

# A Survey on String Constraint Solving

ROBERTO AMADINI, University of Bologna

String constraint solving refers to solving combinatorial problems involving constraints over string variables. String solving approaches have become popular over the last years given the massive use of strings in different application domains like formal analysis, automated testing, database query processing, and cybersecurity.

This paper reports a comprehensive survey on string constraint solving by exploring the large number of approaches that have been proposed over the last decades to solve string constraints.

CCS Concepts: • **Computing methodologies** → **Artificial intelligence**; • **Software and its engineering** → *Constraint and logic languages*; • **Theory of computation** → *Formal languages and automata theory*.

Additional Key Words and Phrases: String Constraint Solving, Constraint Programming, Satisfiability Modulo Theories, Automata Theory, Software Analysis

## ACM Reference Format:

Roberto Amadini. 2021. A Survey on String Constraint Solving. 1, 1 (July 2021), 36 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Strings are everywhere across and beyond computer science. They are a fundamental datatype in all the modern programming languages, and operations over strings frequently occur in disparate fields such as software analysis, model checking, database applications, web security, bioinformatics and so on [5, 10, 20, 26, 35, 36, 60, 71, 80]. Some of the most common uses of strings in modern software development include input sanitization and validation, query generations for databases, automatic generation of code and data, dynamic class loading and method invocation [41].

Reasoning over strings requires handling arbitrarily complex string operations, i.e., relations defined on a number of string variables. In this paper, we refer to these string operations as *string constraints*. Typical examples of string constraints are string length, (dis-)equality, concatenation, substring, regular expression matching.

With the term “*string constraint solving*” (in short, string solving or SCS) we refer to the process of modeling, processing, and solving combinatorial problems involving string constraints. We may see SCS as a declarative paradigm which falls into the intersection between constraint solving and combinatorics on words: the user states a problem with string variables and constraints, and a suitable *string solver* seeks a solution for that problem.

Although works on the combinatorics of words were already published in the 1940s [130], the dawn of a “more general” concept of SCS for handling a variety of different string constraints dates back to the late 1980s in correspondence with the rise of *constraint programming* (CP) [136] and *constraint logic programming* (CLP) [87] paradigms. Pioneers in this

---

Author’s address: Roberto Amadini, roberto.amadini@unibo.it, University of Bologna, Bologna, Italy.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2021 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

1

Listing 1. PHP sanitization code snippet.

```

1 <?php
2   $www = $_GET["www"];
3   $_otherinfo = "URL";
4   $www = preg_replace("/[^\A-Za-z0-9 .-@:\\/]/", "", $www);
5   echo $_otherinfo . ": " . $www;
6 ?>

```

field were for example Trilogly [170], a language providing strings, integer and real constraints, and CLP( $\Sigma^*$ ) [172], an instance of the CLP scheme representing strings as regular sets.

Later in the 1990s and 2000s, string solving has sparked some interest (e.g., [47, 75, 83, 94, 96, 118, 129, 131, 148]) without however leaving a mark. It was only from the 2010s that SCS finally took hold in application domains where string processing plays a central role such as test-case generation, software verification, model checking and web security. In particular, a number of proposals based on the *satisfiability modulo theory* (SMT) [30] paradigm emerged. The increased interest in the SCS field is nowadays witnessed, e.g., by the introduction of a string track starting from the 2018 SMT competition [174], the annual challenge for SMT solvers. In 2019, the first workshop on string constraints and applications (MOSCA) has been organized [54]. In 2020, the SMT-LIB standard [29] officially introduced a theory of Unicode strings and string benchmarks are flourishing [37, 143].

A plausible reason for the growing interest in string solving might be the remarkable performance improvements that constraint solvers have achieved over the last years, both on the CP and the SMT side. It is therefore unsurprising that nowadays CP and SMT are the main paradigms for solving general SCS problems. Arguably, the widespread interest in *cybersecurity* has given strong impulse to SCS because strings are often the silent enabler of software vulnerabilities, and a bad string manipulation can have disastrous effects especially for web applications developed in languages like PHP or JavaScript.

For example, let us have a look at Listing 1 showing a snippet of PHP code taken from [41] and representing a simplified version of a sanitization program taken from the web application *MyEasyMarket* [24]. The program first assigns the query string provided by the user via the `$_GET` array to the `$www` variable (line 2). Then, after assigning the "URL" string to the `$_otherinfo` variable, in line 4 the function `preg_replace` is used to find all the substrings of `$www` matching the pattern `[^\A-Za-z0-9 .-@:\\/]` and to replace all of them with the empty string. In this way, all the matched substrings will be deleted from `$www`. After this sanitization, the concatenation of `$_otherinfo`, `": "`, and `$www` is printed.

The code in Listing 1 contains a subtle bug. Indeed, the goal of the programmer is to delete every character that is *not* in the set `{a, ..., Z, A, ..., Z, 0, ..., 9, , ., -, @, :, /}`. However, the missing escape character `'\'` before character `'-'` in the matching expression will cause the PHP engine to interpret the pattern `.-@` as the *interval* of characters between `'.'` and `'@'` according to the ASCII ordering, i.e., the set `{., /, 0, ..., 9, :, ;, <, =, >, ?, @}`. This oversight might lead to dangerous executions because the sanitization fails to remove the symbol `<`, which can be used to start a potentially malicious script. This is an example of *cross-site scripting* (XSS) vulnerability.

The automated detection of vulnerabilities like XSS or SQL injections would be very hard without suitable string solving procedures. For example, as we shall see in Section 5, one could model the semantics of the `preg_replace` function with a regular expression constraint over `$www` and add a “counterexample constraint” enforcing the occurrence of `'<'` in `$www`. If this constraint is satisfied, the PHP program may not be XSS-safe.

In this survey we classify the large number of different SCS approaches that have emerged over the last decades into three main categories:

- (i) *Automata-based approaches*: mainly relying on finite state automata to represent the domain of string variables and to handle string operations.
- (ii) *Word-based approaches*: algebraic approaches based on systems of *word equations*. They mainly use SMT solvers to tackle string constraints.
- (iii) *Unfolding-based approaches*: explicitly reducing each string into a number of contiguous elements denoting its characters.

The goal of this survey is to provide an overview of SCS that can serve as a solid base for both experienced and novice users. We provide a “high-level” perspective of string solving by focusing in particular on its technological and practical aspects, without however ignoring its theoretical foundations.

*Paper structure.* In Section 2 we give the necessary background notions. In Section 3 we discuss the theoretical foundations of string solving. In Section 4 we provide a detailed review of several SCS approaches grouped by the categories listed above. In Section 5 we focus on the technological aspects, by exploring the SCS modeling languages, benchmarks, and applications that have been developed. In Section 6 we conclude by giving some possible future directions.

## 2 BACKGROUND

In this section we define the basics of string constraint solving. We recall some preliminary notions about strings and automata theory, and we define the notions of string variables and constraints. Then, we show at a high level how SCS problems are handled by CP and SMT paradigms—nowadays the state-of-the-art technologies for string solving.

We assume that the reader is familiar with the basic concepts of first-order logic.

### 2.1 Preliminaries

Let us fix a finite *alphabet*, i.e., a set  $\Sigma = \{a_1, \dots, a_n\}$  of  $n > 1$  symbols also called *characters*. A *string* (or a *word*)  $w$  is a finite sequence of  $|w| \geq 0$  characters of  $\Sigma$ , and  $|w|$  denotes the length of  $w$  (in this work we do not consider infinite-length strings). The empty string is denoted with  $\epsilon$ . The countable set  $\Sigma^*$  of all the strings over  $\Sigma$  is inductively defined as follows: (i)  $\epsilon \in \Sigma^*$ ; (ii) if  $a \in \Sigma$  and  $w \in \Sigma^*$ , then  $wa \in \Sigma^*$ .

The string *concatenation* of  $v, w \in \Sigma^*$  is denoted by  $v \cdot w$  (or simply with  $vw$  when not ambiguous). We denote with  $w^n$  the *iterated concatenation* of  $w$  for  $n$  times, i.e.,  $w^0 = \epsilon$  and  $w^n = ww^{n-1}$  for  $n > 0$ . Analogously, we define the concatenation between *sets of strings*: given  $V, W \subseteq \Sigma^*$ , we denote with  $V \cdot W = \{vw \mid v \in V, w \in W\}$  (or simply with  $VW$ ) their concatenation and with  $W^n$  the iterated concatenation, i.e.,  $W^0 = \{\epsilon\}$  and  $W^n = WW^{n-1}$  for  $n > 0$ . We use the 1-based notation to lookup the symbols in a string:  $w[i]$  is the  $i$ -th symbol of string  $w$ , with  $1 \leq i \leq |w|$ . We use the notation  $w[i..j]$  as a shortcut for the substring  $w[i]w[i+1] \cdots w[j]$ ; if  $i > j$ , we assume  $w[i..j] = \epsilon$ .

A set of strings  $L \subseteq \Sigma^*$  is called a *formal language*. Formal languages are potentially infinite sets that can be recognised by models of computation having different expressiveness. *Finite-state automata* (FSA) are among the best known models of computation for denoting sets of strings. A FSA is a system  $M = \langle \Sigma, Q, \delta, q_0, F \rangle$  with a finite alphabet  $\Sigma$ , a finite set of states  $Q$ , a transition function  $\delta$  over  $Q \times \Sigma$  defining the *state transition* occurring when  $M$  is in a state  $q \in Q$  and a character  $w[i] \in \Sigma$  of an input string  $w \in \Sigma^*$  is read for  $i = 1, \dots, |w|$  (the first character  $w[1]$  is always

read in the initial state  $q_0$ ). A number  $k \geq 0$  of final or accepting states  $F$  determine if  $w$  belongs to the language  $\mathcal{L}(M)$  denoted by  $M$  or not. A language recognized by a FSA is called a *regular language*.

Different variants and extensions of FSA have been proposed. For example, in a *deterministic* FSA (or DFA) the state transition is always deterministic, i.e.,  $\delta : Q \times \Sigma \rightarrow Q$ . Instead, for a *non-deterministic* FSA (or NFA) the transition function is defined as  $\delta : Q \times \Sigma \rightarrow \mathcal{P}(Q)$  where  $\mathcal{P}(Q)$  is the powerset of  $Q$ . A  $\epsilon$ -NFA adds to the NFA formalism the possibility of “empty moves” (or  $\epsilon$ -moves) i.e.,  $\delta : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow \mathcal{P}(Q)$ . Other variants are the *Boolean* FSA (or BFA) where  $\delta : Q \times \Sigma \rightarrow \mathcal{B}^Q$  and  $\mathcal{B}^Q$  is the set of the  $2^{2^{|Q|}}$  Boolean functions over  $Q$ , and the *alternating* FSA (or AFA), a special case of BFA where the initial state is a projection over the states of  $Q$ . It is well-known that all the variants mentioned above have the same expressiveness, i.e., they all denote regular languages. Common extensions of FSA are instead the *push-down automata* (PDA, to recognize *context-free* languages) and the *finite-state transducers* (FST, i.e., FSA with input/output tapes defining relations between sets of strings).

## 2.2 String variables and constraints

A *string variable* is a variable that can only be assigned to a string of  $\Sigma^*$  or, equivalently, whose *domain* is a formal language of  $\Sigma^*$ . We can classify string variables into three hierarchical classes:

- *unbounded-length* variables: they can take any value in  $\Sigma^*$
- *bounded-length* variables: given an integer  $\lambda \geq 0$ , they can only take values in  $\bigcup_{i=1}^{\lambda} \Sigma^i = \{w \in \Sigma^* \mid |w| \leq \lambda\}$
- *fixed-length* variables: given an integer  $\lambda \geq 0$ , they can only take values in  $\Sigma^\lambda = \{w \in \Sigma^* \mid |w| = \lambda\}$

We call *string constraint* a relation over at least a string variable or constant. To avoid confusion, if not better specified, we will use the uppercase notation for variables and the lowercase for constants. In this paper we will only consider *quantifier-free* constraints involving strings and (possibly) integers, e.g., string length or iterated concatenation.

Table 1 provides a high-level overview of the main types of string constraints that one can find in the literature.<sup>1</sup> As we shall see in Section 4.5, a modern SCS approach should be able to properly handle most of them, especially when the domains of the variables involved are finite. If they are infinite, decidability issues may arise as we shall see Section 3.

It is noteworthy that some of the constraints in Table 1 are arbitrarily interchangeable. For example, *replace* can be rewritten in terms of *find*, *length* and *concatenation* constraints as done in [14]. However, e.g., the same constraint can also be rewritten in terms of *replaceAll* and other string constraints:

$$Y = \text{replace}(X, Q, Q') \iff Y = \text{replaceAll}(X[1..N + |Q| - 1], Q, Q') \cdot X[N + |Q|..|X|] \wedge N = \text{find}(Q, X).$$

In the above formulation we assume the character indexing starting from 1, so  $w[1]$  and  $w[|w|]$  are respectively the first and last character of a string  $w$ . This choice is also arbitrary: one SCS approach may index characters starting from zero—as in the case of SMT-LIB standard [29], while another one might prefer a more mathematical 1-based indexing notation—as in the case of MiniZinc [121]. Furthermore, there is no standard naming convention for string constraints. For instance, the  $\text{find}(X, Y)$  constraint defined in [14] is also referred as  $\text{indexOf}(Y, X)$ .

The semantics of *toNum* and *toStr* constraints, respectively handling string-number and number-string conversions, are also tricky. According to SMT-LIB specifications (see Section 5.1),  $N = \text{toNum}(X)$  returns the non-negative integer in base 10 denoted by  $X$ , if  $X \in \{0, \dots, 9\}^*$ ; otherwise  $\text{toNum}(X) = -1$ . The function  $\text{toStr}(N)$  instead returns the string representation of  $N$ , if  $N \in \mathbb{N}$ ; otherwise,  $\text{toStr}(N) = \epsilon$ .<sup>2</sup> This design choice makes totally sense, but it also prevents the conversion of negative numbers. On the other hand, if negative integers are allowed, what value should

<sup>1</sup>The count constraint is also referred as *global cardinality count* (GCC) [11].

<sup>2</sup>In the SMT-LIB 2.6 specifications, the *toNum* and *toStr* operations are respectively called `str.to_int` and `int.to_str`

Table 1. Main string constraints. Variables  $X, Y, Z, Q, Q'$  are string variables, while variables  $I, J, N$  are integer variables.

String constraint	Description
$X = Y, X \neq Y$	equality, inequality
$X < Y, X \leq Y, X \geq Y, X > Y$	lexicographic ordering constraint
$N =  X $	string length
$Z = X \cdot Y, Y = X^N$	concatenation, iterated concatenation $N$ times
$Y = X^{-1}$	string reversal
$Y = X[I..J]$	substring from index $I$ to index $J$
$N = \text{find}(X, Y)$	$N$ is the index of the first occurrence of $X$ in $Y$
$Y = \text{replace}(X, Q, Q')$	$Y$ is obtained by replacing the first occurrence of $Q$ with $Q'$ in $X$
$Y = \text{replaceAll}(X, Q, Q')$	$Y$ is obtained by replacing all the occurrences of $Q$ with $Q'$ in $X$
$N = \text{count}(a, X)$	$N$ is the number of occurrences of character $a$ in $X$
$N = \text{toNum}(X), X = \text{toStr}(N)$	$X$ is the string denoting the number $N$
$X \in \mathcal{L}(\mathcal{R})$	membership of $X$ in regular language denoted by $\mathcal{R}$
$X \in \mathcal{L}(\mathcal{G})$	membership of $X$ in context-free grammar language denoted by $\mathcal{G}$

$\text{toNum}(X)$  return when  $X$  does not denote a valid integer? A possible workaround might be to consider these functions as binary predicates  $\text{toNum}(X, N)$  and  $\text{toStr}(N, X)$  returning false if  $X \notin \{\epsilon, -\} \cdot \{0, \dots, 9\}^*$  or  $N \notin \mathbb{Z}$ ; otherwise, they have the expected semantics. In this way,  $\text{toNum}(X, N) \Leftrightarrow \text{toStr}(N, X)$ . This choice might be however too strict in contexts where it is fine converting non-numeric strings to numbers. For example, in JavaScript the conversion of a not numeric string to a number returns the string “NaN”.

Finally, note that restricting to a core language by pruning redundant constraints may be appropriate for theoretical purposes but it is often counterproductive for practical use. For example, [17] shows that having a dedicated CP propagator for `replace` is more effective than rewriting it into basic string constraints. Hence, formally defining a “good” SCS core language is not a goal of this paper because the definition of “good” is arbitrary and strongly depends on the underlying SCS solving approach and application domain. Nevertheless, we are interested in exploring which types of string constraints are handled by different SCS approaches, as we shall see in Section 4.5.

### 2.3 SCS and CP/SMT solving

The ultimate goal of string constraint solving is to determine whether or not a set of string constraints is feasible. CP and SMT are probably the two main state-of-the-art technologies for solving SCS problems. However, they tackle the (string) constraints in different ways.

**2.3.1 SCS and CP solving.** From a CP perspective, string constraint solving means solving a particular case of *constraint satisfaction problem* (CSP). Formally, a CSP is a triple  $\mathcal{P} = (\mathcal{X}, \mathcal{D}, \mathcal{C})$  where:  $\mathcal{X} = \{X_1, \dots, X_n\}$  are the *variables*;  $\mathcal{D} = \{D(X_1), \dots, D(X_n)\}$  are the *domains*, where for  $i = 1, \dots, n$  each  $D(x_i)$  is a set of values that  $x_i$  can take;  $\mathcal{C} = \{C_1, \dots, C_m\}$  are the *constraints*, i.e., relations over the variables of  $\mathcal{X}$  defining the feasible values for the variables. The goal is to find a *solution* of  $\mathcal{P}$ , which is an assignment  $\sigma : \mathcal{X} \rightarrow \bigcup \mathcal{D}$  such that  $\sigma(X_i) \in D(x_i)$  for  $i = 1, \dots, n$  and  $(\sigma(X_{i_1}), \dots, \sigma(X_{i_k})) \in C$  for each constraint  $C \in \mathcal{C}$  defined over variables  $X_{i_1}, \dots, X_{i_k} \in \mathcal{X}$ . CP solving combines two main techniques: (i) *propagation*, which works on individual constraints trying to prune the domains of the variables involved until a fixpoint is reached, and (ii) *branching*, which aims to actually find a solution via heuristic search.

Fixed an alphabet  $\Sigma$  we call a CSP with strings, or  $\Sigma$ -CSP, a CSP  $\mathcal{P} = (\mathcal{X}, \mathcal{D}, \mathcal{C})$  having  $k > 0$  string variables  $\{W_1, \dots, W_k\} \subseteq \mathcal{X}$  such that  $D(W_i) \subseteq \Sigma^*$  for  $i = 1, \dots, k$ , and a number of constraints in  $\mathcal{C}$  over such variables. To find

a solution, a CP solver can try to compile down a  $\Sigma$ -CSP into a CSP with only integer variables [11], or it can define specialised string propagators and branchers [16, 19, 147]. The latter approach has proved to be much more efficient.

