

A Decision Procedure for Path Feasibility of String Manipulating Programs with Integer Data Type

Taolue Chen¹, Matthew Hague², Jinlong He^{3,6}, Denghang Hu^{3,6},
Anthony Widjaja Lin⁴, Philipp Rümmer⁵, and Zhilin Wu^{3,7,8}

¹ University of Surrey, UK

² Royal Holloway, University of London, UK

³ State Key Laboratory of Computer Science,
Institute of Software, Chinese Academy of Sciences, China

⁴ Technical University of Kaiserslautern, Germany

⁵ Uppsala University, Sweden

⁶ University of Chinese Academy of Sciences, China

⁷ Shanghai Key Laboratory of Trustworthy Computing, East China Normal University, China

⁸ Institute of Intelligent Software, Guangzhou, China

Abstract. In this paper, we propose a decision procedure for a class of string-manipulating programs which includes not only a wide range of string operations such as concatenation, `replaceAll`, `reverse`, and finite transducers, but also those involving the integer data-type such as `length`, `indexOf`, and `substring`. To the best of our knowledge, this represents one of the most expressive string constraint languages that is currently known to be decidable. Our decision procedure is based on a variant of cost register automata. We implement the decision procedure, giving rise to a new solver OSTRICH+. We evaluate the performance of OSTRICH+ on a wide range of existing and new benchmarks. The experimental results show that OSTRICH+ is the first string decision procedure capable of tackling finite transducers and integer constraints, whilst its overall performance is comparable with the state-of-the-art string constraint solvers.

1 Introduction

String-manipulating programs are notoriously subtle, and their potential bugs may bring severe security consequences. A typical example is cross-site scripting (XSS), which is among the OWASP Top 10 Application Security Risks [29]. Integer data type occurs naturally and extensively in string-manipulating programs. An effective and increasingly popular method for identifying bugs, including XSS, is symbolic execution [11]. In a nutshell, this technique analyses static paths through the program being considered. Each of these paths can be viewed as a constraint φ over appropriate data domains, and symbolic execution tools demand fast constraint solvers to check the satisfiability of φ . Such constraint solvers need to support all data-type operations occurring in a program.

Typically, mainstream programming languages provide standard string functions such as concatenation, `replace`, and `replaceAll`. Moreover, Web programming languages usually provide complex string operations (e.g. `htmlEscape` and `trim`), which are conveniently modelled as finite transducers, to sanitise malicious user inputs [19]. Nevertheless, apart from these operations involving only the string data type, functions

such as `length`, `indexOf`, and `substring`, which can convert strings to integers and vice versa, are also heavily used in practice; for instance, it was reported [26] that `length`, `indexOf`, `substring`, and variants thereof, comprise over 80% of string function occurrences in 18 popular JavaScript applications, notably outnumbering concatenation. The introduction of integers exacerbates the intricacy of string-manipulating programs, and poses new theoretical and practical challenges in solver development.

When combining strings and integers, decidability can easily be lost; for instance, the string theory with concatenation and letter counting functions is undecidable [8, 15]. Remarkably, it is still a major open problem whether the string theory with concatenation (arguably the simplest string operation) and length function (arguably the most common string-number function) is decidable [17, 22]. One promising approach to retain decidability is to enforce a syntactic restriction to the constraints. In the literature, these restriction include solved forms [17], acyclicity [5, 2, 3], and straight-line fragment (aka programs in single static assignment form) [21, 12, 14, 18]. On the one hand, such a restriction has led to decidability of string constraint solving with complex string operations (not only concatenation, but also finite transducers) and integer operations (letter-counting, `length`, `indexOf`, etc.); see, e.g., [21]. On the other hand, there is a lot of evidence (e.g. from benchmark) that many practical string constraints do satisfy such syntactic restrictions.

Approaches to building practical string solvers could essentially be classified into two categories. Firstly, one could support as many constraints as possible, but primarily resort to heuristics, offering no completeness/termination guarantee. This is a realistic approach since, as mentioned above, the problem involving both string and integer data types is in general undecidable. Many solvers belong to this category, e.g., CVC4 [20], Z3 [7, 16], Z3-str3 [6], S3(P) [27, 28], Trau [1] (or its variants Trau+ [3] and Z3-Trau [9]), ABC [10], and Slent [32]. Completeness guarantees are, however, valuable since the performance of heuristics can be difficult to predict. The second approach is to develop solvers for decidable fragments supporting both strings and integers (e.g. [17, 5, 2, 3, 21, 12, 14, 18]). Solvers in this category include Norn [2], SLOTH [18], and OSTRICH [14]. The fragment *without* complex string operations (e.g. `replaceAll` and finite transducers, but `length`) can be handled quite well by Norn. The fragment *without* length constraints (but `replaceAll` and finite transducers) can be handled effectively by OSTRICH and SLOTH. Moreover, most existing solvers that belong to the first category do not support complex string operations like `replaceAll` and finite transducers as well. This motivates the following problem: *provide a decision procedure that supports both string and integer data type, with completeness guarantee and meanwhile admitting efficient implementation.*

We argue that this problem is highly challenging. A deeper examination of the algorithms used by OSTRICH and SLOTH reveals that, unlike the case for Norn, it would *not* be straightforward to extend OSTRICH and SLOTH with integer constraints. First and foremost, the complexity of the fragment used by Norn (i.e. without transducers and `replaceAll`) is solvable in exponential time, even in the presence of integer constraints. This is not the case for the straight-line fragments with transducers/`replaceAll`, which require at least double exponential time (regardless of the integer constraints). This unfortunately manifests itself in the size of symbolic representations of the solutions.

SLOTH [18] computes a representation of all solutions “eagerly” as (alternating) finite transducers. Dealing with integer data type requires to compute the Parikh images of these transducers [21], which would result in a quantifier-free linear integer arithmetic formula (LIA for short) of double exponential size, thus giving us a triple exponential time algorithm, since LIA formulas are solved in exponential time (see e.g. [30]). Lin and Barcelo [21] provided a double exponential upper bound in the length of the strings in the solution, and showed that the double exponential time theoretical complexity could be retained. This, however, does not result in a practical algorithm since it requires all strings of double exponential size to be enumerated. OSTRICH [14] adopted a “lazy” approach and computed the pre-images of regular languages step by step, which is more scalable than the “eager” approach adopted by SLOTH and results in a highly competitive solver. It uses *recognisable relations* (a finite union of products of regular languages) as symbolic representations. Nevertheless, extending this approach to integer constraints is not obvious since integer constraints break the independence between different string variables in the recognisable relations.

Contribution. We provide a decision procedure for an expressive class of string constraints involving the integer data type, which includes not only concatenation, replace/replaceAll, reverse, finite transducers, and regular constraints, but also length, indexOf and substring. The decision procedure utilizes a variant of cost-register automata introduced by Alur et al. [4], which are called *cost-enriched finite automata* (CEFA) for convenience. Intuitively, each CEFA records the connection between a string variable and its associated integer variables. With CEFAs, the concept of recognisable relations is then naturally extended to accommodate integers. The integer constraints, however, are detached from CEFAs rather than being part of CEFAs. This allows to preserve the independence of string variables in the recognisable relation. The crux of the decision procedure is to compute the backward images of CEFAs under string functions, where each cost register (integer variable) might be split into several ones, thus extending but still in the same flavour as OSTRICH for string constraints *without* the integer data type [14]. Such an approach is able to treat a wide range of string functions in a generic, and yet simple, way. To the best of our knowledge, the class of string constraints considered in this paper is currently one of the most expressive string theories involving the integer data type known to enjoy a decision procedure.

We implement the decision procedure based on the recent OSTRICH solver [14], resulting in OSTRICH+. We perform experiments on a wide range of benchmark suites, including those where both replace/replaceAll/finite transducers and length/indexOf/substring occur, as well as the well-known benchmarks KALUZA and PyEx. The results show that 1) OSTRICH+ so far is the only string constraint solver capable of dealing with finite transducers and integer constraints, and 2) its overall performance is comparable with the best state-of-the-art string constraint solvers (e.g. CVC4 and Z3-Trau) which are short of completeness guarantees.

