



Solving String Constraints Using SAT

Kevin Lotz¹(✉) , Amit Goel², Bruno Dutertre² , Benjamin Kiesl-Reiter² ,
Soonho Kong² , Rupak Majumdar² , and Dirk Nowotka¹

¹ Department of Computer Science, Kiel University, Kiel, Germany
{kel,dn}@informatik.uni-kiel.de

² Amazon Web Services, Seattle, USA
{amgoel,dutebrun,benkiesl,soonho,rumajumd}@amazon.com

Abstract. String solvers are automated-reasoning tools that can solve combinatorial problems over formal languages. They typically operate on restricted first-order logic formulas that include operations such as string concatenation, substring relationship, and regular expression matching. String solving thus amounts to deciding the satisfiability of such formulas. While there exists a variety of different string solvers, many string problems cannot be solved efficiently by any of them. We present a new approach to string solving that encodes input problems into propositional logic and leverages incremental SAT solving. We evaluate our approach on a broad set of benchmarks. On the logical fragment that our tool supports, it is competitive with state-of-the-art solvers. Our experiments also demonstrate that an eager SAT-based approach complements existing approaches to string solving in this specific fragment.

1 Introduction

Many problems in software verification require reasoning about strings. To tackle these problems, numerous *string solvers*—automated decision procedures for quantifier-free first-order theories of strings and string operations—have been developed over the last years. These solvers form the workhorse of automated-reasoning tools in several domains, including web-application security [19,31,33], software model checking [15], and conformance checking for cloud-access-control policies [2,30].

The general theory of strings relies on deep results in combinatorics on words [5,16,23,29]; unfortunately, the related decision procedures remain intractable in practice. Practical string solvers achieve scalability through a judicious mix of heuristics and restrictions on the language of constraints.

We present a new approach to string solving that relies on an *eager* reduction to the Boolean satisfiability problem (SAT), using incremental solving and unsatisfiable-core analysis for completeness and scalability. Our approach supports a theory that contains Boolean combinations of regular membership constraints and equality constraints on string variables, and captures a large set of practical queries [6].

Our solving method iteratively searches for satisfying assignments up to a length bound on each string variable; it stops and reports unsatisfiability when the search reaches computed upper bounds without finding a solution. Similar to the solver WOORPJE [12], we formulate regular membership constraints as reachability problems in nondeterministic finite automata. By bounding the number of transitions allowed by each automaton, we obtain a finite problem that we encode into propositional logic. To cut down the search space of the underlying SAT problem, we perform an *alphabet reduction* step (SMT-LIB string constraints are defined over an alphabet of $3 \cdot 2^{16}$ letters and a naive reduction to SAT does not scale). Inspired by bounded model checking [8], we iteratively increase bounds and utilize an incremental SAT solver to solve the resulting series of propositional formulas. We perform an unsatisfiable-core analysis after each unsatisfiable incremental call to increase only the bounds of a minimal subset of variables until a theoretical upper bound is reached.

We have evaluated our solver on a large set of benchmarks. The results show that our SAT-based approach is competitive with state-of-the-art SMT solvers in the logical fragment that we support. It is particularly effective on satisfiable instances.

Closest to our work is the WOORPJE solver [12], which also employs an eager reduction to SAT. WOORPJE reduces systems of word equations with linear constraints to a single Boolean formula and calls a SAT solver. An extension can also handle regular membership constraints [21]. However, WOORPJE does not handle the full language of constraints considered here and does not employ the reduction and incremental solving techniques that make our tool scale in practice. More importantly, in contrast to our solver, WOORPJE is not complete—it does not terminate on unsatisfiable instances.

Other solvers such as Hampi [19] and Kaluza [31] encode string problems into constraints on fixed-size bit-vector, which can be solved by reduction to SAT. These tools support expressive constraints but they require a user-provided bound on the length of string variables.

Further from our work are approaches based on the *lazy* SMT paradigm, which tightly integrates dedicated, heuristic, theory solvers for strings using the CDCL(T) architecture (also called DPLL(T) in early papers). Solvers that follow this paradigm include OSTRICH [11], Z3 [25], Z3STR4 [24], CVC5 [3], Z3STR3RE [7], TRAU [1], and CERTISTR [17]. Our evaluation shows that our eager approach is competitive with lazy solvers overall, but it also shows that combining both types of solvers in a portfolio is most effective. Our eager approach tends to perform best on satisfiable instances while lazy approaches work better on unsatisfiable problems.

2 Preliminaries

We assume a fixed alphabet Σ and a fixed set of variables Γ . Words of Σ^* are denoted by w , w' , w'' , etc. Variables are denoted by x, y, z . Our decision procedure supports the theory described in Fig. 1.

$$\begin{aligned}
F &:= F \vee F \mid F \wedge F \mid \neg F \mid \text{Atom} \\
\text{Atom} &:= x \in RE \mid x \doteq y \\
RE &:= RE \cup RE \mid RE \cdot RE \mid RE^* \mid RE \cap RE \mid ? \mid w
\end{aligned}$$

Fig. 1. Syntax: x and y denote string variables and w denotes a word of Σ^* . The symbol $?$ is the wildcard character.

Atoms in this theory include *regular membership constraints* (or *regular constraints* for short) of the form $x \in RE$, where RE is a regular expression, and *variable equations* of the form $x \doteq y$. Concatenation is not allowed in equations.

Regular expressions are defined inductively using union, concatenation, intersection, and the Kleene star. Atomic regular expressions are constant words $w \in \Sigma^*$ and the wildcard character $?$, which is a placeholder for an arbitrary symbol $c \in \Sigma$. All regular expressions are grounded, meaning that they do not contain variables. We use the symbols \nsubseteq and \neq as a shorthand notation for negations of atoms using the respective predicate symbols. The following is an example formula in our language: $\neg(x \in a \cdot ?^* \wedge y \in ?^* \cdot b) \vee x \neq y \vee x \in a \cdot b$.

Using our basic syntax, we can define additional relations, such as *constant equations* $x \doteq w$, and *prefix and suffix constraints*, written $w \sqsubseteq x$ and $w \sqsupseteq x$, respectively. Even though these relations can be expressed as regular constraints (e.g., the prefix constraint $ab \sqsubseteq x$ can be expressed as $x \in a \cdot b \cdot ?^*$), we can generate more efficient reductions to SAT by encoding them explicitly.

This string theory is not as expressive as others, since it does not include string concatenation, but it still has important practical applications. It is used in the ZELKOVA tool described by Backes, et al. [2] to support analysis of AWS security policies. ZELKOVA is a major industrial application of SMT solvers [30].

