

A Decision Procedure for String Constraints with String-Integer Conversion and Flat Regular Constraints

Hao Wu^{1,2}, Yu-Fang Chen³, Zhilin Wu^{1,2} and Naijun Zhan^{1,2*}

^{1*}State Key Laboratory of Computer Science, Institute of Software,
Chinese Academy of Sciences, Beijing, China.

²University of Chinese Academy of Sciences, Beijing, China.

³Institute of Information Science, Academia Sinica, Taiwan, Republic
of China.

*Corresponding author(s). E-mail(s): znj@ios.ac.cn;

Contributing authors: wuhao@ios.ac.cn; yfc@iis.sinica.edu.tw;
wuzl@ios.ac;

Abstract

String constraint solving is the core of various testing and verification approaches for scripting languages. Among algorithms for solving string constraints, flattening is a well-known approach that is particularly useful in handling satisfiable instances. As string-integer conversion is an important function appearing in almost all scripting languages, Abdulla et al. extended the flattening approach to this function recently. However, their approach supports only a special flattening pattern and leaves the support of the general flat regular constraints as an open problem. In this paper, we fill the gap and propose a complete flattening approach for the string-integer conversion. The approach is built upon a quantifier elimination procedure for the linear-exponential arithmetic (namely, the extension of Presburger arithmetic with exponential functions) proposed by Point in 1986. The complexity of the quantifier elimination procedure is analyzed for the first time and is shown to be 3-EXPSPACE if the formula contains only existential quantifiers. While the quantifier elimination procedure by Point is too expensive to be implemented efficiently, we propose various optimizations and provide a prototypical implementation. We evaluate the performance of our implementation on the benchmarks that are generated from the string hash functions as well as randomly. The experimental results show that our implementation outperforms the state-of-the-art solvers.

Keywords: String-integer conversion, Flat regular constraints, Exponential function, Presburger arithmetic, Quantifier elimination

1 Introduction

The emerging of scripting languages boosted the needs of efficient approaches and tools to ensure program quality. Comparing with traditional programming languages, string data type plays a more critical role in its analysis. String constraint solvers are the engine of modern scripting program analysis techniques. Due to the high demand, in recent years, there is a boosting amount of publications on this subject.

However, research progress of string constraint solving has been hampered by many major challenges in both theory and tool implementation aspects (including long-standing open problems). Logical theories over strings have to allow string concatenation, which is arguably the most fundamental operation of strings. The most celebrated result concerning theories of strings is Makanin's result on deciding the satisfiability problem for *word equations* [1]. A simple example of a word equation is $xaby = ybax$, where x, y are variables, and $a, b \in \Sigma$ are constant letters. A word equation is satisfiable if it has a solution, i.e., an assignment that maps variables to strings over the alphabet Σ which equates the left-hand side with the right-hand side of the equation. The correctness proof of Makanin's algorithm is arguably one of the most complex termination proofs in computer science. Makanin's result can be extended to include *regular constraints* (a.k.a. regular expression matching, e.g., $x \in (ba)^*$), and arbitrary Boolean connectives. This extension is called word equations with regular constraints. However, the satisfiability problem of word equations together with length constraints (e.g., $|x| = |y| + 1 \wedge wx = yx$) is still open. The complexity of the satisfiability of word equations with regular constraints was proven to be PSPACE-complete by Plandowski [2], after decades of improvement of the original algorithm by Makanin.

Satisfiability of word equations is a special instance of Hilbert's 10th problem. In the past, the original motivation of studying word equations was to find an undecidability proof of Hilbert's 10th problem. However, the motivation is no longer valid since Makanin finds a decision procedure. Recently, driven by the need for program analysis, people started to revisit the problem and its extensions to describe the complete string library APIs in conventional programming languages. Many highly efficient solvers for string constraints, to name a few, CVC4 [3], Z3 [4], Z3-Str3 [5], S3 [6], Norn [7], Ostrich [8], Sloth [9], ABC [10], Stranger [11], Trau [12] and so on, are developed in the last decade. The satisfiability of *string-integer conversion constraints*, e.g., $wx = yx \wedge |x| > \text{parseInt}(y)$, has been proven undecidable in [13]. However, this kind of constraints is pervasive in scripting language programs. For example, it is common that programs read string inputs from text files and converts a part of the string input to integers. Even more crucially, in many programming languages, the string-integer conversion is a part of the definition of their core semantics [14]. JavaScript, which powers most interactive content on the Web and increasingly server-side code with Node.js, is one of such languages.

Due to the difficulties in solving string constraints and, in practice, satisfiable string constraints are more critical for automatic testing, one idea is to have separate specialized procedures for solving satisfiable sub-problems. Currently, there are two main specialized approaches for proving satisfiability. The first is to consider only strings of bounded length. This approach is taken in the first-generation solvers such as Hampi [15] and Kuluza [16]. Although they are useful in handling many practical cases, they fail to find an answer when the all string solutions exceed the selected bound. For example, a constraint of the form $x.y \neq z \wedge |x| > 2000$ would be quite challenging to handle using those solvers.

One more recent approach is flattening [17–19]. The idea is to restrict the solution space of string variables to *flat languages* (see Section 3). The major benefit of considering this class is two-fold. First, under the restriction, the potential solution space is still infinite, which gives us a higher potential of finding solutions. For instance, we can find a solution for $x.y \neq z \wedge |x| > 2000$ under a very simple restriction: all variables are in a^* , where a is an element of the alphabet. Second, more importantly, because we can convert the membership problem of a flat language to the satisfiability problem of a Presburger arithmetic formula, the class of word equations + flat languages + length constraints is decidable.

The paper of Abdulla et al. [17] has considered adding string-integer conversion constraints to the above class, which proposed an algorithm for a restricted form of flat languages and left the support of general flat languages as an open problem. For string-integer conversion constraints, their approach projects the solution to a finite solution space (in a way similar to the PASS [20] approach).

In this paper, we give a complete solution to this problem. We propose a decision procedure for the class of word equations + flat languages + length constraints + string-integer conversion. The basic idea of our approach can be sketched as follows: we first reduce the satisfiability problem to the corresponding satisfiability problem of the theory of Presburger Arithmetic with exponential functions (denoted by ExpPA), more precisely, the existential fragment of ExpPA; then, according to the decidability of ExpPA, we obtain the decidability of the original satisfiability problem.

The decidability of ExpPA was first shown by Semenöv in [21]. Nevertheless, Semenöv did not provide an explicit decision procedure. To remedy this, in [22], Point presented the first quantifier elimination procedure for the satisfiability of ExpPA. Partially attributed to the fact that Point's procedure in [22] was presented in a mathematical and dense way, this quantifier elimination procedure has mostly eluded the attentions of computer science community¹. To the best of our knowledge, no implementation based on Point's procedure was available up to now.

Aiming at introducing Point's procedure to computer science community, we reformulate (and slightly improve) Point's procedure in a way, which, hopefully, is more accessible for computer science community (Section 4). This can be seen as another contribution of this paper. We also analyze the deterministic upper bound of complexity for the first time: given a ExpPA formula with only existential quantifiers, eliminating all quantifiers has a 3-EXPSpace complexity. Furthermore, we propose

¹There are only two citations in Google Scholar.

various optimizations (Section 6) and achieve the first prototypical implementation of Point’s procedure (Section 7).

In fact, other than the theoretical difficulties, in practice, the string-integer conversion is quite challenging for state-of-the-art solvers. Here we illustrate a toy example that mimic the “mining” step of block-chain construction. Essentially, given a string hash function $\text{hash}(w)$, the goal of the mining step is to find a nonce n such that when inserting n to the text to be protected, say w_1 and w_2 , $\text{hash}(w_1 \cdot n \cdot w_2)$ satisfies a certain pattern, e.g., the last k digits are zeros. If w_1 or w_2 are modified, one needs to compute another n satisfies the desired pattern. Below we consider a simple hash function: $\text{hash}(w) = \sum_{i=1}^n a_i p^{n-i} \bmod m$, where $p, m \in \mathbb{Z}^+$ with $p, m \geq 2$. It is easy to see that $\text{hash}(w)$ can be seen as a generalization of `parseInt` followed by a modulo operation. In particular, if $\Sigma = \Sigma_{\text{num}}$ and $p = 10$, then $\text{hash}(w) = \text{parseInt}(w) \bmod m$. Thus, the problem of finding a suitable input w such that the last k digits of $\text{hash}(w)$ are zeros can be modeled as a string constraint with `parseInt`. Although the example is seemingly simple, it is already challenging for most state-of-the-art solvers, as shown by our experiment results in Section 7. With the optimizations introduced in Section 6, our implementation manages to solve several variants of the string-hash examples as well as some randomly generated arithmetic problem instances better than the state-of-the-art solvers (Section 7).

Structure

After the preliminaries in Section 2, we show how to flatten a string constraint with string-integer conversion to an existential ExpPA formula in Section 3. We describe the quantifier elimination procedure for ExpPA in Section 4 and analyse its complexity in Section 5. Several optimization techniques are presented in Section 6. Finally, we describe the implementation and experiment results in Section 7.

2 Preliminary

In this section, we fix the notations and introduce some basic concepts, including Presburger arithmetic, finite-state automata, and flat languages.

Integers, strings, and languages

Let \mathbb{N} denote the set of natural numbers, \mathbb{Z} denote the set of integers, and \mathbb{Z}^+ denote the set of positive integers. For $n \in \mathbb{Z}^+$, let $[n]$ denote $\{1, \dots, n\}$.

An *alphabet* Σ is a finite set. Each element of Σ is called a *letter*. A *string* w over Σ is a (possibly empty) finite sequence $a_1 \dots a_n$ with $a_i \in \Sigma$ for every $i \in [n]$. Let ε denote the empty string, namely, the empty sequence. Let Σ^* denote the set of all strings over Σ . Let Σ^+ denote the set of nonempty strings over Σ . For convenience, we also use Σ_ϵ to denote $\Sigma \cup \{\epsilon\}$. For a string $w = a_1 \dots a_n \in \Sigma^*$, let $\text{len}(w)$ denote the *length* of w , i.e. n . In particular, $\text{len}(\varepsilon) = 0$. For $w_1 = a_1 \dots a_m, w_2 = b_1 \dots b_n \in \Sigma^*$, let $w_1 \cdot w_2$ denote the *concatenation* of w_1 and w_2 , that is, $a_1 \dots a_m b_1 \dots b_n$. A language L over Σ is a subset of Σ^* .

Presburger Arithmetic

Presburger Arithmetic (PA) is the first-order logic of integers with addition. A term of PA, denoted by \mathfrak{t} , is of the form

$$\mathfrak{t} \doteq c \mid \mathfrak{x} \mid \mathfrak{t} + \mathfrak{t} \mid \mathfrak{t} - \mathfrak{t},$$

where \mathfrak{x} and c represent integer variables and integer constants respectively. A formula of PA, denoted by ϕ , is of the form

$$\phi \doteq \mathfrak{t} \odot \mathfrak{t} \mid c \mid \mathfrak{t} \mid c \nmid \mathfrak{t} \mid \neg \phi \mid \phi \wedge \phi \mid \phi \vee \phi \mid \exists \mathfrak{x}. \phi \mid \forall \mathfrak{x}. \phi,$$

where $\odot \in \{=, <, >, \leq, \geq\}$.

A *quantifier-free* PA (QFPA) formula is a PA formula containing no quantifiers. A PA formula is called *existential* if it contains no occurrences of universal quantifiers. The set of free variables of ϕ is denoted by $\text{Free}(\phi)$. We usually write $\phi(\mathfrak{x}_1, \dots, \mathfrak{x}_k)$ to denote an PA formula ϕ such that $\text{Free}(\phi) \subseteq \{\mathfrak{x}_1, \dots, \mathfrak{x}_k\}$. Given an PA formula ϕ , and an integer interpretation of $\text{Free}(\phi)$, i.e. a function $I : \text{Free}(\phi) \rightarrow \mathbb{Z}$, we denote by $I \models \phi$ that I satisfies ϕ (which is defined in the standard manner, with $+$, $-$, \mid and \nmid interpreted as the integer addition, subtraction, divisibility and undivisibility relation respectively), and call I a *model* of ϕ . We use $\llbracket \phi \rrbracket$ to denote the set of models of ϕ .