For example, consider CSP  $\mathcal{P} = (\{X, Y, N\}, \{\{ab, bc, abcd\}, \Sigma^*, [1, 3]\}, \{N = |X|, X = Y^{-1}\})$  where  $X, Y$  are string variables with associated alphabet  $\Sigma = \{a, b, c, d\}$  and  $N$  is an integer variable. Propagating  $N = |X|$  will exclude the string  $abcd$  from  $D(X)$  because it has length 4, while the domain of  $N$  is the interval  $[1, 3]$ . An optimal propagator for  $X = Y^{-1}$  would narrow the domain of  $Y$  from  $\Sigma^*$  to  $\{ba, cb\}$ . Note that propagation is a compromise between effectiveness (how many values are pruned) and efficiency (the computational cost of pruning), so sometimes it makes sense to settle for efficient but sub-optimal propagators. Then,  $N = |X|$  will narrow the domain of  $N$  to singleton  $\{2\}$ , which actually means assigning the value 2 to  $N$ . At this stage, a *fixpoint* is reached, i.e., no more propagation is possible: we have to branch on  $\{X, Y\}$  to possibly find a solution. Let us suppose that the variable choice heuristics selects variable  $X$  and the value choice heuristics assigns to it the value  $ab$ ; in this case the propagator of  $Y = X^{-1}$  is able to conclude that  $Y = ba$  so a feasible solution for  $\mathcal{P}$  (not the only one) is  $\{X = ab, Y = ba, N = 2\}$ .

Virtually all the CSPs referred in the literature have *finite domains*, i.e., the cardinality of each domain of  $\mathcal{D}$  is bounded. Having finite domains guarantees the decidability of CSPs—which are usually NP-complete problems—by enumeration, but at the same time prevents the use of unbounded-length string variables for  $\Sigma$ -CSPs. As we shall see in Section 4.3, the existing CP approaches for string solving do not handle unbounded-length variables. In fact, although CP provides the  $\text{regular}(X, M)$  constraint—stating that  $X$  must belong to the language denoted by the FSA  $M$ —it is also true that the string variable (or the array of integer variables)  $X$  must have a fixed [124] or bounded [15] length. The C(L)P proposals we are aware handling unbounded-length strings via regular sets are [75, 99, 172].

**2.3.2 SCS and SMT solving.** In a nutshell, satisfiability modulo theories generalises the Boolean satisfiability problem to decide whether a formula in first-order logic is satisfiable w.r.t. some *background theory*  $T$  that fixes the interpretations of predicates and functions [30]. SMT theories can be arbitrarily enriched and also combined together, even though the latter is in general a difficult task. Over the last decades, several *decision procedures* have been developed to tackle the most disparate theories and sub-theories, including the theory of (non-)linear arithmetic, bit-vectors, floating points, arrays, difference logic, and uninterpreted functions. In particular, well-known SMT solvers like, e.g., CVC4 [105] and Z3 [58] implement the *theory of strings* often in conjunction with related theories, such as linear arithmetic for length constraints and regular expressions.

For example, the quantifier-free theory  $T_{SLIA}$  of strings (or *word equations*) and linear arithmetic deals with integers and unbounded-length strings in  $\Sigma^*$ , where  $\Sigma$  is a given alphabet. Its terms are string/integer variables/constants, concatenation and length. The formulas of  $T_{SLIA}$  are equalities between strings and linear arithmetic constraints. For example, the formula  $\phi \equiv X = ab \cdot Z \wedge |X| + |Y| \leq 5 \wedge (abcd \cdot X = Y \vee |X| > 5)$  where  $a, b, c, d \in \Sigma$  and  $X, Y, Z$  are string variables is well-formed for this theory. Unfortunately, the decidability of  $T_{SLIA}$  is still unknown [54].

As we shall see in Section 4.2, over the last years a growing number of modern SMT solvers has integrated the theory of strings. Most of them are based on the  $\text{DPLL}(T)$  [72] procedure.  $\text{DPLL}(T)$  is a general framework extending the original DPLL algorithm (tailored for SAT solving) to deal with an arbitrary theory  $T$  through the interaction between a SAT solver and a solver specific to  $T$ . In a nutshell,  $\text{DPLL}(T)$  lazily decomposes a SMT problem into a SAT formula, which is handled by a DPLL-based SAT solver which in turn interacts with a theory-specific solver for  $T$ , whose job is to check the feasibility of the model returned by the SAT solver.

As an example, let us consider the  $T_{SLIA}$  theory with the above unsatisfiable formula  $\phi$ . That formula is *abstracted* into a Boolean formula  $\phi'$ , obtained by treating the atomic formulas of  $\phi$  as Boolean variables, and handled by a SAT

solver which can return “unsatisfiable” or a satisfying assignment (this however does not imply that the overall formula is satisfiable). In the latter case, the constraints corresponding to such assignments are distributed to the different theories. For instance, if the formula  $X = ab \cdot Z \wedge |X| + |Y| \leq 5 \wedge abcd \cdot X = Y$  is returned, the constraints  $X = ab \cdot Z$  and  $abcd \cdot X = Y$  are delivered to the string solver, while  $|X| + |Y| \leq 5$  will be solved by an arithmetic solver. Now, the theory solvers can either find that the constraints are  $T_{SLIA}$ -satisfiable or return *theory lemmas* to the SAT solver. For example, the string solver might return  $\neg(X = ab \cdot Z) \vee |X| = |Z| + 2$  to the SAT solver, which will add the corresponding clause to its knowledge base. The SAT solver will then produce a new model or return “unsatisfiable”, and the process will be iteratively repeated until either (un-)satisfiability of  $\phi$  is proven or a resource limit is reached. The termination of this procedure is not guaranteed in general, even for decidable theories—often solvers use heuristics to speed-up the search.

### 3 SCS THEORETICAL FOUNDATIONS

String constraint solving lays its theoretical foundations in the theory of automata [53, 74, 85, 138, 171] and combinatorics on words [31, 45, 110–112]. In particular, the concept of *word equation* is central and constitutes the starting block for a number of string constraints. We can define a word equation as a particular string constraint of the form  $L = R$  with  $L, R \in (\Sigma \cup \mathcal{V})^*$  where  $\Sigma$  is an alphabet and  $\mathcal{V}$  is a set of string variables. Algebraically speaking, the structure  $(\Sigma^*, \cdot, \epsilon)$  is a *free monoid* on  $\Sigma$ . The *free semigroup* on  $\Sigma$  is instead the semigroup  $(\Sigma^+, \cdot)$  where  $\Sigma^+ = \Sigma^* - \{\epsilon\}$ . There is therefore a number of theoretical results for free semigroups that can be applied for string solving [59, 116, 166].

In 1946, Quine [130] proved that the first-order theory of word equations (i.e., word equations with Boolean connectives and quantification over variables) is *undecidable* by proving its equivalence with the first-order theory of arithmetic, which is known to be undecidable. However, in 1977 Makanin [116] set a milestone in the history of SCS by proving that the satisfiability problem for the *quantifier-free* theory of word equations is decidable. More than twenty years later, Plandowski—who refers to Makanin’s decision procedure as “*one of the most complicated termination proofs existing in the literature*” [125]—showed that this problem is in PSPACE. Inspired by the work by Plandowski, Jez [88] used a technique called *local recompression* for reducing the space consumption to  $O(n \log n)$ , where  $n$  is the size of the input equation, i.e., the sum of the character count of both sides of the equation.

Makanin’s, Plandowski’s, Jez’s and related techniques can be used to decide the satisfiability of word equations, but their computational complexity is often not suitable for string solvers. Conversely, a main theoretical result used by many word-based string solvers (see Section 4.2) is the *Levi’s lemma* [103], sometimes called *Nielsen’s transformation* by analogy with the Nielsen transformation for groups [62], stating that for all strings  $u, v, x, y \in \Sigma^*$  if  $uw = xy$  then there exists a string  $w \in \Sigma^*$  such that:

$$\begin{aligned} uw = x & \quad \text{if } |u| \leq |x| \\ vw = y & \quad \text{if } |u| \geq |x| \end{aligned}$$

The Levi’s lemma is the main paradigm underlying the Makanin’s method, and it is used by most SMT solvers for solving word equations. It is particularly useful for *quadratic* word equations [135], where each variable may occur at most twice. The decidability in this case can be assessed by incrementally building a *finite* graph of substitutions obtained by repeatedly applying the lemma to reduce the size of the equation. Each node of the graph is a word equation  $L = R$  with  $L, R \in (\Sigma \cup \mathcal{V})^*$ , and edges are possibly added after comparing the prefix of  $L$  against the prefix of  $R$ . For example, a node of the form  $XbaXYZ = YaZbaa$  defines three outgoing edges: one for the case when  $X = Y$ , one for the case where  $X$  prefix of  $Y$ , and one for  $Y$  prefix of  $X$ . In the first case, the next node will be  $baXXZ = aZbaa$ ; in

the other cases, it will be respectively  $baXXYZ = YaZbaa$  and  $XbaYXYZ = aZbaa$ . The key property, which ensures the termination of the algorithm, is that for each edge  $(L = R, L' = R')$  we have  $|L' = R'| \leq |L = R|$ , where  $|\cdot|$  is the number of variables of the equation. At the end, the original equation is satisfiable if and only if from its node we can reach a leaf of the form  $\epsilon = \epsilon$ . This procedure has a polynomial space complexity but, as proven in [135] by reduction from 3-SAT, solving quadratic word equations is in general NP-hard even when a single equation is involved.

Clearly, using word equations only is limiting because this formalism cannot capture the disparate string constraints arising from real-world applications (see Table 1), and many languages cannot be expressed using word equations alone. Because of their widespread use in modern programs, two string constraints in particular can be considered as important as word equations: the string length and the regular expression membership.

Dealing with string length is tricky because this constraint bridges the strings' and the integers' domains. A reasonable approach for handling length constraints is to map them to the *Presburger arithmetic* [77], which is known to be decidable. Unfortunately, reducing to Presburger arithmetic is not always possible. For example, consider the equation  $XabY = YabX$  with  $X, Y$  variables and  $a, b$  characters. The set of pairs of integers:

$$\{(|\sigma(X)|, |\sigma(Y)|) \mid \sigma \text{ is a solution of } XabY = YabX\}$$

is not definable in Presburger arithmetic because, as proven in [108], this set actually corresponds to the set of pairs  $(n, m) \in \mathbb{N} \times \mathbb{N}$  such that:

$$n = m \vee (n = 0 \wedge \text{Even}(m)) \vee (m = 0 \wedge \text{Even}(n)) \vee (n, m > 0 \wedge \text{GCD}(n + 2, m + 2) > 1)$$

where GCD is the greatest common divisor, which is not Presburger-definable.

At present, the decidability of the theory of word equations and arithmetic over length functions is still a major *open problem* for string constraint solving. However, decidability results have been proved for *fragments* of this theory. For example, in [108] Lin et al. defined a class of quadratic word equations with length and regular constraints for which satisfiability is decidable. Their approach extends the Levi-based construction mentioned above with counters addressing the length of each variable. The decidability of the fragment follows from the decidability of the existential theory of Presburger arithmetic with divisibility [77].

In [69] Ganesh et al. defined a *solved form* notion for word equations such that a word equation  $\omega$  has a solved form if there is a finite set  $S$  of formulae logically equivalent to  $\omega$  and: (i) each formula in  $S$  is of the form  $X = t$ , with  $X$  variable and  $t$  finite concatenation of constants; (ii) each variable in  $\omega$  occurs exactly once on the left hand side of an equation in  $S$  and never on the right hand side of an equation in  $S$ . For example, the equations  $Xa = aY \wedge Ya = Xa$  can be rewritten in solved form as  $X = a^i \wedge Y = a^i$  with  $i \geq 0$ . The solved form guarantees that the implied length constraints are linear inequalities, hence their satisfiability problem is decidable. Although not all word equations can be rewritten in solved form, the authors claimed that most word equations encountered in practice, within the context of software testing and verification, are either in solved form or can be converted into one. The authors also constructively show that the satisfiability problem for word equations, string length, and regular expressions is decidable provided that the given word equations have a solved form consisting of regular expressions without unfixed parts.<sup>3</sup>

The works in [32, 56, 67] extend [69] by considering the (un-)decidability of many fragments of the first order theory of word equations and their extensions. For example, the authors prove the undecidability of the theory of word equations with string length, linear arithmetic, and string-number conversion. In [32] Berzish et al. proved that the quantifier-free first-order theory  $T_{LRE,n,c}$  of linear integer arithmetic over string length function ( $L$ ), regex ( $RE$ )

<sup>3</sup>A formula with unfixed parts, also called a *unifier*, is a formula with free variables representing an infinite family of solutions [126].

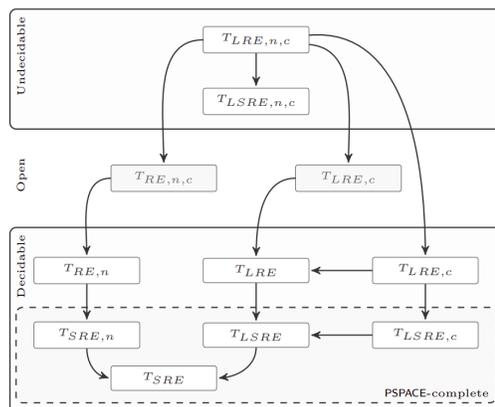


Fig. 1. Decidability of various various fragments of  $T_{LRE}$ .  $SRE$  denotes *simple regular expressions*, i.e., regular expressions where the complement operation is not allowed. Arrows model the relation “is more expressive than”. This picture is taken from [32].

membership predicates, string-number conversion ( $n$ ), and string concatenation ( $c$ ) in undecidable. Note that all of these theories allow equalities between integers but not between strings terms; hence,  $T_{LRE,n,c}$  is not able to express general word equations. The authors also show that several fragments of  $T_{LRE,n,c}$  are decidable and some of them are in PSPACE, as summarized in Fig. 1.<sup>4</sup> Interestingly, the decidability of  $T_{RE,n,c}$  and  $T_{LRE,n}$  is not settled.

In [5], Abdulla et al. used the notion of *acyclic form* to decide the satisfiability of a formula in a fragment of the theory of word equations with arithmetic and regular constraints. Intuitively, acyclicity is a syntactic property guaranteeing the absence of recursive dependencies between variables. In general, acyclicity conditions for a formula  $\varphi$  can be defined starting from the undirected graph  $G(\varphi)$  whose nodes are the variables in  $\varphi$  and there is an edge  $\{X, Y\}$  if there is a constraint in  $\varphi$  involving both  $X$  and  $Y$ . The graph  $G(\varphi)$  is acyclic if and only if  $\varphi$  is in acyclic form. For example, the formula  $Y = aX \wedge Z = Xb \wedge Y = Z$  is not acyclic because the corresponding graph  $(\{X, Y, Z\}, \{\{X, Y\}, \{X, Z\}, \{Y, Z\}\})$  contains a loop. In particular, [5] defines an acyclic form preventing variables to appear twice in (dis-)equations. The authors then proved the completeness of their decision procedure for formulae in this acyclic form.

Other interesting string constraints frequently occurring in problems derived from program analysis are *replace* and its generalization *replaceAll* (see Table 1). These operations are often used, e.g., for sanitizing or transforming a given input string. In [107] the authors show that any non-trivial theory of strings containing the *replaceAll* operation is undecidable unless some kind of *straight-line* restriction is imposed on the formulae. The undecidability follows by considering *replaceAll* constraints as particular cases of finite-state transducers and by reduction from the *Post correspondence problem* [127]. The *straight-line* fragment is based on *straight-line* string constraints, which intuitively correspond to sequences of assignments in *static single assignment* (SSA) form possibly interleaved with assertions of regular properties. More formally, a relational constraint  $\varphi$  is said to be *straight-line* if it can be rewritten into conjunctions  $\bigwedge_{i=1}^m X_i = P_i$  such that  $X_1, \dots, X_m$  are different variables and each predicate  $P_i$  contains either variables in  $\varphi$  or in  $\{X_1, \dots, X_{i-1}\}$ . A string constraint is *straight-line* if it is a conjunction of a *straight-line* relational constraint with a regular constraint.

<sup>4</sup>Not all the results in Fig. 1 were proven in [32]. For example, the decidability of the theory  $T_{LRE,c}$  was already proven in [106].

The straight-line fragment is decidable (the complexity ranges from PSPACE to EXPSPACE) and reasonable when handling constraints generated by (dynamic) symbolic execution. For example, it allows one to express constraints required for analyzing mutated XSS (mXSS) in web applications. Also, this fragment remains decidable in the presence of length, letter-counting, regular, indexOf, and disequality constraints. In [43] Chen et al. also provided a systematic study of the straight-line fragment with replace and regular membership as basic operations.

#### 4 SCS APPROACHES

In this section we provide a comprehensive overview of the main SCS approaches reported in the literature. The classification that we provide into automata-based (Section 4.1), word-based (Section 4.2) and unfolding-based (Section 4.3) approaches follows the categorization provided in [13, 14, 18]. However, similar—if not equivalent—classifications of SCS approaches are reported in many other papers.

For example, back in 2013, [178] presented one of the first approaches based on word-equations (i.e., a word-based approach). The authors stated that: “*Based on the underlining representation, existing string analyses can be roughly categorized into two kinds: automata-based [...] and bit-vector-based [...]*”. We use the term “unfolding-based” to generalize the category of bit-vector based approaches. In an unfolding-based approach, a string can be unfolded into an ordered sequence of elements having arbitrary type (e.g., integers). The unfolding-based definition allows us to capture SCS approaches like, e.g., PASS [104] which models strings with parametric arrays of symbolic length. The same categorization is repeated two years later in [177]. Importantly, in this paper the term “word-based” is introduced for the first time to denote SCS approaches handling word equations without abstractions or representation conversions. In this paper, as also done in [13, 14, 18], we stick to this definition.

Almost at the same time of [178], two other word-based approaches were presented [6, 105]. In [105], the authors argued that (at that time, 2014) most string solvers were “*standalone tools [...] based on reductions to satisfiability problems over other data types [...] or to automata*” [105]. The approaches based on reductions correspond to what we call unfolding-based problems, to emphasize the “expansion” into sequences of elements representing in some way the characters of the string. The same classification of [178] is reported in [6].

More recently, in [173], the authors reported that “*string analysis methods are mainly automata-based or satisfiability-based*”. Their definition of satisfiability-based method includes both bit-vector and SMT-based approaches. SMT-based methods are roughly the word-based methods reported in Section 4.2. However, we argue that our distinction is more general because we decouple a given SCS approach from the solving paradigm used to implement it.

We underline that the classification we propose is not a formal ontology, but the result of a synthesis among the classifications found in the literature. For this reason, as we shall see in Section 4.4, “gray areas” exist corresponding to the hybrid methodologies integrating different SCS approaches. Finally, in Sect. 4.5 we show how the most recent SCS approaches compare to each other in terms of expressiveness and performance.

##### 4.1 Automata-based approaches

In general, the domain of a string variable is a formal language, i.e., a potentially infinite (yet countable) set. A natural way to denote these sets is through (extensions of) finite-state automata. It is therefore unsurprising that a large number of string solving approaches are based on FSA, possibly enriched with other data structures.

