  $[[ \mathbf{s} ]] \text{mm} = \mathcal{L} \in \wp(\Sigma^*)$ ,  $I = [[ \mathbf{a}_1 ]] \text{mm}$ ,  $J = [[ \mathbf{a}_2 ]] \text{mm}$  and  $\text{substr}: \Sigma^* \times \mathbb{N} \times \mathbb{N} \rightarrow \Sigma^*$  corresponds to the concrete semantics of Figure 3. Without loss of generality, we can prove the semantics to be sound when  $[[ \mathbf{a}_1 ]]^\# \text{mm}^\# = [i, i]$  and  $[[ \mathbf{a}_2 ]]^\# \text{mm}^\# = [j, j]$ , with  $i, j \in \mathbb{N}, 0 \leq i \leq j$ , as the abstract semantics lifts such result to



non-singleton intervals applying lub. Let  $[[s]]^{\#mm^{\#}} = A$  s.t.  $\gamma(A) = \mathcal{L}$ , and that  $[[a_1]]^{\#mm^{\#}} = [i,i]$  and  $[[a_2]]^{\#mm^{\#}} = [j,j]$ , with  $i, j \in \mathbb{N}$ . Furthermore, let  $r \equiv A$  be the regex equivalent to  $A$ . We can thus prove completeness of the semantics by proving the following:

$$[[\text{substr}(\mathcal{L}, \{i\}, \{j\})]]^{\gamma(mm^{\#})} = \gamma(\text{Min}(\{\sigma \mid (\sigma, 0, 0) \in \text{Sb}(r, i, j - i)\})).$$

Completeness is proven by structural induction over the structure of the regular expression, referencing the lines of Algorithm 4 that are involved in the computation as  $\$x$ , where  $x$  is the line number. Moreover, when  $\text{Sb}(r, i, j)$  produces the set  $S = \{(\sigma_1, i_1, j_1), \dots, (\sigma_n, i_n, j_n)\}$ , we denote the automaton  $\text{Min}(\{\sigma \mid (\sigma, 0, 0) \in S\})$  as either, abusing notation,  $\text{Min}(\text{Sb}(r, i, j))$  or  $\text{Min}(\{(\sigma_1, i_1, j_1), \dots, (\sigma_n, i_n, j_n)\})$ . With the latter notation, we abuse notation writing  $\sigma_i \notin \text{Sb}$  to denote that  $(\sigma_i, i, j)$  is not in the final result of  $\text{Sb}$ .

#### Base cases

$r = \emptyset$  ( $\mathcal{L}(r) = \emptyset$ ):

$$\begin{aligned} [[\text{substr}(\emptyset, \{i\}, \{j\})]] &= \emptyset \\ &= \gamma(\text{Min}(\emptyset)) && \text{? automata lift ?} \\ &= \gamma(\text{Min}(\text{Sb}(\emptyset, i, j - i))) && \text{? §2 ?} \end{aligned}$$

$r = \sigma \in \Sigma^*$ : here, we identify three cases. If  $i \leq j < |\sigma|$ :

$$\begin{aligned} [[\text{substr}(\{\sigma\}, \{i\}, \{j\})]] &= \{\sigma_i \dots \sigma_j\} \\ &= \gamma(\text{Min}(\{\sigma_i \dots \sigma_j\})) && \text{? automata lift ?} \\ &= \gamma(\text{Min}(\text{Sb}(\{\sigma\}, i, j - i))) && \text{? §6 ?} \end{aligned}$$

Instead, when  $i > |\sigma|$ :

$$\begin{aligned} [[\text{substr}(\{\sigma\}, \{i\}, \{j\})]] &= \emptyset \\ &= \gamma(\text{Min}(\{(\epsilon, i - |\sigma|, j - i)\})) && \text{? } i - |\sigma| > 0 \Rightarrow \epsilon \notin \text{Sb} \text{ ?} \\ &= \gamma(\text{Min}(\text{Sb}(\{\sigma\}, i, j - i))) && \text{? §4 ?} \end{aligned}$$

computing an empty partial substring (i.e., still concretized as the empty set of strings), but taking into account that  $\sigma$  has been read ( $i - |\sigma|$ ) and no character from  $\sigma$  has been taken ( $j - i$ ). Finally, if  $i < |\sigma|$  and  $j > |\sigma|$  (where  $k = j - |\sigma| + i$ ):