We always assume, without loss of expressiveness, that quantifier-free formulas contain no negations, because we can transform a formula into negation normal form and negations before atomic formulas can be absorbed by changing the predicates (for example, $\neg(\mathfrak{x} \geq \mathfrak{y}) \equiv \mathfrak{x} < \mathfrak{y}$). The only time when we deal with negations is to transform an universal quantifier \forall into an existential one $\neg \exists \neg$.

It is well-known that PA admits quantifier elimination, for example, Cooper's algorithm[23]. In Sec 4, we will use this fact and apply Cooper's algorithm as a subprocedure.

Finite state automata

A *finite state automaton* (FA) is a tuple $\mathcal{A} = \langle Q, \Sigma, \Delta, q_{\text{init}}, F \rangle$, where Q is a finite set of states, Σ is a finite alphabet, $\Delta \subseteq Q \times \Sigma_\epsilon \times Q$ is the transition relation, q_{init} is the initial state, $F \subseteq Q$ is the set of accepting states. A *run* of \mathcal{A} on a string $w = a_1 \dots a_n$ is a sequence $q_0 \xrightarrow{b_1} q_1 \xrightarrow{b_2} \dots \xrightarrow{b_{m-1}} q_{m-1} \xrightarrow{b_m} q_m$ such that $q_0 = q_{\text{init}}$, $(q_{i-1}, b_i, q_i) \in \Delta$ for every $i \in [m]$, and $a_1 \dots a_n = b_1 \dots b_m$. We say a run $q_0 \xrightarrow{b_1} q_1 \xrightarrow{b_2} \dots \xrightarrow{b_{m-1}} q_{m-1} \xrightarrow{b_m} q_m$ is *accepting* if $q_m \in F$. A string w is *accepted* by \mathcal{A} if there is an accepting run of \mathcal{A} on w . Let $\mathcal{L}(\mathcal{A})$ denote the set of strings accepted by \mathcal{A} . A language $L \subseteq \Sigma^*$ is *regular* if it can be defined by some FA \mathcal{A} , namely, $L = \mathcal{L}(\mathcal{A})$.

Flat languages

We will present flat languages. This is a high level, language-theoretical view at the flat automata from [17]. A *flat language* (FL) over Σ is the set of strings that conform

to a regular expression of the form $(a_1^1 \dots a_{\ell_1}^1)^* \dots (a_1^k \dots a_{\ell_k}^k)^*$, where $a_j^i \in \Sigma$ for each $i \in [k]$ and $j \in [\ell_i]$. Intuitively, an FL is a sequence of loops, and the body of each loop is a string of the form $a_1^i \dots a_{\ell_i}^i$ with ℓ_i letters. For instance, the language defined by $(ab)^*(a)^*(bb)^*$ is an FL.

3 Flattening string constraints with string-integer conversion

In this section, we first define the class of string constraints with string-integer conversion, denoted by $\text{STR}_{\text{parseInt}}$. Then we define the extension of Presburger arithmetic with exponential functions, denoted by ExpPA . Finally, we show how to flatten the string constraints in $\text{STR}_{\text{parseInt}}$ into the arithmetic constraints in the existential fragment of ExpPA .

3.1 String constraints with string-integer conversion ($\text{STR}_{\text{parseInt}}$)

In the sequel, we shall define $\text{STR}_{\text{parseInt}}$, the class of string constraints with the string-integer conversion function parseInt .

The function parseInt takes a decimal string as input and returns the integer represented by the string². For example, $\text{parseInt}('0123') = \text{parseInt}('123') = 1 * 10^2 + 2 * 10 + 3 = 123$. Note here we use the quotation marks to delimit strings.

Formally, the semantics of the parseInt function is defined as follows. In order to simplify the presentation, we assume all string variables ranging over numerical symbols $\Sigma_{\text{num}} = \{0, 1, \dots, 9\}$. Note that one can easily extend our approach to allow arbitrary finite alphabet. Then $\text{parseInt} : \Sigma_{\text{num}}^+ \mapsto \mathbb{N}$ is recursively defined by, for every $w \in \Sigma_{\text{num}}^+$,

- if $w = 'i'$ for $i \in \Sigma_{\text{num}}$, then $\text{parseInt}('i') = i$;
- for $w = w' \cdot 'i'$ for $i \in \Sigma_{\text{num}}$ with $\text{len}(w') \geq 1$, $\text{parseInt}(w) = 10 * \text{parseInt}(w') + \text{parseInt}('i')$.

Note that parseInt is undefined with ε as the input.

In $\text{STR}_{\text{parseInt}}$, there are two types of variables, i.e. the string variables $x, y, \dots \in \mathcal{X}$ and the integer variables $\mathbb{x}, \mathbb{y}, \dots \in \mathbb{X}$. The syntax of $\text{STR}_{\text{parseInt}}$ is defined as follows: a string term, denoted by t , is of the form

$$t \triangleq a \mid x \mid t \cdot t,$$

an integer term, denoted by \mathbb{t} , is of the form

$$\mathbb{t} \triangleq n \mid \mathbb{x} \mid \text{len}(t) \mid \text{parseInt}(t) \mid \mathbb{t} + \mathbb{t} \mid \mathbb{t} - \mathbb{t},$$

²The parseInt function in scripting languages, e.g. Javascript, is more general in the sense that the base can be a number between 2 and 36. Although our approach works for arbitrary positive bases, we choose to focus on the base 10 in this paper, for simplicity.

and a $\text{STR}_{\text{parseInt}}$ formula, denoted by φ is of the form

$$\varphi \hat{=} t = t \mid x \in \mathcal{A} \mid \mathfrak{t} \odot \mathfrak{t} \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \neg \varphi$$

where $a \in (\Sigma_{\text{num}})_\varepsilon$, $n \in \mathbb{N}$, \mathcal{A} is an FA, $\odot \in \{=, <, >, \leq, \geq\}$ and $\text{len}(t)$ denotes the length of a string t . Let us call $t = t$ as string equality constraints, $x \in \mathcal{A}$ as regular constraints, $\mathfrak{t} \odot \mathfrak{t}$ as arithmetic constraints. Let $\text{SVar}(\varphi)$ and $\text{IVar}(\varphi)$ denote the set of string variables and integer variables occurring in φ respectively.

The $\text{STR}_{\text{parseInt}}$ formulas are interpreted on $I = (I_s, I_i)$ where I_s is a partial function from \mathcal{X} (the set of string variables) to Σ^* and I_i is a partial function from \mathbb{X} (the set of integer variables) to \mathbb{N} . Moreover, it is required that the domains of I_s, I_i are finite. Given $I = (I_s, I_i)$, the interpretations of string terms, integer terms, as well as the formulas φ under I are easy to comprehend, thus omitted to avoid tediousness. For $\varphi \in \text{STR}_{\text{parseInt}}$ and $I = (I_s, I_i)$, I satisfies φ if the interpretation of φ under I is **true**. Let us use $\|\varphi\|$ to denote the set of $I = (I_s, I_i)$ satisfying φ .

Example 1 Given a FA \mathcal{A} , the constraint

$$x \in \mathcal{A} \wedge \text{parseInt}(x) = 109_{\mathbb{X}} \wedge \text{len}(x) < 100$$

is a $\text{STR}_{\text{parseInt}}$ formula.

The satisfiability problem of $\text{STR}_{\text{parseInt}}$ is to decide for a given constraint $\varphi \in \text{STR}_{\text{parseInt}}$, whether $\|\varphi\| \neq \emptyset$.

3.2 An Extension of Presburger Arithmetic with Exponential Functions (ExpPA)

ExpPA extends Presburger arithmetic with two partial functions, the *exponential* function $10^{\mathbb{X}}$ and the *integer logarithmic* function $\ell_{10}(\mathbb{X})$ (cf. [22]). The function $\ell_{10}(\mathbb{X})$ is defined as follows: for $n \geq 1$, $\ell_{10}(n) = m$ if $10^m \leq n < 10^{m+1}$. Note that $10^{\mathbb{X}}$ is undefined for $\mathbb{X} < 0$ and $\ell_{10}(\mathfrak{t})$ is undefined for $\mathfrak{t} \leq 0$.

The syntax of ExpPA is obtained from that of PA by adding $10^{\mathfrak{t}}$ and $\ell_{10}(\mathfrak{t})$ to the definition of terms. An ExpPA term, denoted by \mathfrak{t} , is of the form

$$\mathfrak{t} \hat{=} c \mid \mathbb{X} \mid \mathfrak{t} + \mathfrak{t} \mid \mathfrak{t} - \mathfrak{t} \mid 10^{\mathfrak{t}} \mid \ell_{10}(\mathfrak{t})$$

where \mathbb{X} is an integer variable and $c \in \mathbb{Z}$ is a constant integer.

A formula of ExpPA, denoted by ϕ , is of the form

$$\phi \hat{=} \mathfrak{t} \odot \mathfrak{t} \mid c \mid \mathfrak{t} \mid c \nmid \mathfrak{t} \mid \neg \phi \mid \phi \wedge \phi \mid \phi \vee \phi \mid \exists \mathbb{X}. \phi \mid \forall \mathbb{X}. \phi,$$

where $\odot \in \{=, <, >, \leq, \geq\}$.

The semantics of ExpPA are defined similarly to PA with the only difference that variables are interpreted over \mathbb{N} rather than \mathbb{Z} , that is, an interpretation is a function $I : \text{Free}(\phi) \rightarrow \mathbb{N}$. This restriction ensures that terms $10^{\mathfrak{t}}$ and $\ell_{10}(\mathfrak{t})$ are well-defined: In

the *normalization* step of the quantifier elimination procedure, for each $\ell_{10}(\mathfrak{t})$ term, we will replace all its occurrences by a fresh variable \mathbb{z} , and add a conjunct $10^{\mathbb{z}} \leq \mathfrak{t} < 10 * 10^{\mathbb{z}}$. Then $\mathbb{z} \in \mathbb{N}$ guarantees that $\mathfrak{t} \geq 1$ is always satisfied. Exponential functions are treated similarly (more details can be found in Section 4.1). Note that we can extend the interpretation to \mathbb{Z} since each integer can be written as the difference of two positive integers.

The quantifier-free and existential ExpPA formulas are defined similarly to PA as well. We also assume that quantifier-free ExpPA formulas are free of negations, for the same reason we have explained for PA.

For convenience, when working with exponential and logarithmic functions, we use the notation $\lambda_{10}(\mathfrak{t})$ to denote $10^{\ell_{10}(\mathfrak{t})}$. It is easy to show that for all $n \geq 1$, $\lambda_{10}(n) \leq n < 10\lambda_{10}(n)$ holds.

Example 2 $10^{10^{\mathbb{z}}} + 10^{\mathbb{z}+\mathbb{y}} + 3\mathbb{y} = 2\mathbb{z} + 1$ is an ExpPA formula.

3.3 Flattening STR_{parseInt} into ExpPA

We first recall the flattening approach for string constraints in [17], then show how to extend it to deal with parseInt.

A *flat domain restriction* for a string constraint $\varphi \in \text{STR}_{\text{parseInt}}$ is a function \mathcal{F}_{φ} that maps each string variable $x \in \text{SVar}(\varphi)$ to a flat language $(w_{x,1})^* \cdots (w_{x,k_x})^*$, where $w_{x,i} \in \Sigma_{\text{num}}^+$ for every $i \in [k_x]$. The flattened semantics of $\varphi \in \text{STR}_{\text{parseInt}}$ is defined as $\llbracket \varphi \rrbracket_{\mathcal{F}_{\varphi}} = \{(I_s, I_i) \in \llbracket \varphi \rrbracket \mid \forall x \in \text{SVar}(\varphi). I_s(x) \in \mathcal{F}_{\varphi}(x)\}$.