Given a formula ψ , we denote by $atoms(\psi)$ the set of atoms occurring in ψ , by $V(\psi)$ the set of variables occurring in ψ , and by $\Sigma(\psi)$ the set of constant symbols occurring in ψ . We call $\Sigma(\psi)$ the *alphabet* of ψ . Similarly, given a regular expression R , we denote by $\Sigma(R)$ the set of characters occurring in R . In particular, we have $\Sigma(?) = \emptyset$.

We call a formula *conjunctive* if it is a conjunction of literals and we call it a *clause* if it is a disjunction of literals. We say that a formula is in *normal form* if it is a conjunctive formula without unnegated variable equations. Every conjunctive formula can be turned into normal form by substitution, i.e., by repeatedly rewriting $\psi \wedge x \doteq y$ to $\psi[x := y]$. If ψ is in negation normal form (NNF), meaning that the negation symbol occurs only directly in front of atoms, we denote by $lits(\psi)$ the set of literals occurring in ψ . We say that an atom a occurs with *positive polarity* in ψ if $a \in lits(\psi)$ and that it occurs with *negative polarity* in ψ if $\neg a \in lits(\psi)$; we denote the respective sets of atoms of ψ by $atoms^+(\psi)$ and $atoms^-(\psi)$. The notion of polarity can be extended to arbitrary formulas (not necessarily in NNF), intuitively by considering polarity in a formula's corresponding NNF (see [26] for details).

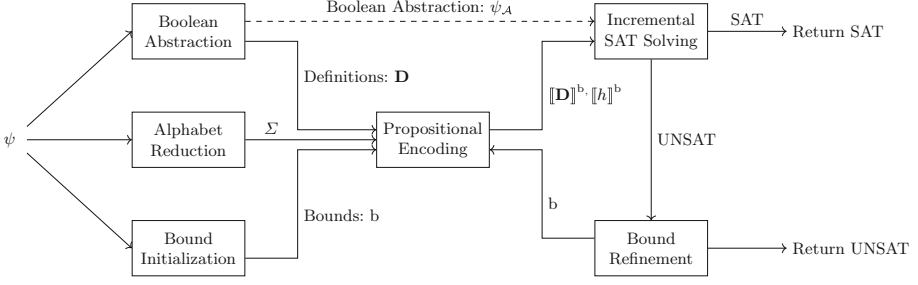


Fig. 2. Overview of the solving process.

The semantics of our language is standard. A regular expression R defines a regular language $\mathcal{L}(R)$ over Σ in the usual way. An *interpretation* is a mapping (also called a *substitution*) $h: \Gamma \rightarrow \Sigma^*$ from string variables to words. Atoms are interpreted as usual, and a *model* (also called a *solution*) is an interpretation that makes a formula evaluate to true under the usual semantics of the Boolean connectives.

3 Overview

Our solving method is illustrated in Fig. 2. It first performs three preprocessing steps that generate a Boolean abstraction of the input formula, reduce the size of the input alphabet, and initialize bounds on the lengths of all string variables. After preprocessing, we enter an encode-solve-and-refine loop that iteratively queries a SAT solver with a problem encoding based on the current bounds and refines the bounds after each unsatisfiable solver call. We repeat this loop until either the propositional encoding is satisfiable, in which case we conclude satisfiability of the input formula, or each bound has reached a theoretical upper bound, in which case we conclude unsatisfiability.

Generating the Boolean Abstraction. We abstract the input formula ψ by replacing each theory atom $a \in \text{atoms}(\psi)$ with a new Boolean variable $\mathbf{d}(a)$, and keep track of the mapping between a and $\mathbf{d}(a)$. This gives us a *Boolean abstraction* $\psi_{\mathcal{A}}$ of ψ and a set \mathbf{D} of definitions, where each definition expresses the relationship between an atom a and its corresponding Boolean variable $\mathbf{d}(a)$. If a occurs with only one polarity in ψ , we encode the corresponding definition as an implication, i.e., as $\mathbf{d}(a) \rightarrow a$ or as $\neg \mathbf{d}(a) \rightarrow \neg a$, depending on the polarity of a . Otherwise, if a occurs with both polarities, we encode it as an equivalence consisting of both implications. This encoding, which is based on ideas behind the well-known *Plaisted-Greenbaum transformation* [28], ensures that the formulas ψ and $\psi_{\mathcal{A}} \wedge \bigwedge_{d \in \mathbf{D}} d$ are equisatisfiable. An example is shown in Fig. 3.

Reducing the Alphabet. In the SMT-LIB theory of strings [4], the alphabet Σ comprises $3 \cdot 2^{16}$ letters, but we can typically use a much smaller alphabet without affecting satisfiability. In Sect. 4, we show that using $\Sigma(\psi)$ and one extra

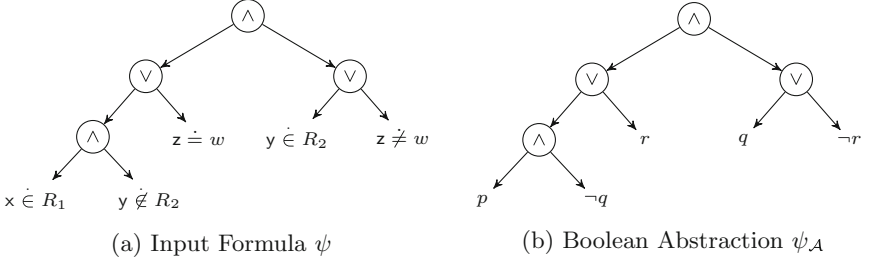


Fig. 3. Example of Boolean abstraction. The formula ψ , whose expression tree is shown on the left, results in the Boolean abstraction illustrated on the right, where p , q , and r are fresh Boolean variables. We additionally get the definitions $p \rightarrow x \in R_1$, $q \leftrightarrow y \in R_2$, and $r \leftrightarrow z \doteq w$. We use an implication (instead of an equivalence) for atom $x \in R_1$ since it occurs only with positive polarity within ψ .

character per string variable is sufficient. Reducing the alphabet is critical for our SAT encoding to be practical.

Initializing Bounds. A model for the original first-order formula ψ is a substitution $h : \Gamma \rightarrow \Sigma^*$ that maps each string variable to a word of arbitrary length such that ψ evaluates to true. As we use a SAT solver to find such substitutions, we need to bound the lengths of strings, which we do by defining a bound function $b : \Gamma \rightarrow \mathbb{N}$ that assigns an upper bound to each string variable. We initialize a small upper bound for each variable, relying on simple heuristics. If the bounds are too small, we increase them in a later refinement step.