$$\begin{aligned} [[\text{substr}(\{\sigma\}, \{i\}, \{j\})]] &= \emptyset \\ &= \gamma(\text{Min}(\{(\sigma_i \dots \sigma_{|\sigma|-1}, 0, k)\})) && \text{? } k > 0 \Rightarrow \sigma_i \dots \sigma_{|\sigma|-1} \notin \text{Sb} \text{ ?} \\ &= \gamma(\text{Min}(\text{Sb}(\{\sigma\}, i, j - i))) && \text{? §5 ?} \end{aligned}$$

computing a partial substring (i.e., still concretized as the empty set of strings) that is a suffix of  $\sigma$ , and noting that  $j - (i - |\sigma_i \dots \sigma_{|\sigma|-1}|)$  characters still have to be read before completing the substring.  $r = T$ :

$$\begin{aligned} [[\text{substr}(\Sigma^*, \{i\}, \{j\})]] &= \{\sigma \mid |\sigma| = j - i\} \\ &= \gamma(\text{Min}(\{(\bullet^{j-i}, 0, 0)\})) && \text{? automata lift ?} \\ &\cup \gamma(\text{Min}(\{(\bullet^l, 0, j - l)\}), l < j - i) && \text{? } j - l > 0 \Rightarrow \bullet^l \notin \text{Sb} \text{ ?} \\ &\cup \gamma(\text{Min}(\{(\epsilon, i - l, j)\}), 0 \leq l \leq i) && \text{? } i - l > 0 \Rightarrow \epsilon \notin \text{Sb} \text{ ?} \\ &= \gamma(\text{Min}(\text{Sb}(T, i, j - i))) && \text{? §8, §9 ?} \end{aligned}$$

where, for the sake of clarity, strings returned by  $\text{Sb}$  are split into three sets, the first  $\{(\bullet^{j-i}, 0, 0)\}$  simulating substrings generated when  $i, j \leq |\sigma|$ , the second one  $\{(\bullet^l, 0, j - l)\}$  representing partial substrings when  $i \geq |\sigma|$ , and the third symbolizing partial substrings generated when  $i < |\sigma| \wedge j \geq |\sigma|$ . Note that only strings from the first set are part of the final concretization, while partial substrings from the second and third automata only serve in computations of successive substrings.

### Inductive steps

$r = r_1 || r_2$ : let  $\mathcal{L}, \mathcal{L}_1, \mathcal{L}_2 \in \wp(\Sigma^*)$  be the languages recognized by  $r$ ,  $r_1$ , and  $r_2$ , respectively. It is easy to see that  $[[ \text{substr}(\mathcal{L}, \{i\}, \{j\}) ]] = [[ \text{substr}(\mathcal{L}_1, \{i\}, \{j\}) ]] \cup [[ \text{substr}(\mathcal{L}_2, \{i\}, \{j\}) ]]$ . We assume  $[[ \text{substr}(\mathcal{L}_1, \{i\}, \{j\}) ]] = \gamma(\text{Min}(\text{Sb}(r_1, i, j - i)))$  and  $[[ \text{substr}(\mathcal{L}_2, \{i\}, \{j\}) ]] = \gamma(\text{Min}(\text{Sb}(r_2, i, j - i)))$  to hold for inductive hypothesis. We then prove soundness with the following:

$$\begin{aligned} [[ \text{substr}(\mathcal{L}, \{i\}, \{j\}) ]] &= [[ \text{substr}(\mathcal{L}_1, \{i\}, \{j\}) ]] \cup [[ \text{substr}(\mathcal{L}_2, \{i\}, \{j\}) ]] \\ &= \gamma(\text{Min}(\text{Sb}(r_1, i, j - i))) \cup \gamma(\text{Min}(\text{Sb}(r_2, i, j - i))) && \text{? ind. hp. } \S \\ &= \gamma(\text{Min}(\text{Sb}(r_1 || r_2, i, j - i))) && \text{? } \S 21 \S. \end{aligned}$$

$r = r_1 r_2$ : let  $\mathcal{L}, \mathcal{L}_1, \mathcal{L}_2 \in \wp(\Sigma^*)$  be the languages recognized by  $r$ ,  $r_1$ , and  $r_2$ , respectively. The concrete semantics is the union of two sets:  $[[ \text{substr}(\mathcal{L}_1, \{i\}, \{j\}) ]]$  (i.e., substrings that are fully contained in  $\mathcal{L}_1$ ), and  $[[ \text{substr}(\mathcal{L}_1 \cdot \mathcal{L}_2, \{i\}, \{j\}) ]]$  (i.e., substrings that straddle  $\mathcal{L}_1$  and  $\mathcal{L}_2$ ). We prove soundness assuming the inductive hypotheses  $[[ \text{substr}(\mathcal{L}_1, \{i\}, \{j\}) ]] = \gamma(\text{Min}(\text{Sb}(r_1, i, j - i)))$  and  $[[ \text{substr}(\mathcal{L}_2, \{i\}, \{j\}) ]] = \gamma(\text{Min}(\text{Sb}(r_2, i, j - i)))$ :