Under a flat domain restriction \mathcal{F}_{φ} , the flattening of $\varphi \in \text{STR}_{\text{parseInt}}$ is an ExpPA formula, denoted by $\text{flatten}_{\mathcal{F}_{\varphi}}(\varphi)$, that encodes the flattened semantics $\llbracket \varphi \rrbracket_{\mathcal{F}_{\varphi}}$. The idea is that x can be expressed as $(w_{x,1})^{\#_{x,1}} \cdots (w_{x,k_x})^{\#_{x,k_x}}$ with $\#_{x,i}$ denoting the number of repetitions of $w_{x,i}$, so we can eliminate the occurrences of x in φ by introduce a set of integer variables $\text{PVar}_{\mathcal{F}_{\varphi}}(x) = \{\#_{x,i} \mid i \in [k_x]\}$. More concretely, $\text{flatten}_{\mathcal{F}_{\varphi}}(\varphi)$ is a formula over the integer variables $\text{IVar}(\varphi)$, and flattening variables $\text{PVar}_{\mathcal{F}_{\varphi}}(\varphi) = \bigcup_{x \in \text{SVar}(\varphi)} \text{PVar}_{\mathcal{F}_{\varphi}}(x)$ plus some other auxiliary variables, such that the interpretations $I_e : \text{IVar}(\varphi) \cup \text{PVar}_{\mathcal{F}_{\varphi}}(\varphi) \rightarrow \mathbb{Z}$ of $\llbracket \text{flatten}_{\mathcal{F}_{\varphi}}(\varphi) \rrbracket$ and interpretations (I_s, I_i) of $\llbracket \varphi \rrbracket_{\mathcal{F}_{\varphi}}$ have the following correspondence

- for every $x \in \text{SVar}(\varphi)$, $I_s(x) = w_{x,1}^{I_e(\#_{x,1})} \cdots w_{x,k_x}^{I_e(\#_{x,k_x})}$,
- for every $\mathbb{x} \in \text{IVar}(\varphi)$, $I_i(\mathbb{x}) = I_e(\mathbb{x})$.

The formula $\text{flatten}_{\mathcal{F}_{\varphi}}(\varphi)$ is constructed inductively on the structure of φ : $\text{flatten}_{\mathcal{F}_{\varphi}}(\varphi_1 \circ \varphi_2) = \text{flatten}_{\mathcal{F}_{\varphi}}(\varphi_1) \circ \text{flatten}_{\mathcal{F}_{\varphi}}(\varphi_2)$, where $\circ \in \{\wedge, \vee\}$, and $\text{flatten}_{\mathcal{F}_{\varphi}}(\neg \varphi_1) = \neg \text{flatten}_{\mathcal{F}_{\varphi}}(\varphi_1)$. Therefore, it is sufficient to show how to construct $\text{flatten}_{\mathcal{F}_{\varphi}}(\varphi)$ for atomic constraints φ . In the sequel, we will show how to construct $\text{flatten}_{\mathcal{F}_{\varphi}}(\mathfrak{t}_1 \odot \mathfrak{t}_2)$ where $\text{parseInt}(t)$ may occur in \mathfrak{t}_1 or \mathfrak{t}_2 . The construction of $\text{flatten}_{\mathcal{F}_{\varphi}}(\varphi)$ for the other atomic constraints is omitted because it is essentially the same as that in [17] (we summarize as Theorem 1).

Theorem 1 [17] *The satisfiability of Boolean combinations of string equality constraints, regular constraints and arithmetic constraints under flat domain restrictions can be reduced to the satisfiability of existential PA formulas, thus is decidable.*

For simplicity, we assume that each occurrence of `parseInt` (resp. `len(t)`) in $\mathbb{t}_1 \odot \mathbb{t}_2$ is of the form `parseInt(x)` (resp. `len(x)`) for a string variable x . (Otherwise, we can introduce a fresh variable x' for t in `parseInt(t)` or `len(t)` and add the constraint $x' = t$.) Then $\text{flatten}_{\mathcal{F}_\varphi}(\mathbb{t}_1 \odot \mathbb{t}_2)$ is obtained from $\mathbb{t}_1 \odot \mathbb{t}_2$ by replacing `parseInt(x)` with $\text{flatten}_{\mathcal{F}_\varphi}(\text{parseInt}(x))$ and `len(x)` with $\text{flatten}_{\mathcal{F}_\varphi}(\text{len}(x))$, where

- $\text{flatten}_{\mathcal{F}_\varphi}(\text{parseInt}(x)) \hat{=} \mathbb{t}_{x,1}$ such that $(\mathbb{t}_{x,i})_{i \in [k_x]}$ are inductively defined as follows:
 - for $i = k_x$,

$$\mathbb{t}_{x,i} = \frac{\text{parseInt}(w_{x,k_x})(10^{\text{len}(w_{x,k_x})\#_{x,k_x}} - 1)}{(10^{\text{len}(w_{x,k_x})} - 1)}$$

- for $i \in [k_x - 1]$,

$$\mathbb{t}_{x,i} = \frac{\text{parseInt}(w_{x,i})(10^{\text{len}(w_{x,i})\#_{x,i}} - 1)10^{(\sum_{i+1 \leq j \leq k_x} \text{len}(w_{x,j})\#_{x,j})}}{(10^{\text{len}(w_{x,i})} - 1)} + \mathbb{t}_{x,i+1}$$

- Notice that here $\text{len}(w_{x,-})$ are constants while $\#_{x,-}$ are (flattening) variables. So we have

$$\text{flatten}_{\mathcal{F}_\varphi}(\text{len}(x)) \hat{=} \sum_{i \in [k_x]} \text{len}(w_{x,i})\#_{x,i}.$$

The following example helps to illustrate the main idea of the flattening technique.

Example 3 Suppose `parseInt(x) = 2x` is an atomic constraint and $\mathcal{F}_\varphi(x) = 1^*2^*$. Then

$$\begin{aligned} \text{flatten}_{\mathcal{F}_\varphi}(\text{parseInt}(x) = 2x) \\ &\hat{=} 1 \frac{10^{\#_{x,1}} - 1}{10 - 1} 10^{\#_{x,2}} + 2 \frac{10^{\#_{x,2}} - 1}{10 - 1} = 2x \\ &\equiv 10^{\#_{x,1} + \#_{x,2}} - 10^{\#_{x,2}} + 2 * 10^{2\#_{x,2}} - 2 = 18x \\ &\equiv 10^{\#_{x,1} + \#_{x,2}} + 10^{\#_{x,2}} = 18x + 2. \end{aligned}$$

By Theorem 1, the above reduction, and the decidability of ExpPA's satisfiability (see Theorem 3), we have

Theorem 2 *The satisfiability of $\text{STR}_{\text{parseInt}}$ under flat domain restrictions can be reduced to the satisfiability of existential ExpPA formulas, and thus is decidable.*

4 Decision procedure for ExpPA

Semënov first proved that ExpPA admits quantifier elimination in [21], thus its satisfiability problem is decidable. However, Semënov did not give a concrete quantifier elimination procedure. Remedying this, Point proposed a quantifier elimination procedure for ExpPA in [22]. In this section, we describe how Point's quantifier elimination procedure works.

Theorem 3 ([22]) *ExpPA admits quantifier elimination.*

Compared to [22], the presentation here is more accessible to computer science researchers. Moreover, the procedure presented here slightly improves Point's procedure, in the following two aspects: 1) DNF (disjunctive normal form) was required in Point's procedure, which is not required here, 2) in Point's procedure, the divisibility constraints produced by the elimination of linear occurrences of a variable should be converted to equality constraints (by introducing fresh variables) before the elimination of exponential occurrences of some other variable, which is unnecessary here, since the divisibility constraints are directly dealt with in the elimination of exponential occurrences of variables.

As $\forall \mathbb{x}. \varphi$ is equivalent to $\neg \exists \mathbb{x}. \neg \varphi$, to prove Theorem 3, we only need to show that every existentially quantified ExpPA formula $\exists \mathbb{x}. \varphi \in \text{ExpPA}$, where φ is quantifier-free, can be transformed into an equivalent quantifier-free formula $\varphi' \in \text{ExpPA}$.

An informal sketch of Point's procedure is presented in Algorithm 1.

Algorithm 1: Quantifier elimination for ExpPA

Input: ExpPA formula $\exists \mathbb{x}_1. \varphi$, where φ is quantifier-free
Output: an equivalent quantifier-free formula
 $\varphi \leftarrow \text{normalize } \varphi \text{ w.r.t } \mathbb{x}_1$ (Subsection 4.1);
 // suppose fresh variables $\mathbb{x}_2, \dots, \mathbb{x}_n$ are introduced
 $\mathcal{S}_n \leftarrow \text{enumerate linear orders among } \mathbb{x}_1, \dots, \mathbb{x}_n$ (Subsection 4.2);
For each $\sigma \in \mathcal{S}_n$ **do**
 $\varphi'_\sigma \leftarrow \exists \vec{\mathbb{x}}. \varphi' \wedge \bigwedge_{i \in [n-1]} \mathbb{x}_{\sigma(i)} \leq \mathbb{x}_{\sigma(i+1)}$;
 $i \leftarrow n$;
 for $i > 0$ **do**
 eliminate exponential occurrences of $\mathbb{x}_{\sigma(i)}$ in φ'_σ (Subsection 4.3);
 eliminate linear occurrences of $\mathbb{x}_{\sigma(i)}$ in φ'_σ (Cooper's QE algorithm);
 // variable $\mathbb{x}_{\sigma(i)}$ is removed
 $i \leftarrow i - 1$;
 end
 // all quantified variables in φ' are removed
 $\varphi''_\sigma \leftarrow \varphi'_\sigma$;
end
return $\bigvee_{\sigma \in \mathcal{S}_n} \varphi''_\sigma$;

4.1 Normalization

The normalization step comprises the following sub-steps.

1. *replace $\ell_{10}(\mathfrak{t})$ terms* Repeat the following procedure, until there are no $\ell_{10}(\mathfrak{t})$ with \mathfrak{x} occurs in \mathfrak{t} ; for each occurrence of $\ell_{10}(\mathfrak{t})$ such that \mathfrak{x} occurs in \mathfrak{t} , introduce a fresh variable, say \mathfrak{z} , and replace all occurrences of $\ell_{10}(\mathfrak{t})$ by \mathfrak{z} , moreover, add the constraint $10^{\mathfrak{z}} \leq \mathfrak{t} < 10^{\mathfrak{z}+1}$ as a conjunct. Note that if \mathfrak{t} contains no variables, then $\ell_{10}(\mathfrak{t})$ is a constant. In this case, we can also assume that \mathfrak{t} contains \mathfrak{x} and perform the same replacements, which helps in the analysis of complexity in 5. Let the resulting formula be φ'' .
2. *flatten $10^{\mathfrak{t}}$ terms* Then repeat the following procedure to φ'' , until for each occurrence of $10^{\mathfrak{t}}$ with \mathfrak{x} occurs in \mathfrak{t} , we have $\mathfrak{t} = \mathfrak{x}$: For each occurrence of the $10^{\mathfrak{t}}$ in φ'' , such that \mathfrak{t} contains \mathfrak{x} but is not equal to \mathfrak{x} , introduce a fresh variable, say \mathfrak{z} , and replace all occurrences of $10^{\mathfrak{t}}$ by $10^{\mathfrak{z}}$, moreover, add the constraint $\mathfrak{z} = \mathfrak{t}$ as a conjunct. Let φ''' denote the resulting formula.
3. *\leq transformation* Do the following replacements to φ''' , so that all the atomic formulas in φ^{\dagger} are of the form $\mathfrak{t}_1 \leq \mathfrak{t}_2$, $c|\mathfrak{t}$ or $c \nmid \mathfrak{t}$: Replace every occurrence of $\mathfrak{t}_1 \geq \mathfrak{t}_2$ with $\mathfrak{t}_2 \leq \mathfrak{t}_1$. Replace every occurrence of $\mathfrak{t}_1 < \mathfrak{t}_2$ (resp. $\mathfrak{t}_1 > \mathfrak{t}_2$) with $\mathfrak{t}_1 \leq \mathfrak{t}_2 - 1$ (resp. $\mathfrak{t}_2 \leq \mathfrak{t}_1 - 1$). Replace ever occurrence of $\mathfrak{t}_1 = \mathfrak{t}_2$ with $\mathfrak{t}_2 \leq \mathfrak{t}_1 \wedge \mathfrak{t}_1 \leq \mathfrak{t}_2$. Let φ^{\dagger} denote the resulting formula.
4. Let $\overline{\mathfrak{z}} = \mathfrak{z}_1, \dots, \mathfrak{z}_n$ be an enumeration of the freshly introduced variables. Then the result of the normalization procedure is $\exists \overline{\mathfrak{z}} \exists \mathfrak{x}. \varphi^{\dagger}$.