The rest of the paper is structured as follows: Section 2 introduces the preliminaries. Section 3 defines the class of string-manipulating programs with integer data type. Section 4 presents the decision procedure. Section 5 presents the benchmarks and experiments for the evaluation. The paper is concluded in Section 6. Missing proofs, implementation details and further examples can be found in the full version [13].

2 Preliminaries

We write \mathbb{N} and \mathbb{Z} for the sets of natural and integer numbers, respectively. For $n \in \mathbb{N}$ with $n \geq 1$, $[n]$ denotes $\{1, \dots, n\}$; for $m, n \in \mathbb{N}$ with $m \leq n$, $[m, n]$ denotes $\{i \in \mathbb{N} \mid m \leq i \leq n\}$. Throughout the paper, Σ is a finite alphabet, ranged over by a, b, \dots .

Strings, languages, and transductions. A string over Σ is a (possibly empty) sequence of elements from Σ , denoted by u, v, w, \dots . An empty string is denoted by ε . We write Σ^* (resp., Σ^+) for the set of all (resp. nonempty) strings over Σ . For a string u , we use $|u|$ to denote the number of letters in u . In particular, $|\varepsilon| = 0$. Moreover, for $a \in \Sigma$, let $|u|_a$ denote the number of occurrences of a in u . Assume $u = a_0 \cdots a_{n-1}$ is nonempty and $i < j \in [0, n-1]$. We let $u[i]$ denote a_i and $u[i, j]$ for the substring $a_i \cdots a_j$.

Let u, v be two strings. We use $u \cdot v$ to denote the *concatenation* of u and v . The string u is said to be a *prefix* of v if $v = u \cdot v'$ for some string v' . In addition, if $u \neq v$, then u is said to be a *strict prefix* of v . If $v = u \cdot v'$ for some string v' , then we use $u^{-1}v$ to denote v' . In particular, $\varepsilon^{-1}v = v$. If $u = a_0 \cdots a_{n-1}$ is nonempty, then we use $u^{(r)}$ to denote the *reverse* of u , that is, $u^{(r)} = a_{n-1} \cdots a_0$.

A *transduction* over Σ is a binary relation over Σ^* , namely, a subset of $\Sigma^* \times \Sigma^*$. We will use T_1, T_2, \dots to denote transductions. For two transductions T_1 and T_2 , we will use $T_1 \cdot T_2$ to denote the *composition* of T_1 and T_2 , namely, $T_1 \cdot T_2 = \{(u, w) \in \Sigma^* \times \Sigma^* \mid \text{there exists } v \in \Sigma^* \text{ s.t. } (u, v) \in T_1 \text{ and } (v, w) \in T_2\}$.

Recognisable relations. We assume familiarity with standard regular language. Recall that a regular language L can be represented by a regular expression $e \in \text{RegExp}$ whereby we usually write $L = \mathcal{L}(e)$.

Intuitively, a *recognisable relation* is simply a finite union of Cartesian products of regular languages. Formally, an r -ary relation $R \subseteq \Sigma^* \times \cdots \times \Sigma^*$ is *recognisable* if $R = \bigcup_{i=1}^n L_1^{(i)} \times \cdots \times L_r^{(i)}$ where $L_j^{(i)}$ is regular for each $j \in [r]$. A *representation* of a recognisable relation $R = \bigcup_{i=1}^n L_1^{(i)} \times \cdots \times L_r^{(i)}$ is $(\mathcal{A}_1^{(i)}, \dots, \mathcal{A}_r^{(i)})_{1 \leq i \leq n}$ such that each $\mathcal{A}_j^{(i)}$ is an NFA with $\mathcal{L}(\mathcal{A}_j^{(i)}) = L_j^{(i)}$. The tuples $(\mathcal{A}_1^{(i)}, \dots, \mathcal{A}_r^{(i)})$ are called the *disjuncts* of the representation and the NFAs $\mathcal{A}_j^{(i)}$ are called the *atoms* of the representation.

Automata models. A (*nondeterministic*) *finite automaton* (NFA) is a tuple $\mathcal{A} = (Q, \Sigma, \delta, I, F)$, where Q is a finite set of states, Σ is a finite alphabet, $\delta \subseteq Q \times \Sigma \times Q$ is the transition relation, $I, F \subseteq Q$ are the set of initial and final states respectively. For readability, we write a transition $(q, a, q') \in \delta$ as $q \xrightarrow[\delta]{a} q'$ (or simply $q \xrightarrow{a} q'$). The *size* of an NFA \mathcal{A} , denoted by $|\mathcal{A}|$, is defined as the number of transitions of \mathcal{A} . A *run* of \mathcal{A} on a string $w = a_1 \cdots a_n$ is a sequence of transitions $q_0 \xrightarrow{a_1} q_1 \cdots q_{n-1} \xrightarrow{a_n} q_n$ with $q_0 \in I$. The run is *accepting* if $q_n \in F$. A string w is accepted by an NFA \mathcal{A} if there is an accepting run of \mathcal{A} on w . In particular, the empty string ε is accepted by \mathcal{A} if $I \cap F \neq \emptyset$. The language of \mathcal{A} , denoted by $\mathcal{L}(\mathcal{A})$, is the set of strings accepted by \mathcal{A} . An NFA \mathcal{A} is said to be *deterministic* if I is a singleton and, for every $q \in Q$ and $a \in \Sigma$, there is at most one state $q' \in Q$ such that $(q, a, q') \in \delta$. It is well-known that finite automata capture regular languages precisely.

A *nondeterministic finite transducer* (NFT) \mathcal{T} is an extension of NFA with outputs. Formally, an NFT \mathcal{T} is a tuple $(Q, \Sigma, \delta, I, F)$, where Q, Σ, I, F are as in NFA and the transition relation δ is a finite subset of $Q \times \Sigma \times Q \times \Sigma^*$. Similarly to NFA, for readability,

we write a transition $(q, a, q', u) \in \delta$ as $q \xrightarrow[\delta]{a,u} q'$ or $q \xrightarrow{a,u} q'$. The *size* of an NFT \mathcal{T} , denoted by $|\mathcal{T}|$, is defined as the sum of the sizes of the transitions of \mathcal{T} , where the size of a transition $q \xrightarrow{a,u} q'$ is defined as $|u| + 3$. A run of \mathcal{T} over a string $w = a_1 \cdots a_n$ is a sequence of transitions $q_0 \xrightarrow{a_1, u_1} q_1 \cdots q_{n-1} \xrightarrow{a_n, u_n} q_n$ with $q_0 \in I$. The run is accepting if $q_n \in F$. The string $u_1 \cdots u_n$ is called the output of the run. The transduction defined by \mathcal{T} , denoted by $\mathcal{T}(\mathcal{T})$, is the set of string pairs (w, u) such that there is an accepting run of \mathcal{T} on w , with the output u . An NFT \mathcal{T} is said to be *deterministic* if I is a singleton, and, for every $q \in Q$ and $a \in \Sigma$ there is at most one pair $(q', u) \in Q \times \Sigma^*$ such that $(q, a, q', u) \in \delta$. In this paper, we are primarily interested in *functional* finite transducers (FFT), i.e., finite transducers that define functions instead of relations. (For instance, deterministic finite transducers are always functional.)

We will also use standard quantifier-free/existential *linear integer arithmetic* (LIA) formulae, which are typically ranged over by ϕ, φ , etc.

3 String-Manipulating Programs with Integer Data Type

In this paper, we consider logics involving two data-types, i.e., the string data-type and the integer data-type. As a convention, u, v, \dots denote string constants, c, d, \dots denote integer constants, x, y, \dots denote string variables, and i, j, \dots denote integer variables.