$$\begin{aligned} [[ \text{substr}(\mathcal{L}, \{i\}, \{j\}) ]] &= [[ \text{substr}(\mathcal{L}_1, \{i\}, \{j\}) ]] \cup [[ \text{substr}(\mathcal{L}_1 \cdot \mathcal{L}_2, \{i\}, \{j\}) ]] \\ &= \gamma(\text{Min}(\text{Sb}(r_1, i, j - i))) \cup \gamma(\text{Min}(\left\{ \left( \sigma_1^1 \cdot \sigma_2^1, i_2^1, j_2^1 \right), \dots, \left( \sigma_1^n \cdot \sigma_2^n, i_2^n, j_2^n \right) \right\})) && \text{? ind. hp. } \S \\ &= \gamma(\text{Min}(\text{Sb}(r_1 r_2, i, j - i))) && \text{? } \S 1, \S 16, \S 18 \S, \end{aligned}$$

where, for the sake of clarity, strings returned by  $\text{Sb}$  are split in two sets, the first  $(\text{Sb}(r_1, i, j - i))$  corresponding to substrings that entirely contained into  $r_1$ , the second one  $\left( \left\{ \left( \sigma_1^1 \cdot \sigma_2^1, i_2^1, j_2^1 \right), \dots, \left( \sigma_1^n \cdot \sigma_2^n, i_2^n, j_2^n \right) \right\} \right)$  that models substrings straddling  $r_1$  and  $r_2$ , where  $\forall i. \left( \sigma_1^i, i_1^i, j_1^i \right) \in \text{Sb}(r_1, i, j - i), j_1^i \neq 0 \wedge \left( \sigma_2^i, i_2^i, j_2^i \right) \in \text{Sb}(r_2, i_1^i, j_1^i)$ . Strings in the latter set are built by offsetting substrings of  $r_2$  by the length of the substrings of  $r_1$ .

$r = (r_1)^*$ . The proof of this case is similar to the one for concatenation, because  $(r_1)^*$  can be seen as an (undefined) concatenation of the regular expression  $r_1$ , and is thus left implicit.  $\square$

**Theorem A10.**  $[[ \text{charAt}(s, a) ]]$  is a sound and complete abstraction of  $[[ \text{charAt}(s, a) ]]$ . Formally,

$$\forall \text{mm}^\# \in \mathbb{M}^\#, \forall s \in \text{SE}, \forall a \in \text{AE}. [[ \text{charAt}(s, a) ]] \gamma(\text{mm}^\#) = \gamma([[ \text{charAt}(s, a) ]])^\# \text{mm}^\#.$$

*Proof.* Because the abstract semantics of  $\text{charAt}$  relies on the one of  $\text{substr}$ , soundness and completeness come from Theorem A9.  $\square$

### A.9 | String equality

We report the soundness and completeness proof of the abstract semantics of string equality. First, we prove the soundness of the  $\text{eq}$  function (whose algorithm is reported in Algorithm 1).

**Theorem A11.** Given  $\sigma_1, \sigma_2 \in \Sigma^* \cup \{\text{T}\}$ ,  $[[ \text{Flat}(\sigma_1) == \text{Flat}(\sigma_2) ]]$  is a sound and complete approximation of  $\text{eq}(\sigma_1, \sigma_2)$ . Formally,

$$\forall \sigma_1, \sigma_2 \in \Sigma^* \cup \{\text{T}\}. \text{Flat}(\sigma_1) == \text{Flat}(\sigma_2) = \gamma(\text{eq}(\sigma_1, \sigma_2)),$$

where we abuse notation denoting by  $==$  also the collecting semantics of string equality.

*Proof.* The proof is done by natural induction over the length of  $\sigma_1$  and  $\sigma_2$ .

#### Base cases

$$|\sigma_1| = 0 \Leftrightarrow \sigma_1 = \epsilon$$

- $|\sigma_2| = 0 \Leftrightarrow \sigma_2 = \epsilon$ . In this case,  $\text{Flat}(\sigma_1) == \text{Flat}(\sigma_2) = \{\text{true}\}$ , that is, equal to the result returned by  $\text{eq}$  in Algorithm 1 when both strings are empty (Lines 1 and 2).