Intuitively, the normalization step first removes all occurrences of $\ell_{10}(\mathfrak{t})$ where \mathfrak{x} occurs in \mathfrak{t} , by encoding them with the exponential function. Moreover, for each occurrence of $10^{\mathfrak{t}}$ such that \mathfrak{x} occurs in \mathfrak{t} , it introduces a fresh variable \mathfrak{z} , replaces $10^{\mathfrak{t}}$ with $10^{\mathfrak{z}}$, and adds the equality $\mathfrak{z} = \mathfrak{t}$. All equalities and inequalities will be rewritten into the form $\mathfrak{t}_1 \leq \mathfrak{t}_2$. Finally, add quantifiers for the introduced fresh variables.

After the normalization, the resulting formula is of the following shape: 1) it contains no occurrences of $\ell_{10}(\mathfrak{t})$ such that \mathfrak{x} occurs in \mathfrak{t} , 2) it contains no occurrences of $10^{\mathfrak{t}}$ such that \mathfrak{x} occurs in \mathfrak{t} , but $\mathfrak{t} \neq \mathfrak{x}$, 3) all the atomic formulas are of the form $\mathfrak{t}_1 \leq \mathfrak{t}_2$, $c|\mathfrak{t}$ or $c \nmid \mathfrak{t}$.

4.2 Enumeration of the variable orders

Suppose $n - 1$ fresh variables are introduced in the normalization procedure, rename the original variable \mathfrak{x} and the $n - 1$ introduced variables as $\mathfrak{x}_i, 1 \leq i \leq n$. Let the output of the normalization procedure be $\exists \overline{\mathfrak{x}}. \varphi'$ with $\overline{\mathfrak{x}} = (\mathfrak{x}_1, \dots, \mathfrak{x}_n)$. We then enumerate all the linear orders of $\{\mathfrak{x}_1, \dots, \mathfrak{x}_n\}$. Each linear order can be represented by a permutation $\sigma \in \mathcal{S}_n$ (where \mathcal{S}_n is the permutation group on $[n]$), with the intention that $\mathfrak{x}_{\sigma(n)} \geq \dots \geq \mathfrak{x}_{\sigma(1)}$.

Assuming a linear order $\sigma \in \mathcal{S}_n$ of $\{\mathfrak{x}_1, \dots, \mathfrak{x}_n\}$, we then consider $\varphi'_{\sigma} = \exists \overline{\mathfrak{x}}. \varphi' \wedge \bigwedge_{i \in [n-1]} \mathfrak{x}_{\sigma(i)} \leq \mathfrak{x}_{\sigma(i+1)}$ and eliminate the quantifiers $\exists \mathfrak{x}_{\sigma(n)}, \dots, \exists \mathfrak{x}_{\sigma(1)}$,

one by one and from $\mathfrak{x}_{\sigma(n)}$ to $\mathfrak{x}_{\sigma(1)}$. To eliminate the i -th quantifier $\exists \mathfrak{x}_{\sigma(i)}$, we first eliminate all its exponential occurrences, i.e. $10^{\mathfrak{x}_{\sigma(i)}}$, using the fact that $\mathfrak{x}_{\sigma(i)}$ is the largest among the remaining quantified variables (see Lemma 1). Then linear terms

of $\mathbb{x}_{\sigma(i)}$ are eliminated further by applying Cooper's algorithm so that the quantifier $\exists \mathbb{x}_{\sigma(i)}$ can be removed. Let φ''_{σ} denote the resulting formula.

Finally, $\exists \mathbb{x}. \varphi'$ is transformed into the quantifier-free formula $\bigvee_{\sigma \in \mathcal{S}_n} \varphi''_{\sigma}$.

4.3 Elimination of exponential occurrences of variables

Let $i \in [n]$ and $\exists \mathbb{x}_{\sigma(1)} \dots \exists \mathbb{x}_{\sigma(i)}. \varphi''_{\sigma,i}(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ be the formula obtained from φ'_{σ} by eliminating the quantifiers $\exists \mathbb{x}_{\sigma(n)}, \dots, \exists \mathbb{x}_{\sigma(i+1)}$. We show how to eliminate the exponential occurrences of $\mathbb{x}_{\sigma(i)}$ in $\varphi''_{\sigma,i}$. The elimination is *local* in the sense that it is applied to the atomic formulas independently.

Recall that after normalization, the atomic formulas are of the form $\mathbb{t}_1 \leq \mathbb{t}_2$, $c \mid \mathbb{t}$ or $c \nmid \mathbb{t}$. Therefore, we can assume that the atomic formulas in $\varphi''_{\sigma,i}$ are of the form $a_i 10^{\mathbb{x}_{\sigma(i)}} + \sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^i b_k \mathbb{x}_{\sigma(k)} \leq \mathbb{t}(\vec{y})$ or $c \mid (a_i 10^{\mathbb{x}_{\sigma(i)}} + \sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^i b_k \mathbb{x}_{\sigma(k)} + \mathbb{t}(\vec{y}))$ (or \nmid), where $\mathbb{t}(\vec{y})$ collects all terms in the formula without \mathbb{x}_i .

4.3.1 Inequality atoms

In the following, we illustrate how to eliminate the exponential occurrences of $\mathbb{x}_{\sigma(i)}$ for inequality atomic formulas of the form

$$\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y}) \triangleq a_i 10^{\mathbb{x}_{\sigma(i)}} + \sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^i b_k \mathbb{x}_{\sigma(k)} \leq \mathbb{t}(\vec{y}) \quad (1)$$

where $a_i \neq 0$, $\mathbb{t}(\vec{y})$ is the sum of all other terms without $\mathbb{x}_{\sigma(i)}$.

The elimination of the exponential occurrences of $\mathbb{x}_{\sigma(i)}$ in $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ relies on the following lemma. Intuitively, The lemma states that when the left-hand-side of the inequality is dominated by the $a_i 10^{\mathbb{x}_{\sigma(i)}}$ term, if $\mathbb{x}_{\sigma(i)} \leq \alpha(\vec{y}) - 1$ or $\mathbb{x}_{\sigma(i)} \geq \alpha(\vec{y}) + 2$, it can be determined directly whether the formula is true or false.

Lemma 1 Let $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ be of form (1) with $a_i \neq 0$. Let $A \triangleq 1 + \sum_{j=1}^{i-1} |a_j|$, $B \triangleq 1 + \sum_{j=1}^i |b_j|$, $\delta \triangleq \ell_{10}(A) + 2$, and $\gamma \triangleq 2\ell_{10}(B) + 3$. When $\mathbb{x}_{\sigma(i)} \geq \gamma$ and $\mathbb{x}_{\sigma(i)} \geq \mathbb{x}_{\sigma(i-1)} + \delta$ are all satisfied,

- for $a_i > 0$, let $\alpha(\vec{y}) \triangleq \ell_{10}(\mathbb{t}(\vec{y})) - \ell_{10}(a_i)$, then
 - if $\mathbb{x}_{\sigma(i)} \leq \alpha(\vec{y}) - 1$, then $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ holds,
 - if $\mathbb{x}_{\sigma(i)} \geq \alpha(\vec{y}) + 2$, then $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ **does not** hold.
- for $a_i < 0$, let $\alpha(\vec{y}) \triangleq \ell_{10}(-\mathbb{t}(\vec{y})) - \ell_{10}(-a_i)$, then
 - if $\mathbb{x}_{\sigma(i)} \leq \alpha(\vec{y}) - 1$, then $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ **does not** hold,
 - if $\mathbb{x}_{\sigma(i)} \geq \alpha(\vec{y}) + 2$, then $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ holds.

We need the following two propositions for the proof of Lemma 1. Proposition 4 shows that when \mathbb{x} is large enough, the linear term $n\mathbb{x}$ can always be bound by the

exponential term $10^{\mathbb{x}}$. Proposition 5 can be seen as a special case of Lemma 1 where the left-hand-side of $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ has only one term $a_i 10^{\mathbb{x}_{\sigma(i)}}$.

Proposition 4 *Given $n \geq 1$, if $\mathbb{x} \geq 2\ell_{10}(n) + 1$, then $n\mathbb{x} \leq 10^{\mathbb{x}}$ holds.*

Proof When $1 \leq n \leq 9$, we have $n\mathbb{x} \leq 10\mathbb{x} \leq 10^{\mathbb{x}}$ for all $\mathbb{x} \in \mathbb{N}$.

For $n \geq 10$, that is, $\ell_{10}(n) \geq 1$, then

$$\begin{aligned} 10^{\mathbb{x}} &\geq 10^{\ell_{10}(n)+1} 10^{\mathbb{x}-\ell_{10}(n)-1} \\ &\geq 10^{\ell_{10}(n)+1} 10^{(\mathbb{x}-\ell_{10}(n)-1)} \\ &\geq 10^{\ell_{10}(n)+1} \mathbb{x} \\ &\geq n\mathbb{x} \end{aligned}$$

□

Proposition 5 *Given $a > 0$, if $\mathbb{x} \leq \ell_{10}(\mathbb{y}) - \ell_{10}(a) - 1$, then $(a + \frac{1}{2}) \cdot 10^{\mathbb{x}} \leq \mathbb{y}$; if $\mathbb{x} \geq \ell_{10}(\mathbb{y}) - \ell_{10}(a) + 2$, then $(a - \frac{1}{2}) \cdot 10^{\mathbb{x}} \geq \mathbb{y}$.*

Proof Suppose that $\mathbb{x} \leq \ell_{10}(\mathbb{y}) - \ell_{10}(a) - 1$, we have

$$(a + \frac{1}{2}) \cdot 10^{\mathbb{x}} \leq (a + \frac{1}{2}) \cdot \frac{\lambda_{10}(\mathbb{y})}{10\lambda_{10}(a)} = \frac{a + \frac{1}{2}}{10\lambda_{10}(a)} \lambda_{10}(\mathbb{y}) \leq \lambda_{10}(\mathbb{y}) \leq \mathbb{y}.$$

Now, suppose that $\mathbb{x} \geq \ell_{10}(\mathbb{y}) - \ell_{10}(a) + 2$,

$$(a - \frac{1}{2}) \cdot 10^{\mathbb{x}} \geq (a - \frac{1}{2}) \frac{100\lambda_{10}(\mathbb{y})}{\lambda_{10}(a)} = \frac{10(a - \frac{1}{2})}{\lambda_{10}(a)} 10\lambda_{10}(\mathbb{y}) \geq 10\lambda_{10}(\mathbb{y}) \geq \mathbb{y}$$

□

Proof of Lemma 1 We only prove for the case $a_i > 0$, the $a_i < 0$ case is symmetric. The proof is twofold.

