

Tarsis: an effective automata-based abstract domain for string analysis

Luca Negrini² | Vincenzo Arceri*¹ | Agostino Cortesi² | Pietro Ferrara²

¹Department of Mathematical, Physical and Computer Sciences, University of Parma, Italy

²Department of Environmental Sciences, Informatics and Statistics, Ca' Foscari University of Venice, Italy

Correspondence

*Vincenzo Arceri

Email: vincenzo.arceri@unipr.it

Present address

Parco Area delle Scienze, 53/A, 43124, Parma, Italy

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.

KEY WORDS

String analysis, Static analysis, Abstract interpretation.

1 | INTRODUCTION

Nowadays, string values play a key role in any modern programming language, since 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^{1,2} 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 code³, reflection^{4,5}, string-to-code statement analysis⁶, and injection vulnerabilities^{7,8}.

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, space) required to retrieve useful information on strings. On the one hand, finite height string abstractions⁹ 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^{10,11}) 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 domain¹⁰. Over-approximating strings into finite state automata 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.

In this paper we introduce TARSIS, a new abstract domain for string values based on finite state automata (FSA). 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 in¹⁰, that we use as reference, either defined directly on the automaton or on its equivalent regular expression.

```

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

TARSIS has been implemented in GoLiSA¹², a static analyzer for Go based on LiSA^{13,12,14}. 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 of¹⁵. Specifically, we completed the previous version of the paper by covering all of the string operations considered in¹⁰, with the addition of the string equality operator. For all supported operations (both newly added and already formalized in¹⁵), we also reported proofs of soundness and completeness (or incompleteness). Moreover, we repeated the original experimental evaluation using GoLiSA¹² instead of the prototypical analyzer used in the original paper. Finally, to clearly assess the performance gain of TARSIS w.r.t. the finite-state automata abstract domain introduced in¹⁰, we extended our evaluation to deeply compare execution times of the two domains.

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

2 | MOTIVATING EXAMPLE

Consider the code of Fig. 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, since it requires to correctly model a set of string operations (namely `len`, `Index`, and substring) and their interaction. State-of-the-art string analyses fail to precisely model most of such operations, since 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 Fig. 1. As we will discuss in Sect. 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.

3 | PRELIMINARIES

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 $s' \curvearrowright_s s$. 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

```

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

two sets S and T , $\wp(S)$ is the powerset of S , $S \setminus T$ 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 , i.e., $S \cdot T = \{s \cdot t \mid s \in S, t \in T\}$. Given a set S and $n \in \mathbb{N}$, S^n is recursively defined as $S^0 \triangleq \{\epsilon\}$, and $S^{n>0} \triangleq S \cdot S^{n-1}$.

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 $\langle L, \leq, \vee, \wedge, \top, \perp \rangle$.

Abstract interpretation. Abstract interpretation^{1,2} 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 $C \xleftrightarrow[\alpha]{\gamma} A$. Given $C \xleftrightarrow[\alpha]{\gamma} A$, a concrete function $f : C \rightarrow C$ is, in general, not computable. Hence, a function $f^\# : A \rightarrow A$ that must *correctly* approximate the function f is needed. If so, we say that the function $f^\#$ is *sound*. Given $C \xleftrightarrow[\alpha]{\gamma} A$ and a concrete function $f : C \rightarrow C$, an abstract function $f^\# : A \rightarrow A$ is sound w.r.t. f if $\forall c \in C. \alpha(f(c)) \leq_A f^\#(\alpha(c))$, or equivalently $\forall a \in A. f(\gamma(a)) \leq_C \gamma(f^\#(a))$. Completeness¹⁶ can be obtained by enforcing the equality of the soundness conditions. Doing so, we obtain two notion of completeness. Given $C \xleftrightarrow[\alpha]{\gamma} A$, a concrete function $f : C \rightarrow C$ and an abstract function $f^\# : A \rightarrow A$, $f^\#$ is *backward complete* w.r.t. f if $\forall c \in C. \alpha(f(c)) = f^\#(\alpha(c))$, and it is *forward complete* w.r.t. f if $\forall a \in A. f(\gamma(a)) = \gamma(f^\#(a))$.

Finite state automata and regular expression notation. We follow the notation reported in¹⁰ for introducing finite state automata. A finite state automaton (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, Σ 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 finite state automaton. The set of all the FSAs is FA. If $\mathcal{L} \subseteq \Sigma^*$ is recognized by a FSA, we say that \mathcal{L} is a regular language. Given $\mathbb{A} \in \text{FA}$, $\mathcal{L}(\mathbb{A})$ 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, $\text{lminPath}(\mathbb{A}) \in \mathbb{N}$ denotes the (unique) length of a minimum path. If \mathbb{A} is a cycle-free automaton, $|\text{lmaxPath}(\mathbb{A})| \in \mathbb{N}$ denotes the (unique) length of a maximum path. Given $\pi = t_0 \dots t_n \in \text{paths}(\mathbb{A})$, σ_{π_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 , $r_1 \parallel r_2$ 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 .

The finite state automata abstract domain. Here, we report the necessary notions about the finite state automata abstract domain presented in¹⁰, over-approximating string properties as the minimum deterministic finite state automaton recognizing

$$\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' \sim_{\mathbf{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 $\sigma = \llbracket \mathbf{s} \rrbracket_{\mathfrak{m}}, \sigma' = \llbracket \mathbf{s}' \rrbracket_{\mathfrak{m}}, \sigma'' = \llbracket \mathbf{s}'' \rrbracket_{\mathfrak{m}}$

them. Given an alphabet Σ , the finite state automata domain is defined as $\langle \text{FA}_{/\equiv}, \sqsubseteq_{\text{FA}}, \sqcup_{\text{FA}}, \sqcap_{\text{FA}}, \text{Min}(\emptyset), \text{Min}(\Sigma^*) \rangle$, where $\text{FA}_{/\equiv}$ is the quotient set of FA w.r.t. the equivalence relation induced by language equality, \sqsubseteq_{FA} is the partial order induced by language inclusion, \sqcup_{FA} and \sqcap_{FA} 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 Σ . We abuse notation by representing equivalence classes in $\text{FA}_{/\equiv}$ by one of its automaton (usually the minimum), i.e., when we write $\mathbb{A} \in \text{FA}_{/\equiv}$ we mean $[\mathbb{A}]_{/\equiv}$. Since $\text{FA}_{/\equiv}$ does not satisfy the Ascending Chain Condition (ACC), i.e., it contains infinite ascending chains, it is equipped with the parametric widening ∇_{FA}^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^{17,18}. For instance, let us consider the automata $\mathbb{A}, \mathbb{A}' \in \text{FA}_{/\equiv}$ recognizing the languages $\mathcal{L} = \{\epsilon, a\}$ and $\mathcal{L}' = \{\epsilon, a, aa\}$, respectively. The result of the application of the widening ∇_{FA}^n , with $n = 1$, is $\mathbb{A} \nabla_{\text{FA}}^n \mathbb{A}' = \mathbb{A}''$ s.t. $\mathcal{L}(\mathbb{A}'') = \{a^n \mid n \in \mathbb{N}\}$.

Core language and semantics. We introduce a core language IMP, whose syntax is reported in Fig. 2. Such language, besides supporting arithmetic expressions (AE) and Boolean expressions (BE), also supports all of the string expressions (SE) discussed in¹⁰, 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 \mathfrak{m} . The concrete semantics of IMP statements is captured by the function $\llbracket \text{st} \rrbracket : \mathbb{M} \rightarrow \mathbb{M}$. The semantics is defined in a standard way and for this reason has been omitted. Such semantics relies on the one of expressions, that we capture, abusing notation, as $\llbracket \text{e} \rrbracket : \mathbb{M} \rightarrow \text{VAL}$. We define the part concerning strings in Fig. 3.

4 | THE TARSIS ABSTRACT DOMAIN

In this section, we recast the original finite state abstract domain working over an alphabet of characters Σ , reported in Sect. 3, to an augmented abstract domain based on finite state automata 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 Σ^* . The main concern is that Σ^* 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 strings. Given an IMP program P , we denote by Σ_P^* any substring of strings appearing in P (the set Σ_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 in¹⁰ with the new alphabet as

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

The alphabet on which finite state automata are defined is $\mathbb{A}_P \triangleq \Sigma_P^* \cup \{\top\}$, where \top is a special symbol that we intend as "any possible string". Let \mathcal{TFA} be the set of any deterministic FSA over the alphabet \mathbb{A}_P . Since we can have more automata recognizing a language, $\mathcal{TFA}_{/\equiv}$ is the quotient set of \mathcal{TFA} w.r.t. the equivalence relation induced by language equality, that is, the elements of domain are equivalence classes. For simplicity, when we write $A \in \mathcal{TFA}_{/\equiv}$, we intend the equivalence class of A . $\sqsubseteq_{\mathcal{T}}$ is the partial order induced by language inclusion, $\sqcup_{\mathcal{T}}$ and $\sqcap_{\mathcal{T}}$ are the lub and the glb over elements of $\mathcal{TFA}_{/\equiv}$, 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 $\text{Min}(\mathbb{A}_P^*)$, namely the automaton recognizing any string over \mathbb{A}_P .

Similarly to the standard FSA domain $FA_{/\equiv}$, also $\mathcal{TFA}_{/\equiv}$ is not a complete lattice and, consequently, it does not form a Galois Connection with the string concrete domain $\wp(\Sigma^*)$. This comes from the non-existence, in general, of the best abstraction of a string set in $\mathcal{TFA}_{/\equiv}$ (e.g., a context-free language has no best abstract element in $\mathcal{TFA}_{/\equiv}$ approximating it). Nevertheless, this is not a concern since weaker forms of abstract interpretation are still possible¹⁹ still guaranteeing soundness relations between concrete and abstract elements (e.g., polyhedra^{20,21}). In particular, we can still ensure soundness comparing the concretizations of our abstract elements (cf. Sect. 8 of¹⁹). Hence, we define the concretization function $\gamma_{\mathcal{T}} : \mathcal{TFA}_{/\equiv} \rightarrow \wp(\Sigma^*)$ as $\gamma_{\mathcal{T}}(A) \triangleq \bigcup_{\sigma \in \mathcal{L}(A)} \text{Flat}(\sigma)$, where Flat converts a string in \mathbb{A}_P^* into a set of strings in Σ^* . For instance, $\text{Flat}(a \top b b c) = \{ a \sigma b b c \mid \sigma \in \Sigma^* \}$. Note that, the language of strings recognized by A corresponds to the concretization function reported above, namely $\mathcal{L}(A) = \gamma_{\mathcal{T}}(A)$.

Widening. Similarly to the standard automata domain $FA_{/\equiv}$, also $\mathcal{TFA}_{/\equiv}$ does not satisfy ACC, meaning that fixpoint computations over $\mathcal{TFA}_{/\equiv}$ may not converge in a finite time. Hence, we need to equip $\mathcal{TFA}_{/\equiv}$ with a widening operator to ensure the convergence of the analysis. We define the widening $\nabla_{\mathcal{T}}^n : \mathcal{TFA}_{/\equiv} \times \mathcal{TFA}_{/\equiv} \rightarrow \mathcal{TFA}_{/\equiv}$, 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 Sect. 3, that, informally speaking, can be seen as a *subset construction* algorithm²² 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 \top (i.e., $\text{Min}(\{\top\})$). The result of the function f is a string of the form $id = \top$, repeated zero or more times. Since 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 , i.e., the value of res at the end of the function. Let A , reported in Fig. 4a, be the automaton abstracting the value of res before starting the second iteration of the loop, and let A' , reported in Fig. 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 $\nabla_{\mathcal{T}}^n$, between A and A' , working as follows. We first compute $A \sqcup_{\mathcal{T}} A'$ (corresponding to the automaton reported in Fig. 4b except that also q_0 is also a final state). On this automaton, we merge any state that recognizes the same \mathbb{A}_P -strings of length n , with $n \in \mathbb{N}$. In our example, let n be 2. The

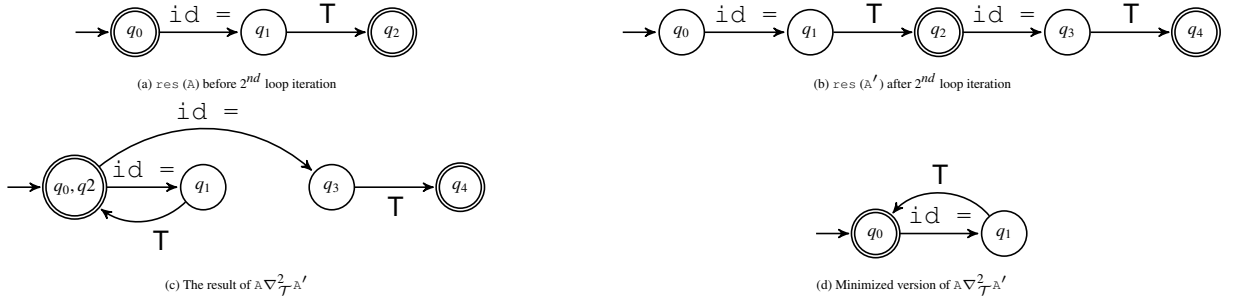


FIGURE 4 Example of widening application

resulting automaton is reported in Fig. 4c, where q_0 and q_2 are put together, and the other states are left as singletons. Fig. 4d depicts the minimized version of Fig. 4c.