Encoding, Solving, and Refining Bounds. Given a bound function b , we build a propositional formula $\llbracket \psi \rrbracket^b$ that is satisfiable if and only if the original formula ψ has a solution h such that $|h(x)| \leq b(x)$ for all $x \in \Gamma$. We encode $\llbracket \psi \rrbracket^b$ as the conjunction $\psi_A \wedge \llbracket \mathbf{D} \rrbracket^b \wedge \llbracket h \rrbracket^b$, where ψ_A is the Boolean abstraction of ψ , $\llbracket \mathbf{D} \rrbracket^b$ is an encoding of the definitions \mathbf{D} , and $\llbracket h \rrbracket^b$ is an encoding of the set of possible substitutions. We discuss details of the encoding in Sect. 5. A key property is that it relies on *incremental SAT solving under assumptions* [13]. Increasing bounds amounts to adding new clauses to the formula $\llbracket \psi \rrbracket^b$ and fixing a set of assumptions, i.e., temporarily fixing the truth values of a set of Boolean variables. If $\llbracket \psi \rrbracket^b$ is satisfiable, we can construct a substitution h from a Boolean model ω of $\llbracket \psi \rrbracket^b$. Otherwise, we examine an unsatisfiable core (i.e., an unsatisfiable subformula) of $\llbracket \psi \rrbracket^b$ to determine whether increasing the bounds may give a solution and, if so, to identify the variables whose bounds must be increased. In Sect. 6, we explain in detail how we analyze unsatisfiable cores, increase bounds, and conclude unsatisfiability.

4 Reducing the Alphabet

In many applications, the alphabet Σ is large—typically Unicode or an approximation of Unicode as defined in the SMT-LIB standard—but formulas use much fewer symbols (less than 100 symbols is common in our experiments). In order to check the satisfiability of a formula ψ , we can restrict the alphabet to the symbols that occur in ψ and add one extra character per variable. This allows us to produce compact propositional encodings that can be solved efficiently in practice.

To prove that such a reduced alphabet A is sufficient, we show that a model $h: \Gamma \rightarrow \Sigma^*$ of ψ can be transformed into a model $h': \Gamma \rightarrow A^*$ of ψ by replacing characters of Σ that do not occur in ψ by new symbols—one new symbol per variable of ψ . For example, suppose $V(\psi) = \{x_1, x_2\}$, $\Sigma(\psi) = \{a, c, d\}$, and h is a model of ψ such that $h(x_1) = abcdef$ and $h(x_2) = abbd$. We introduce two new symbols $\alpha_1, \alpha_2 \in \Sigma \setminus \Sigma(\psi)$, define $h'(x_1) = a\alpha_1cd\alpha_1\alpha_1$ and $h'(x_2) = a\alpha_2\alpha_2d$, and argue that h' is a model as well.

More generally, assume B is a subset of Σ and n is a positive integer such that $|B| \leq |\Sigma| - n$. We can then pick n distinct symbols $\alpha_1, \dots, \alpha_n$ from $\Sigma \setminus B$. Let A be the set $B \cup \{\alpha_1, \dots, \alpha_n\}$. We construct n functions f_1, \dots, f_n from Σ to A by setting $f_i(a) = a$ if $a \in B$, and $f_i(a) = \alpha_i$ otherwise. We extend f_i to words of Σ^* in the natural way: $f_i(\varepsilon) = \varepsilon$ and $f_i(a \cdot w) = f_i(a) \cdot f_i(w)$. This construction satisfies the following property:

Lemma 4.1. *Let f_1, \dots, f_n be mappings as defined above, and let $i, j \in 1, \dots, n$ such that $i \neq j$. Then, the following holds:*

1. *If a and b are distinct symbols of Σ , then $f_i(a) \neq f_j(b)$.*
2. *If w and w' are distinct words of Σ^* , then $f_i(w) \neq f_j(w')$.*

Proof. The first part is an easy case analysis. For the second part, we have that $|f_i(w)| = |w|$ and $|f_j(w')| = |w'|$, so the statement holds if w and w' have different lengths. Assume now that w and w' have the same length and let v be the longest common prefix of w and w' . Since w and w' are distinct, we have that $w = v \cdot a \cdot u$ and $w' = v \cdot b \cdot u'$, where $a \neq b$ are symbols of Σ and u and u' are words of Σ^* . By the first part, we have $f_i(a) \neq f_j(b)$, so $f_i(w)$ and $f_j(w')$ must be distinct. \square

The following lemma can be proved by induction on R .

Lemma 4.2. *Let f_1, \dots, f_n be mappings as defined above and let R be a regular expression with $\Sigma(R) \subseteq B$. Then, for all words $w \in \Sigma^*$ and all $i \in 1, \dots, n$, $w \in \mathcal{L}(R)$ if and only if $f_i(w) \in \mathcal{L}(R)$.*

Given a subset A of Σ , we say that ψ is satisfiable in A if there is a model $h: V(\psi) \rightarrow A^*$ of ψ . We can now prove the main theorem of this section, which shows how to reduce the alphabet while maintaining satisfiability.

Theorem 4.3. *Let ψ be a formula with at most n string variables x_1, \dots, x_n such that $|\Sigma(\psi)| + n \leq |\Sigma|$. Then, ψ is satisfiable if and only if it is satisfiable in an alphabet $A \subseteq \Sigma$ of cardinality $|A| = |\Sigma(\psi)| + n$.*

Proof. We set $B = \Sigma(\psi)$ and use the previous construction. So the alphabet $A = B \cup \{\alpha_1, \dots, \alpha_n\}$ has cardinality $|\Sigma(\psi)| + n$, where $\alpha_1, \dots, \alpha_n$ are distinct symbols of $\Sigma \setminus B$. We can assume that ψ is in disjunctive normal form, meaning that it is a disjunction of the form $\psi = \psi_1 \vee \dots \vee \psi_m$, where each ψ_t is a conjunctive formula. If ψ is satisfiable, then one of the disjuncts ψ_k is satisfiable and we have $\Sigma(\psi_k) \subseteq B$. We can turn ψ_k into normal form by eliminating all variable equalities of the form $x_i \doteq x_j$ from ψ_k , resulting in a conjunction φ_k of literals of the form $x_i \in R$, $x_i \notin R$, or $x_i \neq x_j$. Clearly, for any $A \subseteq \Sigma$, φ_k is satisfiable in A if and only if ψ_k is satisfiable in A .