1) First, we prove that if $\mathbb{x}_{\sigma(i)} \geq \gamma$ and $\mathbb{x}_{\sigma(i)} \geq \mathbb{x}_{\sigma(i-1)} + \delta$, then the left-hand-side of $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ is dominated by $a_i 10^{\mathbb{x}_{\sigma(i)}}$, i.e. ,

$$(a_i - \frac{1}{2})10^{\mathbb{x}_{\sigma(i)}} < a_i 10^{\mathbb{x}_{\sigma(i)}} + \sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^i b_k \mathbb{x}_{\sigma(k)} < (a_i + \frac{1}{2})10^{\mathbb{x}_{\sigma(i)}} \quad (2)$$

To obtain the formula above, we need to bound the absolute values $|\sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}}|$ and $|\sum_{k=1}^i b_k \mathbb{x}_{\sigma(k)}|$ by some fraction of $10^{\mathbb{x}_{\sigma(i)}}$ respectively. Here we choose the bound to be $\frac{1}{10} 10^{\mathbb{x}_{\sigma(i)}}$.

Using $\mathbb{x}_{\sigma(i)} \geq \mathbb{x}_{\sigma(i-1)} + \delta$ and the linear order $\mathbb{x}_{\sigma(i)} \geq \mathbb{x}_{\sigma(i-1)} \geq \dots \geq \mathbb{x}_{\sigma(1)}$, we can prove that

$$\begin{aligned} |\sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}}| &\leq \sum_{j=1}^{i-1} |a_j| 10^{\mathbb{x}_{\sigma(j)}} \\ &< A \cdot 10^{\mathbb{x}_{\sigma(i)} - \delta} \\ &\leq \frac{A}{100\lambda_{10}(A)} 10^{\mathbb{x}_{\sigma(i)}} \\ &< \frac{1}{10} 10^{\mathbb{x}_{\sigma(i)}} \quad (\text{by } \mathbb{x} < 10\lambda_{10}(\mathbb{x})). \end{aligned} \quad (3)$$

Similarly, using $\mathbb{x}_{\sigma(i)} \geq \gamma$, we have

$$\begin{aligned}
 \left| \sum_{k=1}^i b_k \mathbb{x}_{\sigma(k)} \right| &\leq \sum_{k=1}^i |b_k| \mathbb{x}_{\sigma(k)} \\
 &< B \cdot \mathbb{x}_{\sigma(i)} \\
 &= \frac{1}{10} \cdot 10B \cdot \mathbb{x}_{\sigma(i)} \\
 &\leq \frac{1}{10} 10^{\mathbb{x}_{\sigma(i)}} \quad (\text{by Proposition 4, set } n = 10B).
 \end{aligned} \tag{4}$$

Combining formula (3) and (4), we have $|\sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^i b_k \mathbb{x}_{\sigma(k)}| < \frac{3}{10} 10^{\mathbb{x}_{\sigma(i)}} < \frac{1}{2} 10^{\mathbb{x}_{\sigma(i)}}$, then formula (2) can be easily deduced.

2) Then, we give the sufficient conditions for $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ to hold or not by comparing $\mathfrak{t}(\vec{y})$ with $(a_i \pm \frac{1}{2})10^{\mathbb{x}_{\sigma(i)}}$. Note that when $\mathfrak{t}(\vec{y}) \leq 0$, $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ does not hold since $0 < (a_i - \frac{1}{2})10^{\mathbb{x}_{\sigma(i)}}$, so we only need to consider the $\mathfrak{t}(\vec{y}) > 0$ case.

Utilizing Proposition 5, we obtain that

$$\begin{aligned}
 \mathbb{x}_{\sigma(i)} \leq \alpha(\vec{y}) - 1 &\implies (a_i + \frac{1}{2}) \cdot 10^{\mathbb{x}_{\sigma(i)}} \leq \mathfrak{t}(\vec{y}) \\
 &\implies \tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y}) \text{ is satisfied}
 \end{aligned} \tag{5}$$

and

$$\begin{aligned}
 \mathbb{x}_{\sigma(i)} \geq \alpha(\vec{y}) + 2 &\implies (a_i - \frac{1}{2}) \cdot 10^{\mathbb{x}_{\sigma(i)}} \geq \mathfrak{t}(\vec{y}) \\
 &\implies \tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y}) \text{ is unsatisfied}
 \end{aligned} \tag{6}$$

Thus the lemma is proved. \square

For $a_i > 0$, by lemma 1, $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ (for readability, denoted by τ below) is equivalent to following formula without exponential occurrences of $\mathbb{x}_{\sigma(i)}$:

$$\begin{aligned}
 &\bigvee_{p=0}^{\gamma-1} (\mathbb{x}_{\sigma(i)} = p \wedge \tau[p/\mathbb{x}_{\sigma(i)}]) \\
 &\bigvee_{p=0}^{\delta-1} (\mathbb{x}_{\sigma(i)} = \mathbb{x}_{\sigma(i-1)} + p \wedge \tau[\mathbb{x}_{\sigma(i-1)} + p/\mathbb{x}_{\sigma(i)}]) \\
 &\bigvee \left(\mathbb{x}_{\sigma(i)} \geq \gamma \wedge \mathbb{x}_{\sigma(i)} \geq \mathbb{x}_{\sigma(i-1)} + \delta \wedge \mathbb{x}_{\sigma(i)} \leq \alpha(\vec{y}) - 1 \right) \\
 &\bigvee \bigvee_{p=0,1} \left(\mathbb{x}_{\sigma(i)} \geq \gamma \wedge \mathbb{x}_{\sigma(i)} \geq \mathbb{x}_{\sigma(i-1)} + \delta \wedge \mathbb{x}_{\sigma(i)} = \alpha(\vec{y}) + p \wedge \tau[\alpha(\vec{y}) + p/\mathbb{x}_{\sigma(i)}] \right).
 \end{aligned} \tag{7}$$

The elimination of the exponential occurrences of $\mathbb{x}_{\sigma(i)}$ for the case $a_i < 0$ is similar.

We would like to make a few remarks about Lemma 1: (1) When $i = 1$, we can just ignore the condition $\mathbb{x}_{\sigma(i)} \geq \mathbb{x}_{\sigma(i-1)} + \delta$. (2) The constant 1 in the definition of A, B is to ensure that $\ell_{10}(A), \ell_{10}(B)$ are well-defined. (3) The lemma still holds when the base of exponential function is changed to any natural number $n \geq 2$ (constants in δ, γ should be changed accordingly).

4.3.2 Divisibility atoms

For divisibility atoms, we utilize its periodic property to enumerate $\mathbb{x}_{\sigma(i)}$ within a finite range. Consider a divisibility atomic formula

$$\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y}) \hat{=} d \mid \left(a_i 10^{\mathbb{x}_{\sigma(i)}} + \sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^i b_k \mathbb{x}_{\sigma(k)} + \mathfrak{t}(\vec{y}) \right) \quad (8)$$

with $a_i \neq 0, d \geq 1$. The indivisibility case can be treated analogously.

Let $d = 2^{r_1} 5^{r_2} d_0$ such that d_0 is divisible by neither 2 nor 5. Moreover, let $r = \max(r_1, r_2)$. Then $d \mid (10^r d_0)$.

Since 10 and d_0 are relatively prime, according to Euler's theorem (cf. [24]), $10^{\phi(d_0)} \equiv 1 \pmod{d_0}$, where ϕ is the Euler function. Suppose $10^{\phi(d_0)} = kd_0 + 1$ for some $k \in \mathbb{N}$. Then for every $n \in \mathbb{N}$ with $n \geq r$,

$$\begin{aligned} 10^{n+\phi(d_0)} &\equiv 10^{n-r} 10^r (kd_0 + 1) \pmod{d} \\ &\equiv 10^{n-r} (k10^r d_0 + 10^r) \pmod{d} \\ &\equiv 10^{n-r} (0 + 10^r) \pmod{d} \\ &\equiv 10^n \pmod{d}. \end{aligned}$$

Then $\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})$ is equivalent to the following formula

$$\bigvee_{p=0}^{r-1} \tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}, \vec{y})[p/\mathbb{x}_{\sigma(i)}] \bigvee \left(\begin{array}{c} \mathbb{x}_{\sigma(i)} \geq r \wedge \\ \bigvee_{p=0}^{\phi(d_0)-1} \left(d \mid \left(a_i 10^{r+p} + \sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^i b_k \mathbb{x}_{\sigma(k)} + \mathfrak{t}(\vec{y}) \right) \right) \end{array} \right), \quad (9)$$

where the exponential occurrences of $\mathbb{x}_{\sigma(i)}$ disappear. When $\mathbb{x}_{\sigma(i)} > r$, we only need to enumerate $\mathbb{x}_{\sigma(i)}$ in one period, i.e. from r to $r + \phi(d_0) - 1$.

For a special case $d_0 = 1$, since $\phi(1) = 1$, (9) can be simplified to

$$\bigvee_{p=0}^{r-1} \tau[p/\mathbb{x}_{\sigma(i)}] \vee \left(\mathbb{x}_{\sigma(i)} \geq r \wedge d \mid \left(\sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^i b_k \mathbb{x}_{\sigma(k)} + \mathfrak{t}(\vec{y}) \right) \right).$$

5 Complexity Analysis

In this section, we analyse the complexity of Point's procedure applied on an existential quantified ExpPA formula.

Consider a formula of the form

$$\exists \mathbb{x}^1 \dots \exists \mathbb{x}^m. \varphi(\mathbb{x}^1, \dots, \mathbb{x}^m, \vec{y})$$

with free variables \vec{y} . Let N denote the length of the formula. We will prove that the space complexity to eliminate all (existential) quantifiers has a 3-EXPSPACE upper bound. Since Cooper's quantifier elimination algorithm for PA works as a sub-procedure and also has a 3-EXPSPACE upper bound, this bound may not be easily improved.

We first give a brief analysis of the normalization step, in which fresh variables are introduced and the length of the formula grows linearly at most. For the rest of the algorithm, we adopt the strategy of Oppen's analysis of Cooper's algorithm (cf. [25]). The idea is that the upper bound of the formula length can be expressed using the product of the number of atoms, the number of coefficients and the length of the maximum constant. The critical point of our analysis is that only coefficients of linear occurrences of quantified variables need to be considered.

Normalization

In the normalization step, we normalize the formula w.r.t quantified variables from \mathbb{x}^m to \mathbb{x}^1 . Suppose that fresh variables $\mathbb{x}_1^i, \dots, \mathbb{x}_{n_i}^i$ are introduced for each \mathbb{x}^i ($1 \leq i \leq m$), the formula then becomes

$$\exists \mathbb{x}^1 \exists \mathbb{x}_1^1 \dots \exists \mathbb{x}_{n_1}^1 \dots \exists \mathbb{x}^m \exists \mathbb{x}_1^m \dots \exists \mathbb{x}_{n_m}^m. \varphi'(\mathbb{x}^1, \dots, \mathbb{x}_{n_m}^m, \vec{y}).$$

Since the number of newly introduced variables is less than the number of occurrences of exponential and logarithmic functions in the original formula, we have at most $m + \sum_{i=1}^m n_i \leq N$ quantified variables after normalization.

The increase in the length of the formula during normalization comes from two sources, the conjuncts for each introduced variables and the additional formulas for translating $\neq, =$ relations to \leq . It is not hard to see that both these operations will at most increase the length of the formula by a constant factor.

In the sequel, to avoid redundant symbols, we still use m to denote the number of quantified variables and N to denote the length of the formula after normalization. Just keep in mind that $m \leq N$.

Enumeration of linear orders (the outer for-loop)

Denote the normalized formula by $\exists \vec{\mathbb{x}}. \varphi(\vec{\mathbb{x}}, \vec{y})$ with $\vec{\mathbb{x}} = (\mathbb{x}_1, \dots, \mathbb{x}_m)$. According to Point's procedure, we then eliminate the quantified variables one by one: each time select the largest variable among \mathbb{x}_i , $1 \leq i \leq m$, first eliminate the exponential occurrences, then linear occurrences. Each linear order of quantified variables is given by a permutation. For m quantified variables, there are $m!$ possible linear orders, which means the inner for-loop repeats $m!$ times.