The widening ∇_{τ}^n has been proved to meet the widening requirements (i.e., over-approximation of the least upper bounds and convergence on infinite ascending chains) in¹⁸. The parameter n , tuning the widening precision, is arbitrary and can be chosen by the user. As highlighted in¹⁰, the higher n is, the more the corresponding widening operator is precise in over-approximating lub of infinite ascending chains (i.e., in fixpoint computations). A classical improvement on widening-based fixpoint computations is to integrate a threshold²³, 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 ∇_{τ}^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 τ is not automatically inferred, since 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 Sect. 3 over the new string domain $\mathcal{TFA}_{/\equiv}$. Soundness and completeness (or incompleteness) proofs of the TARSIS's abstract semantics are reported in Appendix A. While TARSIS also implements `startsWith` and `endsWith` operators, their abstract semantics is not discussed in this section since we adopt the same ones reported in¹⁰.

Since 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 domain¹ 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_{Intv} and \sqcup_{Bool} the lub between intervals and Booleans, respectively. In particular, we merge such abstract domains in VAL^{\sharp} by the smashed sum abstract domain²⁴ $\text{VAL}^{\sharp} \triangleq \mathcal{TFA}_{/\equiv} \oplus \text{Intv} \oplus \text{Bool}$ 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 $\llbracket \text{st} \rrbracket^{\sharp} : \mathbb{M}^{\sharp} \rightarrow \mathbb{M}^{\sharp}$, relying on the abstract semantics of expressions defined by, abusing notation, $\llbracket e \rrbracket^{\sharp} : \mathbb{M}^{\sharp} \rightarrow \text{VAL}^{\sharp}$. We focus on the abstract semantics of string operations, while the semantics of the other expressions is standard and does not involve strings.

Concat

Given $A, A' \in \mathcal{TFA}_{/\equiv}$, the abstract semantics of `concat` returns a new automaton recognizing the language $\{ \sigma \cdot \sigma' \mid \sigma \in \mathcal{L}(A), \sigma' \in \mathcal{L}(A') \}$, that is, the concatenation between the strings of $\mathcal{L}(A)$ with the strings of $\mathcal{L}(A')$. This is achievable relying on the standard automata concatenation²². Let $\mathbf{s}, \mathbf{s}' \in \text{SE}$ and suppose that $\llbracket \mathbf{s} \rrbracket^{\sharp} = \langle Q, \Delta_P, \delta, q_0, F \rangle \in \mathcal{TFA}_{/\equiv}$, $\llbracket \mathbf{s}' \rrbracket^{\sharp} = \langle Q', \Delta_{P'}, \delta', q'_0, F' \rangle \in \mathcal{TFA}_{/\equiv}$. The abstract semantics of `concat` is:

$$\llbracket \text{concat}(\mathbf{s}, \mathbf{s}') \rrbracket^{\sharp} \triangleq \text{Min}(\langle Q \cup Q', \Delta_P, \delta \cup \delta' \cup \{ (q_f, \epsilon, q'_0) \mid q_f \in F \}, q_0, F' \rangle)$$

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



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

Length

Given $\mathbb{A} \in \mathcal{TFA}_{\equiv}$, 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 finite state automata. Let $\mathbf{s} \in \mathbf{SE}$, supposing that $\llbracket \mathbf{s} \rrbracket^{\sharp m^{\sharp}} = \mathbb{A} \in \mathcal{TFA}_{\equiv}$. The `length` abstract semantics is:

$$\llbracket \text{length}(\mathbf{s}) \rrbracket^{\sharp m^{\sharp}} \triangleq \begin{cases} [\text{lminPath}(\mathbb{A}), +\infty] & \text{if } \text{cyclic}(\mathbb{A}) \vee \text{readsTop}(\mathbb{A}) \\ [\text{lminPath}(\mathbb{A}), \text{lmaxPath}(\mathbb{A})] & \text{otherwise} \end{cases}$$

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 Fig. 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. Since 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 Fig. 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 $[3, 7]$.

Contains

Given $\mathbb{A}, \mathbb{A}' \in \mathcal{TFA}_{\equiv}$, 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 Fig. 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*²⁵, 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 σ_{sp} its longest string. Let $\mathbf{s}, \mathbf{s}' \in \mathbf{SE}$, supposing that $\llbracket \mathbf{s} \rrbracket^{\sharp m^{\sharp}} = \mathbb{A} \in \mathcal{TFA}_{\equiv}$, $\llbracket \mathbf{s}' \rrbracket^{\sharp m^{\sharp}} = \mathbb{A}' \in \mathcal{TFA}_{\equiv}$. The `contains` abstract semantics is:

$$\llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket^{\sharp m^{\sharp}} \triangleq \begin{cases} \{\text{false}\} & \text{if } \mathbb{A}' \sqcap_{\mathcal{T}} \text{FA}(\mathbb{A}) = \text{Min}(\emptyset) \\ \{\text{true}\} & \text{if } \text{singlePath}(\mathbb{A}') \wedge \forall \pi \in \text{paths}(\mathbb{A}^{ac}). \sigma_{\text{sp}} \curvearrow_{\mathbf{s}} \sigma_{\pi} \\ \{\text{true}, \text{false}\} & \text{otherwise} \end{cases}$$

In the first case, we denote by $\text{FA}(\mathbb{A})$ the *factor* automaton of \mathbb{A} , i.e., the automaton recognizing any substring of \mathbb{A} . In particular, if none of \mathbb{A} 's substrings is part of \mathbb{A}' , the abstract semantics safely returns `false` (checking the emptiness of the greatest lower bound between $\text{FA}(\mathbb{A})$ and \mathbb{A}'). Then, if \mathbb{A}' is a single-path automaton, the abstract semantics returns `true` if every path of \mathbb{A}^{ac} reads the longest string of \mathbb{A}' , with \mathbb{A}^{ac} being a copy of \mathbb{A} where all the cycles have been removed. Considering \mathbb{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 $\sigma_{\text{sp}} \curvearrow_{\mathbf{s}} \sigma_{\pi}$ the fact that σ_{sp} is a substring of each string in $\text{Flat}(\sigma_{\pi})$. Otherwise, $\{\text{true}, \text{false}\}$ is returned.

String equality

Given $\mathbb{A}, \mathbb{A}' \in \mathcal{TFA}_{\equiv}$, the abstract semantics of `string equality` returns `true` when \mathbb{A} and \mathbb{A}' recognize a singleton string and they are equal, `false` if no string recognized by \mathbb{A} is equal to any string recognized by \mathbb{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 $\Sigma^* \cup \{\mathbf{T}\}$. For example, the strings $a\mathbf{T}b$ and abb may be equal, while $a\mathbf{T}b$ and abd definitely are not. Alg. 1 defines the function $\text{eq} : \{\Sigma \cup \{\mathbf{T}\}\}^* \times \{\Sigma \cup \{\mathbf{T}\}\}^* \rightarrow \text{Bool}$, working on the expanded alphabet $\{\Sigma \cup \{\mathbf{T}\}\}^*$.

Algorithm 1 eq algorithm

```

  let  $\sigma, \sigma' \in \{\Sigma \cup \{\top\}\}^*$ 
1: if  $\sigma = \sigma' \wedge \sigma = \epsilon$  then
2:   return true
3: else if  $(\sigma = \epsilon \wedge \sigma' = \top) \vee (\sigma' = \epsilon \wedge \sigma = \top)$  then
4:   return {true, false}
5: else if  $(\sigma = \epsilon \wedge \sigma' \neq \top) \vee (\sigma' = \epsilon \wedge \sigma \neq \top)$  then
6:   return false
7: else if  $\sigma = \top \vee \sigma' = \top$  then
8:   return {true, false}
9: else if  $\sigma_0 \neq \sigma'_0 \wedge \sigma_0 \neq \top \wedge \sigma'_0 \neq \top$  then
10:  return false
11: else if  $\sigma_0 = \sigma'_0 \wedge \sigma_0 \neq \top$  then
12:  return eq( $\sigma[1:]$ ,  $\sigma'[1:]$ )
13: else if  $\sigma_0 = \top \vee \sigma'_0 = \top$  then
14:  return false  $\sqcup$  eq( $\sigma, \sigma'[1:]$ )  $\sqcup$  eq( $\sigma[1:]$ ,  $\sigma'[1:]$ )  $\sqcup$  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 \top characters do not appear (lines 2, 6, and 10). Note that, when one of the given strings begins with \top (lines 13-14), the algorithm can only prove inequality.

Let $\mathbf{s}, \mathbf{s}' \in \text{SE}$ and suppose $\llbracket \mathbf{s} \rrbracket^{\#m\#} = \mathbb{A}$ and $\llbracket \mathbf{s}' \rrbracket^{\#m\#} = \mathbb{A}'$. The abstract semantics of string equality is defined as:

$$\llbracket \mathbf{s} == \mathbf{s}' \rrbracket^{\#m\#} \triangleq \begin{cases} \{\text{false}\} & \text{if } \mathbb{A} \sqcap_{\mathcal{T}} \mathbb{A}' = \text{Min}(\emptyset) \\ \{\text{true}, \text{false}\} & \text{if } \text{cyclic}(\mathbb{A}) \vee \text{cyclic}(\mathbb{A}') \\ \bigsqcup_{\sigma \in \mathcal{L}(\mathbb{A}), \sigma' \in \mathcal{L}(\mathbb{A}')} \text{eq}(\sigma, \sigma') & \text{otherwise} \end{cases}$$

In the first case, if the greatest lower bound between \mathbb{A} and \mathbb{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 \mathbb{A} or \mathbb{A}' are cyclic, the abstract semantics of string equality returns $\{\text{true}, \text{false}\}$. Otherwise, we rely on `eq` to compare the \mathbb{A}_P -strings recognized by \mathbb{A} and \mathbb{A}' and we lub the results. To avoid cluttering the notation, the conversion from strings over $\Sigma^* \cup \{\top\}$ to strings over $\{\Sigma \cup \{\top\}\}^*$ when calling the function `eq` is implicit.

IndexOf

Given $\mathbb{A}, \mathbb{A}' \in \mathcal{TFA}_{\equiv}$, the `indexOf` abstract semantics returns an interval of the first indexes of the strings of $\mathcal{L}(\mathbb{A}')$ inside strings of $\mathcal{L}(\mathbb{A})$, recalling that when at least one string of $\mathcal{L}(\mathbb{A}')$ is not a substring of any string of $\mathcal{L}(\mathbb{A})$, the resulting interval must take into account -1 as well. Let $\mathbf{s}, \mathbf{s}' \in \text{SE}$ and suppose $\llbracket \mathbf{s} \rrbracket^{\#m\#} = \mathbb{A}$ and $\llbracket \mathbf{s}' \rrbracket^{\#m\#} = \mathbb{A}'$. The abstract semantics of `indexOf` is defined as:

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

If one of the automata has cycles or the automaton abstracting strings we aim to search for (i.e., \mathbb{A}') has a \top -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 since \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 `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 σ in strings recognized by \mathbb{A} . Since \mathbb{A}' recognizes a finite language, we apply `IO`(\mathbb{A}, σ) to each $\sigma \in \mathcal{L}(\mathbb{A}')$ and to return the least upper bound of the resulting intervals. In particular, the function `IO`(\mathbb{A}, σ) returns an interval $[i, j] \in \text{Intv}$ where, i and j are computed as follows.

$$i = \begin{cases} -1 & \text{if } \exists \pi \in \text{paths}(\mathbb{A}). \sigma \not\curvearrow_{\mathbf{s}} \sigma_{\pi} \\ \min_{\pi \in \text{paths}(\mathbb{A})} \left\{ i \mid \sigma_f \in \text{Flat}(\sigma_{\pi}) \wedge \sigma_{f_i} \dots \sigma_{f_{i+n}} = \sigma \right\} & \text{otherwise} \end{cases}$$

$$j = \begin{cases} -1 & \text{if } \forall \pi \in \text{paths}(\mathbb{A}). \sigma \not\prec_{\mathcal{S}} \sigma_{\pi} \\ +\infty & \text{if } \exists \pi \in \text{paths}(\mathbb{A}). \sigma \prec_{\mathcal{S}} \sigma_{\pi} \wedge \pi \text{ reads } \mathbb{T} \text{ before } \sigma \\ \max_{\pi \in \text{paths}(\mathbb{A})} \left\{ i \mid \sigma_f \in \text{Flat}(\sigma_{\pi}) \wedge \sigma_{f_i} \dots \sigma_{f_{i+n}} = \sigma \wedge \sigma \not\prec_{\mathcal{S}} \sigma_{f_0} \dots \sigma_{f_{i+n-1}} \right\} & \text{otherwise} \end{cases}$$

As for the abstract semantics of `contains`, we abuse notation denoting with $\sigma \prec_{\mathcal{S}} \sigma_{\pi}$ the fact that σ 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 σ 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 π of \mathbb{A} s.t. σ is not recognized by σ_{π} , then the minimal position where σ can be found in \mathbb{A} does not exist and -1 is returned. Otherwise, the minimal position where σ begins across π is returned. Let us consider now the computation of the maximal position. If all paths of the automaton do not recognize σ , then -1 is returned. If there exists a path where σ is recognized but the character \mathbb{T} appears earlier in the path, then $+\infty$ is returned. Otherwise, the maximal index of the first occurrences of σ across the paths of \mathbb{A} is returned.

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 Alg. 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 1–2 handle the case when n is zero and return the automaton recognizing the empty string. Lines 3–8 handle the case when \mathbb{A}' is a cyclic automaton and they build the automaton \mathbb{A}' returned at line 8, corresponding to the n -concatenation of the automaton \mathbb{A} . Otherwise, for each path π of \mathbb{A} , lines 12–15 builds the automaton \mathbb{A}' that corresponds to the n -repetition of σ_{π} , i.e., the string read by the path π . This operation is repeated for each path of the automaton (lines 11–17), and the obtained automata are lubbed together in \mathbb{A}_r at line 16, that is finally returned.

Let $\mathbf{s} \in \text{SE}$ and $\mathbf{a} \in \text{AE}$, supposing that $\llbracket \mathbf{s} \rrbracket^{\#} \mathbf{m}^{\#} = \mathbb{A} \in \mathcal{TFA}_{/\equiv}$ and $\llbracket \mathbf{a} \rrbracket^{\#} \mathbf{m}^{\#} = [i, j] \in \text{Intv}$. 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:

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



FIGURE 6 (a) A , (b) $\llbracket \text{repeat}(A, [2, +\infty]) \rrbracket^\sharp$

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 $i \in \mathbb{N} \setminus \{0\}$, since 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 Fig. 6. Otherwise, i.e., the interval $[i, j]$ is finite, the repeat 's abstract semantics returns the lub of the k -repetition of A , for each $i \leq k \leq j$.

TrimLeft, TrimRight, and Trim

The concrete semantics of trimLeft removes leading whitespaces (i.e., at the begin) from a string. Similarly, its abstract semantics inputs an automaton $A \in \mathcal{TFA}_{\equiv}$ and removes leading whitespaces from the begin of each string recognized by A . Let $A \in \mathcal{TFA}_{\equiv}$ 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.

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

The predicate $\text{readWS}(r)$ holds if the language of r recognizes a whitespace string, i.e., $\text{readWS}(r) \iff \exists ws \in \{_ \}^* . ws \in \mathcal{L}(r)$. If the regex is empty or \top , trimL behaves as the identity function (first case). If the regex is an atom, we rely on the concrete semantics of trimLeft (second case). If the regex is a disjunction, the result is the disjunction of the application of 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 $(_ \parallel a)_ b$. The concrete semantics of trimLeft would return $\{a_b, b\}$. Our abstract semantics, in this case, correctly returns $\text{trimL}((_ \parallel a)(_ b \parallel \text{trimL}(_ b))) = a(_ b \parallel b) = a_b \parallel b$. 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 trimR that removes trailing whitespace from the input regex. The definition of trimR is left implicit since it is analogous to the one of trimL .

Let $s \in SE$, supposing that $\llbracket s \rrbracket^\sharp m^\sharp = A \in \mathcal{TFA}_{\equiv}$ and let r be the regex equivalent to A . The abstract semantics of trimLeft , trimRight , and trim are:

$$\llbracket \text{trimLeft}(s) \rrbracket^\sharp m^\sharp = \text{trimL}(r) \quad \llbracket \text{trimRight}(s) \rrbracket^\sharp m^\sharp = \text{trimR}(r) \quad \llbracket \text{trim}(s) \rrbracket^\sharp m^\sharp = \text{trimL}(\text{trimR}(r))$$

Replace

To give the intuition about how the abstract semantics of replace works, consider three automata $A, A_s, A_r \in \mathcal{TFA}_{\equiv}$. Let us refer to A_s as the *search automaton* and to A_r as the *replace automaton*. Roughly speaking, the abstract semantics of 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



FIGURE 7 Example of may-replacement

of the following example. Consider $A \in \mathcal{TFA}_{\equiv}$ that is depicted in Fig. 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 `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 $\{bbb, cc\}$. Since it is unknown which string *must* be replaced (between bbb and cc), the `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 `replace` only extends A with the replace automaton. An example of may replacement is reported in Fig. 7, where A is the one reported in Fig. 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 `replace`, we define how to replace a string into an automaton's path. In particular, we define algorithm `RP` in Alg. 3, that given a path π of an arbitrary automaton, a replace automaton A^r , and $\sigma^s \in \Sigma^* \cup \{\mathbf{T}\}$ returns a new automaton built starting from the path, but where portions of the path that recognize σ^s have been replaced with A^r .

Algorithm 3 RP algorithm