Let $h: V(\varphi_k) \rightarrow \Sigma^*$ be a model of φ_k and define the mapping $h': V(\varphi_k) \rightarrow A^*$ as $h'(x_i) = f_i(h(x_i))$. We show that h' is a model of φ_k . Consider a literal l of φ_k . We have three cases:

- l is of the form $x_i \in R$ where $\Sigma(R) \subseteq \Sigma(\psi) = B$. Since h satisfies φ_k , we must have $h(x_i) \in \mathcal{L}(R)$ so $h'(x_i) = f_i(h(x_i))$ is also in $\mathcal{L}(R)$ by Lemma 4.2.
- l is of the form $x_i \notin R$ with $\Sigma(R) \subseteq B$. Then, $h'(x_i) \notin \mathcal{L}(R)$ and we can conclude $h'(x_i) \notin \mathcal{L}(R)$ again by Lemma 4.2.
- l is of the form $x_i \neq x_j$. Since h satisfies φ_k , we must have $i \neq j$ and $h(x_i) \neq h(x_j)$, which implies $h'(x_i) = f_i(h(x_i)) \neq f_j(h(x_j)) = h'(x_j)$ by Lemma 4.1.

All literals of φ_k are then satisfied by h' , hence φ_k is satisfiable in A and thus so is ψ_k . It follows that ψ is satisfiable in A . \square

The reduction presented here can be improved and generalized. For example, it can be worthwhile to use different alphabets for different variables or to reduce large character intervals to smaller sets.

5 Propositional Encodings

Our algorithm performs a series of calls to a SAT solver. Each call determines the satisfiability of the propositional encoding $\llbracket \psi \rrbracket^b$ of ψ for some upper bounds b . Recall that $\llbracket \psi \rrbracket^b = \psi_{\mathcal{A}} \wedge \llbracket h \rrbracket^b \wedge \llbracket \mathbf{D} \rrbracket^b$, where $\psi_{\mathcal{A}}$ is the Boolean abstraction of ψ , $\llbracket h \rrbracket^b$ is an encoding of the set of possible substitutions, and $\llbracket \mathbf{D} \rrbracket^b$ is an encoding of the theory-literal definitions, both bounded by b . Intuitively, $\llbracket h \rrbracket^b$ tells the SAT solver to “guess” a substitution, $\llbracket \mathbf{D} \rrbracket^b$ makes sure that all theory literals are assigned proper truth values according to the substitution, and $\psi_{\mathcal{A}}$ forces the evaluation of the whole formula under these truth values.

Suppose the algorithm performs n calls and let $b_k: \Gamma \rightarrow \mathbb{N}$ for $k \in 1, \dots, n$ denote the upper bounds used in the k -th call to the SAT solver. For convenience, we additionally define $b_0(x) = 0$ for all $x \in \Gamma$. In the k -th call, the SAT solver decides whether $\llbracket \psi \rrbracket^{b_k}$ is satisfiable. The Boolean abstraction $\psi_{\mathcal{A}}$, which we already discussed in Sect. 3, stays the same for each call. In the following, we thus discuss the encodings of the substitutions $\llbracket h \rrbracket^{b_k}$ and of the various theory literals $\llbracket a \rrbracket^{b_k}$ and $\llbracket \neg a \rrbracket^{b_k}$ that are part of $\llbracket \mathbf{D} \rrbracket^{b_k}$. Even though SAT solvers expect their input in CNF, we do not present the encodings in CNF to simplify

the presentation, but they can be converted to CNF using simple equivalence transformations.

Most of our encodings are *incremental* in the sense that the formula for call k is constructed by only adding clauses to the formula for call $k - 1$. In other words, for substitution encodings we have $\llbracket h \rrbracket^{b_k} = \llbracket h \rrbracket^{b_{k-1}} \wedge \llbracket h \rrbracket_{b_{k-1}}^{b_k}$ and for literals we have $\llbracket l \rrbracket^{b_k} = \llbracket l \rrbracket^{b_{k-1}} \wedge \llbracket l \rrbracket_{b_{k-1}}^{b_k}$, with the base case $\llbracket h \rrbracket^{b_0} = \llbracket l \rrbracket^{b_0} = \top$. In these cases, it is thus enough to encode the incremental additions $\llbracket l \rrbracket_{b_{k-1}}^{b_k}$ and $\llbracket h \rrbracket_{b_{k-1}}^{b_k}$ for each call to the SAT solver. Some of our encodings, however, introduce clauses that are valid only for a specific bound b_k and thus become invalid for larger bounds. We handle the deactivation of these encodings with *selector variables* as is common in incremental SAT solving.

Our encodings are correct in the following sense.¹

Theorem 5.1. *Let l be a literal and let $b : \Gamma \rightarrow \mathbb{N}$ be a bound function. Then, l has a model that is bounded by b if and only if $\llbracket h \rrbracket^b \wedge \llbracket l \rrbracket^b$ is satisfiable.*

5.1 Substitutions

We encode substitutions by defining for each variable $x \in \Gamma$ the characters to which each of x 's positions is mapped. Specifically, given x and its corresponding upper bound $b(x)$, we represent the substitution $h(x)$ by introducing new variables $x[1], \dots, x[b(x)]$, one for each symbol $h(x)[i]$ of the word $h(x)$. We call these variables *filler variables* and we denote the set of all filler variables by $\tilde{\Gamma}$. By introducing a new symbol $\lambda \notin \Sigma$, which stands for an unused filler variable, we can define h based on a substitution $\tilde{h} : \tilde{\Gamma} \rightarrow \Sigma_\lambda$ over the filler variables, where $\Sigma_\lambda = \Sigma \cup \{\lambda\}$:

$$h(x)[i] = \begin{cases} \varepsilon & \text{if } \tilde{h}(x[i]) = \lambda \\ \tilde{h}(x[i]) & \text{otherwise} \end{cases}$$

We use this representation of substitutions (known as “filling the positions” [18]) because it has a straightforward propositional encoding: For each variable $x \in \Gamma$ and each position $i \in 1, \dots, b(x)$, we create a set $\{h_{x[i]}^a \mid a \in \Sigma_\lambda\}$ of Boolean variables, where $h_{x[i]}^a$ is true if $\tilde{h}(x[i]) = a$. We then use a propositional encoding of an *exactly-one* (EO) constraint (e.g., [20]) to assert that exactly one variable in this set must be true:

$$\llbracket h \rrbracket_{b_{k-1}}^{b_k} = \bigwedge_{x \in \Gamma} \bigwedge_{i=b_{k-1}(x)+1}^{b_k(x)} \text{EO}(\{h_{x[i]}^a \mid a \in \Sigma_\lambda\}) \quad (1)$$

$$\wedge \bigwedge_{x \in \Gamma} \bigwedge_{i=b_{k-1}(x)}^{b_k(x)-1} h_{x[i]}^\lambda \rightarrow h_{x[i+1]}^\lambda \quad (2)$$

¹ Proof is omitted due to space constraints but made available for review purposes.