We consider symbolic execution of string-manipulating programs with numeric conditions (abbreviated as SL_{int}), defined by the following rules,

$$\begin{aligned} S ::= & x := y \cdot z \mid x := \text{replaceAll}_{e,u}(y) \mid x := \text{reverse}(y) \mid x := \mathcal{T}(y) \mid \\ & x := \text{substring}(y, t_1, t_2) \mid \text{assert}(\varphi) \mid S; S, \\ \varphi ::= & x \in \mathcal{A} \mid t_1 \circ t_2 \mid \varphi \vee \varphi \mid \varphi \wedge \varphi, \end{aligned}$$

where e is a regular expression over Σ , $u \in \Sigma^*$, \mathcal{T} is an FFT, \mathcal{A} is an NFA, $\circ \in \{=, \neq, \geq, \leq, >, <\}$, and t_1, t_2 are integer terms defined by the following rules,

$$t ::= i \mid c \mid \text{length}(x) \mid \text{indexOf}_v(x, i) \mid ct \mid t + t, \text{ where } c \in \mathbb{Z}, v \in \Sigma^+.$$

We require that the string-manipulating programs are in **single static assignment (SSA) form**. Note that SSA form imposes restrictions only on the assignment statements, but not on the assertions. A string variable x in an SL_{int} program S is called an *input string variable* of S if it does not appear on the left-hand side of the assignment statements of S . A variable in S is called an *input variable* if it is either an input string variable or an integer variable.

Semantics. The semantics of SL_{int} is explained as follows.

- The assignment $x := y \cdot z$ denotes that x is the concatenation of two strings y and z .
- The assignment $x := \text{replaceAll}_{e,u}(y)$ denotes that x is the string obtained by replacing all occurrences of e in y with u , where the *leftmost and longest* matching of e is used. For instance, $\text{replaceAll}_{(ab)^+,c}(aababaab) = ac \cdot \text{replaceAll}_{(ab)^+,c}(aab) = acac$, since the leftmost and longest matching of $(ab)^+$ in $aababaab$ is $abab$. Here we require that the language defined by e does *not* contain the empty string, in order to avoid the troublesome definition of the semantics of the matching of the empty string. The formal semantics of the `replaceAll` function can be found in [12].

- The assignment $x := \text{reverse}(y)$ denotes that x is the reverse of y .
- The assignment $x := \mathcal{T}(y)$ denotes that $(y, x) \in \mathcal{T}$.
- The assignment $x := \text{substring}(y, t_1, t_2)$ denotes that x is equal to the return value of $\text{substring}(y, t_1, t_2)$, where

$$\text{substring}(y, t_1, t_2) = \begin{cases} \epsilon & \text{if } t_1 < 0 \vee t_1 \geq |y| \vee t_2 = 0 \\ y[t_1, \min\{t_1 + t_2 - 1, |y| - 1\}] & \text{o/w} \end{cases}$$

For instance, $\text{substring}(abaab, -1, 1) = \epsilon$, $\text{substring}(abaab, 3, 0) = \epsilon$, $\text{substring}(abaab, 3, 2) = ab$, and $\text{substring}(abaab, 3, 3) = ab$.

- The conditional statement $\text{assert}(x \in \mathcal{A})$ denotes that x belongs to $\mathcal{L}(\mathcal{A})$.
- The conditional statement $\text{assert}(t_1 \circ t_2)$ denotes that the value of t_1 is equal to (not equal to, ...) that of t_2 , if $\circ \in \{=, \neq, \geq, >, \leq, <\}$.
- The integer term $\text{length}(x)$ denotes the length of x .
- The function $\text{indexOf}_v(x, i)$ returns the starting position of the first occurrence of v in x after the position i , if such an occurrence exists, and -1 otherwise. Note that if $i < 0$, then $\text{indexOf}_v(x, i)$ returns $\text{indexOf}_v(x, 0)$, and if $i \geq \text{length}(x)$, then $\text{indexOf}_v(x, i)$ returns -1 . For instance, $\text{indexOf}_{ab}(aaba, -1) = 1$, $\text{indexOf}_{ab}(aaba, 1) = 1$, $\text{indexOf}_{ab}(aaba, 2) = -1$, and $\text{indexOf}_{ab}(aaba, 4) = -1$.

Path feasibility problem. Given an SL_{int} program S , decide whether there are valuations of the input variables so that S can execute to the end.

4 Decision Procedures for Path Feasibility

In this section, we present a decision procedure for the path feasibility problem of SL_{int} . A distinguished feature of the decision procedure is that it conducts backward computation which is lazy and can be done in a modular way. To support this, we extend a regular language with quantitative information of the strings in the language, giving rise to cost-enriched regular languages and corresponding finite automata (Section 4.1). The crux of the decision procedure is thus to show that the pre-images of cost-enriched regular languages under the string operations in SL_{int} (i.e., concatenation \cdot , $\text{replaceAll}_{e,u}$, reverse, FFTs \mathcal{T} , and substring) are representable by so called cost-enriched recognisable relations (Section 4.2). The overall decision procedure is presented in Section 4.3, supplied by additional complexity analysis.

4.1 Cost-Enriched Regular Languages and Recognisable Relations

Let $k \in \mathbb{N}$ with $k > 0$. A k -cost-enriched string is $(w, (n_1, \dots, n_k))$ where w is a string and $n_i \in \mathbb{Z}$ for all $i \in [k]$. A k -cost-enriched language L is a subset of $\Sigma^* \times \mathbb{Z}^k$. For our purpose, we identify a “regular” fragment of cost-enriched languages as follows.

Definition 1 (Cost-enriched regular languages). Let $k \in \mathbb{N}$ with $k > 0$. A k -cost-enriched language is regular (abbreviated as CERL) if it can be accepted by a cost-enriched finite automaton.

A cost-enriched finite automaton (CEFA) \mathcal{A} is a tuple $(Q, \Sigma, R, \delta, I, F)$ where

- Q, Σ, I, F are defined as in NFAs,
- $R = (r_1, \dots, r_k)$ is a vector of (mutually distinct) cost registers,
- δ is the transition relation which is a finite set of tuples (q, a, q', η) where $q, q' \in Q$, $a \in \Sigma$, and $\eta : R \rightarrow \mathbb{Z}$ is a cost register update function.

For convenience, we usually write $(q, a, q', \eta) \in \Delta$ as $q \xrightarrow{a, \eta} q'$.

A run of \mathcal{A} on a k -cost-enriched string $(a_1 \dots a_m, (n_1, \dots, n_k))$ is a transition sequence $q_0 \xrightarrow{a_1, \eta_1} q_1 \dots q_{m-1} \xrightarrow{a_m, \eta_m} q_m$ such that $q_0 \in I$ and $n_i = \sum_{1 \leq j \leq m} \eta_j(r_i)$ for each $i \in [k]$

(Note that the initial values of cost registers are zero). The run is accepting if $q_m \in F$. A k -cost-enriched string $(w, (n_1, \dots, n_k))$ is accepted by \mathcal{A} if there is an accepting run of \mathcal{A} on $(w, (n_1, \dots, n_k))$. In particular, (ε, n) is accepted by \mathcal{A} if $n = 0$ and $I \cap F \neq \emptyset$. The k -cost-enriched language defined by \mathcal{A} , denoted by $\mathcal{L}(\mathcal{A})$, is the set of k -cost-enriched strings accepted by \mathcal{A} .

The size of a CEFA $\mathcal{A} = (Q, \Sigma, R, \delta, I, F)$, denoted by $|\mathcal{A}|$, is defined as the sum of the sizes of its transitions, where the size of each transition (q, a, q', η) is $\sum_{r \in R} \lceil \log_2(|\eta(r)|) \rceil + 3$. Note here the integer constants in \mathcal{A} are encoded in binary.

Remark 1. CEFAs can be seen as a variant of Cost Register Automata [4], by admitting nondeterminism and discarding partial final cost functions. CEFAs are also closely related to monotonic counter machines [21]. The main difference is that CEFAs discard guards in transitions and allow binary-encoded integers in cost updates, while monotonic counter machines allow guards in transitions but restrict the cost updates to being monotonic and unary, i.e. 0, 1 only. Moreover, we explicitly define CEFAs as language acceptors for cost-enriched languages.

Example 1 (CEFA for length). The string function length can be captured by CEFAs. For any NFA $\mathcal{A} = (Q, \Sigma, \delta, I, F)$, it is not difficult to see that the cost-enriched language $\{(w, \text{length}(w)) \mid w \in \mathcal{L}(\mathcal{A})\}$ is accepted by a CEFA, i.e., $(Q, \Sigma, (r_1), \delta', I, F)$ such that for each $(q, a, q') \in \delta$, we let $(q, a, q', \eta) \in \delta'$, where $\eta(r_1) = 1$.