```

1: let  $\pi = (q_0, \sigma_0, q_1), \dots, (q_{n-1}, \sigma_{n-1}, q_n), A^r = \langle Q^r, \mathbb{A}, \delta^r, q_0^r, F^r \rangle \in \mathcal{TFA}_{\equiv}, \sigma^s \in \Sigma^* \cup \{\mathbf{T}\}$ 
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  $\langle Q^{result}, \mathbb{A}, \delta^{result}, q_0^o, F^o \rangle$ 

```

Alg. 3 searches the given string σ^s across path π , collecting the sequences of transitions that recognize the search string σ^s and extracting them from π (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 ϵ -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 5-6). 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), since 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 \mathbf{T} -transition, we return $\text{Min}(\{\mathbf{T}\})$, namely the automaton recognizing \mathbf{T} . Otherwise, the `replace` abstract semantics is:

$$\llbracket \text{replace}(\mathbf{S}, \mathbf{S}_s, \mathbf{S}_r) \rrbracket^{\#m\#} \triangleq \begin{cases} A & \text{if } \forall \sigma_s \in \mathcal{L}(A_s). \nexists \sigma \in \mathcal{L}(A). \sigma_s \curvearrow_{\mathbf{S}} \sigma \\ \bigsqcup_{\pi \in \text{paths}(A)} \text{RP}(\pi, \sigma_s, A_r) & \text{if } \mathcal{L}(A_s) = \{\sigma_s\} \\ \bigsqcup_{\substack{\sigma \in \mathcal{L}(A_s) \\ \pi \in \text{paths}(A)}} \text{RP}(\pi, \sigma, A_r \sqcup_{\mathcal{T}} \text{Min}(\{\sigma\})) & \text{otherwise} \end{cases}$$

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 σ_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 σ 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.

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  $\{(\sigma_i \dots \sigma_{|\sigma| - 1}, 0, j - |\sigma| + i)\}$ 
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  $r = r_1 \parallel r_2$  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

```

Substr and CharAt

Given $A \in \mathcal{TFA}_{/\equiv}$ 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 $A \in \mathcal{TFA}_{/\equiv}$, in the definition of the `substr` semantics, we rely on the corresponding regex r since the two representations are equivalent and regexes allow us to define a more intuitive formalization of the semantics of `substr`. Let us suppose that $\llbracket s \rrbracket^{\sharp} m^{\sharp} = A \in \mathcal{TFA}_{/\equiv}$ 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 $\llbracket \mathbf{a}_1 \rrbracket^{\sharp m^{\sharp}} = [i, i]$ and $\llbracket \mathbf{a}_2 \rrbracket^{\sharp m^{\sharp}} = [j, j]$, with $i, j \in \mathbb{Z}$. In this case, the abstract semantics is defined as:

$$\llbracket \text{substr}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_2) \rrbracket^{\sharp m^{\sharp}} \triangleq \bigsqcup \text{Min}(\{ \sigma \mid (\sigma, 0, 0) \in \text{Sb}(r, i, j-i) \})$$

where 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, Sb is defined by Alg. 4 and, given a regex r and $i, j \in \mathbb{N}$, it returns a set of triples of the form (σ, n_1, n_2) , such that σ is the *partial substring* that Alg. 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 Alg. 4 still needs to look for to successfully compute a substring. Hence, given $\text{Sb}(r, i, j)$, the result is a set of such triples; note that given an element of the resulting set (σ, n_1, n_2) , $n_2 = 0$ means that no more characters are needed and σ 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 . Alg. 4 is defined by case on the structure of the input regex r :

1. $j = 0$ or $r = \emptyset$ (lines 1-2): \emptyset is returned since 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 σ but ends in subsequent regexes, and we return a singleton set containing the substring of σ from i to its end, with $n_1 = 0$ since we begun collecting characters, and $n_2 = j - |\sigma| + i$ since we collected $|\sigma| - i$ characters; otherwise, the substring is fully inside σ , and we return the substring of σ from i to $i + j$, setting both n_1 and n_2 to 0;
3. $r = \top$ (lines 11-14): since 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. $r = r_1 \parallel r_2$ (lines 26-27): 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 $\llbracket \mathbf{a}_1 \rrbracket^{\sharp m^{\sharp}} = [i, j]$ and $\llbracket \mathbf{a}_2 \rrbracket^{\sharp m^{\sharp}} = [l, k]$, with $i, j, l, k \in \mathbb{Z}$, the abstract semantics of substr is

$$\llbracket \text{substr}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_2) \rrbracket^{\sharp m^{\sharp}} \triangleq \bigsqcup_{a \in [i, j], b \in [l, k], a \leq b} \text{Min}(\{ \sigma \mid (\sigma, 0, 0) \in \text{Sb}(r, a, b-a) \})$$

We do not report the cases when input intervals are unbounded (e.g., $[1, +\infty]$). Nevertheless, these cases have been already considered in¹⁰ 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:

$$\llbracket \text{charAt}(\mathbf{s}, \mathbf{a}_1) \rrbracket^{\sharp m^{\sharp}} \triangleq \llbracket \text{substr}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_1 + 1) \rrbracket^{\sharp m^{\sharp}}$$

5 | EXPERIMENTAL RESULTS

TARSIS has been compared with five string abstract domains, namely the prefix (PR), suffix (SU), char inclusion (C1), bricks (BR) domains (all defined in⁹), and FA_{\equiv} (defined in¹⁰). All domains have been implemented in GoLiSA, that we will briefly introduce before presenting our experimental results. TARSIS and FA_{\equiv} share a common implementation for the automata structure that is

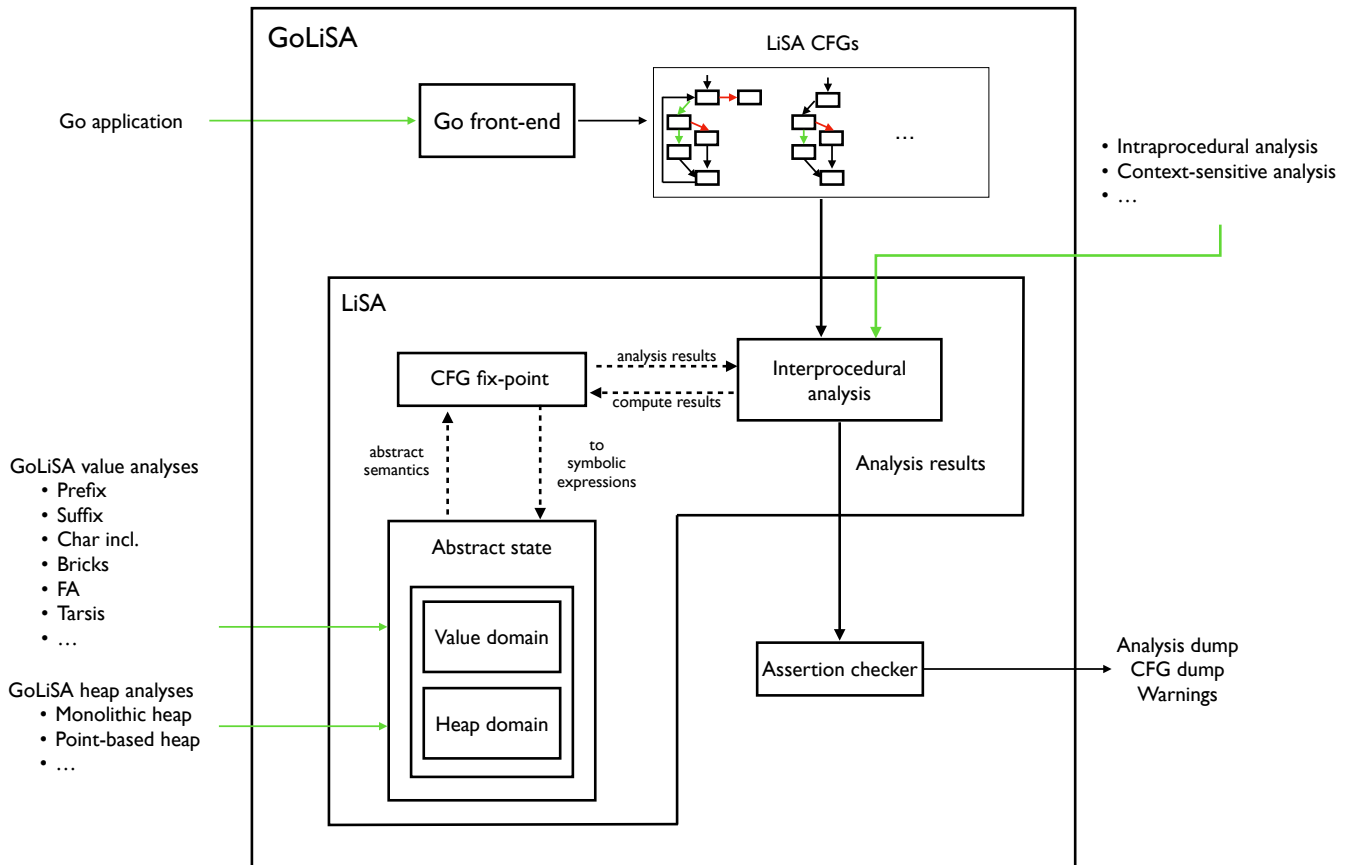


FIGURE 8 Schema of GoLiSA's architecture (taken from¹²)

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 Sect. 4.2 with trace partitioning²⁶ (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 @ 1.8GHz 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: Sect. 5.1 compares them to TARSIS and FA_{\equiv} without expanding the scope of such evaluations. Sect. 5.2 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, Sect. 5.4 focuses on the performance difference between FA_{\equiv} and TARSIS, benchmarking their lattice operations and abstract transformers.