Constraint (2) prevents the SAT solver from considering filled substitutions that are equivalent modulo λ -substitutions—it enforces that if a position i is mapped to λ , all following positions are mapped to λ too. For instance, $ab\lambda\lambda$, $a\lambda b\lambda$, and $\lambda\lambda ab$ all correspond to the same word ab , but our encoding allows only $ab\lambda\lambda$. Thus, every Boolean assignment ω that satisfies $\llbracket h \rrbracket^b$ encodes exactly one substitution h_ω , and for every substitution h (bounded by b) there exists a corresponding assignment ω_h that satisfies $\llbracket h \rrbracket^b$.

5.2 Theory Literals

The only theory literals of our core language are regular constraints ($x \in R$) and variable equations ($x \doteq y$) with their negations. Constant equations ($x \doteq w$) as well as prefix and suffix constraints ($w \sqsubseteq x$ and $w \sqsupseteq x$) could be expressed as regular constraints, but we encode them explicitly to improve performance.

Regular Constraints. We encode a regular constraint $x \in R$ by constructing a propositional formula that is true if and only if the word $h(x)$ is accepted by a specific nondeterministic finite automaton that accepts the language $\mathcal{L}(R)$. Let $x \in R$ be a regular constraint and let $M = (Q, \Sigma, \delta, q_0, F)$ be a nondeterministic finite automaton (with states Q , alphabet Σ , transition relation δ , initial state q_0 , and accepting states F) that accepts $\mathcal{L}(R)$ and that additionally allows λ -self-transitions on every state. Given that λ is a placeholder for the empty symbol, λ -transitions do not change the language accepted by M . We allow them so that M performs exactly $b(x)$ transitions, even for substitutions of length less than $b(x)$. This reduces checking whether the automaton accepts a word to only evaluating the states reached after exactly $b(x)$ transitions.

Given a model $\omega \models \llbracket h \rrbracket^b$, we express the semantics of M in propositional logic by encoding which states are reachable after reading $h_\omega(x)$. To this end, we assign $b(x) + 1$ Boolean variables $\{S_q^0, S_q^1, \dots, S_q^{b(x)}\}$ to each state $q \in Q$ and assert that $\omega_h(S_q^i) = 1$ if and only if q can be reached by reading prefix $h_\omega(x)[1..i]$. We encode this as a conjunction $\llbracket (M; x) \rrbracket = \llbracket I_{(M; x)} \rrbracket \wedge \llbracket T_{(M; x)} \rrbracket \wedge \llbracket P_{(M; x)} \rrbracket$ of three formulas, modelling the semantics of the initial state, the transition relation, and the predecessor relation of M . We assert that the initial state q_0 is the only state reachable after reading the prefix of length 0, i.e., $\llbracket I_{(M; x)} \rrbracket^{b_1} = S_{q_0}^0 \wedge \bigwedge_{q \in Q \setminus \{q_0\}} \neg S_q^0$. The condition is independent of the bound on x , thus we set $\llbracket I_{(M; x)} \rrbracket_{b_{k-1}}^{b_k} = \top$ for all $k > 1$.

We encode the transition relation of M by stating that if M is in some state q after reading $h_\omega(x)[1..i]$, and if there exists a transition from q to q' labelled with an a , then M can reach state q' after $i + 1$ transitions if $h_\omega(x)[i + 1] = a$. This is expressed in the following formula:

$$\llbracket T_{(M; x)} \rrbracket_{b_{k-1}}^{b_k} = \bigwedge_{i=b_{k-1}(x)}^{b_k(x)-1} \bigwedge_{(q,a) \in \text{dom}(\delta)} \bigwedge_{q' \in \delta(q,a)} (S_q^i \wedge h_{x[i+1]}^a) \rightarrow S_{q'}^{i+1}$$

The formula captures all possible forward moves from each state. We must also ensure that a state is reachable only if it has a reachable predecessor, which we encode with the following formula, where $\text{pred}(q') = \{(q, a) \mid q' \in \delta(q, a)\}$:

$$\llbracket P(M; x) \rrbracket_{b_{k-1}}^{b_k} = \bigwedge_{i=b_{k-1}(x)+1}^{b_k(x)} \bigwedge_{q' \in Q} (S_{q'}^i \rightarrow \bigvee_{(q,a) \in \text{pred}(q')} (S_q^{i-1} \wedge h_{x[i]}^a))$$

The formula states that if state q' is reachable after $i \geq 1$ transitions, then there must be a reachable predecessor state $q \in \hat{\delta}(\{q_0\}, h_\omega(x)[1..i-1])$ such that $q' \in \delta(q, h_\omega(x)[i])$.

To decide whether the automaton accepts $h_\omega(x)$, we encode that it must reach an accepting state after $b_k(x)$ transitions. Our corresponding encoding is only valid for the particular bound $b_k(x)$. To account for this, we introduce a fresh selector variable s_k and define $\llbracket \text{accept}_{x \in M} \rrbracket_{b_{k-1}}^{b_k} = s_k \rightarrow \bigvee_{q_f \in F} S_{q_f}^{b_k(x)}$. Analogously, we define $\llbracket \text{reject}_{x \in M} \rrbracket_{b_{k-1}}^{b_k} = s_k \rightarrow \bigwedge_{q_f \in F} \neg S_{q_f}^{b_k(x)}$. In the k -th call to the SAT solver and all following calls with the same bound on x , we solve under the assumption that s_k is true. In the first call k' with $b_k(x) < b_{k'}(x)$, we re-encode the condition using a new selector variable $s_{k'}$ and solve under the assumption that s_k is false and $s_{k'}$ is true. The full encoding of the regular constraint $x \in R$ is thus given by

$$\llbracket x \in R \rrbracket_{b_{k-1}}^{b_k} = \llbracket (M; x) \rrbracket_{b_{k-1}}^{b_k} \wedge \llbracket \text{accept}_{x \in M} \rrbracket_{b_{k-1}}^{b_k}$$

and its negation $x \notin R$ is encoded as

$$\llbracket x \notin R \rrbracket_{b_{k-1}}^{b_k} = \llbracket (M; x) \rrbracket_{b_{k-1}}^{b_k} \wedge \llbracket \text{reject}_{x \in M} \rrbracket_{b_{k-1}}^{b_k}.$$

Variable Equations. Let $x, y \in \Gamma$ be two string variables, let $l = \min(b_{k-1}(x), b_{k-1}(y))$, and let $u = \min(b_k(x), b_k(y))$. We encode equality between x and y with respect to b_k position-wise up to u :

$$\llbracket x \doteq y \rrbracket_{b_{k-1}}^{b_k} = \bigwedge_{i=l+1}^u \bigwedge_{a \in \Sigma_\lambda} (h_{x[i]}^a \rightarrow h_{y[i]}^a).$$

The formula asserts that for each position $i \in l+1, \dots, u$, if $x[i]$ is mapped to a symbol, then $y[i]$ is mapped to the same symbol (including λ). Since our encoding of substitutions ensures that every position in a string variable is mapped to exactly one character, $\llbracket x \doteq y \rrbracket_{b_{k-1}}^{b_k}$ ensures $x[i] = y[i]$ for $i \in l+1, \dots, u$. In conjunction with $\llbracket x \doteq y \rrbracket^{b_{k-1}}$, which encodes equality up to the l -th position, we have symbol-wise equality of x and y up to bound u . Thus, if $b_k(x) = b_k(y)$, then the formula ensures the equality of both variables. If $b_k(x) > b_k(y)$, we add $h_{x[u+1]}^\lambda$ as an assumption to the solver to ensure $x[i] = \lambda$ for $i \in u+1, \dots, b_k(x)$ and, symmetrically, we add the assumption $h_{y[u+1]}^\lambda$ if $b_k(y) > b_k(x)$.