- $|\sigma_2| = 1 \Leftrightarrow \sigma_2 = c \in \Sigma \cup \{T\}$ . We can split the proof in two cases:
  - $c = T$ : the empty string may be equal to T, thus  $\text{Flat}(\sigma_1) == \text{Flat}(\sigma_2) = \{\text{true}, \text{false}\}$ , that is, equal to the result returned by eq in Algorithm 1 when one of the string is empty and the other is equal to T (Lines 3 and 4).
  - $c \neq T$ : the empty string is not equal to any string, thus  $\text{Flat}(\sigma_1) == \text{Flat}(\sigma_2) = \{\text{false}\}$ , that is, equal to the result returned by eq in Algorithm 1 when one of string is empty and the other has a single character not equal to T (Lines 5 and 6).

$|\sigma_1| = 1 \Leftrightarrow \sigma_1 = c \in \Sigma \cup \{T\}$ . We can split the proof in two cases:

- $c = T$ .
  - $|\sigma_2| = 0 \Leftrightarrow \sigma_2 = \epsilon$ . This case is analogous to the second point of the previous base case.
  - $|\sigma_2| = 1 \Leftrightarrow \sigma_2 = c' \in \Sigma \cup \{T\}$ . We have two cases:
    - $c' = T$ . Two strings just made of a singleton T character may be equal, thus  $\text{Flat}(\sigma_1) == \text{Flat}(\sigma_2) = \{\text{true}, \text{false}\}$ , that is, equal to the result returned by eq in Algorithm 1 when both strings have a single character equal to T (Lines 7 and 8).
    - $c' \neq T$ . A string just made of a singleton T may be equal to  $c \in \Sigma \cup \{T\}$ , thus  $\text{Flat}(\sigma_1) == \text{Flat}(\sigma_2) = \{\text{true}, \text{false}\}$ , that is, equal to the result returned by eq in this case (Lines 7 and 8) Algorithm 1.
- $c \neq T$ .
  - $|\sigma_2| = 0 \Leftrightarrow \sigma_2 = \epsilon$ . This case is analogous to the second case of the first base case.
  - $|\sigma_2| = 1 \Leftrightarrow \sigma_2 = c' \in \Sigma \cup \{T\}$ . This case is analogous to the previous case ( $c = T$ ).

### Inductive steps

Let  $n \in \mathbb{N}$  and let  $\sigma_1, \sigma_2 \in (\Sigma \cup \{T\})^*$  such that  $|\sigma_1| \leq n, |\sigma_2| \leq n$ . For inductive hypothesis, the following holds:

$$\text{Flat}(\sigma_1) == \text{Flat}(\sigma_2) \subseteq \gamma(\text{eq}(\sigma_1, \sigma_2)).$$

Given  $\rho_1, \rho_2 \in (\Sigma \cup \{T\})^*$  such that  $|\rho_1| > n, |\rho_2| > n$ , we prove that

$$\text{Flat}(\rho_1) == \text{Flat}(\rho_2) \subseteq \gamma(\text{eq}(\rho_1, \rho_2)).$$

Let us consider  $\rho_1 = c\sigma_1$  and  $\rho_2 = c'\sigma_2$ . We split the proof in the following cases.

$c \neq T, c' \neq T$

$$\begin{aligned} \text{Flat}(\rho_1) == \text{Flat}(\rho_2) &= \\ &= \text{Flat}(c\sigma_1) == \text{Flat}(c'\sigma_2) && \{ \text{def. } \rho_1 \text{ and } \rho_2 \} \\ &= \{ cs == c's' \mid s \in \text{Flat}(\sigma_1), s' \in \text{Flat}(\sigma_2) \} && \{ \text{def. } == \} \end{aligned}$$

We split the proof in the following cases.

- $c = c'$

$$\begin{aligned} \{ cs == c's' \mid s \in \text{Flat}(\sigma_1), s' \in \text{Flat}(\sigma_2) \} &= \\ &= \{ s == s' \mid s \in \text{Flat}(\sigma_1), s' \in \text{Flat}(\sigma_2) \} && \{ c = c' \} \\ &= \gamma(\text{eq}(\sigma_1, \sigma_2)) && \{ \text{ind. hp.} \} \\ &= \gamma(\text{eq}(\rho_1, \rho_2)) && \{ \text{Lines 11 and 12} \}. \end{aligned}$$

- $c \neq c'$

$$\begin{aligned} \{ cs == c's' \mid s \in \text{Flat}(\sigma_1), s' \in \text{Flat}(\sigma_2) \} &= \\ &= \{\text{false}\} && \{ c \neq c' \} \\ &= \gamma(\text{eq}(\rho_1, \rho_2)) && \{ \text{Lines 9 and 10} \}. \end{aligned}$$