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^{13,14}, whose high level infrastructure is visible in Fig. 8.

LiSA (**L**ibrary for **S**tatic **A**nalysis) 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

```

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

Domain	Program SUBS		Program LOOP	
PR	ring test	✗	Repeat:	✗
SU	ϵ	✗	ϵ	✗
CI	[[abdefgilnprstu]	✗	[:aepRt] [!:aepRt T]	✓
BR	[{ring test fai, ring test pas}] (1,1)	✓	[{T}] (0, + ∞)	✗
FA/ \equiv	ring test (pas fai)	✓	Repeat: (T)*	✓
TARISIS	(ring test pas ring test fai)	✓	Repeat: (T!)*	✓

TABLE 1 Values of res at the first assert of each program

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.

Code parsing logic and the definition of language-specific statements are provided by *Frontends* such as GoLiSA, that can also provide implementations for LiSA's components. These constitute 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 \top_{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`.

5.1 | Precision of the various domains on test cases

We start by considering programs SUBS (Fig. 9a) and LOOP (Fig. 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 an user input (i.e., an unknown string), to a constant value. Tab. 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, TARISIS and FA/ \equiv approximations are expressed as regexes.

When analyzing SUBS, both PR and SU lose precision since 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

```

1 func ToString(names []string) {
2   res := "People : {"
3   i := 0;
4   for i < 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 }

```

(a) Program TOSTRING