For the negation $x \neq y$, we encode that $h(x)$ and $h(y)$ must disagree on at least one position, which can happen either because they map to different symbols or because the variable with the higher bound is mapped to a longer word. As for the regular constraints, we again use selector variable s_k to deactivate the encoding for all later bounds, for which it will be re-encoded:

$$\llbracket x \neq y \rrbracket_{b_{k-1}}^{b_k} = \begin{cases} s_k \rightarrow (\bigvee_{i=1}^u \bigvee_{a \in \Sigma_\lambda} (\neg h_{x[i]}^a \wedge h_{y[i]}^a)) & \text{if } b_k(x) = b_k(y) \\ s_k \rightarrow (\bigvee_{i=1}^u \bigvee_{a \in \Sigma_\lambda} (\neg h_{x[i]}^a \wedge h_{y[i]}^a)) \vee \neg h_{y[u+1]}^\lambda & \text{if } b_k(x) < b_k(y) \\ s_k \rightarrow (\bigvee_{i=1}^u \bigvee_{a \in \Sigma_\lambda} (\neg h_{x[i]}^a \wedge h_{y[i]}^a)) \vee \neg h_{x[u+1]}^\lambda & \text{if } b_k(x) > b_k(y) \end{cases}$$

Constant Equations. Given a constant equation $x \doteq w$, if the upper bound of x is less than $|w|$, the atom is trivially unsatisfiable. Thus, for all i such that $b_i(x) < |w|$, we encode $x \doteq w$ with a simple literal $\neg s_{x,w}$ and add $s_{x,w}$ to the assumptions. For $b_k(x) \geq |w|$, the encoding is based on the value of $b_{k-1}(x)$:

$$\llbracket x \doteq w \rrbracket_{b_{k-1}}^{b_k} = \begin{cases} \bigwedge_{i=1}^{|w|} h_{x[i]}^{w[i]} & \text{if } b_{k-1}(x) < |w| = b_k(x) \\ \bigwedge_{i=1}^{|w|} h_{x[i]}^{w[i]} \wedge h_{x[|w|+1]}^\lambda & \text{if } b_{k-1}(x) < |w| < b_k(x) \\ h_{x[|w|+1]}^\lambda & \text{if } b_{k-1}(x) = |w| < b_k(x) \\ \top & \text{if } |w| < b_{k-1}(x) \end{cases}$$

If $b_{k-1}(x) < |w|$, then equality is encoded for all positions $1, \dots, |w|$. Additionally, if $b_k(x) > |w|$, we ensure that the suffix of x is empty starting from position $|w| + 1$. If $b_{k-1}(x) = |w| < b_k(x)$, then only the empty suffix has to be ensured. Lastly, if $|w| < b_{k-1}(x)$, then $\llbracket x \doteq w \rrbracket^{b_{k-1}} \Leftrightarrow \llbracket x \doteq w \rrbracket^{b_k}$.

Conversely, for an inequality $x \neq w$, if $b_k(x) < |w|$, then any substitution trivially is a solution, which we simply encode with \top . Otherwise, we introduce a selector variable $s'_{x,w}$ and define

$$\llbracket x \neq w \rrbracket_{b_{k-1}}^{b_k} = \begin{cases} s'_{x,w} \rightarrow \bigvee_{i=1}^{|w|} \neg h_{x[i]}^{w[i]} & \text{if } b_{k-1}(x) < |w| = b_k(x) \\ \bigvee_{i=1}^{|w|} \neg h_{x[i]}^{w[i]} \vee \neg h_{x[|w|+1]}^\lambda & \text{if } b_{k-1}(x) < |w| < b_k(x) \\ \top & \text{if } |w| < b_{k-1}(x) \leq b_k(x) \end{cases}$$

If $b_k(x) = |w|$, then a substitution h satisfies the constraint if and only if $h(x)[i] \neq w[i]$ for some $i \in 1, \dots, |w|$. If $b_k(x) > |w|$, in addition, h satisfies the constraint if $|h(x)| > |w|$. Thus, if $b_k(x) = |w|$, we perform solver call k under the assumption $s'_{x,w}$, and if $b_k(x) > |w|$, we perform it under the assumption $\neg s'_{x,w}$. Again, if $|w| < b_{k-1}(x)$, then $\llbracket x \neq w \rrbracket^{b_{k-1}} \Leftrightarrow \llbracket x \neq w \rrbracket^{b_k}$.

Prefix and Suffix Constraints. A prefix constraint $w \sqsubseteq x$ expresses that the first $|w|$ positions of x must be mapped exactly onto w . As with equations between a variable x and a constant word w , we could express this as a regular

constraint of the form $x \in w \cdot ?^*$. However, we achieve a more efficient encoding simply by dropping from the encoding of $\llbracket x \doteq w \rrbracket$ the assertion that the suffix of x starting at $|w| + 1$ be empty. Accordingly, a negated prefix constraint $w \not\sqsubseteq x$ expresses that there is an index $i \in 1, \dots, |w|$ such that the i -th position of x is mapped onto a symbol different from $w[i]$, which we encode by repurposing $\llbracket x \neq w \rrbracket$ in a similar manner. Suffix constraints $w \supseteq x$ and $w \not\supseteq x$ can be encoded by analogous modifications to the encodings of $x \doteq w$ and $x \neq w$.

6 Refining Upper Bounds

Our procedure solves a series of SAT problems where the length bounds on string variables increase after each unsatisfiable solver call. The procedure terminates once the bounds are large enough so that further increasing them would be futile. To determine when this is the case, we rely on the upper bounds of a *shortest solution* to a formula ψ . We call a model h of ψ a shortest solution of ψ if ψ has no model h' such that $\sum_{x \in \Gamma} |h'(x)| < \sum_{x \in \Gamma} |h(x)|$. We first establish this bound for conjunctive formulas in normal form, where all literals are of the form $x \neq y$, $x \in R$, or $x \notin R$. Once established, we show how the bound can be generalized to arbitrary formulas.