We can say that a SCS approach is automata-based if the solver primarily relies on automata, i.e., string variables are principally represented by automata and string constraints are mainly mapped to corresponding automata operations. As reported in [105], generally speaking there are two sorts of automata-based SCS approaches: the ones where each

transition in the automaton represents a single character (e.g., [65, 175]), and the “symbolic” ones, where each transition represents a set of characters (e.g., [84, 167]). In the following, we provide a number of approaches—not necessarily string solvers in the strict sense—that may be classified as automata-based.

CLP( $\Sigma^*$ ) was one of the first attempts to incorporate strings in the CLP framework to strengthen the standard string-handling operations such as concatenation and substrings [172]. This approach was further developed by Golden et al. [75] about 15 years later. Their main contribution is to use FSA to represent and manipulate regular sets. The flexibility and expressiveness of using FSA however comes at a price: all the operations discussed by the authors (operations on FSA or string constraints like concatenation, containment, length) are linear or quadratic in the size of the FSA representing the string domain.

More recently, in [99] Krings et al. implemented ConString, a CLP system over strings relying on constraint handling rules (CHR) [64] in order to generate test data. The domains of the string variables of ConString are represented with FSA as done in [75]. ConString is built on top of SWI-Prolog, which does not handle natively FSA. For this reason, FSA are encoded as quaternary Prolog terms `automaton_dom/4` whose arguments are the set of states, the transition relation, the set of initial and final states. However, the authors admit that this representation has several efficiency drawbacks.

In [78] Hansen et al. proposed an approach to solve  $\Sigma$ -CSPs where string constraints are regular expression memberships. The algorithm uses a *Multi-DFA* (MDFA) for joining  $n$  different DFAs without reducing to a single DFA, which size would be too big for practical use. Each MDFA has the form  $M = \langle Q, \Sigma, \delta, q_0, \phi \rangle$  where instead of the set of final states  $F \subseteq Q$  there is a function  $\phi : Q \rightarrow \{0, 1\}^n$  denoting if a state  $q$  is final or not in the  $i$ -th DFA for  $i = 1, \dots, n$ . The authors extend the work of [75] by also using *binary decision diagrams* (BDDs [89]) to handle the interactive configuration of variables—not necessarily string variables—having finite domain.

The `regular(A, M)` global constraint proposed in [124] for solving CP problems treats a fixed-size array  $A$  of integer variables as a *fixed-length* string belonging to the regular language denoted by a given (non-)deterministic FSA  $M$ . This constraint was introduced to solve finite-domains CP problems like rostering and car sequencing, and not targeted to string solving—in fact, it is a useful constraint that has been used in many different CP applications, but its effectiveness diminishes in proportion to the length of  $A$ . The natural extension of `regular` is the grammar constraint [90, 129], where instead of a FSA we have a context-free grammar. However, despite more expressive, the grammar constraint never reached the popularity of `regular` in the CP community.

An interesting paper about automata-based approaches is [82], where Hooimeijer et al. study a comprehensive set of algorithms and data structures for automata operations in order to give a fair comparison between different automata-based SCS frameworks [47, 83, 96, 104, 118, 169, 175]. According to their experiments, the best results were achieved when using BDDs in combination with lazy versions of automata intersection and difference. BDDs are used to represent UTF-16 character sets: each bit of a character representation corresponds to a variable of the BDD. “Lazy version” of language intersection and difference means that the full automaton is not built: disjointness checking  $L(A) \cap L(B) = \emptyset$  and subset checking  $L(A) \subseteq L(B)$  are respectively performed instead.

MONA [96] is a tool developed in the ’90s that acts as a decision procedure for *Monadic Second-Order Logic* (M2L) on finite strings, and as a translator to FSA based on BDDs. M2L on finite strings, or M2L(Str), joins together (subsets of) *positions* in a string, indexed starting from 0, quantification and Boolean connectives. The language  $L(\phi)$  denoted by a M2L(Str) formula  $\phi$  is always regular, so each M2L(Str) formula  $\phi$  can be translated into an equivalent FSA  $M$  such that  $L(\phi) = L(M)$ . The key idea of MONA is to *extend* the original alphabet  $\Sigma$  with bit vectors for encoding the position information. The transition function of the resulting FSA is represented via BDDs to overcome the exponential explosion

of the extended alphabet.<sup>5</sup> Note that MONA is not a string solver in the strict sense, but the FSA transformation makes it possible to solve string constraints as a “side effect”.

FIDO [97] is a domain-specific programming formalism built on top of MONA, designed to get a more high-level and succinct syntax. Its variables can have four different types: pos (positions), set (set of positions), dom (finite domain) and tree ( $k$ -ary trees). FIDO first compiles the source formula into pure M2L, and then transforms it into a FSA through the MONA tool. The FIDO compiler does optimizations at many levels, in most cases relying on the type structure, to detect simple tautologies and eliminate redundant variables and quantifiers.

Another M2L-based approach is PISA [154], a path- and index-sensitive tool for static strings analysis. PISA can encode a string-manipulating method in *static single assignment* (SSA) form into an M2L formula  $\phi$  representing all the possible strings returned by that method. Then, if  $\phi'$  is a M2L formula denoting a set of *unsafe* strings, it checks  $\phi \wedge \phi'$  to verify if the method may return potentially dangerous strings. PISA handles index-sensitive operations (e.g., `indexOf`) as well as string replacement operations. It also employs a simple form of path sensitivity, by encoding the branch conditions for a specific variable into M2L predicates. PISA was integrated into the taint analysis algorithm used by the IBM Rational AppScan [86].

JSA [47] is a framework for the static analysis of Java programs. In order to soundly reason about string values, it builds a flow graph from Java class files (*front-end*) and then it derives a FSA from such graph (*back-end*). The FSA is derived by first constructing a context-free grammar from the flow graphs with a non-terminal symbol for each node. This grammar is then over-approximated into a regular grammar by using a variant of the *Mohri-Nederhof* algorithm, originally targeted for speech recognition [120]. The resulting regular grammar is finally transformed into a well-founded directed acyclic graph of NFA called *multi-level automata* (MLFA). The MLFA allows the efficient extraction of a minimal DFA for particular string expressions of interest called *hotspots*.

Inspired by JSA, in [118] Minamide developed a string analyzer for the PHP scripting language to detect cross-site software vulnerabilities in a server-side program, and to validate web pages dynamically-generated by the program. To statically check the generated pages, the string output of a program is approximated via *transducers* by a context-free grammar denoting (a superset of) all the possible output strings possibly generated by the program. Unlike JSA, this tool does not convert the context free grammar into a regular one.

Stranger [175] is another tool for the static analysis of string-related security vulnerabilities in PHP applications. As JSA [47] and the Minamide’s tool [118], Stranger is not a string solver in the strict sense but it uses string solving techniques to compute the possible values that string expressions can take during program execution. In particular, Stranger implements an automaton-based approach based on symbolic string analysis. It encodes the set of values that string variables can take with a DFA and implements string manipulation functions with a symbolic automata-based representation provided by the MONA automata package. This symbolic encoding also enables Stranger to deal with large alphabets.

Rex [169] is a tool based on the Z3 solver [58] for symbolically expressing and analyzing regular expression constraints. It relies on *symbolic finite-state automata* (SFA) where moves are labeled by formulas representing *sets* of characters instead of individual characters. SFAs are then translated into axioms describing the acceptance conditions, and Z3 is used for their satisfiability checking and to produce witnesses (i.e., models) for satisfiable formulas.

A more recent string analysis approach is SLOG [173]. It is based on a NFA manipulation engine with logic circuit representation to support string and automata operations. The idea is to avoid determinization as much as possible by

<sup>5</sup>As also mentioned in [84], the BDD representation by MONA is different from that used in [82].

performing automata manipulations implicitly via logic circuits inspired by those used for industrial applications in electronic design automation. SLOG also supports symbolic automata and enables the generation of counterexamples. The scalability of SLOG is shown for automata with large alphabets in contrast to BDD-based representations.

DPRLE [83] is a decision procedure for solving systems of equations over regular language variables. The goal here is not to assign satisfying literals to string variables, but instead assigning satisfying *regular languages* to the corresponding variables. These variables are involved in constraints of the form  $e \subseteq L$  where  $L$  is a regular language and  $e$  is an expression defined by the concatenation of variables and constant languages. The authors tackle what they call the *regular matching assignments* (RMA) problem and a subclass, the *concatenation-intersection* (CI) problem. For the RMA problem, DPRLE looks for an assignment that is also *maximal*: for each satisfying variable assignment of the form  $X \leftarrow L$ , the assignment  $X \leftarrow L'$  with  $L \subset L'$  is unsatisfiable. The CI problem is a RMA problem allowing the concatenation of at most two variables, each of which has a subset constraint:  $X \subseteq L \wedge X' \subseteq L' \wedge X \cdot Y \subseteq L''$ .

StrSolve [84] is a decision procedure supporting similar operations to those allowed by DPRLE. However, StrSolve produces single witnesses rather than atomically generating entire sets of satisfying assignments. In other words, StrSolve assigns satisfying literals to string variables, rather than assigning regular languages to corresponding variables. Hence the constraints have the form  $e \in L$  instead of  $e \subseteq L$ . This can speed-up the procedure, although the worst-case performance of StrSolve corresponds to that of DPRLE. The decision procedure of StrSolve is lazy in the sense that string constraints are processed without requiring *a priori* length bounds. This is based on the observation that eager encoding work (i.e., an eager unfolding) is unnecessary if the goal is to find a single solution as quickly as possible.

SUSHI [65] is a string solver based on the *Simple Linear String Equation* (SISE) formalism [66] to represent path conditions and attack patterns. The application domain is the security analysis of web applications. A SISE is a string equation of the form  $L \equiv R$  where  $R$  is a regular expression and  $L$  a string expression obtainable by concatenation, substring, or substitution of a pattern with a string literal. Each string variable can occur at most once in  $L$ . SISE uses an automata-based approach to recursively construct “solution pools” from which “concrete solutions” are derived. As for DPRLE, a concrete solution is a consistent assignment of regular expressions to corresponding variables. To model the semantics of regular substitution, *finite state transducers* (FST) are used.

Luu et al. developed SMC [113], a *model counter* for determining the number of solutions of combinatorial problems involving string variables having unbounded length. Their work is based on *generating functions*, a mathematical tool for reasoning about infinite series that also provides a mechanism to handle the cardinality bounds of string sets. SMC uses FSA to handle conjunctions of regular membership constraints via automata product. SMC can analyze string operators for C and JavaScript programs.

An Automata-Based model Counter for string constraints (ABC) inspired by SMC is implemented in [21]. The ABC solver first constructs an automaton accepting all the feasible solutions of a string constraint, and then counts the total number of solutions within a given length bound. Model counting corresponds to counting the accepting paths of the resulting automaton up to a given length bound  $k$ . This is however performed by producing a generating function with the *transfer matrix method* [63]. In [22] ABC is extended with relational and numeric constraints.

An automata-based approach called Qzy is presented in [52], where *Boolean finite automata* (BFA) are used to represent regular expression operations in SMT problems derived from program analysis, in order to avoid the overhead of NFA operations like determinization, complementation and intersection. The bottleneck here is the emptiness testing: for a NFA is linear, while for a BFA is PSPACE-complete. However, in their experiments the authors empirically verified that the IC3 hardware model checking algorithm [39] is effective to decide the BFA emptiness and, by extension, the satisfiability of SCS problems containing regular expressions.

In [179] Zhu et al. proposed a SCS procedure where atomic string constraints are represented by *streaming string transducers* (SSTs) [7]. In particular, the authors use bounded SSTs to prove that an input straight-line string constraint is satisfiable if and only if the domain of the corresponding SST is not empty. This approach can also solve length constraints by deriving, and solving via a SMT solver, integer constraints from the *Parikh image* of the SST. This is computed by considering the Parikh vector  $\Psi(w) = (c_1, \dots, c_n)$  for each valid output word  $w$  of the SST, where  $c_i$  is the number of occurrences of the  $i$ -th symbol of the output alphabet in  $w$ .

Trau [4] is a SMT string solver based on the flattening technique introduced in [3], where *flat automata* are used to capture simple patterns of the form  $w_1^* w_2^* \dots w_n^*$  where  $w_1, w_2, \dots, w_n$  are finite words. Trau relies on a Counter-Example Guided Abstraction Refinement (CEGAR) framework [48] where an under- and an over-approximation module interact to increase the string solving precision. In addition, Trau implements string transduction by reduction to context-free membership constraints.

Trau has been extended by Z3-Trau [2] to take advantage of the capabilities of the underlying Z3 solver. The main difference between the original Trau and Z3-Trau is that, to handle string-number constraints more efficiently, the latter works fully symbolically with *parametric flat automata* (PFA), i.e., flat automata of the form  $(A, \psi)$  where:  $A$  is an automaton having an alphabet  $V$  of *variables* over an input alphabet  $\Sigma \subseteq \mathbb{N}$ , and  $\psi$  is a linear formula over  $V$ . The parametric words  $v_1 \dots v_k \in V^*$  accepted by a PFA define its semantics, i.e., the set of all the strings  $I(v_1) \dots I(v_k) \in \Sigma^*$  where  $I$  is any interpretation  $I : V \rightarrow \Sigma \cup \{\epsilon\}$  satisfying  $\psi$ .

Sloth [81] is based on the reduction of the satisfiability problem for formulae in the straight-line and the acyclic fragment (see Section 3) to the emptiness problem for *alternating finite-state automata* (AFA). The emptiness problem can be decided by model checking algorithms (e.g., the IC3 algorithm also used by Qzy [52]). In this way Sloth can handle string constraints with concatenation, finite-state transducers (hence, also `replaceAll`), and regular constraints. Sloth was the first solver handling complex string constraints derived from HTML5 applications with sanitisation and implicit browser transductions. Sloth is built on the top of SMT solver Princess [137].

OSTRICH [44] is a string solver implementing a decision procedure for the *path feasibility* of *bounded* programs (i.e., programs without with neither loops nor branching) possibly containing string operations like concatenation, reverse, functional transducers, and `replaceAll`. OSTRICH can be seen as an extension of Sloth, with the main difference that OSTRICH also supports variables in both argument positions of `replaceAll`, while SLOTH only accepts constant strings as the second argument. The empirical evaluation provided by the authors shows that OSTRICH outperforms Sloth over all the benchmarks. Like Sloth, OSTRICH also extends the SMT solver Princess [137].

It is worth noting that approaches like, e.g., Sloth, OSTRICH or (Z3-)Trau are not fully automata-based but rather lie in the intersection between word-based and automata-based methods (see Section 4.2). It would therefore be reasonable to consider them also as word-based methods. We shall further discuss these aspects in Section 4.4.

*Pros and Cons.* Automata enable us to represent infinite sets of strings with finite machines, hence they are a natural and elegant way to represent unbounded-length strings. The theory of automata [171] is well defined and studied, and non-trivial string solving would be unthinkable without making use of automata results.

Unfortunately, the performance of automata-based approaches for string constraint solving has been hampered by two main factors: (i) the possible state explosion due to automata operations (e.g., the determinization or the intersection of FSA); (ii) the difficulty, in terms of both efficiency and expressiveness, of capturing an exhaustive set of string operations with automata only. Think, e.g., to a simple language such as  $\{a^n b^n c^n \mid n \in \mathbb{N}\}$  which is not regular and not even context-free.

For these reasons, the effectiveness of “general” SCS approaches *purely* relying on automata is limited: modern approaches using automata are in some sense “hybrid”, i.e., they embed automata within other solving techniques. However, for some specific applications automata can still be the best choice. For example, the abstract interpretation [51] of programs containing string expressions can take advantage of the expressiveness and the simplicity of automata.

## 4.2 Word-based approaches

We call a SCS approach *word-based* if it *natively* handles theory of word-equations, possibly enriched with other theories (e.g., integers or regular expressions) [177]. Word-based approaches rely on algebraic results (e.g., the Levi’s lemma [103]) for solving string constraints over the theory of unbounded strings, without systematically reducing to other data structures such as bit vectors or automata. However, this does not mean that bit vectors and automata are not internally used by word-based string solvers. The natural candidates for implementing these approaches are the SMT solvers, which can incorporate and integrate the theory of strings in their frameworks.

Back in 2014, Liang et al. [105] argued that: “*Despite their power and success as back-end reasoning engines, general multitheory SMT solvers so far have provided minimal or no native support for reasoning over strings [...] until very recently the available string solvers were standalone tools that [...] imposed strong restrictions on the expressiveness [...] Traditionally, these solvers were based on reductions to satisfiability problems over other data types*”. However, in the following years the situation has changed, and more and more word-based string solving approaches have emerged.

In [105] the authors integrated a word-based SCS approach into the well-known SMT solver CVC4 [28]. They used a DPLL( $T$ ) approach for solving quantifier-free constraints natively over the theory of unbounded strings with length and positive regular language membership. The focus was on solving efficiently string constraints arising from verification and security applications. The authors claimed that CVC4 was the first solver able to reason about a language of mixed constraints including strings together with integers, reals, arrays, and algebraic datatypes. Their decision procedure is based on derivation rules repeatedly applied until unsatisfiability is detected or a “saturated configuration” is found—in the latter case the input constraints are satisfiable. This procedure is sound (when it terminates, the answer is correct) but in general is not refutation-complete (it may not recognize unsatisfiable formulae) and may not terminate due to the unrolling of regular expressions. The work in [105] has been revised and extended in the following years to handle further string functions frequently occurring in security and verification applications such as `contains`, `indexOf`, `replace`, `string to code point conversion` [132–134].

Meanwhile, almost at the same time, also the well-established SMT solver Z3 [58], started to develop string solving capabilities. Z3-str was introduced in [178] as a *plug-in* of Z3 to solve string constraints arising from web application analysis such as `length`, `concatenation`, `substring`, and `replace`. Z3-str treats strings as a primitive type, and solves string constraints with a procedure that systematically breaks down constant strings into substrings and variables into subvariables, until the variables are eventually bounded with constant strings or a conflict is detected. This procedure may not terminate in general, stuck in infinite splitting, when “*overlapping*” string variables occur. This is a well-known problem when solving word equations.

Consider, e.g., the equation  $aX = Xb$  where  $X$  is a string variable and  $a, b$  are distinct literals. Intuitively,  $X$  is overlapping because if we consider the graphical representation where the substrings of the LHS and RHS of the word equation are aligned as in Fig. 2, we can observe an overlapping substring of  $X$ , denoted with  $X_1$  in Fig. 2, that is not equal to  $X$  and is “shared” between the two sides of the equation. This equation is unsatisfiable because we can rewrite  $aX = Xb$  into an equivalent equation  $aaX_1b = aX_1bb$ , then in turn into  $aaaX_2bb = aaX_2bbb$ , and so on until

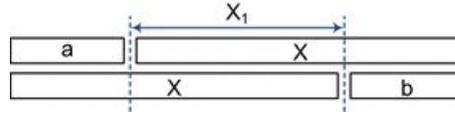


Fig. 2. Example of overlapping string variable  $X$ . Image taken from [177].

we eventually reach an unsatisfiable equation of the form  $a^{k+1}X_k b^k = a^k X_k b^{k+1}$  with  $|X_k| < |a| + |b|$ . However, if the length of  $X$  is unbounded this reasoning cannot be straightforwardly encoded into a terminating procedure.