```

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>0);
10  assert (count==0);
11  assert (count==3);
12 }

```

(b) Program COUNTMATCHES

FIGURE 10 Programs used for assessing domain precision

empty set of included characters, and a set of possibly included characters containing the ones of both strings. Finally, BR, FA_{\equiv} and TARSIS are expressive enough to track any string produced by any concrete execution of SUBS.

When evaluating the assertions of SUBS, a PA should be raised on lines 10 and 11, since "p" or "f" might be in *res*, together with a DA alarm on line 12, since "d" is surely not contained in *res*. No alarm should be raised on line 9 instead, since "g" 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_{\equiv} , or TARSIS. Since 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 (since character "t" is always contained in the resulting string value), and PAs at lines 7 and 8. PR infers as prefix of *res* the string "Repeat: ", 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* since the lub operation occurring at line 3 cannot find a common suffix between "Repeat: " and "!", hence PAs are raised on lines 6-8. Since 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 since the string used in the `contains` call has length 1). The alternation of T and "!" 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_{Bool} on all `contains` calls, resulting in PAs on lines 6-8. The parametric widening of FA_{\equiv} collapses the colon into T. In TARSIS, since the automaton representing *res* grows by two states each iteration, the parametric widening defined in Sect. 4.1 can collapse the whole content of the loop into a 2-states 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_{\equiv} and TARSIS were the two only domains that produced the desired behavior in these rather simple test cases.

5.2 | Evaluation on realistic code samples

In this section, we explore two real world code samples. Method TOSTRING (Fig. 10a) transforms an array of names that come as string values into a single string. While it resembles the code of LOOP in Fig. 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 (Fig. 10b) makes use of `strings.Count` (reported in Sect. 2) to prove properties about its return value. Tab. 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 (since "People" is surely contained in the resulting string), and two PAs at line 13 and 14. PR, SU,

CI and BR raise PAs on all the three assert statements. FA_{\equiv} and TARSIS detect that the assertion at line 12 always holds. Thus, when using them, the analyzer raises PAs on lines 13 and 14 since (i) comma character is part of `res` if the loop is iterated at least once, and (ii) T might match "not".

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. Since 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_{\equiv} 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`, since all possible values of `str` start with `sub`. Since the call to `length` yields $[10, 17]$, all possible substrings from $[2, 2]$ (`idx` plus the length of `sub`) to $[10, 17]$ are computed (namely, "e throat", "is is th", "is is the", ..., "is is the thing"), and the resulting automaton is the one that recognizes all of them. Since 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 (namely "roat") 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_{\equiv} 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.

5.3 | 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_{\equiv} are the only ones able to obtain precise results on them. We now discuss the efficiency of the analyses. Tab. 3 reports the execution times for all the domains on the case 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 in¹⁵. 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_{\equiv} 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_{\equiv} can be accounted on the alphabets underlying the automata. In FA_{\equiv} , 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, since 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_{\equiv} .

TARSIS's alphabet has another peculiarity w.r.t. FA_{\equiv} '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

Domain	Program TOSTRING		Program COUNTMATCHES	
PR	People: {	✗	$[0, +\infty]$	✗
SU	ϵ	✗	$[0, +\infty]$	✗
CI	$\{ \{ : \text{Peopl }] \{ \{ : , \text{Peopl } \text{T} \}$	✗	$[0, +\infty]$	✗
BR	$\{ \{ \text{T} \} (0, +\infty)$	✗	$[0, +\infty]$	✗
FA_{\equiv}	People: $\{ (\text{T}) * \text{T} \}$	✓	$[2, 3]$	✓
TARSIS	People: $\{ \} \parallel \text{People: } \{ (\text{T},) * \text{T} \}$	✓	$[2, 3]$	✓

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

Domain	SUBS	LOOP	TOSTRING	COUNTMATCHES
PR	5ms	28ms	1s 509ms	129ms
SU	5ms	32ms	1s 599ms	110ms
CI	4ms	45ms	1s 594ms	138ms
BR	8ms	78ms	3s 911ms	345ms
FA _{/≡}	11ms	20s 895ms	22m 37s 621ms	9s 175ms
TARSIS	7ms	443ms	12s 437ms	123ms

TABLE 3 Execution times of the domains on each program

without additional tolls on their performances. FA_{/≡}'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_{/≡} 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_{/≡} 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>.

5.4 | Performance benchmark of TARSIS and FA_{/≡}

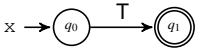
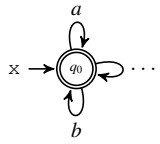
As automata-based abstractions have already been proved to be effective in string program analysis¹⁰, we focus the final part of our evaluation on the difference in resource consumption between TARSIS and FA_{/≡}. In fact, our main goal is proving that the adoption of TARSIS can make automata-based abstractions viable for the analysis of non-trivial code.

To measure the performance of TARSIS and FA_{/≡}, 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, interprocedural analyses), and would hence be affected by their running time and memory consumption. The different semantics of TARSIS and FA_{/≡} 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, i.e., the ones concerning TARSIS and FA_{/≡}, we directly compare each lattice operation and abstract transformer in isolation, benchmarking their execution times. This enables accurate speedup measurements, since such times are not affected by external factors like memory saturation caused by the remaining analysis components.

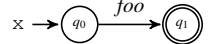
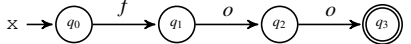
We compare lattice operators \sqsubseteq , \sqcup , \sqcap , and ∇ , together with the semantics of string operators (`concat`, `contains`, `length`, `substr`, `trim`, `repeat`, and `indexOf`); since `trim`'s semantics relies on `trimLeft` and `trimRight`, the comparison of the latter operations have been omitted. Moreover, concerning the `trimLeft`, `trimRight`, and `trim` abstract semantics of FA_{/≡}, we adopted the same one reported in Sect. 4.2 for TARSIS, recasted to standard automata, since a bug was found in the original FA_{/≡} 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_{/≡} returns the automaton recognizing the language $\{ a^n \mid n \in \mathbb{N} \}$, that is an unsound result (e.g., the string `a_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 7 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_{/≡} for the string variable `x` at the end of the fragment.

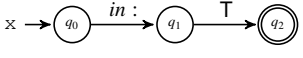
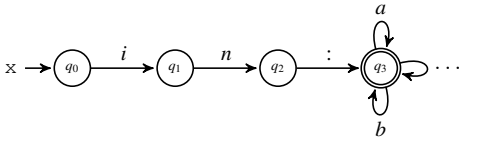
1. Statically unknown strings. Such class contains the abstractions of statically unknown strings, such as an user input. Since both domains have a unique minimum automaton for such strings, this class is made of a single automaton.

Go fragment	TARSIS	FA _{/≡}
<pre>_, err := fmt.Sprintf("%s", &x);</pre>		

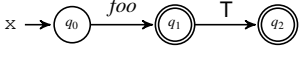
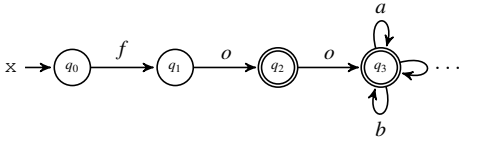
2. *Constant strings.* This class contains automata generated as abstractions of string literals appearing in the program.

Go fragment	TARSIS	FA _{/≡}
<pre>x := "foo";</pre>		

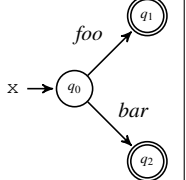
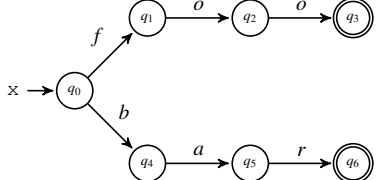
3. *Concatenated strings.* This class contains automata modeling the concatenation of simple strings, with each being either statically unknown or constant. Automata in this class are effective concatenations of ones from the previous classes.

Go fragment	TARSIS	FA _{/≡}
<pre>z := "in:"; _, err := fmt.Sprintf("%s", &y); x := z + y;</pre>		

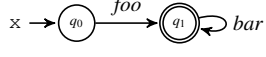
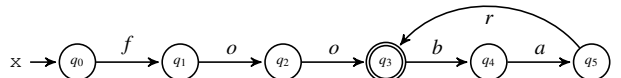
4. *Increasing strings.* The automata in this class model sets of strings built by optionally appending a finite number of strings at the end of an existing one, thus leading to single path automata. Each of such strings can be either statically unknown or constant (i.e., each automaton is part of one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

Go fragment	TARSIS	FA _{/≡}
<pre>x := "foo"; _, err := fmt.Sprintf("%s", &y); if b { x := x + y; }</pre>		

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.

Go fragment	TARSIS	FA _{/≡}
<pre>if b { x := "foo"; } else { x := "bar"; }</pre>		

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.

Go fragment	TARSIS	FA _{/≡}
<pre>x := "foo"; for b { x := x + "bar"; }</pre>		

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.

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

TABLE 4 Lattice operators and abstract transformers benchmark results

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 2 automata A_1 and A_2 from that class, and would then run `concat(A_1, A_2)`). Specifically, 1 round is executed using the single automaton from class 1, 5 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_{/≡} ones, and (iii) the execution of each operation individually, measuring the required time. Each operation’s execution has a 30 seconds 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_{/≡} automata grows fast, we impose some limits on the structure of each generated TARSIS automaton, such that, when we generate the equivalent FA_{/≡} 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;
- automata of class 7 can contain up to 5 states with up to 3 transitions per state; each state has a 25% chance of being an additional final state.

Benchmark results

Tab. 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_{/≡} is determined by only comparing execution times of the rounds where both domains ran successfully, reporting the ratio $t_{FA_{/≡}}/t_{TARSIS}$.

The benchmark confirms that TARSIS is overall reliably faster than FA_{/≡}. TARSIS operations time out significantly less than the ones of FA_{/≡} (2 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 Tab. 4. Moreover, when focusing only on rounds where both domains succeeded within the

required time bound, the speedup (column *Time ratio* of Tab. 4) is still noticeable: even in the worst case of `indexOf`, TARSIS still performs almost twice as fast as FA_{\equiv} .

Let us comment the results concerning `length` and `trim`. For these operations, TARSIS and FA_{\equiv} 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_{\equiv}) and in none when executing `length` (while FA_{\equiv} incurs in 26 timeouts), also enabling significant speedup (more than 2000x for `length` and almost 500x 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^{11,10,27,28,29,30,9}. As already discussed, the TARSIS abstract domain build upon the finite state automata abstract domain defined in¹⁰ 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 statements³¹. As reported by the experimental results in Sect. 5 (Tab. 3 in particular), TARSIS is quite more efficient than FA_{\equiv} 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, 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.

Static analysis of string values: Our experimental results compared TARSIS with several simple abstract domains introduced in^{32,9}. 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 Tabs. 1, 2, and 3, these abstractions are quite more efficient but less precise than TARSIS. Overall, TARSIS exposed execution times comparable to the ones of^{32,9}, and precision similar to^{10,31}.

Static analysis was applied to string values in many different contexts, such as SQL queries programmatically built by code³, reflection^{4,5}, and injection vulnerabilities^{7,8}. Instead, the main goal of our work is to introduce a novel abstract domain (that is, 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.

Automata-based approaches: Several approaches to statically analyze string values through automata (like TARSIS) have been proposed in the literature. For instance, the authors of³³ provided an automata abstraction merged with interval abstractions for analyzing JavaScript arrays and objects. Another interesting automata-based model is symbolic automata³⁴, 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 analysis³⁵, and their usage as program models for mixing syntactic and semantic abstractions over the program³⁶. Generally speaking, automata-based approaches are usually not efficient, since 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.

Regular expressions: A major stream of research focussed on the abstraction of string values through regular expressions. In particular, in²⁸, the authors 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] *)` (that is the string representation of integer values). However, such approximation is strictly more complex and led to less efficient analyses.

Similarly, regular strings³⁷ is an abstraction of the finite state automata 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.

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 formulas³⁸ or proposing new efficient decidable procedures or string constraints representations^{39,40} also based on automata,

such as^{41,42}, or involving type conversion string constraints⁴³. For instance, Z3-str⁴⁴ extended the SMT-solver Z3⁴⁵ 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 slow-downs (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 finite state automata 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 narrowing²³. Hence, a narrowing operator for TARSIS will be studied, in order to get more precise results on loops analyses.

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 domains⁴⁶.

Acknowledgements.

Work partially supported by Bando di Ateneo per la Ricerca 2022, funded by University of Parma, (MUR_DM737_2022_FIL_PROGETTI_B_ARCERI_COFIN, 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.