Let φ be a formula in normal form and let x_1, \dots, x_n be the variables of φ . For each variable x_i , we can collect all the regular constraints on x_i , that is, all the literals of the form $x_i \in R$ or $x_i \notin R$ that occur in φ . We can characterize the solutions to all these constraints by a single nondeterministic finite automaton M_i . If the constraints on x_i are $x_i \in R_1, \dots, x_i \in R_k, x_i \notin R'_1, \dots, x_i \notin R'_l$, then M_i is an NFA that accepts the regular language $\bigcap_{t=1}^k \mathcal{L}(R_t) \cap \bigcap_{t=1}^l \overline{\mathcal{L}(R'_t)}$, where $\overline{\mathcal{L}(R)}$ denotes the complement of $\mathcal{L}(R)$. We say that M_i accepts the regular constraints on x_i in φ . If there are no such constraints on x_i , then M_i is the one-state NFA that accepts the full language Σ^* . Let Q_i denote the set of states of M_i . If we do not take inequalities into account and if the regular constraints on x_i are satisfiable, then a shortest solution h has length $|h(x_i)| \leq |Q_i|$.

Theorem 6.1 gives a bound for the general case with variable inequalities. Intuitively, we prove the theorem by constructing a single automaton \mathcal{P} that takes as input a vector of words $W = (w_1, \dots, w_n)^T$ and accepts W iff the substitution h_W with $h_W(x_i) = w_i$ satisfies φ . To construct \mathcal{P} , we introduce one two-state NFA for each inequality and we then form the product of these NFAs with (slightly modified versions of) the NFAs M_1, \dots, M_n . We can then derive the bound of a shortest solution from the number of states of \mathcal{P} .

Theorem 6.1. *Let φ be a conjunctive formula in normal form over variables x_1, \dots, x_n . Let $M_i = (Q_i, \Sigma, \delta_i, q_{0,i}, F_i)$ be an NFA that accepts the regular constraints on x_i in φ and let k be the number of inequalities occurring in φ . If φ is satisfiable, then it has a model h such that*

$$|h(x_i)| \leq 2^k \times |Q_1| \times \dots \times |Q_n|.$$

Proof. Let λ be a symbol that does not belong to Σ and define $\Sigma_\lambda = \Sigma \cup \{\lambda\}$. As previously, we use λ to extend words of Σ^* by padding. Given a word $w \in \Sigma_\lambda^*$, we denote by \hat{w} the word of Σ^* obtained by removing all occurrences of λ from w . We say that w is well-formed if it can be written as $w = v \cdot \lambda^t$ with $v \in \Sigma^*$ and $t \geq 0$. In this case, we have $\hat{w} = v$. Thus a well-formed word w consists of a prefix in Σ^* followed by a sequence of λ s.

Let Δ be the alphabet Σ_λ^n , i.e., the letters of Δ are the n -letter words over Σ_λ . We can then represent a letter u of Δ as an n -element vector (u_1, \dots, u_n) , and a word W of Δ^t can be written as an $n \times t$ matrix

$$W = \begin{pmatrix} u_{11} & \dots & u_{t1} \\ \vdots & & \vdots \\ u_{1n} & \dots & u_{tn} \end{pmatrix}$$

where $u_{ij} \in \Sigma_\lambda$. Each column of this matrix is a letter in Δ and each row is a word in Σ_λ^t . We denote by $p_i(W)$ the i -th row of this matrix and by $\hat{p}_i(W) = \widehat{p_i(W)}$ the word $p_i(W)$ with all occurrences of λ removed. We say that W is well-formed if the words $p_1(W), \dots, p_n(W)$ are all well-formed. Given a well-formed word W , we can construct a mapping $h_W : \{x_1, \dots, x_n\} \rightarrow \Sigma^*$ by setting $h_W(x_i) = \hat{p}_i(W)$ and we have $|h_W(x_i)| \leq |W| = t$.

To prove the theorem, we build an NFA \mathcal{P} with alphabet Δ such that a well-formed word W is accepted by \mathcal{P} iff h_W satisfies φ . The shortest well-formed W accepted by \mathcal{P} has length no more than the number of states of \mathcal{P} and the bound will follow.

We first extend the NFA $M_i = (Q_i, \Sigma, \delta_i, q_{0,i}, F_i)$ to an automaton M'_i with alphabet Δ . M'_i has the same set of states, initial state, and final states as M_i . Its transition relation δ'_i is defined by

$$\delta'_i(q, u) = \begin{cases} \delta_i(q, u_i) & \text{if } u_i \in \Sigma \\ \{q\} & \text{if } u_i = \lambda \end{cases}$$

One can easily check that M'_i accepts a word W iff M_i accepts $\hat{p}_i(W)$.

For an inequality $x_i \neq x_j$, we construct an NFA $D_{i,j} = (\{e, d\}, \Delta, \delta, e, \{d\})$ with transition function defined as follows:

$$\begin{aligned} \delta(e, u) &= \{e\} & \text{if } u_i = u_j \\ \delta(e, u) &= \{d\} & \text{if } u_i \neq u_j \\ \delta(d, u) &= \{d\}. \end{aligned}$$

This NFA has two states. It starts in state e (for “equal”) and stays in e as long as the characters u_i and u_j are equal. It transitions to state d (for “different”) on the first u where $u_i \neq u_j$ and stays in state d from that point. Since d is the final state, a word W is accepted by $D_{i,j}$ iff $p_i(W) \neq p_j(W)$. If W is well-formed, we also have that W is accepted by $D_{i,j}$ iff $\hat{p}_i(W) \neq \hat{p}_j(W)$.

Let $x_{i_1} \neq x_{j_1}, \dots, x_{i_k} \neq x_{j_k}$ denote the k inequalities of φ . We define \mathcal{P} to be the product of the NFAs M'_1, \dots, M'_n and $D_{i_1, j_1}, \dots, D_{i_k, j_k}$. A well-formed word

W is accepted by \mathcal{P} if it is accepted by all M'_i and all D_{i_t, j_t} , which means that \mathcal{P} accepts a well-formed word W iff h_W satisfies φ .