For later use, we identify a special $\mathcal{A}_{\text{len}} = (\{q_0\}, \Sigma, (r_1), \{(q_0, a, q_0, \eta) \mid \eta(r_1) = 1\}, \{q_0\}, \{q_0\})$. In other words, \mathcal{A}_{len} accepts $\{(w, \text{length}(w)) \mid w \in \Sigma^*\}$.

We can show that the function $\text{indexOf}_v(\cdot, \cdot)$ can be captured by a CEFA as well, in the sense that, for any NFA \mathcal{A} and constant string v , we can construct a CEFA $\mathcal{A}_{\text{indexOf}_v}$ accepting $\{(w, (n, \text{indexOf}_v(w, n))) \mid w \in \mathcal{L}(\mathcal{A}), n \leq \text{indexOf}_v(w, n) < |w|\}$. The construction is slightly technical and can be found in the full version [13].

Note that $\mathcal{A}_{\text{indexOf}_v}$ does not model the corner cases in the semantics of indexOf_v , for instance, $\text{indexOf}_v(w, n) = -1$ if v does not occur after the position n in w .

Given two CEFAs $\mathcal{A}_1 = (Q_1, \Sigma, R_1, \delta_1, I_1, F_1)$ and $\mathcal{A}_2 = (Q_2, \Sigma, \delta_2, R_2, I_2, F_2)$ with $R_1 \cap R_2 = \emptyset$, the product of \mathcal{A}_1 and \mathcal{A}_2 , denoted by $\mathcal{A}_1 \times \mathcal{A}_2$, is defined as $(Q_1 \times Q_2, \Sigma, R_1 \cup R_2, \delta, I_1 \times I_2, F_1 \times F_2)$, where δ comprises the tuples $((q_1, q_2), \sigma, (q'_1, q'_2), \eta)$ such that $(q_1, \sigma, q'_1, \eta_1) \in \delta_1$, $(q_2, \sigma, q'_2, \eta_2) \in \delta_2$, and $\eta = \eta_1 \cup \eta_2$.

For a CEFA \mathcal{A} , we use $R(\mathcal{A})$ to denote the vector of cost registers occurring in \mathcal{A} . Suppose \mathcal{A} is CEFA with $R(\mathcal{A}) = (r_1, \dots, r_k)$ and $\mathbf{i} = (i_1, \dots, i_k)$ is a vector of mutually distinct integer variables such that $R(\mathcal{A}) \cap \mathbf{i} = \emptyset$. We use $\mathcal{A}[\mathbf{i}/R(\mathcal{A})]$ to denote the CEFA obtained from \mathcal{A} by simultaneously replacing r_j with i_j for $j \in [k]$.

Definition 2 (Cost-enriched recognisable relations). Let $(k_1, \dots, k_l) \in \mathbb{N}^l$ with $k_j > 0$ for every $j \in [l]$. A cost-enriched recognisable relation (CERR) $\mathcal{R} \subseteq (\Sigma^* \times \mathbb{Z}^{k_1}) \times \dots \times (\Sigma^* \times \mathbb{Z}^{k_l})$ is a finite union of products of CERLs. Formally, $\mathcal{R} = \bigcup_{i=1}^n L_{i,1} \times \dots \times L_{i,l}$, where for every $j \in [l]$, $L_{i,j} \subseteq \Sigma^* \times \mathbb{Z}^{k_j}$ is a CERL. A CEFA representation of \mathcal{R} is a collection of CEFA tuples $(\mathcal{A}_{i,1}, \dots, \mathcal{A}_{i,l})_{i \in [n]}$ such that $\mathcal{L}(\mathcal{A}_{i,j}) = L_{i,j}$ for every $i \in [n]$ and $j \in [l]$.

4.2 Pre-images of CERLs under string operations

To unify the presentation, we consider string functions $f : (\Sigma^* \times \mathbb{Z}^{k_1}) \times \dots \times (\Sigma^* \times \mathbb{Z}^{k_l}) \rightarrow \Sigma^*$. (If there is no integer input parameter, then k_1, \dots, k_l are zero.)

Definition 3 (Cost-enriched pre-images of CERLs). Suppose that $f : (\Sigma^* \times \mathbb{Z}^{k_1}) \times \dots \times (\Sigma^* \times \mathbb{Z}^{k_l}) \rightarrow \Sigma^*$ is a string function, $L \subseteq \Sigma^* \times \mathbb{Z}^{k_0}$ is a CERL defined by a CEFA $\mathcal{A} = (Q, \Sigma, R, \delta, I, F)$ with $R = (r_1, \dots, r_{k_0})$. Then the R -cost-enriched pre-image of L under f , denoted by $f_R^{-1}(L)$, is a pair $(\mathcal{R}, \mathbf{t})$ such that

- $\mathcal{R} \subseteq (\Sigma^* \times \mathbb{Z}^{k_1+k_0}) \times \dots \times (\Sigma^* \times \mathbb{Z}^{k_l+k_0})$;
- $\mathbf{t} = (t_1, \dots, t_{k_0})$ is a vector of linear integer terms where for each $i \in [k_0]$, t_i is a term whose variables are from $\{r_i^{(1)}, \dots, r_i^{(l)}\}$ which are fresh cost registers and are disjoint from R in \mathcal{A} ;
- L is equal to the language comprising the k_0 -cost-enriched strings

$$(w_0, t_1 [d_1^{(1)}/r_1^{(1)}, \dots, d_1^{(l)}/r_1^{(l)}], \dots, t_{k_0} [d_{k_0}^{(1)}/r_{k_0}^{(1)}, \dots, d_{k_0}^{(l)}/r_{k_0}^{(l)}]),$$

such that

$$w_0 = f((w_1, \mathbf{c}_1), \dots, (w_l, \mathbf{c}_l)) \text{ for some } ((w_1, (\mathbf{c}_1, \mathbf{d}_1)), \dots, (w_l, (\mathbf{c}_l, \mathbf{d}_l))) \in \mathcal{R},$$

where $\mathbf{c}_j \in \mathbb{Z}^{k_j}$, $\mathbf{d}_j = (d_1^{(j)}, \dots, d_{k_0}^{(j)}) \in \mathbb{Z}^{k_0}$ for $j \in [l]$.

The R -cost-enriched pre-image of L under f , say $f_R^{-1}(L) = (\mathcal{R}, \mathbf{t})$, is said to be CERR-definable if \mathcal{R} is a CERR.

Definition 3 is essentially a semantic definition of the pre-images. For the decision procedure, one desires an effective representation of a CERR-definable $f_R^{-1}(L) = (\mathcal{R}, \mathbf{t})$ in terms of CEFAs. Namely, a CEFA representation of $(\mathcal{R}, \mathbf{t})$ (where t_j is over $\{r_j^{(1)}, \dots, r_j^{(l)}\}$ for $j \in [k_0]$) is a tuple $((\mathcal{A}_{i,1}, \dots, \mathcal{A}_{i,l})_{i \in [n]}, \mathbf{t})$ such that $(\mathcal{A}_{i,1}, \dots, \mathcal{A}_{i,l})_{i \in [n]}$ is a CEFA representation of \mathcal{R} , where $R(\mathcal{A}_{i,j}) = (r'_{j,1}, \dots, r'_{j,k_j}, r_1^{(j)}, \dots, r_{k_0}^{(j)})$ for each $i \in [n]$ and $j \in [l]$. (The cost registers $r'_{1,1}, \dots, r'_{1,k_1}, \dots, r'_{l,1}, \dots, r'_{l,k_l}$ are mutually distinct and freshly introduced.)