Elimination of quantifiers (the inner for-loop)

Suppose the specified linear order is $\mathbb{x}_m \geq \dots \geq \mathbb{x}_1$. Let $\exists \mathbb{x}_1 \dots \exists \mathbb{x}_{m-k}. \varphi_k$ denote the formula obtained by eliminating quantifiers $\exists \mathbb{x}_m, \dots, \exists \mathbb{x}_{m-k+1}$ from $\exists \vec{\mathbb{x}}. \varphi(\vec{\mathbb{x}}, \vec{y})$. Let φ'_k denote the formula obtained by further eliminating exponential occurrences of \mathbb{x}_{m-k+1} (from φ_k). Let φ_0 denote φ .

Let c_k be the number of distinct divisor d in atoms of the form $d \mid \mathfrak{t}$ and $d \nmid \mathfrak{t}$ plus the number of distinct coefficients of *linear occurrences* of quantified variables in φ_k . Let s_k be the largest constant (including coefficients) and a_k be the number of atomic formulas in φ_k . Similarly we define c'_k , s'_k and a'_k for φ'_k .

First we analyse the sub-procedure to eliminate exponential occurrences of the inner most quantified variable \mathfrak{x}_m . We prove the following lemma.

Lemma 2

$$\begin{aligned} c'_0 &\leq c_0^2 \\ s'_0 &\leq ms_0^2 \\ a'_0 &\leq s_0 a_0 \end{aligned}$$

Proof The analysis is divided into two cases by assuming all atomic formulas are of the same form (inequalities atoms or divisibility atoms).

If all atoms are inequalities of the form in Lemma 1. We know that each atomic formula τ with exponential occurrence of \mathfrak{x}_m is replaced by a new formula. Note that coefficients of linear occurrences of \mathfrak{x}_i , for $i \leq m-2$, remain changed throughout the substitutions, only the coefficients of linear occurrences of \mathfrak{x}_m and \mathfrak{x}_{m-1} will be changed: constant 1 is introduced as a coefficient for \mathfrak{x}_m , and if we substitute \mathfrak{x}_m by $\mathfrak{x}_{m-1} + p$ for some constant p , coefficient of linear occurrence of \mathfrak{x}_{m-1} will become $b_m + b_{m-1}$ (b_m, b_{m-1} are coefficients for \mathfrak{x}_m and \mathfrak{x}_{m-1} in τ , see Lemma 1). Since the new coefficient is obtained by adding two linear coefficient together, we have $c'_1 \leq c_0^2$. Note that δ and B in Lemma 1 are at most $\ell_{10}(ms_0)$. When we substitute \mathfrak{x}_m by $\mathfrak{x}_{m-1} + \delta - 1$ or by $B - 1$, the largest constant in the formula becomes at most $s_0 \cdot 10^{\ell_{10}(ms_0)} \leq ms_0^2$. And an inequality is replaced by at most $4\ell_{10}(ms_0)$ atomic formulas, so $a'_1 \leq 4\ell_{10}(ms_0)a_0$.

If all atoms are divisibility atomic formulas of the form $d \mid \mathfrak{t}$ or $d \nmid \mathfrak{t}$. We have $c'_1 < 2c_0$ because a divisibility atomic formula will produce at most two forms of atomic formulas $d \mid \mathfrak{t}$ and $\phi(d) \mid \mathfrak{t}$. In a divisibility atom, any constant in the dividend \mathfrak{t} , say l , can be replaced by $(l \bmod d)$. So we have $s'_1 \leq s_0$. When d is a large prime number, $\phi(d) = d - 1$, a divisibility atomic formula is replaced by roughly d atomic formulas, so $a'_1 \leq s_0 a_0$.

Choose larger upper bounds for c'_1 , s'_1 and a'_1 respectively, then the lemma is proved. \square

When φ'_0 is transformed into φ_1 by eliminating linear occurrences of \mathfrak{x}_m , Oppen's analysis gives the following lemma.

Lemma 3 [25]

$$\begin{aligned} c_1 &\leq c_0'^4 \\ s_1 &\leq s_0'^{4c'_0} \\ a_1 &\leq a_0'^4 s_0'^{2c'_0} \end{aligned}$$

Combining Lemma 2 and Lemma 3, we have

$$c_1 \leq c_0^8$$

$$s_1 \leq (ms_0^2)^{4c_0^2}$$

$$a_1 \leq (s_0a_0)^4(m_0s_0^2)^{2c_0^2}$$

By induction on k and assuming $m \leq N$, we get

Lemma 4

$$c_k \leq c_0^{8^k}$$

$$s_k \leq N^{(4c_0)^{s^k}} s_0^{(8c_0)^{s^k}}$$

$$a_k \leq a_0^{4^k} n^{(4c_0)^{s^k}} s_0^{(8c_0)^{s^k}}$$

If we adopt the assumptions in [25], by assuming $c_0 \leq N$, $a_0 \leq N$ and $s_0 \leq N$, the space required to store the quantifier free formula φ_k , is bounded by the product of the number of linear orders $m!$, the number of atoms a_k , the maximum number of constants $2m + 2$ per atom, the maximum amount of space s_k to store each constant and some constant q . So the space complexity is $O(m! \cdot a_k \cdot (2m + 2) \cdot s_k) = O(2^{2^{2^{pn \log n}}})$, which belongs to 3-EXPSpace.

The above analysis is focused on the existential ExpPA formulas and is sufficient for our string constraints solving context. But what about general ExpPA formulas that contain alternating quantifiers? Unfortunately, directly extending our analysis on alternating quantifiers by transforming $\forall \mathbb{x}. \varphi$ into $\neg \exists \mathbb{x}. \neg \varphi$ will not give an elementary upper bound. We will explain the main idea without going into details on this issue. To eliminate quantifiers in a formula, say $\forall \mathbb{x}_2 \exists \mathbb{x}_1. \varphi(\mathbb{x}_1, \mathbb{x}_2, \vec{y})$, we first eliminate $\exists \mathbb{x}_1$, treating \mathbb{x}_2 and \vec{y} both as free variables. The subprocedure to eliminate exponential occurrences of \mathbb{x}_1 is at the cost of introducing logarithmic expressions of \mathbb{x}_2 and \vec{y} , which in the worst case increases *exponentially* the number of fresh variables in the normalization step w.r.t \mathbb{x}_2 . We conjecture that the space complexity of quantifier elimination procedure over ExpPA formulas still has an elementary upper bound, which may require a more detailed analysis or improvements of the algorithm.

6 Optimizations

In the last section we showed the complexity of Point's procedure over existential ExpPA formulas to be 3-EXPSpace, which is quite expensive and a faithful implementation would not scale³. Note that this high complexity holds even for quantifier-free ExpPA formulas: For a quantifier-free formula φ , we solve its satisfiability problem by adding the existential quantifiers for all the variables occurring in φ , then eliminate the quantifiers one by one, resulting into `true` or `false` in the end. The original formula φ is satisfiable if `true` is obtained in the end.

In this section, we focus on quantifier-free ExpPA formulas (or existential ExpPA formulas since they are satisfiability-equivalent), and present various optimizations

³We did implement Point's algorithm and discovered that the implementation could only solve formulas of very small size.

of the quantifier elimination procedure for ExpPA, aiming at an efficient implementation. The focus on quantifier-free ExpPA formulas is motivated by the following two facts: 1) the flattening of $\text{STR}_{\text{parseInt}}$ constraints results into such formulas, 2) these formulas are already challenging for state-of-the-art SMT solvers (with exponential functions defined as recursive functions).

Let φ be a quantifier-free ExpPA formula in the remainder of this section. Moreover, we assume that φ is normalized since the optimizations presented in the sequel are for normalized formulas. Furthermore, for technical convenience, we assume that all the inequality atomic formulas are of the form $\sum_{j=1}^n a_j 10^{\mathbb{x}_j} + \sum_{k=1}^n b_k \mathbb{x}_k \leq c$, where c is an integer constant. (Implicitly, we assume that there are no free variables and all the variables are existentially quantified.)

In the sequel, we will explain two major optimizations in 6.1 and 6.2. Additional optimizations are listed in 6.3.

6.1 Reduce the number of enumerated variable orders by over approximation

Recall that in Point's procedure for ExpPA, after the normalization, the variable orders are enumerated and for each order, the exponential and linear occurrences of variables are eliminated. Since the quantifier elimination is expensive and applied to each possible order of variables, if we could reduce the candidate variable orders in the very beginning, it would facilitate considerable speed-up for the decision procedure.

Our main idea is to consider an over approximation of φ , which is a PA formula φ' , and use φ' to remove the infeasible candidate variable orders.

Note that all the exponential terms in φ is of the form $10^{\mathbb{z}}$ for some integer variable \mathbb{z} . The over approximation is based on the observation that $10^n \geq 9n + 1$ for every $n \in \mathbb{N}$. Then we obtain the over approximation φ' from φ by replacing each exponential term $10^{\mathbb{z}}$ with a fresh variable \mathbb{z}' and add $\mathbb{z}' \geq 9\mathbb{z} + 1$ as a conjunct.

Then during the enumeration of the linear orders for the variables $\mathbb{x}_1, \dots, \mathbb{x}_n$, we can quickly remove those infeasible candidates σ such that $\varphi' \wedge \bigwedge_{i \in [n-1]} \mathbb{x}_{\sigma(i)} \leq$

$\mathbb{x}_{\sigma(i+1)}$ is unsatisfiable. A special case is that if φ' is unsatisfiable, then we can directly conclude that original formula φ is unsatisfiable.

6.2 Avoid the elimination of linear occurrences of variables by under approximation

The decision procedure of ExpPA in Section 4 requires the elimination of both exponential and linear occurrences of variables. Considering the fact that PA formulas can be solved efficiently by the state-of-the-art solvers, e.g. CVC4 and Z3, one natural idea is to try to only eliminate the exponential occurrences, but not the linear occurrences, of variables, and obtain the PA formulas in the end, which can then be solved by the state-of-the-art solvers.

Recall that Lemma 1 enables us to eliminate the exponential occurrences of $\mathbb{x}_{\sigma(i)}$ from an atomic formula τ of the form

$$\tau(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}) \hat{=} a_i 10^{\mathbb{x}_{\sigma(i)}} + \sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^i b_k \mathbb{x}_{\sigma(k)} \leq c. \quad (10)$$

Actually, Lemma 1 does more in the sense that all occurrences of $\mathbb{x}_{\sigma(i)}$, including the linear ones, are eliminated from the atomic formulas resulted from τ , e.g. $\tau[\mathbb{x}_{\sigma(i-1)} + p/\mathbb{x}_{\sigma(i)}]$. Then we can continue eliminating the exponential occurrences of $\mathbb{x}_{\sigma(i-1)}$ from $\tau[\mathbb{x}_{\sigma(i-1)} + p/\mathbb{x}_{\sigma(i)}]$, provided that the coefficient of $\mathbb{x}_{\sigma(i-1)}$ therein is nonzero. Iterating this process would produce a PA formula eventually.

Nevertheless, the condition of Lemma 1, namely $a_i \neq 0$, undermines the aforementioned natural idea. If $a_i = 0$, but $b_i \neq 0$, then we are unable to utilize Lemma 1 to eliminate the linear occurrences of $\mathbb{x}_{\sigma(i)}$ from τ . In this case, the quantifier elimination algorithm of PA has to be applied to eliminate $\mathbb{x}_{\sigma(i)}$, so that later on, we can eliminate the exponential occurrences of $\mathbb{x}_{\sigma(i-1)}$, which requires that $\mathbb{x}_{\sigma(i-1)}$ is the maximum variable in the left-hand side of the inequality.