Z3-str is the progenitor of a number of different string solvers such as Z3str2, Z3strBV, Z3str3, Z3str3RE, and Z3str4. Z3str2 [176] extends Z3-str by including overlapping variables detection and new search heuristics. Overlapping detection allows to individuate cases like that in Fig. 2 to avoid common cases of non-termination. However, Z3str2 may still not terminate due to the unrolling of regular expressions or the interaction between the integer and string components of its theory. The new search heuristics of Z3str2 allows instead the solver to (i) prune the search space via bi-directional integration between the string and integer theories, and (ii) improve the efficiency of string length queries with a binary search based approach.

Z3strBV aims at the software verification, testing, and security analysis of C/C++ programs. Instead of mapping strings to bit-vectors, or representing bit-vectors as natural numbers, Z3strBV defines a decision procedure combining together the theory of string equations, string lengths represented as bit-vectors, and bitvector arithmetic. This procedure is similar to the one of Z3str2: the main differences are the reasoning on length constraints (Z3strBV uses bit-vectors to model underflow/overflow, while Z3str2 uses integers) and the search strategy for length assignments (Z3strBV uses binary search, while Z3str2 uses linear search).

Z3str3 [33] extends Z3str2 by adding a technique called *theory-aware branching* to expose the structure of the theory literals to the underlying SAT solver. The Z3 core solver is modified to take into account the structure of the theory literals underlying the Boolean abstraction of the input formula and prioritize simpler branches over more complex ones. For instance, consider the equation  $XY = AB$  with  $X, Y, A, B$  string variables. Z3str3, as well as Z3str2, handles it by exploring three possible arrangements: (a)  $X = A, Y = B$ , (b)  $X = AX', X'Y = B$ , and (c)  $XX'' = A, Y = X''B$  where  $X', X''$  are fresh variables. But in this case, Z3str3 chooses branch (a) because no variable is introduced.

Z3str3RE [34] is an extension of Z3str3 handling regular expressions without reducing to word equations. Instead, it uses a length-aware, automata-based approach which is sound and complete for the quantifiers-free theory of regular expressions and linear integer arithmetic over string lengths. The basic idea is to take advantage of implicit and explicit length constraints to prune the search space when using automata-based representations. Z3str3RE is based on Z3str3 so it also supports string-number conversion, concatenation and other string operations via reduction to word equations. This solver is a clear example of hybrid approach, because it joins together automata-based, word-based and also unfolding-based techniques as we shall see in Section 4.4.

Z3str4 [158] is, at the moment, the last solver of the Z3str\* saga. It differs from previous approaches because it is actually a meta-solver or *portfolio solver*. Indeed, apart from the above mentioned Z3str3, Z3str4 also includes in its portfolio two additional solvers: LAS, a CEGAR-style length abstraction solver, and Z3seq [149], a Z3-based solver built on top of Z3 (see Section 4.3). Z3str4 uses algorithm selection [98] techniques to select which sequence of algorithms or “arm” is supposedly better for a given unforeseen problem instance. The selection is performed according to some static features extracted from that instance.

Another SMT string solver that came around the same time of CVC4 and Z3-str is Norn [6]. Unlike those solvers, Norn relies on a decision procedure that assumes the *acyclicity* of word equations in order to guarantee its termination (see Section 3). Starting from a conjunction of (dis-)equalities, membership constraints, and arithmetic constraints Norn performs a depth-first search based on rules applied in the following order: (i) arithmetic constraints checking (ii) compound disequalities elimination (iii) complex equalities splitting (iv) complex membership constraints splitting (v) satisfiability check of remaining constraints.

S3 is a *symbolic* string solver [162] motivated by the analysis of web programs inputs. It is called symbolic because it is based on the symbolic representation of string constraints derived from the symbolic execution of an input program. S3 extends Z3-str [178] to handle regular expressions and other high-level operations such as search, replaceAll, regex match, split, test, exec. In particular, S3 added the definition of recursive functions, lazily unfolded during the process of incremental solving. Clearly, this implies that in general the termination of its decision procedure is not guaranteed.

The successor of S3, called S3P [163], uses a progressive search algorithm that mitigates the non-termination issues of recursively defined functions and guides the search towards a “minimal solution”, i.e., a solution having variables with minimal length. S3P tries to detect if a formula is not progressing towards a target solution by pruning the subtrees rooted in those nodes corresponding to formulas that do not contain the minimal solution of the input formula. S3P also implements a conflict clause learning technique to prune the search space. The overall procedure of S3 is still not complete, but in practice works better than S3.

The latest version of S3, called S3# [164], implements instead an algorithm for counting the models of a formula involving the string constraints handled by S3P. S3# uses the reduction rules of S3P but, because a model counter needs to count *all* the solutions, it cannot apply directly the same algorithm. Inspired by model counters SMC [113] and ABC [21] (see Section 4.1), S3# exhaustively builds a reduction tree in a S3P fashion, but with the difference that each node is now associated with a *generating function* representing its count (i.e., the number of solutions of the corresponding formula).

*Pros and Cons.* State-of-the-art word-based approaches are more expressive and flexible than “pure” automata-based approaches when it comes to solve generic SCS problems, possibly involving non-regular languages. Moreover, word-based approaches do not directly suffer from the state-explosion issues of automata. Word-based approaches are naturally built on the top of well-known SMT solvers and can natively handle unbounded-length strings.

One the downside, a fundamental problem of word-based string solvers is that some SCS problems are undecidable or still open. For example, at present it is still unclear if the satisfiability problem for the quantifier-free theory of word equations, regular-expression membership predicate and length function is decidable in general. This makes SMT solving inevitably incomplete unless some restrictions over the input language are imposed (e.g., things become decidable when considering bounded-length variables).

SMT solvers may also suffer from the performance issues due to the disjunctive reasoning of the splitting rule of the underlying DPLL(T) paradigm [72]. Some experimental evaluations [14, 16, 19] showed that these solvers may encounter difficulties when dealing with string variables or literals having big length.

### 4.3 Unfolding-based approaches

An intuitive way of solving string constraints is to reduce them into other well-known data types—for which well-established constraint solving techniques already exist—such as Boolean, integers or bit-vectors. We call a SCS approach

unfolding-based if each string variable is unfolded into an homogeneous sequence of  $k \geq 0$  variables of a type  $T$ , and each string constraint is accordingly mapped into a corresponding operation over  $T$ .

An unfolding approach inherently needs an *upper bound*  $\lambda$  on the string length. So, in general these approaches can handle fixed-length or bounded-length string variables, but cannot deal with unbounded-length variables. One can think to unfolding as a sort of *under-approximation*, where models are bounded by  $\lambda$ . Once again, we underline that hybrid approaches exist. For example, automata-based or word-based solvers can internally use unfolding techniques (e.g., by bounding the string length) to boost the search for solutions. In this scenario, a proper choice of  $\lambda$  is clearly crucial. If  $\lambda$  is too small, one cannot capture solutions having not-small-enough string length. On the other hand, too large a value for  $\lambda$  can significantly worsen the SCS performance even for trivial problems.

Another important choice is whether to unfold *eagerly* (i.e., statically, before the actual solving process) or *lazily* (i.e., dynamically, during the solving process). A hybrid approach typically performs the unfolding lazily, by deferring it after applying some higher-level reasoning.

Hampi [93, 94] was probably the first SMT-based approach unfolding string constraints into constraints over bit-vectors, solved by the underlying STP solver [68]. Its first version [94] only allowed one fixed-length string variable. Its subsequent version [93] added a number of optimisations, e.g., it provided the support for multiple and bounded-length string variables, word equations and substring extraction.

Kaluza [141] was the back-end solver used by Kudzu, a symbolic execution framework for the JavaScript code analysis. Similarly to Hampi, Kaluza dealt with string constraints over bounded-length variables by translating them into bit-vector constraints solved with the STP solver [68]. In fact, Kaluza can be seen as an extension of the first version of Hampi [94] to support multiple, bounded-length string variables.

PASS [104] is a string solver using parameterized arrays of symbolic length as the main data structure to model strings. The indices and the elements of the array can be symbolic too. String constraints are converted into quantified expressions handled by an iterative quantifier elimination algorithm generating equisatisfiable un-quantified constraints, that are then solved by a SMT solver. Moreover, PASS also uses an automaton to handle regular expressions. So, PASS can be seen as an hybrid approach that combines both unfolding and automata methods.

Mapping strings into bit-vectors is suitable for software analysis applications, especially when it comes to precisely handling overflows via *wrapped* integer arithmetic. Plenty of SMT solvers for the quantifier-free bit-vector formulas exist (often relying on *bit blasting*, i.e., by converting bit-vector formulae to SAT problems) while the CP support for bit-vectors appears limited [117].

Woorpje [55] uses a back-end SAT solver to solve bounded-length word equations. However, the unfolding performed by Woorpje is not the naive encoding mapping each string position to a fixed number of Boolean variables. Woorpje instead first encodes with Boolean variables the automaton denoting the solutions of an equation, and then it solves the resulting problem with a SAT solver. This process is refined by considering the length abstraction of word equations, which allows Woorpje to narrow the upper bound of string lengths. The length abstraction is also used to guide the search in the automaton.

As an alternative to the Boolean encoding of strings, a straightforward approach is to reason on sequences of characters (or, equivalently, integers). For example, Z3seq [149] is a Z3-based approach—not belonging to the Z3str\* family—supporting a general theory of sequences over arbitrary datatypes. Hence, Z3seq can be used as a string solver by considering strings as sequences of characters. Z3seq uses *symbolic Boolean derivatives* to reduce regex constraints without constructing automata. In a nutshell, the symbolic derivative of a regex  $R$  is a regex  $\delta(R)$  possibly enriched

with if-then-else conditionals to algebraically manipulate complementation and intersection operations. This enable lazy unfolding: the symbolic conditionals directly map to the underlying character theory.

Mapping sting variables to integer variables is straightforward for C(L)P solvers [11]. For example, in [42] the authors describe a lightweight solver relying on CLP and *Constraint Handling Rules* (CHR) paradigm [64] in order to generate large solutions for tractable string constraints in model finding. This approach unfolds string variables by first labelling their lengths and domains, and then their characters.

CP solvers can use fixed-length or bounded-length arrays of integer variables to deal with regular, grammar and other string constraints [11, 79, 91, 115, 124, 129] without a native support for string variables. However, as shown in [11, 16, 19, 147], having dedicated *propagators* for string variables can make a difference. As explained in Sect. 2, string propagation enables the solver to remove inconsistent values from the domain of string variables by performing a proper high-level reasoning for each type of string constraint occurring in the problem.

In [145] Scott et al. presented a prototypical bounded-length approach based on the *affix* domain to natively handle string variables. This domain allows one to reason about the content of string suffixes even when the length is unknown by using a padding symbol at the end of the string.

The approach in [145] has been subsequently improved [144, 146, 147] with a new structured variable type for strings called *Open-Sequence Representation*, for which suitable propagators were defined. If  $\lambda$  is the maximum string length, each string variable  $X$  is modeled with a pair  $(A_X, N_X)$  where  $A_X$  is an array of integer variables such that  $A[i] \in \Sigma \cup \{\epsilon\}$  denotes the  $i$ -th character of  $X$ , and  $N_X$  is an integer variable denoting the length of  $X$ . The invariant  $A[i] = \epsilon \Leftrightarrow A[j]_{i \leq j \leq \lambda} = \epsilon \Leftrightarrow N_X < i$  is enforced for  $i = 1, \dots, \lambda$ . This approach has been implemented in the Gecode solver [73] and in [11] (and following works) has been referred as `GECODE+S`.

To mitigate the dependency on  $\lambda$  and enable a lazier unfolding, a different CP approach based on *dashed strings* has been introduced [18]. Dashed strings are concatenations of distinct set of strings (called *blocks*) used to represent in a compact way the domain of string variables with potentially very big length. More formally, given an alphabet  $\Sigma$  and a length bound  $\lambda$ , a dashed string has the form  $S_1^{l_1, u_1} \dots S_n^{l_n, u_n}$  with  $0 \leq l_i \leq u_i \leq \lambda$  and  $S_i \subseteq \Sigma$  for  $i = 1, \dots, n$ . Each block  $S_i^{l_i, u_i}$  denotes the language  $\gamma(S_i^{l_i, u_i}) = \{w \in S_i^* \mid l_i \leq |w| \leq u_i\}$ , and the whole dashed string denotes the set of strings  $\{w \in \gamma(S_1^{l_1, u_1}) \cdot \dots \cdot \gamma(S_n^{l_n, u_n}) \mid |w| \leq \lambda\}$ .

Dashed string propagators have been defined for a number of string constraints (e.g., concatenation, length, find/replace, regular expression membership [14–17, 19]) and implemented in the G-Strings solver [9]. These propagators are mainly based on the notion of dashed string *equation*, which can be seen as a semantic unification between dashed strings. For a detailed description of the dashed string formalism, we refer the interested reader to [18].

*Pros and cons.* The unfolding-based approaches allows one to take advantage of already defined theories and propagators without explicitly implementing support for strings. Experimental results show that unfolding approaches, and in particular the CP-based dashed string approach, can be quite effective—especially for satisfiable SCS problems involving long, fixed strings (see, e.g., [14, 16, 19]).

However, unfolding approaches also have a number of limitations. The most obvious one is the impossibility of handling unbounded-length strings. This can be negligible, provided that a good value of  $\lambda$  is chosen. Unfortunately, deciding a good value of  $\lambda$  may not be always trivial. CP solvers may fail on unsatisfiable problems with large domains (because they rely on systematic branching over domain values) and on problems with a lot of logical disjunctions (typically dealt with *reification*). For example, as shown in [18], the “overlapping” problem  $a \cdot X = X \cdot b$  mentioned in Section 4.2 is typically hard for a CP solver because it has to test all the possible values of the domain of  $X$  before

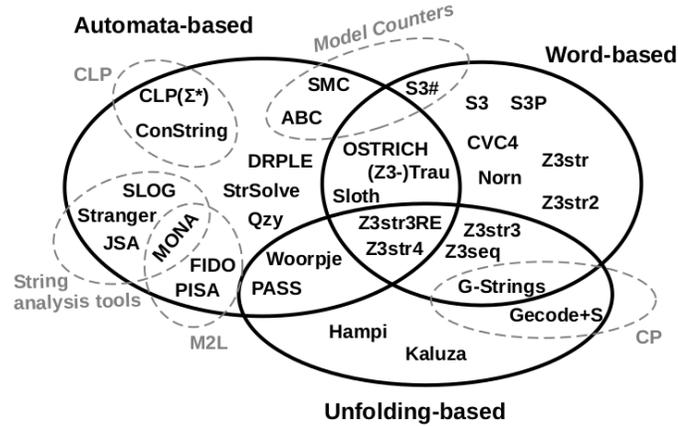


Fig. 3. Graphical representation of SCS approaches.

detecting unsatisfiability. Solving strategies based on *clause learning* [122] would be very helpful in these cases, but unfortunately the learning support for string variables is still an open issue.

#### 4.4 Hybrid approaches

In Section 4.1, 4.2 and 4.3 we provided an overview of several SCS approaches, classified according to the main techniques used to tackle string constraints. This categorization, as mentioned at the beginning of Section 4, has been presented—sometimes under different names—in other works. However, the dividing line between these approaches is not so clear-cut. The more string constraints a SCS approach handles, and the more likely is that different SCS techniques are employed to tackle them. This means that some SCS approaches can be referred as *hybrid*, meaning that orthogonal techniques (e.g., automata in combination with unfolding) are combined together to solve string constraints.

The Venn diagram in Fig. 3 illustrates the relationships between the SCS approaches mentioned in Sections 4.1–4.3 by showing, in the overlapping regions of the diagram, the composition of some hybrid approaches.

For example, the PASS approach based on symbolic arrays presented in [104] is compared against automaton based models (i.e., what we call automata-based approaches) and bit-vector based models (i.e., instances of what we call unfolding-based approaches). The authors make it clear that PASS does not fall in the above categories. Nevertheless, the authors also explicitly say that a contribution of the paper is using interval automata to handle regular expression membership. Hence, we argue that PASS can be considered as a hybrid SCS approach lying at the intersection between automata-based and unfolding-based approaches.

Also the Woorpje solver described in [55] may stand in the middle between automata-based and unfolding-based approaches: it solves word equations by constructing automata that are successively encoded into Boolean formulas tackled by a SAT solver. In [57] Woorpje has been extended with a transformation-system approach using CVC4, Z3str3 and Z3seq as assisting string constraint solvers, which are called according to different heuristics.

As mentioned in Section 4.1, in the intersection between automata-based and word-based approaches we can find Sloth, Ostrich, and (Z3-)Trau. Sloth [81] reduces to alternating finite-state automata to handle string constraints like, e.g., `replaceAll`. OSTRICH [44] is an extension of Sloth that uses FSA to compute the pre-image of string constraints or to handle the `replaceAll` constraint. Both Sloth and OSTRICH however follow the SMT solving approach and are built

on top SMT solver Princess [137]. Analogously, Trau and its extension Z3-Trau can be considered as SMT string solvers, but they also rely on (parametric) flat automata to handle string constraints.

We put G-Strings [18] halfway between unfolding-based approaches and word-based approaches because, even though it relies on lazy unfolding and it does not actually solve word equations in the strict sense, it is strongly based on the notion of dashed string equation. Roughly, each block  $S^{l,u}$  of a dashed string can be seen either as a literal (if the block is fixed, i.e.,  $l = u$  and  $|S| = 1$ ) or a string variable (otherwise). Equating dashed strings does not mean finding a solution, but narrowing the non-fixed blocks by pruning the most inconsistent values.

Z3seq [149] is a word-based approach that may also be considered unfolding-based because it does not implement a theory of strings, but it reduces to the theory of sequences and characters. We cannot consider it an automata-based approach, because it uses symbolic Boolean derivatives to handle regular expressions without actually building automata.

Another solver falling in the word-based/unfolding-based intersection is Z3str3 [33]. In fact, in [158] the authors state that: “*The hybrid approach we use in Z3str3 combines the efficiency of an unfolding-based strategy (reduction of fixed-length word equations to bit-vectors) with the ability of a word-based strategy to reason about string terms of unbounded length (the arrangement method).*”

An interesting case is Z3str3RE [34], an approach that we can consider orthogonal to all classes. Indeed, it is word-based and unfolding-based as explained above. However, it is also automata-based because it uses a length-aware, automata-based approach (see Sect. 4.2).

Finally, we can say that Z3str4 [158] is hybrid by definition, being actually a portfolio joining together different SCS solvers (including Z3str3RE).

#### 4.5 Comparison

Let us conclude Section 4 by giving a closer look on how the SCS techniques mentioned in Sect. 4.1–4.3 compare to each other. These approaches, and many others, have been tested on several string benchmarks—especially coming from the software verification world. To give an idea of the comparisons that have been performed between different string solvers, we considered all the approaches listed in Sections 4.1–4.3 and published from 2015 onward, namely: [2, 4, 6, 14–19, 21, 22, 33, 34, 42, 44, 55, 55, 81, 99, 132–134, 146, 147, 152, 158, 163, 164, 173, 176, 179]. Table 2 shows how many times, in the papers above, the approach on that row has been compared against the approach on that column. We excluded from the count the self-comparisons (i.e., comparisons between different variants of the solver presented in the paper) and the approaches based on model counters, because they only compare to each others. Also, we did not report the approaches in [42, 99, 147] because they did not compare against any other approach in Table 2. We grouped together “families” of solvers: Z3str\* indicates all the Z3str-based solvers seen in Section 4.2, S3\* denotes both S3 and S3P, Trau\* refers to both Trau and Z3-Trau. The entry [179] denotes the unnamed solver described by Zhu et al. in [179], where it is simply referred as “Our Solver”.