Example 2 (substring $_{\bar{R}}^{-1}(L)$). Let $\Sigma = \{a\}$ and $L = \{(w, |w|) \mid w \in \mathcal{L}((aa)^*)\}$. Evidently L is a CERL defined by a CEFA $\mathcal{A} = (Q, \Sigma, R, \delta, \{q_0\}, \{q_0\})$ with $Q = \{q_0, q_1\}$, $R = (r_1)$ and $\delta = \{(q_0, a, q_1), (q_1, a, q_0)\}$. Since substring is from $\Sigma^* \times \mathbb{Z}^2$ to Σ^* , substring $_{\bar{R}}^{-1}(L)$, the R -cost-enriched pre-image of L under substring, is the pair $(\mathcal{R}, \mathbf{t})$, where $\mathbf{t} = r_1^{(1)}$ (note that in this case $l = 1$, $k_0 = 1$, and $k_1 = 2$) and

$$\mathcal{R} = \{(w, n_1, n_2, n_2) \mid w \in \mathcal{L}(a^*), n_1 \geq 0, n_2 \geq 0, n_1 + n_2 \leq |w|, n_2 \text{ is even}\},$$

which is represented by $(\mathcal{A}', \mathbf{t})$ such that $\mathcal{A}' = (Q', \Sigma, R', \delta', I', F')$, where

- $Q' = Q \times \{p_0, p_1, p_2\}$, (Intuitively, p_0 , p_1 , and p_2 denote that the current position is before the starting position, between the starting position and ending position, and after the ending position of the substring respectively.)
- $R' = (r'_{1,1}, r'_{1,2}, r_1^{(1)})$,
- $I' = \{(q_0, p_0)\}$, $F' = \{(q_0, p_2), (q_0, p_0)\}$ (where (q_0, p_0) is used to accept the 3-cost-enriched strings $(w, n_1, 0, 0)$ with $0 \leq n_1 \leq |w|$), and
- δ' is

$$\left\{ \begin{array}{l} (q_0, p_0) \xrightarrow{a, \eta_1} (q_0, p_0), (q_0, p_0) \xrightarrow{a, \eta_2} (q_1, p_1), (q_1, p_1) \xrightarrow{a, \eta_2} (q_0, p_1), \\ (q_0, p_1) \xrightarrow{a, \eta_2} (q_1, p_1), (q_1, p_1) \xrightarrow{a, \eta_2} (q_0, p_2), (q_0, p_2) \xrightarrow{a, \eta_3} (q_0, p_2) \end{array} \right\},$$

where $\eta_1(r'_{1,1}) = 1$, $\eta_1(r'_{1,2}) = 0$, $\eta_1(r_1^{(1)}) = 0$, $\eta_2(r'_{1,1}) = 0$, $\eta_2(r'_{1,2}) = 1$, and $\eta_2(r_1^{(1)}) = 1$, $\eta_3(r'_{1,1}) = 0$, $\eta_3(r'_{1,2}) = 0$, and $\eta_3(r_1^{(1)}) = 0$.

Therefore, $\text{substring}_R^{-1}(L)$ is CERR-definable.

It turns out that for each string function f in the assignment statements of SL_{int} , the cost-enriched pre-images of CERLs under f are CERR-definable.

Proposition 1. *Let L be a CERL defined by a CEFA $\mathcal{A} = (Q, \Sigma, R, \delta, I, F)$. Then for each string function f ranging over \cdot , $\text{replaceAll}_{e,u}$, reverse , FFTs \mathcal{T} , and substring , $f_R^{-1}(L)$ is CERR-definable. In addition,*

- a CEFA representation of $\cdot_R^{-1}(L)$ can be computed in time $O(|\mathcal{A}|^2)$,
- a CEFA representation of $\text{reverse}_R^{-1}(L)$ (resp. $\text{substring}_R^{-1}(L)$) can be computed in time $O(|\mathcal{A}|)$,
- a CEFA representation of $(\mathcal{T}(\mathcal{T}))_R^{-1}(L)$ can be computed in time polynomial in $|\mathcal{A}|$ and exponential in $|\mathcal{T}|$,
- a CEFA representation of $(\text{replaceAll}_{e,u})_R^{-1}(L)$ can be computed in time polynomial in $|\mathcal{A}|$ and exponential in $|e|$ and $|u|$.

The proof of Proposition 1 is given in the full version [13].

4.3 The Decision Procedure

Let S be an SL_{int} program. Without loss of generality, we assume that for every occurrence of assignments of the form $y := \text{substring}(x, t_1, t_2)$, it holds that t_1 and t_2 are integer variables. This is not really a restriction, since, for instance, if in $y := \text{substring}(x, t_1, t_2)$, neither t_1 nor t_2 is an integer variable, then we introduce fresh integer variables i and j , replace t_1, t_2 by i, j respectively, and add $\text{assert}(i = t_1); \text{assert}(j = t_2)$ in S . We present a decision procedure for the path feasibility problem of S which is divided into five steps.

Step I: Reducing to atomic assertions.

Note first that in our language, each assertion is a positive Boolean combination of atomic formulas of the form $x \in \mathcal{A}$ or $t_1 \circ t_2$ (cf. Section 3). Nondeterministically choose, for each assertion $\text{assert}(\varphi)$ of S , a set of atomic formulas $\Phi_\varphi = \{\alpha_1, \dots, \alpha_n\}$ such that φ holds when atomic formulas in Φ_φ are true.

Then each assertion $\text{assert}(\varphi)$ in S with $\Phi_\varphi = \{\alpha_1, \dots, \alpha_n\}$ is replaced by $\text{assert}(\alpha_1); \dots; \text{assert}(\alpha_n)$, and thus S constrains atomic assertions only.

Step II: Dealing with the case splits in the semantics of indexOf_v and substring .

For each integer term of the form $\text{indexOf}_v(x, i)$ in S , nondeterministically choose one of the following five options (which correspond to the semantics of indexOf_v in Section 3).

- (1) Add $\text{assert}(i < 0)$ to S , and replace $\text{indexOf}_v(x, i)$ with $\text{indexOf}_v(x, 0)$ in S .
- (2) Add $\text{assert}(i < 0); \text{assert}(x \in \mathcal{A}_{\Sigma^* \setminus v\Sigma^*})$ to S ; replace $\text{indexOf}_v(x, i)$ with -1 in S .
- (3) Add $\text{assert}(i \geq \text{length}(x))$ to S , and replace $\text{indexOf}_v(x, i)$ with -1 in S .
- (4) Add $\text{assert}(i \geq 0); \text{assert}(i < \text{length}(x))$ to S .
- (5) Add

$$\begin{aligned} & \text{assert}(i \geq 0); \text{assert}(i < \text{length}(x)); \text{assert}(j = \text{length}(x) - i); \\ & y := \text{substring}(x, i, j); \text{assert}(y \in \mathcal{A}_{\Sigma^* \setminus v\Sigma^*}) \end{aligned}$$

to S , where y is a fresh string variable, j is a fresh integer variable, and $\mathcal{A}_{\Sigma^* \setminus v\Sigma^*}$ is an NFA defining the language $\{w \in \Sigma^* \mid v \text{ does not occur as a substring in } w\}$. Replace $\text{indexOf}_v(x, i)$ with -1 in S .

For each assignment $y := \text{substring}(x, i, j)$, nondeterministically choose one of the following three options (which correspond to the semantics of substring in Section 3).

- (1) Add the statements $\text{assert}(i \geq 0); \text{assert}(i + j \leq \text{length}(x))$ to S .
- (2) Add the statements $\text{assert}(i \geq 0); \text{assert}(i \leq \text{length}(x)); \text{assert}(i + j > \text{length}(x)); \text{assert}(i' = \text{length}(x) - i)$ to S , and replace $y := \text{substring}(x, i, j)$ with $y := \text{substring}(x, i, i')$, where i' is a fresh integer variable.
- (3) Add the statement $\text{assert}(i < 0); \text{assert}(y \in \mathcal{A}_\varepsilon)$ to S , and remove $y := \text{substring}(x, i, j)$ from S , where \mathcal{A}_ε is the NFA defining the language $\{\varepsilon\}$.

Step III: Removing length and indexOf .

For each term $\text{length}(x)$ in S , we introduce a *fresh* integer variable i , replace every occurrence of $\text{length}(x)$ by i , and add the statement $\text{assert}(x \in \mathcal{A}_{\text{len}}[i/r_1])$ to S . (See Example 1 for the definition of \mathcal{A}_{len} .)