To avoid applying the quantifier elimination algorithm of PA, we consider the following under approximation of τ , namely, we additionally assume that $\mathbb{x}_{\sigma(i)} \leq 10^u$ for some constant bound $u \in \mathbb{N}$. Then τ can be rewritten as

$$\tau'(\mathbb{x}_{\sigma(i)}, \dots, \mathbb{x}_{\sigma(1)}) \hat{=} \sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^{i-1} b_k \mathbb{x}_{\sigma(k)} \leq c - b_i \mathbb{x}_{\sigma(i)}.$$

Let us assume $a_{i-1} > 0$. Define c_1, c_2 as follows: If $b_i > 0$, then $c_1 = c - b_i 10^u$ and $c_2 = c$, otherwise, $c_1 = c$ and $c_2 = c - b_i 10^u$. It is easy to observe that $c_1 \leq c - b_i \mathbb{x}_{\sigma(i)} \leq c_2$. Then we can apply Lemma 1 to the following two inequalities to eliminate $10^{\mathbb{x}_{\sigma(i-1)}}$,

$$\sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^{i-1} b_k \mathbb{x}_{\sigma(k)} \leq c_1 \quad (11)$$

and

$$\sum_{j=1}^{i-1} a_j 10^{\mathbb{x}_{\sigma(j)}} + \sum_{k=1}^{i-1} b_k \mathbb{x}_{\sigma(k)} \leq c_2. \quad (12)$$

Let $\alpha_1 = \ell_{10}(c_1) - \ell_{10}(a_{i-1})$ and $\alpha_2 = \ell_{10}(c_2) - \ell_{10}(a_{i-1})$. Then from Lemma 1, let $A \hat{=} 1 + \sum_{j=1}^{i-2} |a_j|$, $B \hat{=} 1 + \sum_{j=1}^{i-1} |b_j|$, $\delta \hat{=} \ell_{10}(A) + 2$, and $\gamma \hat{=} 2\ell_{10}(B) + 3$.

- if $\mathbb{x}_{\sigma(i-1)} \geq \gamma$, $\mathbb{x}_{\sigma(i-1)} \leq \alpha_1 - 1$, and $\mathbb{x}_{\sigma(i-1)} \geq \mathbb{x}_{\sigma(i-2)} + \delta$, then inequality (11), thus also τ' , is evaluated to **true**,
- if $\mathbb{x}_{\sigma(i-1)} \geq \gamma$, $\mathbb{x}_{\sigma(i-1)} \geq \alpha_2 + 2$, and $\mathbb{x}_{\sigma(i-1)} \geq \mathbb{x}_{\sigma(i-2)} + \delta$, then inequality (12), thus also τ' , is evaluated to **false**.

Therefore, τ' and τ are equivalent to

$$\begin{aligned}
 & \bigvee_{p=0}^{\gamma-1} (\mathbb{X}_{\sigma(i-1)} = p \wedge \tau[p/\mathbb{X}_{\sigma(i-1)}]) \\
 & \vee \bigvee_{p=0}^{\delta-1} \left(\mathbb{X}_{\sigma(i-1)} = \mathbb{X}_{\sigma(i-2)} + p \wedge \tau[\mathbb{X}_{\sigma(i-2)} + p/\mathbb{X}_{\sigma(i-1)}] \wedge \mathbb{X}_{\sigma(i-1)} \leq \alpha_1 - 1 \right) \\
 & \vee \bigvee_{p=0}^{\delta-1} \left(\mathbb{X}_{\sigma(i-1)} = \mathbb{X}_{\sigma(i-2)} + p \wedge \tau[\mathbb{X}_{\sigma(i-2)} + p/\mathbb{X}_{\sigma(i-1)}] \wedge \mathbb{X}_{\sigma(i-1)} \geq \alpha_2 + 2 \right) \\
 & \vee \left(\mathbb{X}_{\sigma(i-1)} \geq \gamma \wedge \mathbb{X}_{\sigma(i-1)} \geq \mathbb{X}_{\sigma(i-2)} + \delta \wedge \mathbb{X}_{\sigma(i-1)} \leq \alpha(\vec{\gamma}) - 1 \right) \\
 & \vee \bigvee_{p=\alpha_1}^{\alpha_2+1} \left(\mathbb{X}_{\sigma(i-1)} \geq \gamma \wedge \mathbb{X}_{\sigma(i-1)} \geq \mathbb{X}_{\sigma(i-2)} + \delta \wedge \mathbb{X}_{\sigma(i-1)} = p \wedge \tau[p/\mathbb{X}_{\sigma(i-1)}] \right),
 \end{aligned} \tag{13}$$

where all exponential occurrences of $\mathbb{X}_{\sigma(i-1)}$ are eliminated.

Similarly, we can eliminate the exponential occurrences of $\mathbb{X}_{\sigma(i-2)}$ from $\tau'[\mathbb{X}_{\sigma(i-2)} + p/\mathbb{X}_{\sigma(i-1)}]$ as well as $\tau'[p/\mathbb{X}_{\sigma(i-1)}]$, and so on. Eventually, we obtain a PA formula.

6.3 Additional optimization techniques

Eliminate exponential occurrences in atomic formulas simultaneously

Although Lemma 1 is stated for a single atomic formula, the elimination of the exponential occurrences of the same variable in different atomic formulas can actually be conducted simultaneously. That is, let $\alpha_1^\tau, \alpha_2^\tau, \gamma^\tau, \delta^\tau$ be the constants as stated in the aforementioned under-approximation of an inequality τ , define $\alpha_1^{\min}, \alpha_2^{\max}, B^{\max}, \delta^{\max}$ as the minimum of α_1^τ , the maximum of α_2^τ , the maximum of B^τ , and the maximum of δ^τ respectively with τ ranging over the inequalities of φ . Then we can use the same constants $\alpha_1^{\min}, \alpha_2^{\max}, B^{\max}, \delta^{\max}$ for different inequalities when eliminating the exponential occurrences of the same variable.

Avoid the formula-size blow-up by depth-first search

The PA formula resulting from the elimination of exponential occurrences is essentially a disjunction of large number of disjuncts of small size. If we store this large formula naively, then its size quickly blows up and exhausts the memory. Alternatively, we choose to do a depth-first search (DFS) and consider the disjuncts, which are of small sizes, one by one, and solve the satisfiability problem for these disjuncts. When a satisfiable disjunct is found, then the search terminates and “SAT” is reported.

Preprocess with small upper bound

We believe that if a quantifier-free ExpPA formula is satisfiable, then more likely it may be satisfied with an assignment in which all variables are uniformly bounded.

Consequently, as a preprocessing step, we put a small upper bound, e.g. the biggest constant occurring in the formula, on the values of variables, and perform a depth-first search, so that a model can be quickly found if there is any. If this preprocessing is unsuccessful, then we continue the search with the greater upper bound 10^u for some proper predefined $u \in \mathbb{N}$.

7 Implementation and Experiments

7.1 Implementation

We implement the decision procedure in Wolfram Mathematica, called the ExpPA-solver, which is able to solve the satisfiability of ExpPA formulas.

The ExpPA-solver takes a quantifier-free ExpPA formula as input. Moreover, it allows specifying a upper bound 10^u to uniformly bound the values of variables. Given a upper bound 10^u , the problem becomes to decide whether there is an assignment in which all variable values are bounded by 10^u satisfying the given formula. Outputs of the ExpPA-solver are either “SAT”, “UNSAT”, “B-UNSAT”, or “TIME-OUT”, standing for the given formula satisfiable, unsatisfiable, unsatisfiable up to 10^u , or timeout if the solver does not terminate within the time limit. If the output is “SAT”, then a model (namely, an assignment) is returned.

7.2 Benchmarks

To evaluate the performance of the ExpPA-solver, we design two benchmark suites, ARITHMETIC and STRINGHASH⁴.

The ARITHMETIC benchmark suite

This suite comprises three groups of randomly generated ExpPA formulas. Each group is characterized by four parameters (EV , LV , EI , LI), where EV , LV represent the number of variables with exponential occurrences and with only linear occurrences respectively, and EI , LI represent the number of inequalities with exponential terms and with only linear terms respectively. We consider three parameter classes, $(2, 3, 3, 4)$, $(3, 4, 4, 5)$, and $(4, 5, 5, 6)$. Each group of the benchmark suite consists of 200 randomly generated problem instances. The coefficients of exponential terms are randomly selected from the interval $[-10^2, 10^2]$ and the other coefficients/constants are randomly selected from $[-10^5, 10^5]$. The two intervals are chosen with the intention that the coefficients of exponential terms are smaller so that they do not always dominate the left-hand side of the inequalities. Moreover, aiming at a better coverage of the syntactical ingredients of ExpPA, we randomly choose some problem instances and replace the \leq symbol of their first inequalities by $=$. The constant upper bound for the values of variables is set to be 10^{20} , as the largest 64-bit integer is less than 10^{20} . We also create an SMTLib2 file for each problem instance, to facilitate the comparison with the state-of-the-art of SMT solvers CVC4 and Z3. Because neither CVC4 nor Z3 supports the exponential functions directly, in

⁴The benchmarks are available at <https://github.com/EcstasyH/PAexp-Solver/tree/main/Benchmark>

the SMTLib2 files, we encode 10^x as a recursive function $f(x)$ defined by: $f(0) = 1$ and $f(n + 1) = 10 * f(n)$.

The STRINGHASH benchmark suite

This suite comprises two groups of string constraints generated from the string hash functions $\text{hash}(w)$ encoded by `parseInt`. We restrict the nonce strings in one group conforming some flat pattern, while for the other one, we allow any word from Σ_{num}^* to be used as the nonce.

The string constraints in the STRINGHASH benchmark suite are of the form $x \in \mathcal{A} \wedge (\text{parseInt}(x) \bmod m) \bmod m' = 0 \wedge \text{len}(x) < 100$, where \mathcal{A} is an FA, $m, m' \geq 2$. The two groups of string constraints are characterized by flat and non-flat regular constraints respectively. The flat group comprises 300 problem instances, where the flat languages are of the form $12345w_1^+w_2^+$, $12345w_1^+w_2^+6789$, or $w_1^+w_2^+6789$, with $w_1, w_2 \in \Sigma_{\text{num}}^+$, where 12345 and 6789 are the text to be protected, and $w_1^+w_2^+$ is the pattern for nonce string. The non-flat group comprises 300 problem instances, where the non-flat languages are of the form $12345\Sigma_{\text{num}}^*$, $12345\Sigma_{\text{num}}^*6789$, or $\Sigma_{\text{num}}^*6789$. Moreover, the number m is a randomly chosen prime number in the interval $[10^2, 10^5]$ and $1 \leq m' < m$ (m' is not necessarily a prime number). We generate the SMTLib2 files for these string constraints, as inputs to the string constraint solvers. On the other hand, for the ExpPA-solver, we do the following:

- For flat instances, we generate ExpPA formulas corresponding to the string constraints, as inputs to the ExpPA-solver.
- For non-flat instances, we use flat languages $a^*(b_1 \dots b_k)$ to under-approximate Σ_{num}^* , where $a, b_1, \dots, b_k \in \Sigma_{\text{num}}$. We iterate the following procedure until a model is found or timeout: Initially, set $k = 1$ and iterate by assigning $0, \dots, 9$ to a . For each assignment, we encode the resulting string constraint into an ExpPA formula with only one exponential variable. If the resulting ExpPA formula is unsatisfiable, then we increase k by 1 and repeat this process.

We would like to remark that the flattening strategy for non-flat regular constraints here is a strict generalization of that in [17]: Patterns of the form $0^*(b_1 \dots b_k)$ were considered therein and PA formulas are sufficient to encode such patterns. On the other hand, we consider patterns of the form e.g. $(a)^*(b_1 \dots b_k)$ (where $a \in \Sigma_{\text{num}}$ can be nonzero), which requires ExpPA formulas to encode in general.

7.3 Experiments

We compare the ExpPA-solver with the state-of-the-art SMT solvers on the generated benchmarks. Specifically,

- for the ARITHMETIC benchmark suite, we compare the ExpPA-solver against CVC4 (version 1.8) and Z3 (version 4.8.10),
- for the STRINGHASH benchmark suite, we compare the ExpPA-solver against CVC4, Z3, and Trau⁵.