Each entry on the last row denotes how many *distinct* approaches have been compared against the approach on that column; each entry on the last column denotes how many distinct approaches the approach on that row has been compared to. As we can see, most of the comparisons are between SMT-based solvers. Among them, Z3str\* and CVC4 solver are certainly the most popular. This is not surprising given that Z3 and CVC4 are well-established SMT solvers providing a solid base for building SCS capabilities by extension. A curious case is represented by G-Strings and Woorpje: they are compared against different string solvers, but no paper reports a comparison with them. If for the latter the reason for this asymmetry might be that Woorpje is quite a recent approach, for G-Strings the main reason is

Table 2. Number of times the approach on that row has been compared against the approach on that column. Cells with number greater than zero have gray background.

	Slog [179]	Qzy	Trau*	Sloth	Ostrich	CVC4	Z3str*	Norn	S3*	Z3seq	G-Strings	Woorpje	#dist.
Slog [179]	–	0	0	0	0	1	1	1	0	0	0	0	3
Qzy	0	–	0	0	1	0	0	0	0	0	0	0	2
Trau*	0	0	–	0	0	2	2	0	1	1	0	0	4
Sloth	0	0	0	–	0	1	0	0	1	0	0	0	2
Ostrich	0	0	0	0	–	1	1	0	0	0	0	0	3
CVC4	0	0	0	0	1	–	1	0	0	3	0	0	3
Z3str*	0	0	0	1	0	1	–	1	2	1	0	0	6
Norn	0	0	0	0	0	1	1	–	1	0	0	0	3
S3*	0	0	0	0	0	1	2	1	–	0	0	0	3
Z3seq	0	0	0	1	0	1	1	0	0	–	0	0	4
G-Strings	0	0	0	0	0	1	1	1	0	1	–	0	4
Woorpje	0	0	0	0	1	1	1	1	0	1	0	–	5
#dist.	0	0	0	2	3	4	10	9	6	4	5	0	0

likely the absence of a standard interface for modeling SCS problems. On the other hand, as we shall see in Section 5.1, SMT solvers can rely on the SMT-LIB standard.

4.5.1 *Expressiveness.* Let us now focus on the expressiveness of the SCS approaches in Table 2.

Slog [173] handles operations on automata, including the replace operation, for acyclic constraints. However, a big limitation of Slog is that it does not support string length operations. As reported in [44], also Sloth [81] does not support length constraints. OSTRICH extends Sloth by extending the class of transducers it supports. In particular, OSTRICH allows variables in both argument positions of `replaceAll`, while Sloth only accepts constant strings as the second argument. According to [179], also the OSTRICH support of integer constraints is limited since “*in general these constraints do not follow its restrictions*”. However, this issue is not detailed in the paper, so it is hard to understand what restrictions they exactly refer to, or what they mean by “in general”. Qzy [52] supports Perl-compatible regular expressions by encoding Boolean combination of regular language membership constraints with Boolean finite automata. Besides regular expression membership, there is no mention of other string operations. Trau [4] supports word equations, length constraints, context-free membership queries, and transducer constraints. Z3-Trau extends Trau by also adding the support for string-number constraints.

Norn [6] solver supports word equations, length constraints, and regular membership operations but it does not handle constraints like `indexOf`, `replace` and `replaceAll`. These operations are instead supported by the solvers of the S3\* family [162–164], which also handle high-level string operations such as `search`, `match`, `split`, `test`, `exec`. Woorpje [55] instead currently only supports word equations with linear length constraints.

CVC4 and Z3\*-based are mature solvers that nowadays support plenty of string constraints. They started with word equations, string length and regular expression memberships, and incrementally added the support for new constraints such as, e.g., `indexOf`, `replace`, and the lexicographic ordering.<sup>6</sup> They support most of the constraints listed in Table 1 of Sect. 2.1 apart from string reversal, characters count and context-free grammar language membership. However, it is worth noting that among the SCS approaches of Table 2, only Trau can handle context-free grammar language membership, and only G-Strings can deal with characters count.

<sup>6</sup>More details here: <http://cvc4.cs.stanford.edu/wiki/Strings> and here: <https://rise4fun.com/z3/tutorialcontent/sequences>.

G-Strings can also handle most of the constraints of Table 1—all of them except the context-free grammar constraints. Adding the support for a new string constraint is easier for a CP-based string solver because, unlike SMT-based solvers, this does not require the definition of a new decision procedure or an axiomatization based on previously defined theories. The only thing required is the definition of a new propagator for possibly pruning the domains of the variables involved in the constraint.

Summarizing, we can say that nowadays a modern “general-purpose” string solver must be able to handle string (in-)equality, length, concatenation, and regular expression membership constraints. With this in mind, closing the problem of the decidability of theory of word equations and arithmetic over length functions would be a significant breakthrough for string constraint solving, at least from the theoretical perspective. In addition, handling constraints like `indexOf`, `replace`, `replaceAll` and string-number conversion is highly recommended given their frequent use in modern software applications. Other constraints like, e.g., lexicographic ordering, string reverse or characters count are probably less used in practice but they can be surely useful for some specific applications.

*4.5.2 Performance.* We conclude Section 4 by discussing the performance of the SCS solvers in Table 2. Here we rely on what is reported in the corresponding papers. However, it is important to note that the reported evaluations might be influenced by factors such as the choice of the benchmarks, the encoding of the problems, the version and parameters configuration of the solvers. We report here the words of [102]: “[SCS approaches] mostly rely on hand crafted suites exhibiting that one tool is better than the others on this particular set of benchmark”. It is outside the scope of this survey to report an empirical evaluation between different SCS approaches.

Slog [173] outperformed CVC4, Z3str2 and Norn on a benchmark of 20386 string analysis instances generated from real web applications via Stranger tool [175]. These instances contain regular expression operations such as union, concatenation, and replacement. Curiously, no other approach from Table 2 has been compared against Slog. In [81] it is mentioned that “We have not experimented with other semi-decision procedures, such as [...] Slog, since they [...] often are not able to process input in the SMT-LIBv2 format, which would complicate the experiments.” This is probably the reason why OSTRICH [44] was evaluated on Slog benchmarks without including the Slog approach itself in the evaluation.

The approach in [179] is compared against Sloth and OSTRICH. The authors show that this approach is not better than Sloth and OSTRICH when it comes to solve SCS problems with constraints they both support. However, [179] can be more efficient for some particular classes of SCS problems where either Sloth or OSTRICH struggle (e.g., those involving the string reversal constraint).

OSTRICH [44] extends and outperforms Sloth [81]. As Sloth, it can handle constraints (e.g., string transducers) that neither CVC4 nor Z3str3 can solve. However, in the Slog benchmarks without `replaceAll` and in the Kaluza benchmarks CVC4 was more effective.

Qzy [52] compared against Norn but excluded CVC4 and Z3str2 “due to their lack of support for negation of regular language membership”. To the best of our knowledge, already at that time both CVC4 and Z3str2 supported negative regular expressions. However, it is possible that the support for this constraint was still prototypical.

Z3-Trau [2] outperformed CVC4, Z3seq and Z3str especially on string-number conversion benchmarks mainly collected with the Py-Conbyte [128] concolic testing tool for Python. Some benchmarks also came from the symbolic execution of JavaScript programs. Curiously, Z3-Trau was not compared against its predecessor Trau [4].

In [132–134] CVC4 was compared against Z3str2, Z3seq and OSTRICH on different benchmarks, often achieving the best results. Exceptions are, e.g., the Slog benchmark (where OSTRICH achieved the peak performance) or the

unsatisfiable instances of the TermEq and Aplas benchmarks of [132] where Z3seq was better. The comparisons between CVC4 and Z3str\* in these papers is limited to Z3str2, because at that time Z3str3 was still unstable.

Nowadays, Z3 solver [157] offers two main ways to solve string constraints: the Z3str3 solver, using the theory of strings, and Z3seq, employing the theory of sequences. Z3seq is the default solver, while Z3-str and Z3str2 are now outdated. In [34] the authors show that Z3str3RE performs better than Z3seq overall, while in [149] Z3seq outperforms Z3str3. This is due to the fact that the benchmarks used are different, although other factors may have been affected the solvers' performance.

A promising approach is Z3str4 [158], a portfolio solver not currently part of Z3, which selects the SCS approach to run based on algorithm selection techniques. The Z3str4's website [159] reports an evaluation over 20 different benchmarks. Overall, Z3str4 was better than Z3str3 and Z3seq but slightly worse than CVC4. The interesting thing here is that the results are not homogeneous along all the benchmarks. For example, on Kausler benchmark Z3str4 is far better than all the other approaches, while on BanditFuzz problems CVC4 dominates all the others.

Norn [6] and S3\* variants [162–164] are probably a bit outdated nowadays. Norn is however still usable and publicly available at [150], although its last update dates back to 2015. At [165] one can find the binaries of all the three versions of S3 string solver (viz. S3, S3P, and S3#). The sources are available only for S3 and no longer updated since 2014. However, it is noteworthy to quote the words of [149]: “*Many of the advances previously developed in S3 are now integrated within Z3's default string solver*”.

Woorpje [55] is a recent SAT-based SCS approach that showed promising results against CVC4, Z3str3, Z3seq, Norn and Sloth. It was evaluated on a benchmark consisting of five handcrafted tracks containing randomly-generated word equations and length constraints. As the authors underline, these benchmarks are interesting because they challenge established solvers. However, it would be also interesting to assess the performance of Woorpje over other well-established benchmarks (see Section 5.2).

G-Strings [18] is currently the state-of-the-art solver for CP problems with string variables and constraints. It improved GECODE+S with the dashed string abstraction, enabling a lazier unfolding and a higher-level reasoning. Several empirical evaluations, on both handcrafted and existing benchmarks, showed its effectiveness: G-Strings was often able to outperform CVC4, Z3str3, Z3seq and Norn [14–19]. Unfortunately, outside the CP world no other SCS approach has been compared against G-Strings. As we shall see in Section 5, the main reason is arguably the lack of a standard encoding for SCS problems.

Only a few of the approaches not shown in the Table 2 are still available or maintained. The automata-based approaches MONA [156], JSA [155], the PHP string analyzer [119] described in [118], and Rex [168] are still available, but among them the only approach that looks still maintained is the MONA project. Kaluza solver is still available at [140] but no longer maintained. Analogously, GECODE+S is on-line [142] but no longer developed.

## 5 TECHNOLOGICAL ASPECTS

In this Section we consider the more practical aspects of string solving, by focusing in particular on the SCS modeling languages, benchmarks and applications that have been developed.

### 5.1 Modeling

As mentioned in Section 4.5, there is no *lingua franca* for SCS problems even though a number of domain-specific languages have been proposed (e.g., [8, 97]). SMT solvers share a common, standard language called *SMT-LIB* [29] to encode (also) SCS problems. Conversely, there is no standard modeling language for G-Strings and CP solvers in

general, although *MiniZinc* [121] can be considered nowadays a de-facto standard. This lack has certainly hampered the cross-comparisons between heterogeneous SCS approaches, hence making difficult a rigorous and comprehensive evaluation between them.

A good news for SCS modeling came in February 2020: based on an initial proposal by Nikolaj Bjørner, Vijay Ganesh, Raphaël Michel and Margus Veanes, a theory of Unicode-based strings and regular expressions became officially part of SMT-LIB standard. This theory, which undoubtedly represents a milestone for SCS, includes the definition of string constants, variables, and functions such as, e.g., concatenation, length, lexicographic ordering, find/replace, regular expressions manipulation, converting from/to integers, and so on.<sup>7</sup>

Let us see a practical example of a SMT-LIB 2.6 specification involving the theory of strings. Consider again the PHP code in Listing 1. An automated program analysis tool may want to prevent the occurrence of the '`<`' character in the `$www` variable. To do so, the tool may collect the constraints encountered during the (dynamic) symbolic execution of the program and then add an additional “counterexample constraint” imposing the occurrence of '`<`' in `$www`. In this way, if the resulting problem is satisfiable we have found a witness telling us that the PHP program may be vulnerable.

Listing 2 shows a SMT-LIB 2.6 instance corresponding to the SCS problem described above. Lines 1 and 2 define two string variables `www_0` and `www_1` used to model the PHP variable `$www` respectively before and after the statement `$www = preg_replace("/[^\A-Za-z0-9 .-@:\\/]"/, "", $www)` at line 4 of Listing 1.

Lines 3–5 capture the semantics of `preg_replace` through the SMT-LIB 2.6 function `replace_re_all`, which replaces, left-to right, each shortest non-empty match of the given regular expression in `www_0` by the empty string. The resulting string is then equalized to the `www_1` variable.

The input pattern `[^\A-Za-z0-9 .-@:\\/]` of `preg_replace` is modelled by composing the regular expression functions `re.inter`, `re.allchar`, `re.comp`, `re.union`, `re.range`, and `str.to_re`. The resulting regular expression denotes the restricted alphabet  $\Sigma - \{a, \dots, Z, A, \dots, Z, 0, \dots, 9, \dots, @, :, /\}$ , where  $\Sigma$  is the set of Unicode characters.

Finally, at line 6 the `indexof` function is used to assert the presence of '`<`' in `www_1` by imposing that its first occurrence starts at an index greater or equal than zero (note that string indexes are 0-based in SMT-LIB, so `indexof` returns `-1` if the searched string does not occur in the target string).

State-of-the-art solvers like CVC4 and Z3 can parse the code in Listing 2. CVC4 can also instantaneously solve this problem, while Z3 cannot support the `replace_re_all` operator, a new entry in SMT-LIB. This operator is hard to handle because, in addition to the already difficult `replaceAll` constraint (see Table 1 of Section 2), one has to handle also the search for a regular expression pattern in a target string. However, in this particular case we can reformulate the problem as shown in Listing 3. The trick here is the transformation of `replace_re_all` into a logically equivalent regular membership constraint. In practice, instead of removing every string matching `[^\A-Za-z0-9 .-@:\\/]` from `www` we simply impose the membership of `www` in `[A-Za-z0-9 .-@:\\/]*`. With this encoding, also Z3 can instantaneously find a solution.

Let us now switch to the CP side. G. Gange developed `smt2mzn-str` [70], a prototypical compiler from SMT-LIB to MiniZinc able to process string variables and constraints. Unfortunately, the SMT-LIB syntax it supports is earlier than version 2.6. Having an up-to-date compiler is important to compare CP and SMT solvers. In particular, in this way the CP solvers could be submitted to the SMT competition.

Listing 4 shows how `smt2mzn-str` translates the SMT-LIB instance of Listing 3 to MiniZinc (Listing 2 cannot be translated because `replace_re_all` is not supported). Line 2 of Listing 4 encodes the regular expression membership

<sup>7</sup>The full specification is available at <http://smtlib.cs.uiowa.edu/theories-UnicodeStrings.shtml>

Listing 2. SMT-LIB 2.6 encoding of the SCS problem.

```

1 (declare-fun www_0 () String)
2 (declare-fun www_1 () String)
3 (assert (= www_1 (str.replace_re_all www_0 (re.inter re.allchar (re.comp (re.union
4   (re.range "A" "Z") (re.range "a" "z") (re.range "0" "9") (str.to_re " ")
5   (re.range "." "@" (str.to_re ":" (str.to_re "/" )))) " " )))
6 (assert (>= (str.indexof www_1 "<" 0) 0) )
7 (check-sat)

```

Listing 3. Reformulation of Listing 2 without `str.replace_re_all`.

```

1 (declare-fun www () String)
2 (assert (str.in_re www (re.* (re.union
3   (re.range "A" "Z") (re.range "a" "z") (re.range "0" "9") (str.to_re " ")
4   (re.range "." "@" (str.to_re ":" (str.to_re "/" )))))
5 (assert (>= (str.indexof www "<" 0) 0) )
6 (check-sat)

```

Listing 4. MiniZinc model translated from Listing 3 with `smt2mzn-str`.

```

1 var string: www;
2 constraint str_reg(www, "([A-Z]|([a-z])|([0-9])|( )|([.-@])|(:)|(/))*");
3 constraint ((str_find_offset("<", www, (0)+(1)))-(1)) >= (0);
4 solve satisfy;

```

while line 3 encodes the occurrence constraint. Note that here the name “find” is preferred to “indexof” [14] and that the string indexes are 1-based, so in general `str_find_offset(X, Y, K)` returns 0 if and only if `X` does not occur in `Y[K]..Y[|Y|]`.

The MiniZinc model in Listing 4 enables a CP string solver like G-Strings to quickly solve the SCS problem. Unfortunately, this encoding is not standard because strings are not (yet) part of the official MiniZinc release, although the unofficial extension introduced in [11] is currently usable to model and solve string constraints. In fact, the G-Strings solver has an interface to process and solve MiniZinc models with strings.

## 5.2 Benchmarks

The use of compilers as a bridge between SMT-LIB and MiniZinc may greatly help the definition of standard benchmarks processable by both SMT and CP solvers. This may be a first step towards a closer and more fruitful interaction between these communities. For example, Bofill et al. [38] defined a compiler from FlatZinc—the low-level language derived from MiniZinc—to SMT-LIB that was used to run the Yices SMT solver over the CP problems of the MiniZinc Challenges 2010–2012 [151].

Unsurprisingly, apart from some handcrafted MiniZinc models defined in [14, 16, 17], virtually all the string benchmarks that can be found in the literature are encoded in SMT-LIB language. A nice framework for the SCS benchmarks generation and evaluation is *StringFuzz* [37], a modular SMT-LIB problem instance transformer and generator for string solvers. A repository of several SMT-LIB 2.0/2.5 problem instances generated and transformed with *StringFuzz* is available online<sup>8</sup>. Table 3 summarizes the nature of these instances, grouped into twelve different classes.

<sup>8</sup><http://stringfuzz.dmitryblotsky.com/benchmarks/>

Table 3. StringFuzz benchmarks composition

Class	Description	Quantity
<i>Concats-{Small,Big}</i>	Right-heavy, deep tree of concats	120
<i>Concats-Balanced</i>	Balanced, deep tree of concats	100
<i>Concats-Extracts-{Small,Big}</i>	Single concat tree, with character extractions	120
<i>Lengths-{Long,Short}</i>	Single, large length constraint on a variable	200
<i>Lengths-Concats</i>	Tree of fixed-length concats of variables	100
<i>Overlaps-{Small,Big}</i>	Word equations $aX = Xb$ with $X$ variable and $a, b$ literals	80
<i>Regex-{Small,Big}</i>	Complex regex membership test	120
<i>Many-Regexes</i>	Multiple random regex membership tests	40
<i>Regex-Deep</i>	Regex membership test with many nested operators	45
<i>Regex-Pair</i>	Test for membership in one regex, but not another	40
<i>Regex-Lengths</i>	Regex membership test, and a length constraint	40
<i>Different-Prefix</i>	Equality of two deep concats with different prefixes	60

StringFuzz also reports the performance of Z3str3, CVC4, Z3seq, and Norn on such instances. Only these SMT solvers were used because other SCS approaches were either unstable or could not properly process the proposed SMT-LIB syntax. In [18], G-Strings has been compared against, and often outperformed, the solvers above after translating the generated StringFuzz instances into MiniZinc with the `smt2mzn-str` compiler.