For each term $\text{indexOf}_v(x, i)$ occurring in S , introduce two fresh integer variables i_1 and i_2 , replace every occurrence of $\text{indexOf}_v(x, i)$ by i_2 , and add the statements $\text{assert}(i = i_1); \text{assert}(x \in \mathcal{A}_{\text{indexOf}_v}[i_1/r_1, i_2/r_2])$ to S .

Step IV: Removing the assignment statements backwards.

Repeat the following procedure until S contains no assignment statements.

Suppose $y := f(x_1, \mathbf{i}_1, \dots, x_l, \mathbf{i}_l)$ is the *last* assignment of S , where $f : (\Sigma^* \times \mathbb{Z}^{k_1}) \times \dots \times (\Sigma^* \times \mathbb{Z}^{k_l}) \rightarrow \Sigma^*$ is a string function and $\mathbf{i}_j = (i_{j,1}, \dots, i_{j,k_j})$ for each $j \in [l]$.

Let $\{\mathcal{A}_1, \dots, \mathcal{A}_s\}$ be the set of all CEFA's such that $\text{assert}(y \in \mathcal{A}_j)$ occurs in S for every $j \in [s]$. Let $j \in [s]$ and $R(\mathcal{A}_j) = (r_{j,1}, \dots, r_{j,\ell_j})$. Then from Proposition 1, a CEFA representation of $f_{R(\mathcal{A}_j)}^{-1}(\mathcal{L}(\mathcal{A}_j))$, say $\left((\mathcal{B}_{j,j'}^{(1)}, \dots, \mathcal{B}_{j,j'}^{(l)})_{j' \in [m_j]}, \mathbf{t} \right)$, can be effectively computed from \mathcal{A} and f , where we write

$$R(\mathcal{B}_{j,j'}^{(l)}) = \left((r')_j^{(j',1)}, \dots, (r')_j^{(j',k_{j'})}, r_{j,1}^{(j')}, \dots, r_{j,\ell_j}^{(j')} \right)$$

for each $j' \in [m_j]$ and $j'' \in [l]$, and $t = (t_1, \dots, t_{\ell_j})$. Note that the cost registers $(r')_j^{(1,1)}, \dots, (r')_j^{(1,k_1)}, \dots, (r')_j^{(l,1)}, \dots, (r')_j^{(l,k_l)}, r_{j,1}^{(1)}, \dots, r_{j,\ell_j}^{(1)}, \dots, r_{j,1}^{(l)}, \dots, r_{j,\ell_j}^{(l)}$ are mutually distinct and freshly introduced, moreover, $R(\mathcal{B}_{j,j'_1}^{(j'')}) = R(\mathcal{B}_{j,j'_2}^{(j'')})$ for distinct $j'_1, j'_2 \in [m_j]$. Remove $y := f(x_1, \mathbf{i}_1, \dots, x_l, \mathbf{i}_l)$, as well as all the statements $\text{assert}(y \in \mathcal{A}_1), \dots, \text{assert}(y \in \mathcal{A}_s)$ from S . For every $j \in [s]$, nondeterministically choose $j' \in [m_j]$, and add the following statements to S ,

$$\text{assert}(x_1 \in \mathcal{B}_{j,j'}^{(1)}); \dots; \text{assert}(x_l \in \mathcal{B}_{j,j'}^{(l)}); S_{j,j',i_1,\dots,i_l}; S_{j,t}$$

where

$$S_{j,j',i_1,\dots,i_l} \equiv \text{assert}(i_{1,1} = (r')_{j,j'}^{(1,1)}); \dots; \text{assert}(i_{1,k_1} = (r')_{j,j'}^{(1,k_1)}); \\ \dots \\ \text{assert}(i_{l,1} = (r')_{j,j'}^{(l,1)}); \dots; \text{assert}(i_{l,k_l} = (r')_{j,j'}^{(l,k_l)})$$

and

$$S_{j,t} \equiv \text{assert}(r_{j,1} = t_1); \dots, \text{assert}(r_{j,\ell_j} = t_{\ell_j}).$$

Step V: Final satisfiability checking.

In this step, S contains no assignment statements and only assertions of the form $\text{assert}(x \in \mathcal{A})$ and $\text{assert}(t_1 \circ t_2)$ where \mathcal{A} are CEFAs and t_1, t_2 are linear integer terms. Let X denote the set of string variables occurring in S . For each $x \in X$, let $\Lambda_x = \{\mathcal{A}_x^1, \dots, \mathcal{A}_x^{s_x}\}$ denote the set of CEFAs \mathcal{A} such that $\text{assert}(x \in \mathcal{A})$ appears in S . Moreover, let ϕ denote the conjunction of all the LIA formulas $t_1 \circ t_2$ occurring in S . It is straightforward to observe that ϕ is over $R' = \bigcup_{x \in X, j \in [s_x]} R(\mathcal{A}_x^j)$. Then the path feasibility of S is reduced to *the satisfiability problem of LIA formulas w.r.t. CEFAs (abbreviated as SAT_{CEFA}[LIA] problem)* which is defined as

deciding whether ϕ is satisfiable w.r.t. $(\Lambda_x)_{x \in X}$, namely, whether there are an assignment function $\theta : R' \rightarrow \mathbb{Z}$ and strings $(w_x)_{x \in X}$ such that $\phi[\theta(R')/R']$ holds and $(w_x, \theta(R(\mathcal{A}_x^j))) \in \mathcal{L}(\mathcal{A}_x^j)$ for every $x \in X$ and $j \in [s_x]$.

This SAT_{CEFA}[LIA] problem is decidable and PSPACE-complete; The proof can be found in the full version [13].

Proposition 2. SAT_{CEFA}[LIA] is PSPACE-complete.

An example to illustrate the decision procedure can be found in the full version [13]. *Complexity analysis of the decision procedure.* Step I and Step II can be done in nondeterministic linear time. Step III can be done in linear time. In Step IV, for each input string variable x in S , at most exponentially many CEFAs can be generated for x , each of which is of at most exponential size. Therefore, Step IV can be done in nondeterministic exponential space. By Proposition 2, Step V can be done in exponential space. Therefore, we conclude that the path feasibility problem of SL_{int} programs is in NEXPSPACE, thus in EXPSPACE by Savitch's theorem [23].

Remark 2. In this paper, we focus on functional finite transducers (cf. Section 2). Our decision procedure is applicable to general finite transducers as well with minor adaptation. However, the `EXPSpace` complexity upper-bound does not hold any more, because the distributive property $f^{-1}(L_1 \cap L_2) = f^{-1}(L_1) \cap f^{-1}(L_2)$ for regular languages L_1, L_2 only holds for functional finite transducers f .

5 Evaluations

We have implemented the decision procedure presented in the preceding section based on the recent string constraint solver `OSTRICH` [14], resulting in a new solver `OSTRICH+`. `OSTRICH` is written in Scala and based on the SMT solver `Princess` [25]. `OSTRICH+` reuses the parser of `Princess`, but replaces the NFAs from `OSTRICH` with CEFA. Correspondingly, in `OSTRICH+`, the pre-image computation for concatenation, `replaceAll`, `reverse`, and finite transducers is reimplemented, and a new pre-image operator for substring is added. `OSTRICH+` also implements CEFA constructions for `length` and `indexOf`. More details can be found in the full version [13].

We have compared `OSTRICH+` with some of the state-of-the-art solvers on a wide range of benchmarks. We discuss the benchmarks in Section 5.1 and present the experimental results in Section 5.2.

5.1 Benchmarks

Our evaluation focuses on problems that combine string with integer constraints. To this end, we consider the following four sets of benchmarks, all in `SMT-LIB 2` format.