REFERENCES

1. 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, January 1977* ACM 1977:238–252
2. 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, January 1979* ACM Press 1979:269–282
3. 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
4. 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, May 20-28, 2017* IEEE / ACM 2017:507–518
5. 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, May 21-28, 2011* ACM 2011:241–250
6. Arceri V, Mastroeni I. A sound abstract interpreter for dynamic code. In: Hung C, Cerný T, Shin D, Bechini A., eds. *SAC '20: The 35th ACM/SIGAPP Symposium on Applied Computing, online event, [Brno, Czech Republic], March 30 - April 3, 2020* ACM 2020:1979–1988
7. Fu X, Lu X, Peltsverger B, Chen S, Qian K, Tao L. A Static Analysis Framework For Detecting SQL Injection Vulnerabilities. In: IEEE Computer Society 2007:87–96
8. 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, July 31 - August 5, 2005* USENIX Association 2005.
9. 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. doi: 10.1002/SPE.2218
10. Arceri V, Mastroeni I, Xu S. Static Analysis for ECMAScript String Manipulation Programs. *Appl. Sci.*. 2020;10:3525. doi: 10.3390/app10103525
11. 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, August 29-31, 2018* IEEE Computer Society 2018:1–8
12. Olivieri L, Negrini L, Arceri V, et al. Information Flow Analysis for Detecting Non-Determinism in Blockchain. In: Ali K, Salvaneschi G., eds. *37th European Conference on Object-Oriented Programming, ECOOP 2023, July 17-21, 2023, Seattle, Washington, United States*. 263 of *LIPICs*. Schloss Dagstuhl - Leibniz-Zentrum für Informatik 2023:23:1–23:25

13. 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, 22 June, 2021* ACM 2021:1–6
14. Negrini L, Ferrara P, Arceri V, Cortesi A. *LiSA: A Generic Framework for Multilanguage Static Analysis*:19–42; Singapore: Springer Nature Singapore . 2023
15. Negrini L, Arceri V, Ferrara P, Cortesi A. Twinning Automata and Regular Expressions for String Static Analysis. In: Henglein F, Shoham S, Vizel Y., eds. *Verification, Model Checking, and Abstract Interpretation - 22nd International Conference, VMCAI 2021, Copenhagen, Denmark, January 17-19, 2021, Proceedings*. 12597 of *Lecture Notes in Computer Science*. Springer 2021:267–290
16. Giacobazzi R, Ranzato F, Scozzari F. Making abstract interpretations complete. *J. ACM*. 2000;47(2):361–416. doi: 10.1145/333979.333989
17. Bartzis C, Bultan T. Widening Arithmetic Automata. In: Alur R, Peled DA., eds. *Computer Aided Verification, 16th International Conference, CAV 2004, Boston, MA, USA, July 13-17, 2004, Proceedings*. 3114 of *Lecture Notes in Computer Science*. Springer 2004:321–333
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. doi: 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. doi: 10.1016/J.SCICO.2007.08.001
21. Becchi A, Zaffanella E. PPLite: Zero-overhead encoding of NNC polyhedra. *Inf. Comput.*. 2020;275:104620. doi: 10.1016/J.IC.2020.104620
22. Davis MD, Sigal R, Weyuker EJ. *Computability, Complexity, and Languages: Fund. of Theor. CS*. 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. doi: 10.1016/J.CL.2010.09.001
24. Arceri V, Maffeis S. Abstract Domains for Type Juggling. *Electron. Notes Theor. Comput. Sci.*. 2017;331:41–55. doi: 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: Hierons RM, Mosbah M., eds. *Theoretical Aspects of Computing - ICTAC 2019 - 16th International Colloquium, Hammamet, Tunisia, October 31 - November 4, 2019, Proceedings*. 11884 of *Lecture Notes in Computer Science*. Springer 2019:255–272
26. Rival X, Mauborgne L. The Trace Partitioning Abstract Domain. *ACM Trans. Program. Lang. Syst.*. 2007;29(5):26–es. doi: 10.1145/1275497.1275501
27. Park C, Im H, Ryu S. Precise and scalable static analysis of jQuery using a regular expression domain. In: Ierusalimsky R., ed. *Proceedings of the 12th Symposium on Dynamic Languages, DLS 2016, Amsterdam, The Netherlands, November 1, 2016* ACM 2016:25–36
28. Christensen AS, Møller A, Schwartzbach MI. Precise Analysis of String Expressions. In: Cousot R., ed. *Static Analysis, 10th International Symposium, SAS 2003, San Diego, CA, USA, June 11-13, 2003, Proceedings*. 2694 of *Lecture Notes in Computer Science*. Springer 2003:1–18
29. Madsen M, Andreasen E. String Analysis for Dynamic Field Access. In: Cohen A., ed. *Compiler Construction - 23rd International Conference, CC 2014, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2014, Grenoble, France, April 5-13, 2014, Proceedings*. 8409 of *Lecture Notes in Computer Science*. Springer 2014:197–217
30. Abdulla PA, Atig MF, Chen Y, et al. String Constraints for Verification. In: Biere A, Bloem R., eds. *Computer Aided Verification - 26th International Conference, CAV 2014, Held as Part of the Vienna Summer of Logic, VSL 2014, Vienna, Austria, July 18-22, 2014, Proceedings*. 8559 of *Lecture Notes in Computer Science*. Springer 2014:150–166
31. 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. doi: 10.1145/3426470
32. Costantini G, Ferrara P, Cortesi A. Static Analysis of String Values. In: Qin S, Qiu Z., eds. *Formal Methods and Software Engineering - 13th International Conference on Formal Engineering Methods, ICFEM 2011, Durham, UK, October 26-28, 2011, Proceedings*. 6991 of *Lecture Notes in Computer Science*. Springer 2011:505–521
33. Almashfi N, Lu L. Precise String Domain for Analyzing JavaScript Arrays and Objects. In: May Huang SZ., ed. *3rd International Conference on Information and Computer Technologies, ICICT 2020, San Jose, CA, USA, March 9-12, 2020* IEEE 2020:17–23
34. D’Antoni L, Veanes M. Minimization of symbolic automata. In: Jagannathan S, Sewell P., eds. *The 41st Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL ’14, San Diego, CA, USA, January 20-21, 2014* ACM 2014:541–554
35. Veanes M. Applications of Symbolic Finite Automata. In: Konstantinidis S., ed. *Implementation and Application of Automata - 18th International Conference, CIAA 2013, Halifax, NS, Canada, July 16-19, 2013, Proceedings*. 7982 of *Lecture Notes in Computer Science*. Springer 2013:16–23
36. Preda MD, Giacobazzi R, Lakhotia A, Mastroeni I. Abstract Symbolic Automata: Mixed syntactic/semantic similarity analysis of executables. In: Rajamani SK, Walker D., eds. *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2015, Mumbai, India, January 15-17, 2015* ACM 2015:329–341
37. Choi T, Lee O, Kim H, Doh K. A Practical String Analyzer by the Widening Approach. In: Kobayashi N., ed. *Programming Languages and Systems, 4th Asian Symposium, APLAS 2006, Sydney, Australia, November 8-10, 2006, Proceedings*. 4279 of *Lecture Notes in Computer Science*. Springer 2006:374–388
38. Abdulla PA, Atig MF, Diep BP, Holík L, Janku P. Chain-Free String Constraints. In: Chen Y, Cheng C, Esparza J., eds. *Automated Technology for Verification and Analysis - 17th International Symposium, ATVA 2019, Taipei, Taiwan, October 28-31, 2019, Proceedings*. 11781 of *Lecture Notes in Computer Science*. Springer 2019:277–293
39. 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(POPL):49:1–49:30. doi: 10.1145/3290362
40. Amadini R, Gange G, Stuckey PJ. Dashed strings for string constraint solving. *Artif. Intell.*. 2020;289:103368. doi: 10.1016/J.ARTINT.2020.103368
41. Wang H, Chen S, Yu F, Jiang JR. A symbolic model checking approach to the analysis of string and length constraints. In: Huchard M, Kästner C, Fraser G., eds. *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering, ASE 2018, Montpellier, France, September 3-7, 2018* ACM 2018:623–633
42. 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. doi: 10.1007/s10703-013-0189-1
43. Abdulla PA, Atig MF, Chen Y, et al. Efficient handling of string-number conversion. In: Donaldson AF, Torlak E., eds. *Proceedings of the 41st ACM SIGPLAN International Conference on Programming Language Design and Implementation, PLDI 2020, London, UK, June 15-20, 2020* ACM 2020:943–957

44. Zheng Y, Zhang X, Ganesh V. Z3-str: a z3-based string solver for web application analysis. In: Meyer B, Baresi L, Mezini M., eds. *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, August 18-26, 2013* ACM 2013:114–124
45. Moura dLM, Bjørner NS. Z3: An Efficient SMT Solver. In: Ramakrishnan CR, Rehof J., eds. *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, Budapest, Hungary, March 29-April 6, 2008. Proceedings.* 4963 of *Lecture Notes in Computer Science*. Springer 2008:337–340
46. Arceri V, Oliario M, Cortesi A, Ferrara P. Relational String Abstract Domains. In: Finkbeiner B, Wies T., eds. *Verification, Model Checking, and Abstract Interpretation - 23rd International Conference, VMCAI 2022, Philadelphia, PA, USA, January 16-18, 2022, Proceedings.* 13182 of *Lecture Notes in Computer Science*. Springer 2022:20–42

□

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 $\text{VAL}^\sharp \triangleq \mathcal{TFA}_{/\equiv} \oplus \text{Intv} \oplus \text{Bool}$, we compare concretizations of its elements with a concrete smashed sum $\overline{\text{VAL}} \triangleq \wp(\Sigma^*) \cup \wp(\mathbb{Z}) \cup \wp(\{\text{true}, \text{false}\})$, that is defined as a collecting semantics. We abuse notation denoting with $\mathbb{M} : \text{ID} \rightarrow \overline{\text{VAL}}$ the set of collecting memories, ranging over meta-variable m , that associate each identifier to a collecting value. The concrete expression semantics of such domain is defined as $\llbracket e \rrbracket : \mathbb{M} \rightarrow \overline{\text{VAL}}$, evaluating e and returning its possible values. Such semantics is defined as the additive lift of the one in Fig. 3. Function $\gamma_{\text{VAL}^\sharp} : \text{VAL}^\sharp \rightarrow \overline{\text{VAL}}$ is the smashed sum concretization and it is defined as:

$$\gamma_{\text{VAL}^\sharp}(a) \triangleq \begin{cases} \emptyset & \text{if } a = \perp \\ \gamma_{\text{Intv}}(a) & \text{if } a \in \text{Intv} \\ \gamma_{\text{Bool}}(a) & \text{if } a \in \text{Bool} \\ \gamma_{\mathcal{T}}(a) & \text{if } a \in \mathcal{TFA}_{/\equiv} \\ \overline{\text{VAL}} & \text{otherwise} \end{cases}$$

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 $\gamma : \mathbb{M}^\sharp \rightarrow \mathbb{M}$ as $\gamma(m^\sharp) \triangleq \{ (x, \gamma_{\text{VAL}^\sharp}(m^\sharp(x))) \mid x \in \text{dom}(m^\sharp) \}$. With this setup, we prove the abstract semantics to be sound by proving that $\forall m^\sharp \in \mathbb{M}^\sharp. \llbracket e \rrbracket \gamma(m^\sharp) \subseteq \gamma_{\text{VAL}^\sharp}(\llbracket e \rrbracket^\sharp(m^\sharp))$. 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 γ to avoid cluttering the notation, since 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 (that is, 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 Sect. 3.

A.1 Concat

Theorem 1. $\llbracket \text{concat}(\mathbf{s}, \mathbf{s}') \rrbracket^\sharp$ is a sound and complete abstraction of $\llbracket \text{concat}(\mathbf{s}, \mathbf{s}') \rrbracket$. Formally:

$$\forall m^\sharp \in \mathbb{M}^\sharp, \forall \mathbf{s}, \mathbf{s}' \in \text{SE}. \llbracket \text{concat}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(m^\sharp) = \gamma(\llbracket \text{concat}(\mathbf{s}, \mathbf{s}') \rrbracket^\sharp m^\sharp)$$

Proof. Soundness and completeness follow from the fact that finite state automata and regular languages are closed under finite concatenation²². □

A.2 Length

Theorem 2. $\llbracket \text{length}(\mathbf{s}) \rrbracket^\sharp$ is a sound but not complete abstraction of $\llbracket \text{length}(\mathbf{s}) \rrbracket$. Formally:

$$\forall m^\sharp \in \mathbb{M}^\sharp, \forall \mathbf{s} \in \text{SE}. \llbracket \text{length}(\mathbf{s}) \rrbracket \gamma(m^\sharp) \subsetneq \gamma(\llbracket \text{length}(\mathbf{s}) \rrbracket^\sharp m^\sharp)$$

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

$$\begin{aligned} & \llbracket \text{length}(\mathbf{s}) \rrbracket \gamma(m^\sharp) \\ &= I && \{\text{def. } \mathbb{E}\} \\ &\subseteq \gamma([\min I, +\infty]) && \{\text{def. min, } \gamma\} \\ &= \gamma([\text{lminPath}(A), +\infty]) && \{\text{def. minPath}\} \\ &= \gamma(\llbracket \text{length}(\mathbf{s}) \rrbracket^\sharp m^\sharp) && \{\text{def. } \mathbb{E}^\sharp, 1^{\text{st}} \text{ case}\} \end{aligned}$$

Otherwise, since \mathcal{L} is a finite language, soundness is proven as:

$$\begin{aligned} & \llbracket \text{length}(\mathbf{s}) \rrbracket \gamma(m^\sharp) \\ &= I && \{\text{def. } \mathbb{E}\} \\ &\subseteq \gamma([\min I, \max I]) && \{\text{def. max, } \gamma\} \\ &= \gamma([\text{lminPath}(A), \text{lmaxPath}(A)]) && \{\text{def. minPath, maxPath}\} \\ &= \gamma(\llbracket \text{length}(\mathbf{s}) \rrbracket^\sharp m^\sharp) && \{\text{def. } \mathbb{E}^\sharp, 2^{\text{nd}} \text{ case}\} \end{aligned}$$

As a counterexample for completeness, let us consider $\llbracket \mathbf{s} \rrbracket^\sharp m^\sharp = A \in \mathcal{TFA}_{/\equiv}$ s.t. $\gamma(A) = \{a, aaa\}$.

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

□

A.3 Contains

Theorem 3. $\llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket^\sharp$ is a sound but not complete abstraction of $\llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket$. Formally:

$$\forall m^\sharp \in \mathbb{M}^\sharp, \forall \mathbf{s}, \mathbf{s}' \in \text{SE}. \llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(m^\sharp) \subseteq \gamma(\llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket^\sharp m^\sharp)$$

Proof. The collecting semantics of `contains` is defined as the additive lift of the concrete one, that is $\llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket m = \{ b \mid b = \text{contains}(\sigma, \sigma'), \sigma \in \mathcal{L}, \sigma' \in \mathcal{L}' \}$, where $\llbracket \mathbf{s} \rrbracket m = \mathcal{L} \in \wp(\Sigma^*)$, $\llbracket \mathbf{s}' \rrbracket m = \mathcal{L}' \in \wp(\Sigma^*)$ and `contains` : $\Sigma^* \times \Sigma^* \rightarrow \{\text{true}, \text{false}\}$ corresponds to the concrete semantics of Fig. 3. Let $\llbracket \mathbf{s} \rrbracket^\sharp m^\sharp = A \in \mathcal{TFA}_{/\equiv}$ s.t. $\gamma(A) = \mathcal{L} \in \wp(\Sigma^*)$ and $\llbracket \mathbf{s}' \rrbracket^\sharp m^\sharp = A' \in \mathcal{TFA}_{/\equiv}$ s.t. $\gamma(A') = \mathcal{L}' \in \wp(\Sigma^*)$. We split the proof following possible values produced by the concrete semantics.