⁵https://github.com/guluchen/z3/tree/new_trau

Table 1 Experimental Results, left: results for both experiments, right: more detailed results of STRINGHASH benchmark suite. O: Output, S: SAT, U: UNSAT, B: Bounded UNSAT, F: Fail, #: number of problems, T : average time in seconds

Group	O	Z3		CVC4		Trau		ExpPA	
		#	T	#	T	#	T	#	T
(2,3,3,4)	S	56	0.4	42	2.3	-	-	64	0.4
	U	69	0.1	72	0.1	-	-	89	0.1
	B	-	-	-	-	-	-	47	9.5
(3,4,4,5)	F	75	-	86	-	-	-	0	-
	S	33	1.4	25	2.9	-	-	52	3.3
	U	59	0.1	60	0.1	-	-	88	0.1
(4,5,5,6)	B	-	-	-	-	-	-	1	54.0
	F	108	-	115	-	-	-	59	-
	S	35	1.8	19	6.6	-	-	47	22.4
Flat	U	36	0.3	39	0.4	-	-	72	0.1
	B	-	-	-	-	-	-	0	-
	F	129	-	142	-	-	-	81	-
Non-flat	S	34	19.0	88	12.7	5	0.1	115	12.3
	U	0	-	1	4.0	182	2.5	182	47.7
	F	266	-	211	-	113	-	3	-
	S	210	7.8	144	4.9	55	5.9	300	16.7
	U	0	-	0	-	0	-	0	-
	F	90	-	156	-	245	-	0	-

Group	O	Z3		CVC4		Trau		ExpPA	
		#	T	#	T	#	T	#	T
$12345(w_1)^+(w_2)^+$	S	5	14.0	29	8.5	3	0.1	37	9.9
	U	0	-	0	-	60	1.3	60	47.2
	F	95	-	71	-	37	-	3	-
$12345(w_1)^+(w_2)^+6789$	S	11	13.0	29	12.0	0	-	37	10.6
	U	0	-	0	-	63	1.2	63	50.0
	F	89	-	71	-	37	-	0	-
$(w_1)^+(w_2)^+6789$	S	18	24.0	30	9.3	2	0.1	41	16.1
	U	0	-	1	4.0	59	2.5	59	45.8
	F	82	-	69	-	39	-	0	-
$12345\Sigma_{\text{num}}^*$	S	82	8.7	100	2.2	28	5.9	100	18.5
	U	0	-	0	-	0	-	0	-
	F	18	-	0	-	72	-	0	-
$12345\Sigma_{\text{num}}^*6789$	S	60	9.3	17	7.8	3	0.3	100	16.0
	U	0	-	0	-	0	-	0	-
	F	40	-	83	-	97	-	0	-
$\Sigma_{\text{num}}^*6789$	S	68	5.5	27	13.0	24	9.0	100	15.7
	U	0	-	0	-	0	-	0	-
	F	32	-	73	-	76	-	0	-

All the experiments are run on a lap-top with the Intel i5 1.4GHz CPU and 8GB memory. We set the time limit as 60 seconds per problem instance.

The experiment results are summarized in Table 1, where “Fail” means either timeout, unknown, or wrong answers.

For the ARITHMETIC benchmark suite, the ExpPA-solver solves around 20%-60% more instances than Z3, and 30%-100% more instances than CVC4. Moreover, the gap becomes bigger as the sizes of the formulas increase, which demonstrates that the ExpPA-solver is more efficient in solving formulas of greater sizes. The average time of the ExpPA-solver is comparable with Z3 and CVC4. The ExpPA-solver reports “B-UNSAT” for 47 instances of the (2, 3, 3, 4)-group, while it does not report “B-UNSAT” (except one) for the other two groups. If more time is allowed, the ExpPA-solver is able to report “B-UNSAT” for the “TIMEOUT” instances.

For the STRINGHASH benchmark suite, in overall, the ExpPA-solver solves significantly more instances, especially those satisfiable instances, than Z3, CVC4, and Trau. For instance, for flat regular constraints, the ExpPA-solver solves almost all 300 problem instances, except 3 of them⁶, while Z3, CVC4, Trau solve only 34, 89, 187 instances respectively. Trau gets wrong answers for some problem instances, e.g. it

⁶These three instances can actually be solved in 70 seconds.

reports “UNSAT” for some satisfiable instances. From the results, we can see that the ExpPA-solver achieves a good tradeoff between precision and efficiency, although it is slower than the other solvers.

8 Conclusion

In this paper, we proposed a complete flattening approach for string constraints with string-integer conversion and flat regular constraints, based on a quantifier elimination procedure by Point in 1986, for the extension of Presburger arithmetic with exponential functions. We gave a more accessible reformulation of Point’s procedure and for the first time a complexity upper bound. Moreover, we proposed various optimizations and achieved the first prototypical implementation of Point’s procedure. We also did extensive experiments to evaluate the performance of the implementation. The experiment results show the efficacy of our implementation, compared with the state-of-the-art solvers.

References

- [1] Makanin, G.S.: The problem of solvability of equations in a free semigroup. *Matematicheskii Sbornik* **145**(2), 147–236 (1977)
- [2] Plandowski, W.: Satisfiability of word equations with constants is in pspace. In: 40th Annual Symposium on Foundations of Computer Science (Cat. No. 99CB37039), pp. 495–500 (1999). IEEE
- [3] Barrett, C., Conway, C.L., Deters, M., Hadarean, L., Jovanović, D., King, T., Reynolds, A., Tinelli, C.: CVC4. In: International Conference on Computer Aided Verification, pp. 171–177 (2011). Springer
- [4] De Moura, L., Bjørner, N.: Z3: An efficient smt solver. In: International Conference on Tools and Algorithms for the Construction and Analysis of Systems, pp. 337–340 (2008). Springer
- [5] Zheng, Y., Zhang, X., Ganesh, V.: Z3-str: A Z3-based string solver for web application analysis. In: Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering, pp. 114–124 (2013). ACM
- [6] Trinh, M.-T., Chu, D.-H., Jaffar, J.: S3: A symbolic string solver for vulnerability detection in web applications. In: Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security, pp. 1232–1243 (2014). ACM
- [7] Abdulla, P.A., Atig, M.F., Chen, Y.-F., Holík, L., Rezine, A., Rümmer, P., Stenman, J.: Norn: An SMT solver for string constraints. In: International Conference on Computer Aided Verification, pp. 462–469 (2015). Springer

- [8] Chen, T., Chen, Y., Hague, M., Lin, A.W., Wu, Z.: What is decidable about string constraints with the replaceall function. *Proceedings of the ACM on Programming Languages* **2**(POPL), 3 (2017)
- [9] Holík, L., Janků, P., Lin, A.W., Rümmer, P., Vojnar, T.: String constraints with concatenation and transducers solved efficiently. *Proceedings of the ACM on Programming Languages* **2**(POPL), 4 (2017)
- [10] Aydin, A., Eiers, W., Bang, L., Brennan, T., Gavrilov, M., Bultan, T., Yu, F.: Parameterized model counting for string and numeric constraints. In: *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pp. 400–410 (2018). ACM
- [11] Yu, F., Alkhalaf, M., Bultan, T.: Stranger: An automata-based string analysis tool for php. In: *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pp. 154–157 (2010). Springer
- [12] Abdulla, P.A., Atig, M.F., Chen, Y.-F., Diep, B.P., Holík, L., Rezine, A., Rümmer, P.: Trau: SMT solver for string constraints. In: *2018 Formal Methods in Computer Aided Design (FMCAD)*, pp. 1–5 (2018). IEEE
- [13] Day, J.D., Ganesh, V., He, P., Manea, F., Nowotka, D.: The satisfiability of word equations: Decidable and undecidable theories. In: Potapov, I., Reynier, P. (eds.) *Reachability Problems - 12th International Conference, RP 2018, Marseille, France, September 24–26, 2018, Proceedings. Lecture Notes in Computer Science*, vol. 11123, pp. 15–29. Springer, ??? (2018). https://doi.org/10.1007/978-3-030-00250-3_2. https://doi.org/10.1007/978-3-030-00250-3_2
- [14] ECMA Script, E., Association, E.C.M., et al.: Ecma script language specification (2019). <https://www.ecma-international.org/ecma-262/>
- [15] Kiezun, A., Ganesh, V., Artzi, S., Guo, P.J., Hooimeijer, P., Ernst, M.D.: HAMPI: A solver for word equations over strings, regular expressions, and context-free grammars. *ACM Trans. Softw. Eng. Methodol.* **21**(4), 25–12528 (2012). <https://doi.org/10.1145/2377656.2377662>
- [16] Saxena, P., Akhawe, D., Hanna, S., Mao, F., McCamant, S., Song, D.: A symbolic execution framework for JavaScript. In: *31st IEEE Symposium on Security and Privacy, S&P 2010, 16–19 May 2010, Berkeley/Oakland, California, USA*, pp. 513–528. IEEE Computer Society, ??? (2010). <https://doi.org/10.1109/SP.2010.38>. <https://doi.org/10.1109/SP.2010.38>
- [17] Abdulla, P.A., Atig, M.F., Chen, Y., Diep, B.P., Dolby, J., Janku, P., Lin, H., Holík, L., Wu, W.: Efficient handling of string-number conversion. In: *Donaldson, A.F., 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, pp. 943–957. ACM, ??? (2020). <https://doi.org/10.1145/3385412.3386034>. <https://doi.org/10.1145/3385412.3386034>
- [18] Abdulla, P.A., Atig, M.F., Chen, Y., Diep, B.P., Holík, L., Rezine, A., Rümmer, P.: Trau: SMT solver for string constraints. In: Bjørner, N., Gurfinkel, A. (eds.) 2018 Formal Methods in Computer Aided Design, FMCAD 2018, Austin, TX, USA, October 30 - November 2, 2018, pp. 1–5. IEEE, ??? (2018). <https://doi.org/10.23919/FMCAD.2018.8602997>. <https://doi.org/10.23919/FMCAD.2018.8602997>
- [19] Abdulla, P.A., Atig, M.F., Chen, Y., Diep, B.P., Holík, L., Rezine, A., Rümmer, P.: Flatten and conquer: a framework for efficient analysis of string constraints. In: Cohen, A., Vechev, M.T. (eds.) Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2017, Barcelona, Spain, June 18-23, 2017, pp. 602–617. ACM, ??? (2017). <https://doi.org/10.1145/3062341.3062384>. <https://doi.org/10.1145/3062341.3062384>
- [20] Li, G., Ghosh, I.: PASS: string solving with parameterized array and interval automaton. In: Bertacco, V., Legay, A. (eds.) Hardware and Software: Verification and Testing - 9th International Haifa Verification Conference, HVC 2013, Haifa, Israel, November 5-7, 2013, Proceedings. Lecture Notes in Computer Science, vol. 8244, pp. 15–31. Springer, ??? (2013)
- [21] Semënov, A.L.: Logical theories of one-place functions on the set of natural numbers. *Mathematics of the USSR-Izvestiya* **22**(3), 587–618 (1984). <https://doi.org/10.1070/im1984v022n03abeh001456>
- [22] Point, F.: On the expansion $(\mathbb{N}, +, 2^x)$ of Presburger arithmetic. In: Proceedings of the Fourth Easter Conference on Model Theory, Gross Koris, pp. 17–34 (1986)
- [23] Cooper, D.C.: Theorem proving in arithmetic without multiplication. *Machine intelligence* **7**(91-99), 300 (1972)
- [24] G. H. Hardy, E.M.W.: *An Introduction to the Theory of Numbers* (5th Edition). Oxford University Press, ??? (1980)
- [25] Oppen, D.C.: Elementary bounds for Presburger arithmetic. In: Proceedings of the Fifth Annual ACM Symposium on Theory of Computing. STOC '73, pp. 34–37. Association for Computing Machinery, New York, NY, USA (1973)