`TRANSDUCER+` is derived from the `TRANSDUCER` benchmark suite of `OSTRICH` [14]. The `TRANSDUCER` suite involves seven transducers: `toUpper` (replacing all lowercase letters with their uppercase ones) and its dual `toLower`, `htmlEscape` and its dual `htmlUnescape`, `escapeString`, `addslashes`, and `trim`. These transducers are collected from `Stranger` [33] and `SLOTH` [18]. Initially none of the benchmarks involved integers. In `TRANSDUCER+`, we encode four security-relevant properties of transducers [19], with the help of the functions `charAt` and `length`:

- idempotence: given \mathcal{T} , whether $\forall x. \mathcal{T}(\mathcal{T}(x)) = \mathcal{T}(x)$;
- duality: given \mathcal{T}_1 and \mathcal{T}_2 , whether $\forall x. \mathcal{T}_2(\mathcal{T}_1(x)) = x$;
- commutativity: given \mathcal{T}_1 and \mathcal{T}_2 , whether $\forall x. \mathcal{T}_2(\mathcal{T}_1(x)) = \mathcal{T}_1(\mathcal{T}_2(x))$;
- equivalence: given \mathcal{T}_1 and \mathcal{T}_2 , whether $\forall x. \mathcal{T}_1(x) = \mathcal{T}_2(x)$.

For instance, we encode the non-idempotence of \mathcal{T} into the path feasibility of the `SLint` program $y := \mathcal{T}(x); z := \mathcal{T}(y); S_{y \neq z}$, where y and z are two fresh string variables, and $S_{y \neq z}$ is the `SLint` program encoding $y \neq z$ (see the full version [13] for the details). We also include in `TRANSDUCER+` three instances generated from a program to sanitize URLs against XSS attacks (see the full version [13] for the details), where $\mathcal{T}_{\text{trim}}$ is used. In total, we obtain 94 instances for the `TRANSDUCER+` suite.

`SLOG+` is adapted from the `SLOG` benchmark suite [31], containing 3,511 instances about strings only. We obtain `SLOG+` by choosing a string variable x for each instance, and adding the statement `assert(length(x) < 2 * indexOfa(x, 0))` for some $a \in \Sigma$. As in [14], we split `SLOG+` into `SLOG+(REPLACE)` and `SLOG+(REPLACEALL)`, comprising

Benchmark	Output	CVC4	Z3-str3	Z3-Trau	OSTRICH ⁽¹⁾	OSTRICH ⁽²⁾	OSTRICH+
TRANSDUCER+ Total: 94	sat	–	–	–	0	0	84
	unsat	–	–	–	1	1	4
	inconcl.	–	–	–	93	93	6
SLOG+(REPLACEALL) Total: 120	sat	104	–	–	0	0	98
	unsat	11	–	–	7	5	12
	inconcl.	5	–	–	113	115	10
SLOG+(REPLACE) Total: 3,391	sat	1,309	878	–	0	169	584
	unsat	2,082	2,066	–	2,079	2,075	2,082
	inconcl.	0	447	–	1,312	1,147	725
PyEx-td Total: 5,569	sat	4,224	4,068	4,266	68	96	4,141
	unsat	1,284	1,289	1,295	95	93	1,203
	inconcl.	61	212	8	5,406	5,380	225
PyEx-z3 Total: 8,414	sat	6,346	6,040	7,003	76	100	5,489
	unsat	1,358	1,370	1,394	61	53	1,239
	inconcl.	710	1,004	17	8,277	8,261	1,686
PyEx-zz Total: 11,438	sat	10,078	8,804	10,129	71	98	9,033
	unsat	1,204	1,207	1,222	91	61	868
	inconcl.	156	1,427	87	11,276	11,279	1,537
KALUZA Total: 47,284	sat	35,264	33,438	34,769	23,397	28,522	27,962
	unsat	12,014	11,799	12,014	10,445	10,445	9,058
	inconcl.	6	2,047	501	13,442	8,317	10,264
Total: 76,310	solved	75,278	70,959	72,092	36,391	41,718	61,857
	unsolved	1,032	5,351	4,218	39,919	34,592	14,453

Table 1. Experimental results on different benchmark suites. ‘–’ means that the tool is not applicable to the benchmark suite, and ‘inconclusive’ means that a tool gave up, timed out, or crashed.

3,391 and 120 instances respectively. In addition to the `indexOf` and `length` functions, the benchmarks use regular constraints and concatenation; `SLOG+(REPLACE)` also contains the `replace` function (replacing the first occurrence), while `SLOG+(REPLACEALL)` uses the `replaceAll` function (replacing all occurrences).

`PyEx` [24] contains 25,421 instances derived by the `PyEx` tool, a symbolic execution engine for Python programs. The `PyEx` suite was generated by the CVC4 group from four popular Python packages: `httplib2`, `pip`, `pymongo`, and `requests`. These instances use regular constraints, concatenation, `length`, `substring`, and `indexOf` functions. Following [24], the `PyEx` suite is further divided into three parts: `PyEx-td`, `PyEx-z3` and `PyEx-zz`, comprising 5,569, 8,414 and 11,438 instances, respectively.

`KALUZA` [26] is the most well-known benchmark suite in literature, containing 47,284 instances with regular constraints, concatenation, and the `length` function. The 47,284 benchmarks include 28,032 satisfiable and 9,058 unsatisfiable problems in SSA form.

5.2 Experiments

We compare `OSTRICH+` to CVC4 [20], Z3-str3 [34], and Z3-Trau [9], as well as two configurations of `OSTRICH` [14] with standard NFAs. The configuration `OSTRICH(1)`

is a direct implementation of the algorithm in [14], and does not support integer functions. In OSTRICH⁽²⁾, we integrated support for the length function as in Norn [2], based on the computation of length abstractions of regular languages, and handle `indexOf`, `substring`, and `charAt` via an encoding to word equations. The experiments are executed on a computer with an Intel Xeon Silver 4210 2.20GHz and 2.19GHz CPU (2-core) and 8GB main memory, running 64bit Ubuntu 18.04 LTS OS and Java 1.8. We use a timeout of 30 seconds (wall-clock time), and report the number of satisfiable and unsatisfiable problems solved by each of the systems. Table 1 summarises the experimental results. We did not observe incorrect answers by any tool.

There are two additional state-of-the-art solvers Slent and Trau+ which were not included in the evaluation. We exclude Slent [32] because it uses its own input format `laut`, which is different from the SMT-LIB 2 format used for our benchmarks; also, TRANSUCER+ is beyond the scope of Slent. Trau+ [3] integrates Trau with Sloth to deal with both finite transducers and integer constraints. We were unfortunately unable to obtain a working version of Trau+, possibly because Trau requires two separate versions of Z3 to run. In addition, the algorithm in [3] focuses on length-preserving transducers, which means that TRANSUCER+ is beyond the scope of Trau+.

OSTRICH+ and OSTRICH are the only tools applicable to the problems in TRANSUCER+. With a timeout of 30s, OSTRICH+ can solve 88 of the benchmarks, but this number rises to 94 when using a longer timeout of 600s. Given the complexity of those benchmarks, this is an encouraging result. OSTRICH can only solve one of the benchmarks, because the encoding of `charAt` in the benchmarks using equations almost always leads to problems that are not in SSA form.

On SLOG+(REPLACEALL), OSTRICH+ and CVC4 are very close: OSTRICH+ solves 98 satisfiable instances, slightly less than the 104 instances solved by CVC4, while OSTRICH+ solves one more unsatisfiable instance than CVC4 (12 versus 11). The suite is beyond the scope of Z3-str3 and Z3-Trau, which do not support `replaceAll`.

On SLOG+(REPLACE), OSTRICH+, CVC4, and Z3-str3 solve a similar number of unsatisfiable problems, while CVC4 solves the largest number of satisfiable instances (1,309). The suite is beyond the scope of Z3-Trau which does not support `replace`.

On the three PyEx suites, Z3-Trau consistently solves the largest number of instances by some margin. OSTRICH+ solves a similar number of instances as Z3-str3. Interpreting the results, however, it has to be taken into account that PyEx includes 1,334 instances that are *not* in SSA form, which are beyond the scope of OSTRICH+.

The KALUZA problems can be solved most effectively by CVC4. OSTRICH+ can solve almost all of the around 80% of the benchmarks that are in SSA form, however.

OSTRICH+ consistently outperforms OSTRICH⁽¹⁾ and OSTRICH⁽²⁾ in the evaluation, except for the KALUZA benchmarks. For OSTRICH⁽¹⁾, this is expected because most benchmarks considered here contain integer functions. For OSTRICH⁽²⁾, it turns out that the encoding of `indexOf`, `substring`, and `charAt` as word equations usually leads to problems that are not in SSA form, and therefore are beyond the scope of OSTRICH.