A more recent fuzzer for SMT solvers is *BanditFuzz* [143]. *BanditFuzz* uses the abstract syntax tree generation procedure of StringFuzz for generating random SMT instances, and also extends it for handling floating point arithmetic. The key feature of *BanditFuzz* is the isolation of the cause for a performance issue, encoded in the form of grammatical constructs such as, e.g., predicates or functions. To learn which constructs are most likely to cause performance issues, *reinforcement learning* methods are used—in particular the problem of how to optimally mutate an input is reduced to the multi-arm bandit problem [153].

SMT-LIB string benchmarks are also available on the website of the SMT competition [174], that since 2018 has also a single-query string track. In 2020 for the first time two different solvers (viz. Z3str4 and CVC4, with the latter achieving better results) have entered the competition, while in 2018 and 2019 editions only CVC4 participated.

A relevant set of SMT-LIB benchmarks have been collected in [102]. This set includes the following benchmarks:

- *Kaluza*: 47284 instances generated from the dynamic symbolic execution of JavaScript by the Kudzu tool [141]
- *PyEx*: 8414 instances generated by PyEx [23], a symbolic executor for Python programs, by Reynolds et al [134]
- *PISA*: 12 instances derived from real-world Java sanitizer methods [154]
- *AppScan*: 8 instances collected by using the output of security warnings generated by IBM Security AppScan [178]
- *StringFuzz*: 1065 instances generated with StringFuzz tool [37]
- *Norn*: 1027 instances coming from queries generated during verification of string-processing programs [6]
- *LightTrau*: 100 instances generated by Abdulla et al. [4] for testing unsatisfiable formulae
- *Woorpje*: 809 instances handcrafted by Day et al. [55] to evaluate Woorpje tool
- *Joaco*: 94 instances based on 11 open-source Java Web applications and security benchmarks [160]
- *Kausler*: 120 instances derived from 8 Java programs via dynamic symbolic execution by Kausler et al. [92]
- *Cashew*: 394 instances from Kaluza benchmark set normalised via Cashew tool [40]
- *Stranger*: 4 instances manually translated from 4 PHP real-world web applications.

The benchmarks’ collection is not the only contribution of [102], which also introduces the ZaligVinder framework to facilitate the performance evaluation and comparison of SMT-based string solvers over those benchmarks.

In addition to those of ZaligVinder, the Z3str4 website [159] offers other benchmarks such as Leetcode Strings, Sloth, Z3str3 Regression, BanditFuzz. Unfortunately, there is no description of them.

### 5.3 Applications

The main application field for string solving is undoubtedly the area of software verification and testing. In fact, most of the SCS approaches presented in Section 4 were developed to facilitate the automated detection of web vulnerabilities.

Kausler et al. [92] performed an evaluation of string constraint solvers in the context of *symbolic execution* [95]. What they found out is that, as one can expect, one solver might be more appropriate than another depending on the input program. This is also pointed out in [10], where Amadini et al. presented a multi-solver tool for the dynamic symbolic execution of JavaScript built on the top of CVC4, G-Strings and Z3 solvers. These observations probably underlie the development of Z3str4 [158].

In [25], a combined approach based on symbolic execution, string analysis and model counting is used to detect and quantify side-channels vulnerabilities in Java programs. This framework is built on the top of Z3 solver [58] and ABC [21] model counter.

In [139], the PHPRepair tool is used for automatically repairing HTML generation errors in PHP via string constraint solving. The property that all tests of a suite should produce their expected output is modelled with string constraints encoded into the language of Kodkod [161], a SAT-based constraint solver.

ACO-Solver [160] is a tool for the detection of injections and XSS vulnerabilities for Java Web applications that uses a hybrid procedure based on the *ant colony optimization* meta-heuristic. ACO-Solver is a meta-solver that complements the support for string operations provided by existing string solvers. JOACO-CS is an extension of ACO-Solver publicly available at [61].

Abstracting a set of strings with a finite formalism is not merely a string solving affair. For example, the well-known *Abstract Interpretation* [51] framework may require the sound approximation of sets of strings—i.e., all the possible “concrete” values that a string variable of the input program can take—with an abstract counterpart.

Several abstract domains have been proposed to approximate sets of strings. These domains vary according to the properties that one needs to capture (e.g., the string length, the prefix or suffix, the characters occurring in a string). Madsen et al. [114] proposed and evaluated a suite of twelve string domains for the static analysis of dynamic field access. Additional string domains are also discussed in [50].

In [49], a refined segmentation abstract domain called M-String is introduced for the static analysis of strings in the C programming language. M-String is based the parametric segmentation domain introduced by P. Cousot for the representation of arrays.

Choi et al. [46] used restricted regular expressions as an abstract domain for strings in the context of Java analysis. Park et al. [123] use a stricter variant of this idea, with a more clearly defined string abstract domain.

Amadini et al. [20] provided an evaluation on the combination via direct product of different string abstract domains in order to improve the precision of JavaScript static analysis. In [12] this combination is achieved via reduced product by using the set of regular languages as a *reference domain* for the other string domains.

Finally, we mention that string solving techniques might be useful in the context of *Bioinformatics*, where a number of CP techniques has been already applied (see, e.g., the works by Barahona et al. [26, 100, 101]).

## 6 CONCLUSIONS

In this work we provided a comprehensive survey on the various aspects of string constraint solving (SCS), an emerging important field orthogonal to combinatorics on words and constraint solving. We focused in particular on the technological aspects of string solving, by grouping the main SCS approaches we are aware, from the early proposals to the state-of-the-art approaches, into three main categories: automata-based, word-based and unfolding-based.

After exploring several approaches of the SCS galaxy, we can safely say that the main application for string solving is the field of software testing and verification, with particular emphasis on the detection of vulnerabilities conveyed by improper use of strings. Another conclusion we can draw is that “general-purpose” SCS approaches nowadays should probably take advantage of the advances that SAT/SMT and CP technologies made in this field over the last years to handle the string constraints most frequently occurring in modern programs.

Among the future directions for string constraint solving we mention the four challenges reported in [13], namely:

- *Extend* the SCS capabilities to properly handle complex string operations, frequently occurring in web programming [109], such as *back-references*, *lookaheads/lookbehinds* or *greedy matching*. These operations significantly extend the expressiveness of regular expressions in a way that string solvers may have difficulty with—e.g., it is known that regular expressions with backreferences can describe non-regular languages.
- *Improve* the efficiency of SCS solvers with new algorithms and search heuristics. At present, some empirical evidences (e.g. [14, 16, 19]) show that SMT solvers tend to fail when dealing with string literals, while CP strings solvers may struggle to prove unsatisfiability because they do not have (yet) a support for *clause learning*—probably the key to success for CP solvers over integers like Chuffed [1] or OR-Tools [76].
- *Combine* SCS solvers with a *portfolio approach* [98] in order to exploit their different nature and uneven performance across different problem instances. Recent experiments for the dynamic symbolic execution of JavaScript already show some potential for this approach, that can be seen as an instance of the *algorithm selection* problem [98]. The novel Z3str4 solver [159] also includes an algorithm selection architecture choosing the best algorithms to run based on their expected performance.
- *Use* SCS solvers and related tools in different fields. The best candidates are probably software verification and testing, model checking and cybersecurity. In these areas SCS plays a crucial role, given the massive use of string operations in modern (web) applications.

Finally, we hope that the research in SCS will encourage a closer and more fruitful collaboration between the CP and the SAT/SMT communities. A step in this direction was taken by Bardin et al. in a 2019 Dagstuhl seminar [27].

## ACKNOWLEDGMENTS

We would like to thank the anonymous reviewers for their informative and constructive comments. Many thanks also to Murphy Berzish, Pierre Flener, Vijay Ganesh, Andrew Reynolds, Peter J. Stuckey, and Cesare Tinelli for their help and feedback.

## REFERENCES

- [1] [n.d.]. Chuffed, a lazy clause generation solver. Available at <https://github.com/chuffed/chuffed>.
- [2] Parosh Aziz Abdulla, Mohamed Faouzi Atig, Yu-Fang Chen, Bui Phi Diep, Julian Dolby, Petr Janku, Hsin-Hung Lin, Lukás Holik, and Wei-Cheng Wu. 2020. Efficient handling of string-number conversion. In *Proceedings of the 41st ACM SIGPLAN International Conference on Programming Language Design and Implementation, PLDI 2020, London, UK, June 15-20, 2020*, Alastair F. Donaldson and Emina Torlak (Eds.). ACM, 943–957.

- [3] Parosh Aziz Abdulla, Mohamed Faouzi Atig, Yu-Fang Chen, Bui Phi Diep, Lukás Holík, Ahmed Rezine, and Philipp Rümmer. 2017. Flatten and Conquer: A Framework for Efficient Analysis of String Constraints. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2017*. 602–617.
- [4] Parosh Aziz Abdulla, Mohamed Faouzi Atig, Yu-Fang Chen, Bui Phi Diep, Lukás Holík, Ahmed Rezine, and Philipp Rümmer. 2018. Trau: SMT solver for string constraints. In *2018 Formal Methods in Computer Aided Design, FMCAD 2018, Austin, TX, USA, October 30 - November 2, 2018*, Nikolaj Bjørner and Arie Gurfinkel (Eds.). IEEE, 1–5.
- [5] Parosh Aziz Abdulla, Mohamed Faouzi Atig, Yu-Fang Chen, Lukás Holík, Ahmed Rezine, Philipp Rümmer, and Jari Stenman. 2014. String Constraints for Verification. In *Computer Aided Verification - 26th International Conference, CAV 2014, Held as Part of the Vienna Summer of Logic, VSL 2014, Vienna, Austria, July 18-22, 2014. Proceedings (Lecture Notes in Computer Science)*, Armin Biere and Roderick Bloem (Eds.), Vol. 8559. Springer, 150–166.
- [6] Parosh Aziz Abdulla, Mohamed Faouzi Atig, Yu-Fang Chen, Lukás Holík, Ahmed Rezine, Philipp Rümmer, and Jari Stenman. 2015. Norn: An SMT Solver for String Constraints. In *CAV (LNCS)*, Vol. 9206. Springer, 462–469.
- [7] Rajeev Alur. 2011. Streaming String Transducers. In *Logic, Language, Information and Computation*, Lev D. Beklemishev and Ruy de Queiroz (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 1–1.
- [8] Rajeev Alur, Loris D’Antoni, and Mukund Raghothaman. 2015. DRex: A Declarative Language for Efficiently Evaluating Regular String Transformations. In *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2015, Mumbai, India, January 15-17, 2015*, Sriram K. Rajamani and David Walker (Eds.). ACM, 125–137.
- [9] Roberto Amadini. 2020. G-Strings: Gecode with (dashed) string variables. Available at <https://github.com/ramadini/gecode>.
- [10] Roberto Amadini, Mak Andrlon, Graeme Gange, Peter Schachte, Harald Søndergaard, and Peter J. Stuckey. 2019. Constraint Programming for Dynamic Symbolic Execution of JavaScript. In *Integration of Constraint Programming, Artificial Intelligence, and Operations Research - 16th International Conference, CPAIOR 2019, Thessaloniki, Greece, June 4-7, 2019, Proceedings (Lecture Notes in Computer Science)*, Louis-Martin Rousseau and Kostas Stergiou (Eds.), Vol. 11494. Springer.
- [11] Roberto Amadini, Pierre Flener, Justin Pearson, Joseph D. Scott, Peter J. Stuckey, and Guido Tack. 2017. MiniZinc with strings. In *LOPSTR 2016: Revised Selected Papers (LNCS)*, Manuel Hermenegildo and Pedro López-García (Eds.), Vol. 10184. Springer, 59–75.
- [12] Roberto Amadini, Graeme Gange, François Gauthier, Alexander Jordan, Peter Schachte, Harald Søndergaard, Peter J. Stuckey, and Chenyi Zhang. 2018. Reference Abstract Domains and Applications to String Analysis. *Fundam. Inform.* 158, 4 (2018), 297–326.
- [13] Roberto Amadini, Graeme Gange, Peter Schachte, Harald Søndergaard, and Peter J. Stuckey. 2020. String Constraint Solving: Past, Present and Future. In *ECAI 2020 - 24th European Conference on Artificial Intelligence, 29 August-8 September 2020, Santiago de Compostela, Spain, August 29 - September 8, 2020 - Including 10th Conference on Prestigious Applications of Artificial Intelligence (PAIS 2020) (Frontiers in Artificial Intelligence and Applications)*, Giuseppe De Giacomo, Alejandro Catalá, Bistra Dilkina, Michela Milano, Senén Barro, Alberto Bugarín, and Jérôme Lang (Eds.), Vol. 325. IOS Press, 2875–2876.
- [14] Roberto Amadini, Graeme Gange, and Peter J. Stuckey. 2018. Propagating *Lex*, *Find* and *Replace* with Dashed Strings. In *Proc. 15th Int. Conf. Integration of Artificial Intelligence and Operations Research Techniques in Constraint Programming (LNCS)*, W.-J. van Hove (Ed.), Vol. 10848. Springer.
- [15] Roberto Amadini, Graeme Gange, and Peter J. Stuckey. 2018. Propagating Regular Membership with Dashed Strings. In *Proc. 24th Conf. Principles and Practice of Constraint Programming (LNCS)*, J. Hooker (Ed.), Vol. 11008. Springer, 13–29.
- [16] Roberto Amadini, Graeme Gange, and Peter J. Stuckey. 2018. Sweep-Based Propagation for String Constraint Solving. In *Proc. 32nd AAAI Conf. Artificial Intelligence*. AAAI Press, 6557–6564.
- [17] Roberto Amadini, Graeme Gange, and Peter J. Stuckey. 2020. Dashed Strings and the Replace-(all) Constraint. In *Principles and Practice of Constraint Programming - 26th International Conference, CP 2020, Louvain-la-Neuve, Belgium, September 7-11, 2020, Proceedings (Lecture Notes in Computer Science)*, Helmut Simonis (Ed.), Vol. 12333. Springer, 3–20.
- [18] Roberto Amadini, Graeme Gange, and Peter J. Stuckey. 2020. Dashed Strings for String Constraint Solving. *Artif. Intell.* 289 (2020), 103368.
- [19] Roberto Amadini, Graeme Gange, Peter J. Stuckey, and Guido Tack. 2017. A Novel Approach to String Constraint Solving. In *Proc. 23rd Int. Conf. Principles and Practice of Constraint Programming (LNCS)*, J. C. Beck (Ed.), Vol. 10416. Springer, 3–20.
- [20] Roberto Amadini, Alexander Jordan, Graeme Gange, François Gauthier, Peter Schachte, Harald Søndergaard, Peter J. Stuckey, and Chenyi Zhang. 2017. Combining String Abstract Domains for JavaScript Analysis: An Evaluation. In *Tools and Algorithms for the Construction and Analysis of Systems - 23rd International Conference, TACAS 2017, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2017, Uppsala, Sweden, April 22-29, 2017, Proceedings, Part I*. 41–57.
- [21] Abdulkali Aydin, Lucas Bang, and Tevfik Bultan. 2015. Automata-Based Model Counting for String Constraints. In *Computer Aided Verification - 27th International Conference, CAV 2015, San Francisco, CA, USA, July 18-24, 2015, Proceedings, Part I (Lecture Notes in Computer Science)*, Daniel Kroening and Corina S. Pasareanu (Eds.), Vol. 9206. Springer, 255–272.
- [22] Abdulkali Aydin, William Eiers, Lucas Bang, Tegan Brennan, Miroslav Gavrilov, Tevfik Bultan, and Fang Yu. 2018. Parameterized model counting for string and numeric constraints. In *Proceedings of the 2018 ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/SIGSOFT FSE 2018, Lake Buena Vista, FL, USA, November 04-09, 2018*, Gary T. Leavens, Alessandro Garcia, and Corina S. Pasareanu (Eds.). ACM, 400–410.