$c = T, c' \neq T$

$$\begin{aligned}
\text{Flat}(\rho_1) == \text{Flat}(\rho_2) &= \\
&= \text{Flat}(\mathbf{T}\sigma_1) == \text{Flat}(c'\sigma_2) && \text{? def. } \rho_1 \text{ and } \rho_2 \text{ ?} \\
&= \{ ts == c's' \mid s \in \text{Flat}(\sigma_1), s' \in \text{Flat}(\sigma_2), t \in \Sigma^* \} && \text{? def. } == \text{ ?} \\
&= \{ s == c's' \mid s \in \text{Flat}(\sigma_1), s' \in \text{Flat}(\sigma_2), t \in \Sigma^* \} \\
&\cup \{ cs == c's' \mid s \in \text{Flat}(\sigma_1), s' \in \text{Flat}(\sigma_2) \} \\
&\cup \{ ts == c's' \mid s \in \text{Flat}(\sigma_1), s' \in \text{Flat}(\sigma_2), t \in \Sigma^* \{ \epsilon, c \} \} && \text{? set union def. ?} \\
&= \gamma(\text{eq}(\rho_1, \rho_2[1:])) && \text{? ind. hp. ?} \\
&\sqcup \gamma(\text{eq}(\rho_1[1:], \rho_2[1:])) && \text{? ind. hp. ?} \\
&\sqcup \gamma(\text{eq}(\rho_1, \rho_2[1:]) \sqcup \text{false}) && \text{? ind. hp., } t = \epsilon \vee t \neq c \text{ ?} \\
&= \gamma(\text{eq}(\rho_1, \rho_2)) && \text{? Lines 13 and 14 ?}
\end{aligned}$$

$c = T, c' = T$ . The proof is analogous to the previous case.

$c \neq T, c' = T$ . The proof is analogous to the previous cases.  $\square$

**Theorem A12.**  $[[s == s']]^\sharp$  is a sound and complete abstraction of  $[[s == s']]$ . Formally,

$$\forall \text{mm}^\sharp \in \mathbb{M}^\sharp, \forall s, s' \in \text{SE}. [[s == s']] \gamma(\text{mm}^\sharp) \subseteq \gamma([[s == s']]^\sharp \text{mm}^\sharp).$$

*Proof.* In the first case of the abstract semantics of string equality, soundness and completeness are trivially met, while soundness and completeness of the third case follow from Theorem A11. Let us focus on the second case, let  $s, s' \in \text{SE}$  and suppose  $[[s]]^\sharp \text{mm}^\sharp = \mathbb{A} \in \mathcal{TF}_{\mathbb{A}/\equiv}$ ,  $[[s']]^\sharp \text{mm}^\sharp = \mathbb{A}' \in \mathcal{TF}_{\mathbb{A}/\equiv}$ . We prove that if either  $\mathbb{A}$  or  $\mathbb{A}'$  are cyclic,  $[[s == s']] \gamma(\text{mm}^\sharp) = \{\text{true}, \text{false}\}$ , the same result returned by  $\gamma([[s == s']]^\sharp \text{mm}^\sharp)$ , proving completeness. Note that  $[[s == s']] \gamma(\text{mm}^\sharp)$  cannot be  $\{\text{false}\}$  because this case is treated in the first case of the abstract semantics of string equality. By contradiction, let either  $\mathbb{A}$  or  $\mathbb{A}'$  be cyclic and let us suppose that  $[[s == s']] \gamma(\text{mm}^\sharp) = \{\text{true}\}$ .

$$\begin{aligned}
[[s == s']] \gamma(\text{mm}^\sharp) = \{\text{true}\} & \\
\Leftrightarrow \forall \sigma \in \gamma(\mathbb{A}) \forall \sigma' \in \gamma(\mathbb{A}'). \sigma == \sigma' &&& \text{? def. } [[s == s']] \text{ ?} \\
\Rightarrow |\gamma(\mathbb{A})| = |\gamma(\mathbb{A}')| = 1 &&& \text{? set theory, def. } == \text{ ?}
\end{aligned}$$

We supposed that either  $\mathbb{A}$  or  $\mathbb{A}'$  were cyclic, reaching a contradiction. Thus, if either  $\mathbb{A}$  or  $\mathbb{A}'$  are cyclic,  $[[s == s']] \gamma(\text{mm}^\sharp) = \{\text{false}, \text{true}\} = \gamma([[s == s']]^\sharp \text{mm}^\sharp)$ , proving completeness.  $\square$