In summary, we observe that OSTRICH+ is competitive with other solvers, while is able to handle benchmarks that are beyond the scope of the other tools due to the combination of string functions (in particular transducers) and integer constraints. Interestingly, the experiments show that OSTRICH+, at least in its current state, is better

at solving unsatisfiable problems than satisfiable problems; this might be an artefact of the use of nuXmv for analysing products of CEFAs. We expect that further optimisation of our algorithm will lead to additional performance improvements. For instance, a natural optimisation that is to be included in our implementation is to use standard finite automata, as opposed to CEFAs, for simpler problems such as the KALUZA benchmarks. Such a combination of automata representations is mostly an engineering effort.

6 Conclusion

In this paper, we have proposed an expressive string constraint language which can specify constraints on both strings and integers. We provided an automata-theoretic decision procedure for the path feasibility problem of this language. The decision procedure is simple, generic, and amenable to implementation, giving rise to a new solver OSTRICH+. We have evaluated OSTRICH+ on a wide range of existing and newly created benchmarks, and have obtained very encouraging results. OSTRICH+ is shown to be the first solver which is capable of tackling finite transducers and integer constraints with completeness guarantees. Meanwhile, it demonstrates competitive performance against some of the best state-of-the-art string constraint solvers.

Acknowledgements. T. Chen and Z. Wu are supported by Guangdong Science and Technology Department grant (No. 2018B010107004); T. Chen is also supported by Overseas Grant (KFKT2018A16) from the State Key Laboratory of Novel Software Technology, Nanjing University, China and Natural Science Foundation of Guangdong Province, China (No. 2019A1515011689). M. Hague is supported by EPSRC [EP/T00021X/1]; A. Lin is supported by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement no 759969). P. Rümmer is supported by the Swedish Research Council (VR) under grant 2018-04727, and by the Swedish Foundation for Strategic Research (SSF) under the project WebSec (Ref. RIT17-0011). Z. Wu is supported by the Open Project of Shanghai Key Laboratory of Trustworthy Computing (No. 07dz22304201601), the NSFC grants (No. 61872340), and the INRIA-CAS joint research project VIP.

References

1. P. A. Abdulla, M. F. Atig, Y. Chen, B. P. Diep, L. Holík, A. Rezine, and P. Rümmer. Flatten and conquer: a framework for efficient analysis of string constraints. In *PLDI*, pages 602–617, 2017.
2. P. A. Abdulla, M. F. Atig, Y. Chen, L. Holík, A. Rezine, P. Rümmer, and J. Stenman. String constraints for verification. In *CAV*, pages 150–166, 2014.
3. P. A. Abdulla, M. F. Atig, B. P. Diep, L. Holík, and P. Janku. Chain-free string constraints. In *ATVA*, pages 277–293, 2019.
4. R. Alur, L. D’Antoni, J. Deshmukh, M. Raghothaman, and Y. Yuan. Regular functions and cost register automata. In *LICS*, pages 13–22. IEEE Computer Society, 2013.
5. P. Barceló, D. Figueira, and L. Libkin. Graph logics with rational relations. *Logical Methods in Computer Science*, 9(3), 2013.
6. M. Berzish, V. Ganesh, and Y. Zheng. Z3str3: A string solver with theory-aware heuristics. In *FMCAD*, pages 55–59, 2017.
7. N. Bjørner, N. Tillmann, and A. Voronkov. Path feasibility analysis for string-manipulating programs. In *TACAS*, pages 307–321, 2009.

8. J. R. Büchi and S. Senger. Definability in the existential theory of concatenation and undecidable extensions of this theory. In *Collected Works of J. R. Büchi*, pages 671–683. 1990.
9. D. Bui and contributors. Z3-trau, 2019.
10. T. Bultan and contributors. Abc string solver, 2015.
11. C. Cadar and K. Sen. Symbolic execution for software testing: Three decades later. *Commun. ACM*, 56(2):82–90, Feb. 2013.
12. T. Chen, Y. Chen, M. Hague, A. W. Lin, and Z. Wu. What is decidable about string constraints with the replaceall function. *PACMPL*, 2(POPL):3:1–3:29, 2018.
13. T. Chen, M. Hague, J. He, D. Hu, A. W. Lin, P. Rümmer, and Z. Wu. A decision procedure for path feasibility of string manipulating programs with integer data type (full version). 2020. <http://arxiv.org/abs/2007.06913>.
14. T. Chen, M. Hague, A. W. Lin, P. Rümmer, and Z. Wu. Decision procedures for path feasibility of string-manipulating programs with complex operations. *PACMPL*, 3(POPL), 2019.
15. J. D. Day, V. Ganesh, P. He, F. Manea, and D. Nowotka. RP. pages 15–29, 2018.
16. L. de Moura and N. Bjørner. Z3: an efficient SMT solver. In *TACAS*, pages 337–340, 2008.
17. V. Ganesh, M. Minnes, A. Solar-Lezama, and M. C. Rinard. Word equations with length constraints: What’s decidable? In *HVC 2012*, pages 209–226, 2012.
18. L. Holík, P. Janku, A. W. Lin, P. Rümmer, and T. Vojnar. String constraints with concatenation and transducers solved efficiently. *PACMPL*, 2(POPL):4:1–4:32, 2018.
19. P. Hooimeijer, B. Livshits, D. Molnar, P. Saxena, and M. Veanes. Fast and precise sanitizer analysis with BEK. In *USENIX Security Symposium*, 2011.
20. T. Liang, A. Reynolds, C. Tinelli, C. Barrett, and M. Deters. A DPLL(T) theory solver for a theory of strings and regular expressions. In *CAV*, pages 646–662, 2014.
21. A. W. Lin and P. Barceló. String solving with word equations and transducers: Towards a logic for analysing mutation XSS. In *POPL*, pages 123–136. ACM, 2016.
22. A. W. Lin and R. Majumdar. Quadratic word equations with length constraints, counter systems, and presburger arithmetic with divisibility. In *ATVA*, pages 352–369, 2018.
23. C. H. Papadimitriou. *Computational complexity*. Addison-Wesley, 1994.
24. A. Reynolds, M. Woo, C. Barrett, D. Brumley, T. Liang, and C. Tinelli. Scaling up DPLL(T) string solvers using context-dependent simplification. In *CAV*, pages 453–474, 2017.
25. P. Rümmer. A constraint sequent calculus for first-order logic with linear integer arithmetic. In *LPAR*, pages 274–289, 2008.
26. P. Saxena, D. Akhawe, S. Hanna, F. Mao, S. McCamant, and D. Song. A symbolic execution framework for javascript. In *S&P*, pages 513–528, 2010.
27. M. Trinh, D. Chu, and J. Jaffar. S3: A symbolic string solver for vulnerability detection in web applications. In *CCS*, pages 1232–1243, 2014.
28. M. Trinh, D. Chu, and J. Jaffar. Progressive reasoning over recursively-defined strings. In *CAV*, pages 218–240. Springer, 2016.
29. A. van der Stock, B. Glas, N. Smithline, and T. Gigler. OWASP Top 10 – 2017, 2017.
30. K. N. Verma, H. Seidl, and T. Schwentick. On the complexity of equational horn clauses. In *CADE*, pages 337–352, 2005.
31. H. Wang, T. Tsai, C. Lin, F. Yu, and J. R. Jiang. String analysis via automata manipulation with logic circuit representation. In *CAV*, pages 241–260, 2016.
32. H.-E. Wang, S.-Y. Chen, F. Yu, and J.-H. R. Jiang. A symbolic model checking approach to the analysis of string and length constraints. In *ASE*, page 623633. ACM, 2018.
33. F. Yu, M. Alkhalaf, T. Bultan, and O. H. Ibarra. Automata-based symbolic string analysis for vulnerability detection. *Form. Methods Syst. Des.*, 44(1):44–70, 2014.
34. Y. Zheng, X. Zhang, and V. Ganesh. Z3-str: a Z3-based string solver for web application analysis. In *ESEC/SIGSOFT FSE*, pages 114–124, 2013.