▷ $\llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(m^\sharp) = \{\text{false}\}$. Thus, no substring of the strings in \mathcal{L} is in \mathcal{L}' .

$$\begin{aligned} & \llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(m^\sharp) = \{\text{false}\} \xrightarrow{\text{def. } \mathbb{E}} \forall \sigma \in \mathcal{L} \forall \sigma' \in \mathcal{L}'. \sigma' \not\prec_s \sigma \xrightarrow{\text{def. FA}} \mathcal{L}(\text{FA}(A)) \cap \mathcal{L}' = \emptyset \\ & \xrightarrow{\text{automata lift}} \text{FA}(A) \sqcap_{\mathcal{T}} A' = \text{Min}(\emptyset) \xrightarrow{\text{def. } \mathbb{E}^\sharp, 1^{\text{st}} \text{ case}} \llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket^\sharp m^\sharp = \{\text{false}\} \end{aligned}$$

$\triangleright \llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(\mathfrak{m}^\sharp) = \{\text{true}\}$. Thus, all strings in \mathcal{L} contain all the strings of \mathcal{L}' . This invalidates the first case of our abstract semantics, as $\exists \sigma \in \mathcal{L}(\text{FA}(\mathbb{A})) . \sigma \in \mathcal{L}'$. If $\text{singlePath}(\mathbb{A}')$ holds, our semantics matches the concrete one:

$$\begin{aligned} \llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(\mathfrak{m}^\sharp) = \{\text{true}\} \wedge \text{singlePath}(\mathbb{A}') &\xleftrightarrow{\text{def. } \mathbb{E}, \text{singlePath}} \forall \sigma \in \mathcal{L} . \sigma_{\text{sp}} \curvearrow_{\mathbf{s}} \sigma \xleftrightarrow{\mathcal{L}'(\mathbb{A}^{ac}) \subseteq \mathcal{L}} \forall \sigma \in \mathcal{L}(\mathbb{A}^{ac}) . \sigma_{\text{sp}} \curvearrow_{\mathbf{s}} \sigma \\ &\xleftrightarrow{\text{automata lift}} \forall \pi \in \text{paths}(\mathbb{A}^{ac}) . \sigma_{\text{sp}} \curvearrow_{\mathbf{s}} \sigma_\pi \xleftrightarrow{\text{def. } \mathbb{E}^\sharp, 2^{\text{nd}} \text{ case}} \llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket^\sharp \mathfrak{m}^\sharp = \{\text{true}\} \end{aligned}$$

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

$\triangleright \llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(\mathfrak{m}^\sharp) = \{\text{true}, \text{false}\}$. In this case, soundness is met as none of the conditions appearing in the definition of $\llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket^\sharp$ are satisfied, and the latter returns $\{\text{true}, \text{false}\}$ as well (3rd case).

As a counterexample for the completeness of contains , let us consider $\llbracket \mathbf{s} \rrbracket^\sharp \mathfrak{m}^\sharp = \mathbb{A} \in \mathcal{TFA}_{/\equiv}$ s.t. $\gamma(\mathbb{A}) = \{ab, bba\}$ and $\llbracket \mathbf{s}' \rrbracket^\sharp \mathfrak{m}^\sharp = \mathbb{A}' \in \mathcal{TFA}_{/\equiv}$ s.t. $\gamma(\mathbb{A}') = \{a, b\}$.

$$\llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(\mathfrak{m}^\sharp) = \{\text{true}\} \subsetneq \{\text{true}, \text{false}\} = \gamma(\llbracket \text{contains}(\mathbf{s}, \mathbf{s}') \rrbracket^\sharp \mathfrak{m}^\sharp)$$

□

A.4 IndexOf

Theorem 4. $\llbracket \text{indexOf}(\mathbf{s}, \mathbf{s}') \rrbracket^\sharp$ is a sound but not complete abstraction of $\llbracket \text{indexOf}(\mathbf{s}, \mathbf{s}') \rrbracket$. Formally:

$$\forall \mathfrak{m}^\sharp \in \mathbb{M}^\sharp, \forall \mathbf{s}, \mathbf{s}' \in \text{SE} . \llbracket \text{indexOf}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(\mathfrak{m}^\sharp) \subsetneq \gamma(\llbracket \text{indexOf}(\mathbf{s}, \mathbf{s}') \rrbracket^\sharp \mathfrak{m}^\sharp)$$

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

$\triangleright \llbracket \text{indexOf}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(\mathfrak{m}^\sharp) = \{-1\}$. No string of \mathcal{L}' is contained in any string of \mathcal{L} . Formally:

$$\begin{aligned} \llbracket \text{indexOf}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(\mathfrak{m}^\sharp) = \{-1\} &\xleftrightarrow{\text{necessary condition}} \forall \sigma' \in \mathcal{L}' \nexists \sigma \in \mathcal{L} . \sigma' \curvearrow_{\mathbf{s}} \sigma \\ &\xleftrightarrow{\text{def. } \mathbb{E}^\sharp, 2^{\text{nd}} \text{ case}} \llbracket \text{indexOf}(\mathbf{s}, \mathbf{s}') \rrbracket^\sharp \mathfrak{m}^\sharp = [-1, -1] = \end{aligned}$$

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

$\triangleright \llbracket \text{indexOf}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(\mathfrak{m}^\sharp) = I \subseteq \{n \in \mathbb{Z} \mid n \geq -1\}$. $\exists i \in I . i \geq 0$. This implies $\exists \sigma \in \mathcal{L}(\mathbb{A}), \sigma' \in \mathcal{L}(\mathbb{A}') . \sigma' \curvearrow_{\mathbf{s}} \sigma$, 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 σ' 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 γ_{Intv} . We first prove $i \leq m$, identifying two cases. If $m = -1$:

$$-1 \in I_{\sigma'} \xleftrightarrow{\text{necessary condition}} \exists \sigma \in \mathcal{L} . \sigma' \not\curvearrow_{\mathbf{s}} \sigma \xleftrightarrow{\text{automata lift}} \exists \pi \in \text{paths}(\mathbb{A}) . \sigma' \not\curvearrow_{\mathbf{s}} \sigma_\pi \xleftrightarrow{\text{def. } i, 1^{\text{st}} \text{ case}} i = -1$$

Instead, if $m > -1$:

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

$$\xleftrightarrow{\text{def. } i, 2^{\text{nd}} \text{ case}} i = \min k = m$$

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

$$I_{\sigma'} = \{-1\} \xleftrightarrow{\text{necessary condition}} \forall \sigma \in \mathcal{L} . \sigma' \not\prec_{\mathcal{S}} \sigma \xleftrightarrow{\text{automata lift}} \forall \pi \in \text{paths}(\mathbb{A}) . \sigma' \not\prec_{\mathcal{S}} \sigma_{\pi} \xleftrightarrow{\text{def. } j, 1^{\text{st}} \text{ case}} j = 1$$

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

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

Finally, if $M > -1$ and $\exists \pi \in \text{paths}(\mathbb{A}) . \pi$ reads \top before σ , $j = +\infty$ by the 2nd 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 $\llbracket \mathbf{s} \rrbracket^{\#} \mathfrak{m}^{\#} = \mathbb{A} \in \mathcal{TFA}_{/\equiv}$ s.t. $\gamma(\mathbb{A}) = \{a, bba\}$ and $\llbracket \mathbf{s}' \rrbracket^{\#} \mathfrak{m}^{\#} = \mathbb{A}' \in \mathcal{TFA}_{/\equiv}$ s.t. $\gamma(\mathbb{A}') = \{a\}$.

$$\llbracket \text{indexOf}(\mathbf{s}, \mathbf{s}') \rrbracket \gamma(\mathfrak{m}^{\#}) = \{0, 2\} \subsetneq \gamma([0, 2]) = \gamma(\llbracket \text{indexOf}(\mathbf{s}, \mathbf{s}') \rrbracket^{\#} \mathfrak{m}^{\#})$$

A.5 Repeat

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

Theorem 5. *Given $\mathbb{A} \in \mathcal{TFA}_{/\equiv}$, $i \in \mathbb{N}$, `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 $\text{repeat}(\mathcal{L}, i) \triangleq \{ \sigma^i \mid \sigma \in \mathcal{L} \}$.

Proof. We split the proof in the following cases.

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

$$\text{repeat}(\mathcal{L}(\mathbb{A}), i) = \{ \sigma^0 \mid \sigma \in \mathcal{L}(\mathbb{A}) \} = \{\epsilon\} = \gamma(\text{Min}(\{\epsilon\})) \quad \{\text{lines 1-2 of Alg. 2}\}$$

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

$$\{ \sigma^i \mid \sigma \in \mathcal{L} \} \subseteq \mathcal{L}^i \tag{A1}$$

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

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 ϵ . These operations are the ones performed at lines 3–8 of Alg. 2.

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

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

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

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

Theorem 6. $\llbracket \text{repeat}(\mathbf{s}, \mathbf{a}) \rrbracket^\sharp$ is a sound but not complete abstraction of $\llbracket \text{repeat}(\mathbf{s}, \mathbf{a}) \rrbracket$. Formally:

$$\forall m^\sharp \in M^\sharp, \forall s \in SE, \mathbf{a} \in AE. \llbracket \text{repeat}(\mathbf{s}, \mathbf{a}) \rrbracket \gamma(m^\sharp) \subsetneq \gamma(\llbracket \text{repeat}(\mathbf{s}, \mathbf{a}) \rrbracket^\sharp m^\sharp)$$

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

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

$$\begin{aligned}
\llbracket \text{repeat}(\mathbf{s}, \mathbf{a}) \rrbracket \gamma(m^\sharp) &= \\
&= \{ \sigma^i \mid \sigma \in \mathcal{L}(A), i \in \mathbb{N} \} && \{\text{def. repeat}(\mathbf{s}, \mathbf{a})\} \\
&\subseteq \bigcup_{i \in \mathbb{N}} \mathcal{L}(A)^i && \{\text{Eq. A1}\} \\
&= \gamma(\text{Kleene}(A)) && \{\text{def. Kleene}\}
\end{aligned}$$

$\triangleright i = j \wedge i \in \mathbb{N}$. Soundness follows from Thm. 5.

$\triangleright j = +\infty$.

$$\begin{aligned}
\llbracket \text{repeat}(\mathbf{s}, \mathbf{a}) \rrbracket \gamma(m^\sharp) &= \\
&= \{ \sigma^k \mid \sigma \in \mathcal{L}(A), k \geq i \} && \{\text{def. repeat}(\mathbf{s}, \mathbf{a})\} \\
&= \{ \sigma^i \sigma^l \mid \sigma \in \mathcal{L}(A), l \in \mathbb{N} \} && \{\text{def. } \geq, \text{ string concat.}\} \\
&= \{ \sigma^i \mid \sigma \in \mathcal{L}(A) \} \cdot \{ \sigma^l \mid \sigma \in \mathcal{L}(A), l \in \mathbb{N} \} && \{\text{def. lang. concat.}\} \\
&\subseteq \gamma(\llbracket \text{concat}(\text{repeat}(A, i), \text{Kleene}(A)) \rrbracket^\sharp) && \{\text{def. Kleene, repeat, } \gamma\}
\end{aligned}$$

$\triangleright i, j \in \mathbb{N}$. Soundness follows from Thm. 5 and the definition of $\sqcup_{\mathcal{T}}$.

As a counterexample for completeness, let $\llbracket \mathbf{s} \rrbracket^\sharp m^\sharp = A \in \mathcal{TFA}_{/\equiv}$ s.t. $\gamma(A) = \{ a^n \mid n \in \mathbb{N} \} \cup \{ b \}$ and $\llbracket \mathbf{a} \rrbracket^\sharp m^\sharp = [2, 2]$.

$$\llbracket \text{repeat}(\mathbf{s}, \mathbf{a}) \rrbracket \gamma(m^\sharp) = \{ a^n \mid n \in \mathbb{N} \} \cup \{ bb \} \subsetneq \{ a^n b a^m b \mid n, m \in \mathbb{N} \} = \gamma(\llbracket \text{repeat}(\mathbf{s}, \mathbf{a}) \rrbracket^\sharp m^\sharp)$$

\square

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 7. Given $A \in \mathcal{TF}_{A/\exists}$, 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 since $\text{trimL}(\mathcal{L}(\emptyset)) = \emptyset = \gamma(\text{trimL}(\emptyset))$ (1st case).

▷ $r = \top$. Soundness holds since $\text{trimL}(\mathcal{L}(\top)) = \Sigma^* = \gamma(\text{trimL}(\top))$ (1st case).

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

Inductive cases

▷ $r = r_1 \parallel r_2$.

$$\begin{aligned} \text{trimL}(\gamma(r_1 \parallel r_2)) &= \\ &= \text{trimL}(\gamma(r_1) \cup \gamma(r_2)) && \text{\{def. } \gamma \text{\}} \\ &= \text{trimL}(\gamma(r_1)) \cup \text{trimL}(\gamma(r_2)) && \text{\{distrib. of trimL\}} \\ &= \gamma(\text{trimL}(r_1)) \cup \gamma(\text{trimL}(r_2)) && \text{\{ind. hp.\}} \\ &= \gamma(\text{trimL}(r_1) \cup \text{trimL}(r_2)) && \text{\{def. } \cup, \gamma \text{\}} \\ &= \gamma(\text{trimL}(r_1 \parallel r_2)) && \text{\{distrib. of trimL, 3rd case\}} \end{aligned}$$

▷ $r = r_1 r_2$. For regex concatenation, we split the proof in three sub-cases.

- $\gamma(r_1) \subseteq \{_ \}^* \implies \text{trimL}(\gamma(r_1)) = \epsilon$

$$\begin{aligned} \text{trimL}(\gamma(r_1 r_2)) &= \\ &= \left\{ \sigma' \mid \begin{array}{l} \sigma \in \gamma(r_1 r_2), \sigma = \psi\sigma', \\ \psi = \max\{ \psi' \in \{_ \}^* \mid \sigma = \psi'\sigma'' \} \end{array} \right\} && \text{\{def. trimL\}} \\ &= \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \gamma(r_1), \sigma_2 \in \gamma(r_2), \sigma_1\sigma_2 = \psi\sigma', \\ \psi = \max\{ \psi' \in \{_ \}^* \mid \sigma_1\sigma_2 = \psi'\sigma'' \} \end{array} \right\} && \text{\{def. } \gamma(r_1 r_2) \text{\}} \\ &= \left\{ \sigma' \mid \begin{array}{l} \sigma_2 \in \gamma(r_2), \sigma_2 = \psi\sigma', \\ \psi = \max\{ \psi' \in \{_ \}^* \mid \sigma_2 = \psi'\sigma'' \} \end{array} \right\} && \text{\{\text{trimL}(\gamma(r_1)) = \epsilon \implies \sigma_1 \in \psi\}} \\ &= \text{trimL}(\gamma(r_2)) && \text{\{def. trimL\}} \\ &\subseteq \gamma(\text{trimL}(r_2)) && \text{\{ind. hp.\}} \\ &= \gamma(\text{trimL}(r_1 r_2)) && \text{\{4th case\}} \end{aligned}$$

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

$$\begin{aligned} \text{trimL}(\gamma(r_1 r_2)) &= \\ &= \left\{ \sigma' \mid \begin{array}{l} \sigma \in \gamma(r_1 r_2), \sigma = \psi\sigma', \\ \psi = \max\{ \psi' \in \{_ \}^* \mid \sigma = \psi'\sigma'' \} \end{array} \right\} && \text{\{def. trimL\}} \\ &= \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \gamma(r_1), \sigma_2 \in \gamma(r_2), \sigma_1\sigma_2 = \psi\sigma', \\ \psi = \max\{ \psi' \in \{_ \}^* \mid \sigma_1\sigma_2 = \psi'\sigma'' \} \end{array} \right\} && \text{\{def. } \gamma(r_1 r_2) \text{\}} \\ &= \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \mathcal{L}^{\text{ws}}, \sigma_2 \in \gamma(r_2), \sigma_1\sigma_2 = \psi\sigma', \\ \psi = \max\{ \psi' \in \{_ \}^* \mid \sigma_1\sigma_2 = \psi'\sigma'' \} \end{array} \right\} \\ &\cup \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \mathcal{L}^{-\text{ws}}, \sigma_2 \in \gamma(r_2), \sigma_1\sigma_2 = \psi\sigma', \\ \psi = \max\{ \psi' \in \{_ \}^* \mid \sigma_1\sigma_2 = \psi'\sigma'' \} \end{array} \right\} && \text{\{def. } \mathcal{L}^{\text{ws}}, \mathcal{L}^{-\text{ws}} \text{\}} \end{aligned}$$

$$\begin{aligned}
&= \left\{ \sigma' \mid \begin{array}{l} \sigma_2 \in \gamma(r_2), \sigma_2 = \psi\sigma', \\ \psi = \max\{ \psi' \in \{_ \}^* \mid \sigma_2 = \psi'\sigma'' \} \end{array} \right\} \\
&\cup \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \mathcal{L}^{-ws}, \sigma_2 \in \gamma(r_2), \sigma_1\sigma_2 = \psi\sigma', \\ \psi = \max\{ \psi' \in \{_ \}^* \mid \sigma_1\sigma_2 = \psi'\sigma'' \} \end{array} \right\} && \{\sigma_1 \in \{_ \}^*\} \\
&= \left\{ \sigma' \mid \begin{array}{l} \sigma_2 \in \gamma(r_2), \sigma_2 = \psi\sigma', \\ \psi = \max\{ \psi' \in \{_ \}^* \mid \sigma_2 = \psi'\sigma'' \} \end{array} \right\} \\
&\cup \left\{ \sigma' \mid \begin{array}{l} \sigma_1 \in \mathcal{L}^{-ws}, \sigma_2 \in \gamma(r_2), \sigma_1\sigma_2 = \psi\sigma', \\ \psi = \max\{ \psi' \in \{_ \}^* \mid \sigma_1\sigma_2 = \psi'\sigma'' \} \end{array} \right\} && \{\sigma_1 \notin \{_ \}^*, \psi \in \sigma_1\} \\