- [23] Thomas Ball and Jakob Daniel. 2015. Deconstructing Dynamic Symbolic Execution. In *Dependable Software Systems Engineering*, Maximilian Irlbeck, Doron A. Peled, and Alexander Pretschner (Eds.). NATO Science for Peace and Security Series, D: Information and Communication Security, Vol. 40. IOS Press, 26–41.
- [24] Davide Balzarotti, Marco Cova, Viktoria Felmetsger, Nenad Jovanovic, Engin Kirda, Christopher Kruegel, and Giovanni Vigna. 2008. Saner: Composing Static and Dynamic Analysis to Validate Sanitization in Web Applications. In *2008 IEEE Symposium on Security and Privacy (S&P 2008), 18-21 May 2008, Oakland, California, USA*. IEEE Computer Society, 387–401.
- [25] Lucas Bang, Abdulkaki Aydin, Quoc-Sang Phan, Corina S. Pasareanu, and Tevfik Bultan. 2016. String analysis for side channels with segmented oracles. In *Proceedings of the 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering, FSE 2016, Seattle, WA, USA, November 13-18, 2016*, Thomas Zimmermann, Jane Cleland-Huang, and Zhendong Su (Eds.). ACM, 193–204.
- [26] Pedro Barahona and Ludwig Krippahl. 2008. Constraint Programming in Structural Bioinformatics. *Constraints* 13, 1-2 (2008), 3–20.
- [27] Sébastien Bardin, Nikolaj Bjørner, and Cristian Cadar. 2019. Bringing CP, SAT and SMT together: Next Challenges in Constraint Solving (Dagstuhl Seminar 19062). *Dagstuhl Reports* 9, 2 (2019), 27–47. <https://doi.org/10.4230/DagRep.9.2.27>
- [28] Clark Barrett, Christopher L. Conway, Morgan Deters, Liana Hadarean, Dejan Jovanović, Tim King, Andrew Reynolds, and Cesare Tinelli. 2011. CVC4. In *Computer Aided Verification*, Ganesh Gopalakrishnan and Shaz Qadeer (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 171–177.
- [29] Clark Barrett, Pascal Fontaine, and Cesare Tinelli. 2016. The Satisfiability Modulo Theories Library (SMT-LIB). Available at <https://www.SMT-LIB.org>.
- [30] Clark W. Barrett and Cesare Tinelli. 2018. Satisfiability Modulo Theories. In *Handbook of Model Checking*. 305–343.
- [31] Jean Berstel and Dominique Perrin. 2007. The origins of combinatorics on words. *European Journal of Combinatorics* 28, 3 (2007), 996–1022.
- [32] Murphy Berzish, Joel D. Day, Vijay Ganesh, Mitja Kulczynski, Florin Manea, Federico Mora, and Dirk Nowotka. 2021. String Theories Involving Regular Membership Predicates: From Practice to Theory and Back. *CoRR abs/2105.07220* (2021). <https://arxiv.org/abs/2105.07220>
- [33] Murphy Berzish, Vijay Ganesh, and Yunhui Zheng. 2017. Z3str3: A String Solver with Theory-Aware Heuristics. In *Proc. 17th Conf. Formal Methods in Computer-Aided Design*, D. Stewart and G. Weissenbacher (Eds.). FMCAD Inc, 55–59.
- [34] Murphy Berzish, Mitja Kulczynski, Federico Mora, Florin Manea, Joel D. Day, Dirk Nowotka, and Vijay Ganesh. 2020. A Length-aware Regular Expression SMT Solver. *CoRR abs/2010.07253* (2020). arXiv:2010.07253 <https://arxiv.org/abs/2010.07253>
- [35] Prithvi Bisht, Timothy L. Hinrichs, Nazari Skrupsky, and V. N. Venkatakrishnan. 2011. WAPTEC: Whitebox analysis of web applications for parameter tampering exploit construction. In *Proceedings of ACM Conference on Computer and Communications Security*. ACM, 575–586.
- [36] Nikolaj Bjørner, Nikolai Tillmann, and Andrei Voronkov. 2009. Path Feasibility Analysis for String-Manipulating Programs. In *Proceedings of the International Conference on Tools and Algorithms for the Construction and Analysis of Systems (LNCS)*, Vol. 5505. Springer, 307–321.
- [37] Dmitry Blotsky, Federico Mora, Murphy Berzish, Yunhui Zheng, Ifaz Kabir, and Vijay Ganesh. 2018. StringFuzz: A Fuzzer for String Solvers. In *Computer Aided Verification*, Hana Chockler and Georg Weissenbacher (Eds.). Springer International Publishing, Cham, 45–51.
- [38] Miquel Bofill, Josep Suy, and Mateu Villaret. 2010. A system for solving constraint satisfaction problems with SMT. In *SAT (LNCS)*, Vol. 6175. Springer, 300–305.
- [39] Aaron R. Bradley. 2011. SAT-Based Model Checking without Unrolling. In *Verification, Model Checking, and Abstract Interpretation - 12th International Conference, VMCAI 2011, Austin, TX, USA, January 23-25, 2011. Proceedings (Lecture Notes in Computer Science)*, Ranjit Jhala and David A. Schmidt (Eds.), Vol. 6538. Springer, 70–87.
- [40] Tegan Brennan, Nestan Tsiskaridze, Nicolás Rosner, Abdulkaki Aydin, and Tevfik Bultan. 2017. Constraint normalization and parameterized caching for quantitative program analysis. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering, ESEC/FSE 2017, Paderborn, Germany, September 4-8, 2017*, Eric Bodden, Wilhelm Schäfer, Arie van Deursen, and Andrea Zisman (Eds.). ACM, 535–546.
- [41] Tevfik Bultan, Fang Yu, Muath Alkhalaf, and Abdulkaki Aydin. 2017. *String Analysis for Software Verification and Security*. Springer. <https://doi.org/10.1007/978-3-319-68670-7>
- [42] Fabian Büttner and Jordi Cabot. 2015. Lightweight string reasoning in model finding. *Software and Systems Modeling* 14, 1 (2015), 413–427.
- [43] Taolue Chen, Yan Chen, Matthew Hague, Anthony W. Lin, and Zhilin Wu. 2018. What is decidable about string constraints with the ReplaceAll function. *PACMPL* 2, POPL (2018), 3:1–3:29.
- [44] Taolue Chen, Matthew Hague, Anthony W. Lin, Philipp Rümmer, and Zhilin Wu. 2019. Decision procedures for path feasibility of string-manipulating programs with complex operations. *PACMPL* 3, POPL (2019), 49:1–49:30.
- [45] Christian Hoffrut and Juhani Karhumäki. 1997. Combinatorics of words. In *Handbook of formal languages*. Springer, 329–438.
- [46] Tae-Hyoung Choi, Oukseh Lee, Hyunha Kim, and Kyung-Goo Doh. 2006. A Practical String Analyzer by the Widening Approach. In *Programming Languages and Systems, 4th Asian Symposium, APLAS 2006, Sydney, Australia, November 8-10, 2006, Proceedings (Lecture Notes in Computer Science)*, Naoki Kobayashi (Ed.), Vol. 4279. Springer, 374–388.
- [47] Aske Simon Christensen, Anders Møller, and Michael I. Schwartzbach. 2003. Precise Analysis of String Expressions. In *SAS (LNCS)*, Vol. 2694. Springer, 1–18.
- [48] Edmund Clarke, Orna Grumberg, Somesh Jha, Yuan Lu, and Helmut Veith. 2000. Counterexample-Guided Abstraction Refinement. In *Computer Aided Verification*, E. Allen Emerson and Aravinda Prasad Sistla (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 154–169.
- [49] Agostino Cortesi and Martina Olliaro. 2018. M-String Segmentation: A Refined Abstract Domain for String Analysis in C Programs. In *2018 International Symposium on Theoretical Aspects of Software Engineering, TASE 2018, Guangzhou, China, August 29-31, 2018*, Jun Pang, Chenyi Zhang, Jifeng He, and Jian Weng (Eds.). IEEE Computer Society, 1–8.

- [50] Giulia Costantini, Pietro Ferrara, and Agostino Cortesi. 2015. A suite of abstract domains for static analysis of string values. *Software: Practice and Experience* 45, 2 (2015), 245–287.
- [51] Patrick Cousot and Radhia Cousot. 1977. Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints. In *Proceedings of the Fourth ACM Symposium on Principles of Programming Languages*. ACM, 238–252.
- [52] Arlen Cox and Jason Leasure. 2017. Model Checking Regular Language Constraints. *CoRR* abs/1708.09073 (2017). arXiv:1708.09073 <http://arxiv.org/abs/1708.09073>
- [53] M. Dal Cin. 1980. *The Algebraic Theory of Automata*. Vieweg+Teubner Verlag, Wiesbaden, 348–361.
- [54] Loris D’Antoni, Anthony W. Lin, and Philipp Rümmer. 2019. Meeting on String Constraints and Applications. <https://mosca19.github.io/>.
- [55] Joel D. Day, Thorsten Ehlers, Mitja Kulczynski, Florin Manea, Dirk Nowotka, and Danny Bøgsted Poulsen. 2019. On Solving Word Equations Using SAT. In *Reachability Problems - 13th International Conference, RP 2019, Brussels, Belgium, September 11-13, 2019, Proceedings (Lecture Notes in Computer Science)*, Emmanuel Filiot, Raphaël M. Jungers, and Igor Potapov (Eds.), Vol. 11674. Springer, 93–106.
- [56] Joel D. Day, Vijay Ganesh, Paul He, Florin Manea, and Dirk Nowotka. 2018. The Satisfiability of Word Equations: Decidable and Undecidable Theories. In *Reachability Problems - 12th International Conference, RP 2018, Marseille, France, September 24-26, 2018, Proceedings (Lecture Notes in Computer Science)*, Igor Potapov and Pierre-Alain Reynier (Eds.), Vol. 11123. Springer, 15–29.
- [57] Joel D. Day, Mitja Kulczynski, Florin Manea, Dirk Nowotka, and Danny Bøgsted Poulsen. 2020. Rule-based Word Equation Solving. In *FormaliSE@ICSE 2020: 8th International Conference on Formal Methods in Software Engineering, Seoul, Republic of Korea, July 13, 2020*. ACM, 87–97. <https://doi.org/10.1145/3372020.3391556>
- [58] Leonardo Mendonça de Moura and Nikolaj Bjørner. 2008. Z3: An Efficient SMT Solver. In *Tools and Algorithms for the Construction and Analysis of Systems, 14th International Conference, TACAS 2008, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2008, Budapest, Hungary, March 29-April 6, 2008. Proceedings*. 337–340.
- [59] Valery G. Durnev. 1995. Undecidability of the positive  $\forall\exists$ -theory of a free semigroup. *Siberian Mathematical Journal* 36 (1995), 917–929.
- [60] Michael Emmi, Rupak Majumdar, and Koushik Sen. 2007. Dynamic Test Input Generation for Database Applications. In *Proceedings of the ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*. ACM, 151–162.
- [61] Julian Thomé et al. [n.d.]. Joaco. Available at <https://sites.google.com/site/joacosite/>.
- [62] B. Fine, G. Rosenberger, and Michael Stille. 1995. Nielsen Transformations and Applications: A Survey.
- [63] Philippe Flajolet and Robert Sedgewick. 2009. *Analytic combinatorics*. Cambridge University press.
- [64] Thom W. Frühwirth. 1998. Theory and Practice of Constraint Handling Rules. *J. Log. Program.* 37, 1-3 (1998), 95–138.
- [65] Xiang Fu and Chung-Chih Li. 2010. A String Constraint Solver for Detecting Web Application Vulnerability. In *Proceedings of the 22nd International Conference on Software Engineering & Knowledge Engineering (SEKE’2010), Redwood City, San Francisco Bay, CA, USA, July 1 - July 3, 2010*. Knowledge Systems Institute Graduate School, 535–542.
- [66] Xiang Fu, Michael C. Powell, Michael Bantegui, and Chung-Chih Li. 2013. Simple linear string constraints. *Formal Asp. Comput.* 25, 6 (2013), 847–891. <https://doi.org/10.1007/s00165-011-0214-3>
- [67] Vijay Ganesh and Murphy Berzish. 2016. Undecidability of a Theory of Strings, Linear Arithmetic over Length, and String-Number Conversion. *CoRR* abs/1605.09442 (2016).
- [68] Vijay Ganesh and David L. Dill. 2007. A Decision Procedure for Bit-Vectors and Arrays. In *Computer Aided Verification, 19th International Conference, CAV 2007, Berlin, Germany, July 3-7, 2007, Proceedings (Lecture Notes in Computer Science)*, Werner Damm and Holger Hermanns (Eds.), Vol. 4590. Springer, 519–531.
- [69] Vijay Ganesh, Mia Minnes, Armando Solar-Lezama, and Martin C. Rinard. 2012. Word Equations with Length Constraints: What’s Decidable?. In *Hardware and Software: Verification and Testing - 8th International Haifa Verification Conference, HVC 2012, Haifa, Israel, November 6-8, 2012. Revised Selected Papers (Lecture Notes in Computer Science)*, Armin Biere, Amir Nahir, and Tanja E. J. Vos (Eds.), Vol. 7857. Springer, 209–226.
- [70] Graeme Gange. [n.d.]. smt2mzn-str. Available at <https://bitbucket.org/gkgange/smt2mzn-str>.
- [71] Graeme Gange, Jorge A. Navas, Peter J. Stuckey, Harald Søndergaard, and Peter Schachte. 2013. Unbounded Model-Checking with Interpolation for Regular Language Constraints. In *Proceedings of the International Conference on Tools and Algorithms for the Construction and Analysis of Systems (LNCS)*, Vol. 7795. Springer, 277–291.
- [72] H. Ganzinger, G. Hagen, R. Nieuwenhuis, A. Oliveras, and C. Tinelli. 2004. DPLL(T): Fast Decision Procedures. In *Computer Aided Verification: Proc. 16th Int. Conf. (LNCS)*, R. Alur and D. A. Peled (Eds.), Vol. 3114. Springer, 175–188.
- [73] Gecode Team. 2016. Gecode: Generic Constraint Development Environment. Available at <http://www.gecode.org>.
- [74] Victor M. Glushkov. 1961. The Abstract Theory of Automata. *Russian Mathematical Surveys* 16, 5 (1961), 1–53.
- [75] Keith Golden and Wanlin Pang. 2003. Constraint Reasoning over Strings. In *Proceedings of the Ninth International Conference on Principles and Practice of Constraint Programming*. 377–391.
- [76] Google. [n.d.]. OR-Tools. Available at <https://developers.google.com/optimization>.
- [77] Christoph Haase. 2018. A survival guide to Presburger arithmetic. *ACM SIGLOG News* 5, 3 (2018), 67–82.
- [78] Esben Rune Hansen and Henrik Reif Andersen. 2007. Interactive Configuration with Regular String Constraints. In *Proceedings of the Twenty-Second AAAI Conference on Artificial Intelligence, July 22-26, 2007, Vancouver, British Columbia, Canada*. AAAI Press, 217–223.
- [79] Jun He, Pierre Flener, Justin Pearson, and Wei Ming Zhang. 2013. Solving String Constraints: The Case for Constraint Programming. In *Principles and Practice of Constraint Programming - 19th International Conference, CP 2013, Uppsala, Sweden, September 16-20, 2013. Proceedings*, Christian

- Schulte (Ed.), Vol. 8124. Springer, 381–397.
- [80] Hossein Hojjat, Philipp Rümmer, and Ali Shamakhi. 2019. On Strings in Software Model Checking. In *Programming Languages and Systems - 17th Asian Symposium, APLAS 2019, Nusa Dua, Bali, Indonesia, December 1–4, 2019, Proceedings (Lecture Notes in Computer Science)*, Anthony Widjaja Lin (Ed.), Vol. 11893. Springer, 19–30.
- [81] Lukás Holík, Petr Janku, Anthony W. Lin, Philipp Rümmer, and Tomáš Vojnar. 2018. String constraints with concatenation and transducers solved efficiently. *PACMPL* 2, POPL (2018), 4:1–4:32.
- [82] Pieter Hooimeijer and Margus Veanes. 2011. An Evaluation of Automata Algorithms for String Analysis. In *VMCAI (LNCS)*, Vol. 6538. Springer, 248–262.
- [83] Pieter Hooimeijer and Westley Weimer. 2009. A decision procedure for subset constraints over regular languages. In *Proceedings of the 2009 ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2009, Dublin, Ireland, June 15–21, 2009*, Michael Hind and Amer Diwan (Eds.). ACM, 188–198.
- [84] Pieter Hooimeijer and Westley Weimer. 2012. StrSolve: Solving string constraints lazily. *Automated Software Engineering* 19, 4 (2012), 531–559.
- [85] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. 2007. *Introduction to automata theory, languages, and computation, 3rd Edition*. Addison-Wesley.
- [86] IBM. 2018. Security AppScan. Available at <https://www.ibm.com/security/application-security/appscan>.
- [87] Joxan Jaffar and Michael J. Maher. 1994. Constraint Logic Programming: A Survey. *J. Log. Program.* 19/20 (1994), 503–581. [https://doi.org/10.1016/0743-1066\(94\)90033-7](https://doi.org/10.1016/0743-1066(94)90033-7)
- [88] Artur Jez. 2016. Recompression: A Simple and Powerful Technique for Word Equations. *J. ACM* 63, 1 (2016), 4:1–4:51. <https://doi.org/10.1145/2743014>
- [89] Sheldon B. Akers Jr. 1978. Binary Decision Diagrams. *IEEE Trans. Computers* 27, 6 (1978), 509–516. <https://doi.org/10.1109/TC.1978.1675141>
- [90] Serdar Kadioglu and Meinolf Sellmann. 2010. Grammar constraints. *Constraints* 15, 1 (2010), 117–144.
- [91] Serdar Kadioglu and Meinolf Sellmann. 2010. Grammar constraints. *Constraints* 15, 1 (2010), 117–144. <https://doi.org/10.1007/s10601-009-9073-4>
- [92] Scott Kausler and Elena Sherman. 2014. Evaluation of String Constraint Solvers in the Context of Symbolic Execution. In *Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering*. Association for Computing Machinery, New York, NY, USA, 259–270.
- [93] Adam Kiežun, Vijay Ganesh, Shay Artzi, Philip J. Guo, Pieter Hooimeijer, and Michael D. Ernst. 2012. HAMPI: A Solver for Word Equations over Strings, Regular Expressions, and Context-Free Grammars. *ACM Trans. Software Engineering and Methodology* 21, 4 (2012), article 25.
- [94] Adam Kiežun, Vijay Ganesh, Philip J. Guo, Pieter Hooimeijer, and Michael D. Ernst. 2009. HAMPI: A Solver for String Constraints. In *Proc. 18th Int. Symp. Software Testing and Analysis*. ACM, 105–116.
- [95] James C. King. 1976. Symbolic Execution and Program Testing. *Commun. ACM* 19, 7 (1976), 385–394.
- [96] Nils Klarlund. 1998. Mona & Fido: The logic-automaton connection in practice. In *Computer Science Logic*, Mogens Nielsen and Wolfgang Thomas (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 311–326.
- [97] N. Klarlund and M. I. Schwartzbach. 1999. A domain-specific language for regular sets of strings and trees. *IEEE Transactions on Software Engineering* 25, 3 (1999), 378–386.
- [98] Lars Kotthoff. 2016. Algorithm Selection for Combinatorial Search Problems: A Survey. In *Data Mining and Constraint Programming*. LNAI, Vol. 10101. Springer, 149–190.
- [99] Sebastian Krings, Joshua Schmidt, Patrick Skowronek, Jannik Dunkelau, and Dierk Ehmk. 2019. Towards Constraint Logic Programming over Strings for Test Data Generation. In *Declarative Programming and Knowledge Management - Conference on Declarative Programming, DECLARE 2019, Unifying INAP, WLP, and WFLP, Cottbus, Germany, September 9–12, 2019, Revised Selected Papers (Lecture Notes in Computer Science)*, Petra Hofstedt, Salvador Abreu, Ulrich John, Herbert Kuchen, and Dietmar Seipel (Eds.), Vol. 12057. Springer, 139–159.
- [100] Ludwig Krippahl and Pedro Barahona. 2016. Constraining Redundancy to Improve Protein Docking. In *Principles and Practice of Constraint Programming - 22nd International Conference, CP 2016, Toulouse, France, September 5–9, 2016, Proceedings*. 721–732.
- [101] Ludwig Krippahl, Fábio Madeira, and Pedro Barahona. [n.d.]. Constraining Protein Docking with Coevolution Data for Medical Research. In *Artificial Intelligence in Medicine - 14th Conference on Artificial Intelligence in Medicine, AIME 2013, Murcia, Spain, May 29 - June 1, 2013, Proceedings (Lecture Notes in Computer Science)*, Niels Peek, Roque Marín Morales, and Mor Peleg (Eds.), Vol. 7885. Springer, 110–114.
- [102] Mitja Kulczynski, Florin Manea, Dirk Nowotka, and Danny Bøgsted Poulsen. 2020. The Power of String Solving: Simplicity of Comparison. In *AST@ICSE 2020: IEEE/ACM 1st International Conference on Automation of Software Test, Seoul, Republic of Korea, 15–16 July, 2020*. ACM, 85–88.
- [103] Frank W Levi. 1944. On semigroups. *Bull. Calcutta Math. Soc* 36, 141–146 (1944), 82.
- [104] Guodong Li and Indradeep Ghosh. 2013. PASS: String Solving with Parameterized Array and Interval Automaton. In *Proc. 9th Int. Haifa Verification Conf. (LNCS)*, V. Bertacco and A. Legay (Eds.), Vol. 8244. Springer, 15–31.
- [105] Tianyi Liang, Andrew Reynolds, Cesare Tinelli, Clark Barrett, and Morgan Deters. 2014. A DPLL(T) Theory Solver for a Theory of Strings and Regular Expressions. In *Computer Aided Verification: Proc. 26th Int. Conf. (LNCS)*, A. Biere and R. Bloem (Eds.), Vol. 8559. Springer, 646–662.
- [106] Tianyi Liang, Nestan Tsiskaridze, Andrew Reynolds, Cesare Tinelli, and Clark W. Barrett. 2015. A Decision Procedure for Regular Membership and Length Constraints over Unbounded Strings. In *Frontiers of Combining Systems - 10th International Symposium, FroCoS 2015, Wrocław, Poland, September 21–24, 2015, Proceedings (Lecture Notes in Computer Science)*, Carsten Lutz and Silvio Ranise (Eds.), Vol. 9322. Springer, 135–150.
- [107] Anthony Widjaja Lin and Pablo Barceló. 2016. String solving with word equations and transducers: towards a logic for analysing mutation XSS. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2016, St. Petersburg, FL, USA, January 20 - 22, 2016*, Rastislav Bodik and Rupak Majumdar (Eds.). ACM, 123–136.