&= \text{trimL}(\gamma(r_2)) \cup \text{trimL}(\gamma(r_1))\gamma(r_2) && \{\text{def. trimL}\} \\
&\subseteq \gamma(\text{trimL}(r_2)) \cup \gamma(\text{trimL}(r_1)r_2) && \{\text{ind. hp., def. } \gamma\} \\
&= \gamma(\text{trimL}(r_2) \parallel \text{trimL}(r_1)r_2) && \{\text{def. } \gamma(r_1 \parallel r_2)\} \\
&= \gamma(\text{trimL}(r_1r_2)) && \{5^{\text{th}} \text{ case}\}
\end{aligned}$$

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

$\triangleright r = (r_1)^*$. 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 $r = r_1r_2$. \square

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

$$\llbracket \text{trimLeft}(\mathbf{s}) \rrbracket \gamma(m^{\#}) = \{ \sigma ab \mid \sigma \in \Sigma^* \setminus \{_ \}^* \} \subsetneq \gamma(\text{Min}(\{\text{Tab}\})) = \gamma(\llbracket \text{trimLeft}(\mathbf{s}) \rrbracket^{\#} m^{\#})$$

Consequently, also the abstract semantics of `trimRight` and `trim` are not complete.

A.7 Replace

Theorem 8. $\llbracket \text{replace}(\mathbf{s}, \mathbf{s}_s, \mathbf{s}_r) \rrbracket^{\#}$ is a sound but not complete abstraction of $\llbracket \text{replace}(\mathbf{s}, \mathbf{s}_s, \mathbf{s}_r) \rrbracket$. Formally:

$$\forall m^{\#} \in M^{\#}, \forall \mathbf{s}, \mathbf{s}_s, \mathbf{s}_r \in \text{SE}. \llbracket \text{replace}(\mathbf{s}, \mathbf{s}_s, \mathbf{s}_r) \rrbracket \gamma(m^{\#}) \subsetneq \gamma(\llbracket \text{replace}(\mathbf{s}, \mathbf{s}_s, \mathbf{s}_r) \rrbracket^{\#} m^{\#})$$

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

$$\llbracket \text{replace}(\mathcal{L}, \mathcal{L}_s, \mathcal{L}_r) \rrbracket = \mathcal{L} \xleftrightarrow{\text{necessary condition}} \forall \sigma_s \in \mathcal{L}_s \nexists \sigma \in \mathcal{L}. \sigma_s \not\sim_{\mathbf{s}} \sigma \xleftrightarrow{\text{def. } \mathbb{E}^{\#}, 1^{\text{st}} \text{ case}} \llbracket \text{replace}(A, A_s, A_r) \rrbracket = A$$

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 , that can be thoroughly explored since \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 $\forall \sigma \in \mathcal{L}, \forall \sigma_s \in \mathcal{L}_s, \forall \pi \in \text{paths}(A). \sigma_{\pi} = \sigma$:

$$\llbracket \text{replace}(\{\sigma\}, \{\sigma_s\}, \mathcal{L}_r) \rrbracket \subseteq \gamma(\text{RP}(\pi, \sigma_s, A_r))$$

Specifically, `RP` removes every occurrence of σ_s in π (lines 7–8 of Alg. 3, where states and transitions composing σ_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 ϵ -transitions. This means that all $\sigma' \sim_{\mathbf{s}} \sigma_{\pi}. \sigma' = \sigma_s$ are substituted with *all* the strings recognized by A_r . We can then characterize the language of the automaton returned by `RP` as $\{ \sigma_{\pi}[\sigma_s/\sigma_r] \mid \sigma_r \in \mathcal{L}(A_r) \}$. Soundness is thus ensured:

$$\begin{aligned}
\llbracket \text{replace}(\{\sigma\}, \{\sigma_s\}, \mathcal{L}_r) \rrbracket &= \{ \sigma[\sigma_s/\sigma_r] \mid \sigma_r \in \mathcal{L}_r \} \\
&\subseteq \{ \sigma_{\pi}[\sigma_s/\sigma_r] \mid \sigma_r \in \mathcal{L}(A_r) \} && \{\text{automata lift}\}
\end{aligned}$$

$$= \gamma(\text{RP}(\pi, \sigma_s, A_r)) \quad \{\text{def. RP}\}$$

Soundness is thus proven as the result on individual strings can be lifted to languages, and since the 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 $\text{Min}(\{\text{T}\})$ when either the input automaton or the search automaton contain cycles or read T.

As a counterexample for completeness, let $\llbracket \mathbf{s} \rrbracket^{\#m\#} = A \in \mathcal{TFA}_{/\equiv}$, $\llbracket \mathbf{s}^s \rrbracket^{\#m\#} = A^s \in \mathcal{TFA}_{/\equiv}$, $\llbracket \mathbf{s}^r \rrbracket^{\#m\#} = A^r \in \mathcal{TFA}_{/\equiv}$ s.t. $\gamma(A) = \{abc\}$, $\gamma(A^s) = \{a, z\}$ and $\gamma(A^r) = \{r\}$.

$$\llbracket \text{replace}(\mathbf{s}, \mathbf{s}^s, \mathbf{s}^r) \rrbracket^{\#m\#} = \{rbc\} \subsetneq \{abc, rbc\} = \gamma(\llbracket \text{replace}(\mathbf{s}, \mathbf{s}^s, \mathbf{s}^r) \rrbracket^{\#m\#})$$

□

A.8 Substring and CharAt

Theorem 9. $\llbracket \text{substr}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_2) \rrbracket^{\#}$ is a sound and complete abstraction of $\llbracket \text{substr}(\mathbf{s}, \mathbf{a}_1, \mathbf{a}_2) \rrbracket$. Formally:

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

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

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

Completeness is proven by structural induction over the structure of the regular expression, referencing the lines of Alg. 4 that are involved in the computation as $\S x$, 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 (σ_i, i_i, j_i) is not in the final result of Sb .

Base cases

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

$$\begin{aligned} \llbracket \text{substr}(\emptyset, \{i\}, \{j\}) \rrbracket^{\#} &= \emptyset \\ &= \gamma(\text{Min}(\emptyset)) && \{\text{automata lift}\} \\ &= \gamma(\text{Min}(\text{Sb}(\emptyset, i, j-i))) && \{\S 2\} \end{aligned}$$

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

$$\begin{aligned} \llbracket \text{substr}(\{\sigma\}, \{i\}, \{j\}) \rrbracket^{\#} &= \{\sigma_i \dots \sigma_j\} \\ &= \gamma(\text{Min}(\{\sigma_i \dots \sigma_j\})) && \{\text{automata lift}\} \\ &= \gamma(\text{Min}(\text{Sb}(\{\sigma\}, i, j-i))) && \{\S 6\} \end{aligned}$$

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

$$\begin{aligned} \llbracket \text{substr}(\{\sigma\}, \{i\}, \{j\}) \rrbracket^{\#} &= \emptyset \\ &= \gamma(\text{Min}(\{\epsilon, i-|\sigma|, j-i\})) && \{i-|\sigma| > 0 \implies \epsilon \notin \text{Sb}\} \\ &= \gamma(\text{Min}(\text{Sb}(\{\sigma\}, i, j-i))) && \{\S 4\} \end{aligned}$$

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

$$\begin{aligned} \llbracket \text{substr}(\{\sigma\}, \{i\}, \{j\}) \rrbracket &= \emptyset \\ &= \gamma(\text{Min}(\{(\sigma_i \dots \sigma_{|\sigma|}, 0, k)\})) && \{k > 0 \implies \sigma_i \dots \sigma_{|\sigma|} \notin \text{Sb}\} \\ &= \gamma(\text{Min}(\text{Sb}(\{\sigma\}, i, j - i))) && \{ \S 5 \} \end{aligned}$$

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

$\triangleright r = T$:

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

where, for the sake of clarity, strings returned by **Sb** are split into three sets, the first $(\{(\bullet^{i-i}, 0, 0)\})$ 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

$\triangleright r = r_1 \parallel r_2$: 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 $\llbracket \text{substr}(\mathcal{L}, \{i\}, \{j\}) \rrbracket = \llbracket \text{substr}(\mathcal{L}_1, \{i\}, \{j\}) \rrbracket \cup \llbracket \text{substr}(\mathcal{L}_2, \{i\}, \{j\}) \rrbracket$. We assume $\llbracket \text{substr}(\mathcal{L}_1, \{i\}, \{j\}) \rrbracket = \gamma(\text{Min}(\text{Sb}(r_1, i, j - i)))$ and $\llbracket \text{substr}(\mathcal{L}_2, \{i\}, \{j\}) \rrbracket = \gamma(\text{Min}(\text{Sb}(r_2, i, j - i)))$ to hold for inductive hypothesis. We then prove soundness with the following:

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

$\triangleright r = r_1 r_2$: 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: $\llbracket \text{substr}(\mathcal{L}_1, \{i\}, \{j\}) \rrbracket$ (i.e., substrings that are fully contained in \mathcal{L}_1), and $\llbracket \text{substr}(\mathcal{L}_1 \cdot \mathcal{L}_2, \{i\}, \{j\}) \rrbracket$ (i.e., substrings that straddle \mathcal{L}_1 and \mathcal{L}_2). We prove soundness assuming the inductive hypotheses $\llbracket \text{substr}(\mathcal{L}_1, \{i\}, \{j\}) \rrbracket = \gamma(\text{Min}(\text{Sb}(r_1, i, j - i)))$ and $\llbracket \text{substr}(\mathcal{L}_2, \{i\}, \{j\}) \rrbracket = \gamma(\text{Min}(\text{Sb}(r_2, i, j - i)))$:

$$\begin{aligned} \llbracket \text{substr}(\mathcal{L}, \{i\}, \{j\}) \rrbracket &= \llbracket \text{substr}(\mathcal{L}_1, \{i\}, \{j\}) \rrbracket \cup \llbracket \text{substr}(\mathcal{L}_1 \cdot \mathcal{L}_2, \{i\}, \{j\}) \rrbracket \\ &= \gamma(\text{Min}(\text{Sb}(r_1, i, j - i))) \cup \gamma(\text{Min}(\{(\sigma_1^1 \cdot \sigma_2^1, i_1^1, j_1^1), \dots, (\sigma_1^n \cdot \sigma_2^n, i_1^n, j_1^n)\})) && \{ \text{ind. hp.} \} \\ &= \gamma(\text{Min}(\text{Sb}(r_1 r_2, i, j - i))) && \{ \S 1, \S 16, \S 18 \} \end{aligned}$$

where, for the sake of clarity, strings returned by **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 $(\{(\sigma_1^1 \cdot \sigma_2^1, i_1^1, j_1^1), \dots, (\sigma_1^n \cdot \sigma_2^n, i_1^n, j_1^n)\})$ that models substrings straddling r_1 and r_2 , where $\forall i. (\sigma_1^i, i_1^i, j_1^i) \in \text{Sb}(r_1, i, j - i), j_1^i \neq 0 \wedge (\sigma_2^i, i_2^i, j_2^i) \in \text{Sb}(r_2, i_1^i, j_1^i)$. Strings in the latter set are built by offsetting substrings of r_2 by the length of the substrings of r_1 .

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

Theorem 10. $\llbracket \text{charAt}(\mathbf{s}, \mathbf{a}) \rrbracket^\#$ is a sound and complete abstraction of $\llbracket \text{charAt}(\mathbf{s}, \mathbf{a}) \rrbracket$. Formally:

$$\forall m^\# \in \mathbb{M}^\#, \forall \mathbf{s} \in \text{SE}, \forall \mathbf{a} \in \text{AE}. \llbracket \text{charAt}(\mathbf{s}, \mathbf{a}) \rrbracket \gamma(m^\#) = \gamma(\llbracket \text{charAt}(\mathbf{s}, \mathbf{a}) \rrbracket^\# m^\#)$$

Proof. Since the abstract semantics of `charAt` relies on the one of `substr`, soundness and completeness come from Thm 9. \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 `eq` function (whose algorithm is reported in Alg. 1).

Theorem 11. *Given $\sigma_1, \sigma_2 \in \Sigma^* \cup \{\top\}$, $\llbracket \text{Flat}(\sigma_1) == \text{Flat}(\sigma_2) \rrbracket$ is a sound and complete approximation of `eq`(σ_1, σ_2). Formally:*

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

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 σ_1 and σ_2 .

Base cases

$$\triangleright |\sigma_1| = 0 \iff \sigma_1 = \epsilon$$

- $|\sigma_2| = 0 \iff \sigma_2 = \epsilon$. In this case, $\text{Flat}(\sigma_1) == \text{Flat}(\sigma_2) = \{\text{true}\}$ that is equal to the result returned by `eq` in Alg. 1 when both strings are empty (lines 1-2).
- $|\sigma_2| = 1 \iff \sigma_2 = c \in \Sigma \cup \{\top\}$. We can split the proof in two cases:
 - $c = \top$: the empty string may be equal to \top , thus $\text{Flat}(\sigma_1) == \text{Flat}(\sigma_2) = \{\text{true}, \text{false}\}$, that is equal to the result returned by `eq` in Alg. 1 when one of the string is empty and the other is equal to \top (lines 3-4).
 - $c \neq \top$: 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 Alg. 1 when one of string is empty and the other has a single character not equal to \top (lines 5-6).

$$\triangleright |\sigma_1| = 1 \iff \sigma_1 = c \in \Sigma \cup \{\top\}. \text{ We can split the proof in two case}$$

- $c = \top$.
 - $|\sigma_2| = 0 \iff \sigma_2 = \epsilon$. This case is analogous to the second point of the previous base case.
 - $|\sigma_2| = 1 \iff \sigma_2 = c' \in \Sigma \cup \{\top\}$. We have two cases
 - * $c' = \top$. Two strings just made of a singleton \top 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 Alg. 1 when both strings have a single character equal to \top (lines 7-8).
 - * $c' \neq \top$. A string just made of a singleton \top may be equal to $c \in \Sigma \cup \{\top\}$, 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 7-8) Alg. 1
- $c \neq \top$.
 - $|\sigma_2| = 0 \iff \sigma_2 = \epsilon$. This case is analogous to the second case of the first base case.
 - $|\sigma_2| = 1 \iff \sigma_2 = c' \in \Sigma \cup \{\top\}$. This case is analogous to the previous case ($c = \top$).

Inductive steps Let $n \in \mathbb{N}$ and let $\sigma_1, \sigma_2 \in (\Sigma \cup \{\top\})^*$ 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 $\rho_1, \rho_2 \in (\Sigma \cup \{\top\})^*$ 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.

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

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

We split the proof in the following cases.

- $c = c'$

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

- $c \neq c'$

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

$\triangleright c = T, c' \neq T$

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

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

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

Theorem 12. $\llbracket \mathbf{s} == \mathbf{s}' \rrbracket^\#$ is a sound and complete abstraction of $\llbracket \mathbf{s} == \mathbf{s}' \rrbracket$. Formally:

$$\forall m^\# \in \mathbb{M}^\#, \forall \mathbf{s}, \mathbf{s}' \in \text{SE}. \llbracket \mathbf{s} == \mathbf{s}' \rrbracket \gamma(m^\#) \subseteq \gamma(\llbracket \mathbf{s} == \mathbf{s}' \rrbracket^\# m^\#)$$

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 Thm. 11. Let us focus on the second case, let $\mathbf{s}, \mathbf{s}' \in \text{SE}$ and suppose $\llbracket \mathbf{s} \rrbracket^\# m^\# = A \in \mathcal{TFA}_{\neq}, \llbracket \mathbf{s}' \rrbracket^\# m^\# = A' \in \mathcal{TFA}_{\equiv}$. We prove that if either A or A' are cyclic, $\llbracket \mathbf{s} == \mathbf{s}' \rrbracket \gamma(m^\#) = \{\text{true}, \text{false}\}$, the same result returned by $\gamma(\llbracket \mathbf{s} == \mathbf{s}' \rrbracket^\# m^\#)$, proving completeness. Note that $\llbracket \mathbf{s} == \mathbf{s}' \rrbracket \gamma(m^\#)$ cannot be $\{\text{false}\}$ since this case is treated in the first case of the abstract semantics of string equality. By contradiction, let either A or A' be cyclic and let us suppose that $\llbracket \mathbf{s} == \mathbf{s}' \rrbracket \gamma(m^\#) = \{\text{true}\}$.

$$\llbracket \mathbf{s} == \mathbf{s}' \rrbracket \gamma(m^\#) = \{\text{true}\}$$

$$\begin{aligned} &\iff \forall \sigma \in \gamma(A) \forall \sigma' \in \gamma(A') . \sigma == \sigma' && \{ \text{def. } \llbracket \mathbf{s} == \mathbf{s}' \rrbracket \} \\ &\implies |\gamma(A)| = |\gamma(A')| = 1 && \{ \text{set theory, def. } == \} \end{aligned}$$

We supposed that either A or A' were cyclic, reaching a contradiction. Thus, if either A or A' are cyclic, $\llbracket \mathbf{s} == \mathbf{s}' \rrbracket_{\gamma(\mathfrak{m}^\#)} = \{\text{false}, \text{true}\} = \gamma(\llbracket \mathbf{s} == \mathbf{s}' \rrbracket^\#_{\mathfrak{m}^\#})$, proving completeness. \square