- [108] Anthony W. Lin and Rupak Majumdar. 2018. Quadratic Word Equations with Length Constraints, Counter Systems, and Presburger Arithmetic with Divisibility. In *Automated Technology for Verification and Analysis - 16th International Symposium, ATVA 2018, Los Angeles, CA, USA, October 7-10, 2018, Proceedings (Lecture Notes in Computer Science)*, Shuvendu K. Lahiri and Chao Wang (Eds.), Vol. 11138. Springer, 352–369.
- [109] Blake Loring, Duncan Mitchell, and Johannes Kinder. 2019. Sound Regular Expression Semantics for Dynamic Symbolic Execution of JavaScript. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2019*, Kathryn S. McKinley and Kathleen Fisher (Eds.). ACM, 425–438.
- [110] M\_ Lothaire. 1997. *Combinatorics on words*. Vol. 17. Cambridge university press.
- [111] M Lothaire. 2005. *Applied combinatorics on words*. Vol. 105. Cambridge University Press.
- [112] Monsieur Lothaire and M Lothaire. 2002. *Algebraic combinatorics on words*. Vol. 90. Cambridge university press.
- [113] Loi Luu, Shweta Shinde, Prateek Saxena, and Brian Demsky. 2014. A model counter for constraints over unbounded strings. In *ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '14, Edinburgh, United Kingdom - June 09 - 11, 2014*, Michael F. P. O'Boyle and Keshav Pingali (Eds.). ACM, 565–576.
- [114] Magnus Madsen and Esben Andreasen. 2014. String Analysis for Dynamic Field Access. In *23rd International Conference on Compiler Construction (LNCS)*, Vol. 8409. Springer, 197–217.
- [115] Michael J. Maher. 2009. Open Constraints in a Boundable World. In *Integration of AI and OR Techniques in Constraint Programming for Combinatorial Optimization Problems, 6th International Conference, CPAIOR 2009, Pittsburgh, PA, USA, May 27-31, 2009, Proceedings (Lecture Notes in Computer Science)*, Willem Jan van Hoeve and John N. Hooker (Eds.), Vol. 5547. Springer, 163–177.
- [116] Gennadiy Semenovich Makanin. 1977. The problem of solvability of equations in a free semigroup. *Matematicheskii Sbornik* 145, 2 (1977), 147–236.
- [117] Laurent D. Michel and Pascal Van Hentenryck. 2012. Constraint Satisfaction over Bit-Vectors. In *Proceedings of the 18th International Conference on Principles and Practice of Constraint Programming (LNCS)*, Vol. 7514. Springer, 527–543.
- [118] Yasuhiko Minamide. 2005. Static approximation of dynamically generated Web pages. In *WWW*. ACM, 432–441.
- [119] Yasuhiko Minamide and Nobuo Otoi. [n.d.]. PHP String Analyzer. Available at <https://sv.c.titech.ac.jp/minamide/phpsa/>.
- [120] Mehryar Mohri and Mark-Jan Nederhof. 2001. *Regular Approximation of Context-Free Grammars through Transformation*. Springer Netherlands, Dordrecht, 153–163.
- [121] Nicholas Nethercote, Peter J. Stuckey, Ralph Becket, Sebastian Brand, Gregory J. Duck, and Guido Tack. 2007. MiniZinc: Towards a Standard CP Modelling Language. In *Proceedings of the 13th International Conference on Principles and Practice of Constraint Programming (LNCS)*, Vol. 4741. Springer, 529–543.
- [122] Olga Ohrimenko, Peter J. Stuckey, and Mike Codish. 2009. Propagation via Lazy Clause Generation. *Constraints* 14, 3 (2009), 357–391.
- [123] Changhee Park, Hyeonseung Im, and Sukeyoung Ryu. 2016. Precise and scalable static analysis of jQuery using a regular expression domain. In *Proceedings of the 12th Symposium on Dynamic Languages, DLS 2016, Amsterdam, The Netherlands, November 1, 2016*, Roberto Ierusalimsky (Ed.). ACM, 25–36.
- [124] G. Pesant. 2004. A Regular Language Membership Constraint for Finite Sequences of Variables. In *Proceedings of the 10th International Conference on Principles and Practice of Constraint Programming (LNCS)*, M. Wallace (Ed.), Vol. 3258. Springer-Verlag, 482–495.
- [125] Wojciech Plandowski. 2004. Satisfiability of word equations with constants is in PSPACE. *J. ACM* 51, 3 (2004), 483–496.
- [126] Wojciech Plandowski. 2006. An efficient algorithm for solving word equations. In *Proceedings of the 38th Annual ACM Symposium on Theory of Computing, Seattle, WA, USA, May 21-23, 2006*, Jon M. Kleinberg (Ed.). ACM, 467–476.
- [127] Emil L Post. 1946. A variant of a recursively unsolvable problem. *Bull. Amer. Math. Soc.* 52, 4 (1946), 264–268.
- [128] Py-Conbyte Team. 2021. Py-Conbyte: A Python concolic testing tool running on bytecode level. Available at <https://github.com/spencerwuwu/py-conbyte>.
- [129] Claude-Guy Quimper and Toby Walsh. 2006. Global Grammar Constraints. In *Principles and Practice of Constraint Programming - CP 2006, 12th International Conference, CP 2006, Nantes, France, September 25-29, 2006, Proceedings (Lecture Notes in Computer Science)*, Frédéric Benhamou (Ed.), Vol. 4204. Springer, 751–755.
- [130] Willard Van Orman Quine. 1946. Concatenation as a Basis for Arithmetic. *J. Symb. Log.* 11, 4 (1946), 105–114. <https://doi.org/10.2307/2268308>
- [131] Arcot Rajasekar. 1994. Applications in constraint logic programming with strings. In *Principles and Practice of Constraint Programming*, Alan Borning (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 109–122.
- [132] Andrew Reynolds, Andres Nötzli, Clark Barrett, and Cesare Tinelli. 2019. High-Level Abstractions for Simplifying Extended String Constraints in SMT. In *Computer Aided Verification*, Isil Dillig and Serdar Tasiran (Eds.). Springer International Publishing, Cham, 23–42.
- [133] Andrew Reynolds, Andres Nötzli, Clark W. Barrett, and Cesare Tinelli. 2020. A Decision Procedure for String to Code Point Conversion. In *Automated Reasoning - 10th International Joint Conference, IJCAR 2020, Paris, France, July 1-4, 2020, Proceedings, Part I (Lecture Notes in Computer Science)*, Nicolas Peltier and Viorica Sofronie-Stokkermans (Eds.), Vol. 12166. Springer, 218–237.
- [134] Andrew Reynolds, Maverick Woo, Clark Barrett, David Brumley, Tianyi Liang, and Cesare Tinelli. 2017. Scaling Up DPLL(T) String Solvers Using Context-Dependent Simplification. In *Computer Aided Verification*, Rupak Majumdar and Viktor Kunčák (Eds.). Springer International Publishing, Cham, 453–474.
- [135] John Michael Robson and Volker Diekert. 1999. On Quadratic Word Equations. In *STACS 99, 16th Annual Symposium on Theoretical Aspects of Computer Science, Trier, Germany, March 4-6, 1999, Proceedings (Lecture Notes in Computer Science)*, Christoph Meinel and Sophie Tison (Eds.), Vol. 1563. Springer, 217–226. [https://doi.org/10.1007/3-540-49116-3\\_20](https://doi.org/10.1007/3-540-49116-3_20)

- [136] F. Rossi, P. van Beek, and T. Walsh (Eds.). 2006. *Handbook of Constraint Programming*. Elsevier.
- [137] Philipp Rümmer. 2008. A Constraint Sequent Calculus for First-Order Logic with Linear Integer Arithmetic. In *Logic for Programming, Artificial Intelligence, and Reasoning, 15th International Conference, LPAR 2008, Doha, Qatar, November 22–27, 2008. Proceedings (Lecture Notes in Computer Science)*, Iliano Cervesato, Helmut Veith, and Andrei Voronkov (Eds.), Vol. 5330. Springer, 274–289.
- [138] Arto Salomaa and Ian N. Sneddon. 1969. *Theory of Automata*. Pergamon Press Reprint.
- [139] H. Samimi, M. Schäfer, S. Artzi, T. Millstein, F. Tip, and L. Hendren. 2012. Automated repair of HTML generation errors in PHP applications using string constraint solving. In *2012 34th International Conference on Software Engineering (ICSE)*. 277–287.
- [140] Prateek Saxena and Devdatta Akhawe. [n.d.]. Kaluza String Solver. Available at <http://webblaze.cs.berkeley.edu/2010/kaluza/>.
- [141] Prateek Saxena, Devdatta Akhawe, Steve Hanna, Feng Mao, Stephen McCamant, and Dawn Song. 2010. A Symbolic Execution Framework for JavaScript. In *Proc. 2010 IEEE Symp. Security and Privacy*. IEEE Comp. Soc., 513–528.
- [142] Joseph Scott. [n.d.]. Prototype implementation of a bounded string module for the Gecode CP library. Available at <https://github.com/jossc/gecode-string>.
- [143] Joseph Scott, Federico Mora, and Vijay Ganesh. 2020. BanditFuzz: A Reinforcement-Learning Based Performance Fuzzer for SMT Solvers. In *Software Verification - 12th International Conference, VSTTE 2020, and 13th International Workshop, NSV 2020, Los Angeles, CA, USA, July 20–21, 2020, Revised Selected Papers (Lecture Notes in Computer Science)*, Maria Christakis, Nadia Polikarpova, Parasara Sridhar Duggirala, and Peter Schrammel (Eds.), Vol. 12549. Springer, 68–86.
- [144] Joseph D. Scott. 2016. *Other Things Besides Number: Abstraction, Constraint Propagation, and String Variable Types*. Ph.D. Dissertation. Department of Information Technology, Uppsala University, Sweden. Available at <http://urn.kb.se/resolve?urn=urn:nbn:se:uu:diva-273311>.
- [145] Joseph D. Scott, Pierre Flener, and Justin Pearson. 2013. Bounded Strings for Constraint Programming. In *25th IEEE International Conference on Tools with Artificial Intelligence, ICTAI 2013, Herndon, VA, USA, November 4–6, 2013*. IEEE Computer Society, 1036–1043.
- [146] Joseph D. Scott, Pierre Flener, and Justin Pearson. 2015. Constraint Solving on Bounded String Variables. In *Twelfth International Conference on Integration of Artificial Intelligence and Operations Research techniques in Constraint Programming (LNCS)*, Vol. 9075. Springer, 375–392.
- [147] Joseph D. Scott, Pierre Flener, Justin Pearson, and Christian Schulte. 2017. Design and Implementation of Bounded-Length Sequence Variables. In *Proc. 14th Int. Conf. Integration of Artificial Intelligence and Operations Research Techniques in Constraint Programming (LNCS)*, M. Lombardi and D. Salvagnin (Eds.), Vol. 10335. Springer, 51–67.
- [148] David B. Searls. 1995. String Variable Grammar: A Logic Grammar Formalism for the Biological Language of DNA. *J. Log. Program.* 24, 1&2 (1995), 73–102.
- [149] Caleb Stanford, Margus Veanes, and Nikolaj Bjørner. 2020. *Symbolic Boolean Derivatives for Efficiently Solving Extended Regular Expression Constraints*. Technical Report MSR-TR-2020-25. Microsoft. <https://www.microsoft.com/en-us/research/publication/symbolic-boolean-derivatives-for-efficiently-solving-extended-regular-expression-constraints/> Updated November 2020.
- [150] Jari Stenman. [n.d.]. Norn, a solver for string constraints. Available at <http://user.it.uu.se/~jarst116/norn/>.
- [151] Peter J. Stuckey, Thibaut Feydy, Andreas Schutt, Guido Tack, and Julien Fischer. 2014. The MiniZinc Challenge 2008–2013. *AI Magazine* 2 (2014), 55–60.
- [152] Sanu Subramanian, Murphy Berzish, Yunhui Zheng, Omer Tripp, and Vijay Ganesh. 2016. A Solver for a Theory of Strings and Bit-vectors. *CoRR* abs/1605.09446 (2016). <http://arxiv.org/abs/1605.09446>
- [153] Csaba Szepesvári. 2010. *Algorithms for Reinforcement Learning*. Morgan & Claypool Publishers.
- [154] Takaaki Tateishi, Marco Pistoia, and Omer Tripp. 2013. Path- and Index-Sensitive String Analysis Based on Monadic Second-Order Logic. *ACM Trans. Software Engineering Methodology* 22, 4 (2013), article 33.
- [155] BRICS team. [n.d.]. Java String Analyzer. Available at <https://www.brics.dk/JSA>.
- [156] BRICS team. [n.d.]. The MONA Project. Available at <https://www.brics.dk/mona>.
- [157] Z3 team. [n.d.]. Z3 solver. Available at <https://github.com/Z3Prover/Z3>.
- [158] Z3 team. [n.d.]. Z3str4 Description. Available at <https://z3str4.github.io/smtcomp.pdf>.
- [159] Z3 team. [n.d.]. Z3str4 String Solver. Available at <https://z3str4.github.io/>.
- [160] Julian Thomé, Lwin Khin Shar, Domenico Bianculli, and Lionel C. Briand. 2017. Search-driven string constraint solving for vulnerability detection. In *Proceedings of the 39th International Conference on Software Engineering, ICSE 2017, Buenos Aires, Argentina, May 20–28, 2017*. 198–208.
- [161] Emina Torlak and Daniel Jackson. 2007. Kodkod: A Relational Model Finder. In *Tools and Algorithms for the Construction and Analysis of Systems, 13th International Conference, TACAS 2007, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2007 Braga, Portugal, March 24 - April 1, 2007, Proceedings (Lecture Notes in Computer Science)*, Orna Grumberg and Michael Huth (Eds.), Vol. 4424. Springer, 632–647.
- [162] Minh-Thai Trinh, Duc-Hiep Chu, and Joxan Jaffar. 2014. S3: A Symbolic String Solver for Vulnerability Detection in Web Applications. In *SIGSAC*. ACM, 1232–1243.
- [163] Minh-Thai Trinh, Duc-Hiep Chu, and Joxan Jaffar. 2016. Progressive Reasoning over Recursively-Defined Strings. In *Computer Aided Verification - 28th International Conference, CAV 2016, Toronto, ON, Canada, July 17–23, 2016, Proceedings, Part I (Lecture Notes in Computer Science)*, Swarat Chaudhuri and Azadeh Farzan (Eds.), Vol. 9779. Springer, 218–240.
- [164] Minh-Thai Trinh, Duc-Hiep Chu, and Joxan Jaffar. 2017. Model Counting for Recursively-Defined Strings. In *Computer Aided Verification - 29th International Conference, CAV 2017, Heidelberg, Germany, July 24–28, 2017, Proceedings, Part II*. 399–418.

- [165] Min-Thai Trinh. [n.d.]. S3: An Efficient String Solver and Model Counter. Available at <https://trinhmt.github.io/home/S3/>.
- [166] Ju M Važenin and B V Rozenblat. 1983. DECIDABILITY OF THE POSITIVE THEORY OF A FREE COUNTABLY GENERATED SEMIGROUP. *Mathematics of the USSR-Sbornik* 44, 1 (feb 1983), 109–116. <https://doi.org/10.1070/sm1983v044n01abeh000954>
- [167] Margus Veanes. 2013. Applications of Symbolic Finite Automata. In *Implementation and Application of Automata - 18th International Conference, CIAA 2013, Halifax, NS, Canada, July 16-19, 2013. Proceedings (Lecture Notes in Computer Science)*, Stavros Konstantinidis (Ed.), Vol. 7982. Springer, 16–23.
- [168] Margus Veanes and Peli de Halleux. [n.d.]. Rex - Regular Expression Exploration. Available at <https://rise4fun.com/rex>.
- [169] Margus Veanes, Peli de Halleux, and Nikolai Tillmann. 2010. Rex: Symbolic Regular Expression Explorer. In *Third International Conference on Software Testing, Verification and Validation, ICST 2010, Paris, France, April 7-9, 2010*. IEEE Computer Society, 498–507.
- [170] Pavol Voda. 1988. The constraint language trilogy: Semantics and computations. Technical report, Complete Logic Systems, North Vancouver, BC, Canada.
- [171] J. von Neumann. 1951. The general and logical theory of automata. In *Cerebral Mechanisms in Behaviour*, L. A. Jeffress (Ed.). Wiley.
- [172] Clifford Walinsky. 1989. CLP(Sigma\*): Constraint Logic Programming with Regular Sets. In *Logic Programming, Proceedings of the Sixth International Conference, Lisbon, Portugal, June 19-23, 1989*, Giorgio Levi and Maurizio Martelli (Eds.). MIT Press, 181–196.
- [173] Hung-En Wang, Tzung-Lin Tsai, Chun-Han Lin, Fang Yu, and Jie-Hong R. Jiang. 2016. String Analysis via Automata Manipulation with Logic Circuit Representation. In *Computer Aided Verification - 28th International Conference, CAV 2016, Toronto, ON, Canada, July 17-23, 2016, Proceedings, Part I (Lecture Notes in Computer Science)*, Swarat Chaudhuri and Azadeh Farzan (Eds.), Vol. 9779. Springer, 241–260.
- [174] Tjark Weber, Sylvain Conchon, David Déharbe, Matthias Heizmann, Aina Niemetz, and Giles Reger. 2019. The SMT Competition 2015-2018. *J. Satisf. Boolean Model. Comput.* 11, 1 (2019), 221–259. <https://doi.org/10.3233/SAT190123>
- [175] Fang Yu, Muath Alkhalaf, and Tefvik Bultan. 2010. Stranger: An Automata-Based String Analysis Tool for PHP. In *TACAS (LNCS)*, Vol. 6015. Springer, 154–157.
- [176] Yunhui Zheng, Vijay Ganesh, Sanu Subramanian, Omer Tripp, Murphy Berzish, Julian Dolby, and Xiangyu Zhang. 2017. Z3str2: an efficient solver for strings, regular expressions, and length constraints. *Formal Methods in System Design* 50, 2-3 (2017), 249–288.
- [177] Yunhui Zheng, Vijay Ganesh, Sanu Subramanian, Omer Tripp, Julian Dolby, and Xiangyu Zhang. 2015. Effective Search-Space Pruning for Solvers of String Equations, Regular Expressions and Length Constraints. In *CAV (LNCS)*, Vol. 9206. Springer, 235–254.
- [178] Yunhui Zheng, Xiangyu Zhang, and Vijay Ganesh. 2013. Z3-str: A Z3-Based String Solver for Web Application Analysis. In *Proc. 9th Joint Meeting on Foundations of Software Engineering*. ACM, 114–124.
- [179] Qizhen Zhu, Hitoshi Akama, and Yasuhiko Minamide. 2019. Solving String Constraints with Streaming String Transducers. *JIP* 27 (2019), 810–821.