Let P be the set of states of \mathcal{P} . We then have $|P| \leq 2^k \times |Q_1| \times \dots \times |Q_n|$. Assume φ is satisfiable, so \mathcal{P} accepts a well-formed word W . The shortest well-formed word accepted by \mathcal{P} has an accepting run that does not visit the same state twice. So the length of this well-formed word W is no more than $|P|$. The mapping h_W satisfies φ and for every x_i , it satisfies $|h_W(x_i)| = |\hat{p}_i(W)| \leq |W| \leq |P| \leq 2^k \times |Q_1| \times \dots \times |Q_n|$. \square

The bound given by Theorem 6.1 holds if φ is in normal form but it also holds for a general conjunctive formula ψ . This follows from the observation that converting conjunctive formulas to normal form preserves the length of solutions. In particular, we convert $\psi \wedge x \doteq y$ to formula $\psi' = \psi[x := y]$ so x does not occur in ψ' , but clearly, a bound for y in ψ' gives us the same bound for x in ψ .

In practice, before we apply the theorem we decompose the conjunctive formula φ into subformulas that have disjoint sets of variables. We write φ as $\varphi_1 \wedge \dots \wedge \varphi_m$ where the conjuncts have no common variables. Then, φ is satisfiable if each conjunct φ_t is satisfiable and we derive upper bounds on the shortest solution for the variables of φ_t , which gives more precise bounds than deriving bounds from φ directly. In particular, if a variable x_i of ψ does not occur in any inequality, then the bound on $|h(x_i)|$ is $|Q_i|$.

Theorem 6.1 only holds for conjunctive formulas. For an arbitrary (non-conjunctive) formula ψ , a generalization is to convert ψ into disjunctive normal form. Alternatively, it is sufficient to enumerate the subsets of $\text{lits}(\psi)$. Given a subset A of $\text{lits}(\psi)$, let us denote by d_A a mapping that bounds the length of solutions to A , i.e., any solution h to A satisfies $|h(x)| \leq d_A(x)$. This mapping d_A can be computed from Theorem 6.1. The following property gives a bound for ψ .

Proposition 6.2. *If ψ is satisfiable, then it has a model h such that for all $x \in \Gamma$, it holds that $|h(x)| \leq \max\{d_A(x) \mid A \subseteq \text{lits}(\psi)\}$.*

Proof. We can assume that ψ is in negation normal form. We can then convert ψ to disjunctive normal form $\psi \Leftrightarrow \psi_1 \vee \dots \vee \psi_n$ and we have $\text{lits}(\psi_i) \subseteq \text{lits}(\psi)$. Also, ψ is satisfiable if and only if at least one ψ_i is satisfiable and the proposition follows. \square

Since there are $2^{|\text{lits}(\psi)|}$ subsets of $\text{lits}(\psi)$, a direct application of Proposition 6.2 is rarely feasible in practice. Fortunately, we can use unsatisfiable cores to reduce the number of subsets to consider.

6.1 Unsatisfiable-Core Analysis

Instead of calculating the bounds upfront, we use the unsatisfiable core produced by the SAT solver after each incremental call to evaluate whether the upper

bounds on the variables exceed the upper bounds of the shortest solution. If $\llbracket \psi \rrbracket^b$ is unsatisfiable for bounds b , then it has an unsatisfiable core

$$C = C_{\mathcal{A}} \wedge C_h \wedge \bigwedge_{a \in \text{atoms}^+(\psi)} C_a \wedge \bigwedge_{a \in \text{atoms}^-(\psi)} C_{\bar{a}}$$

with (possibly empty) subsets of clauses $C_{\mathcal{A}} \subseteq \psi_{\mathcal{A}}$, $C_h \subseteq \llbracket h \rrbracket^b$, $C_a \subseteq (\mathbf{d}(a) \rightarrow \llbracket a \rrbracket^b)$, and $C_{\bar{a}} \subseteq (\neg \mathbf{d}(a) \rightarrow \llbracket \neg a \rrbracket^b)$. Here we implicitly assume $\psi_{\mathcal{A}}$, $\mathbf{d}(a) \rightarrow \llbracket a \rrbracket^b$, and $\neg \mathbf{d}(a) \rightarrow \llbracket \neg a \rrbracket^b$ to be in CNF. Let $\mathcal{C}^+ = \{a \mid C_a \neq \emptyset\}$ and $\mathcal{C}^- = \{\neg a \mid C_{\bar{a}} \neq \emptyset\}$ be the sets of literals whose encodings contain at least one clause of the core C . Using these sets, we construct the formula

$$\psi^{\mathcal{C}} = \psi_{\mathcal{A}} \wedge \bigwedge_{a \in \mathcal{C}^+} \mathbf{d}(a) \rightarrow a \wedge \bigwedge_{\neg a \in \mathcal{C}^-} \neg \mathbf{d}(a) \rightarrow \neg a,$$

which consists of the conjunction of the abstraction and the definitions of the literals that are contained in \mathcal{C}^+ , respectively \mathcal{C}^- . Recall that ψ is equisatisfiable to the conjunction $\psi_{\mathcal{A}} \wedge \bigwedge_{d \in \mathbf{D}} d$ of the abstraction and all definitions in \mathbf{D} . Let ψ' denote this formula, i.e.,

$$\psi' = \psi_{\mathcal{A}} \wedge \bigwedge_{a \in \text{atoms}^+(\psi)} \mathbf{d}(a) \rightarrow a \wedge \bigwedge_{\neg a \in \text{atoms}^-(\psi)} \neg \mathbf{d}(a) \rightarrow \neg a.$$

The following proposition shows that it suffices to refine the bounds according to $\psi^{\mathcal{C}}$.

Proposition 6.3. *Let ψ be unsatisfiable with respect to b and let C be an unsatisfiable core of $\llbracket \psi \rrbracket^b$. Then, $\psi^{\mathcal{C}}$ is unsatisfiable with respect to b and $\psi' \models \psi^{\mathcal{C}}$.*

Proof. By definition, we have $\llbracket \psi^{\mathcal{C}} \rrbracket^b = \psi_{\mathcal{A}} \wedge \llbracket h \rrbracket^b \wedge \bigwedge_{a \in \mathcal{C}^+} \mathbf{d}(a) \rightarrow \llbracket a \rrbracket^b \wedge \bigwedge_{\neg a \in \mathcal{C}^-} \neg \mathbf{d}(a) \rightarrow \neg \llbracket \neg a \rrbracket^b$. This implies $C \subseteq \llbracket \psi^{\mathcal{C}} \rrbracket^b$ and, since C is an unsatisfiable core, $\llbracket \psi^{\mathcal{C}} \rrbracket^b$ is unsatisfiable. That is, $\psi^{\mathcal{C}}$ is unsatisfiable with respect to b . We also have $\psi' \models \psi^{\mathcal{C}}$ since $\mathcal{C}^+ \subseteq \text{atoms}^+(\psi)$ and $\mathcal{C}^- \subseteq \text{atoms}^-(\psi)$. \square

Applying Proposition 6.2 to $\psi^{\mathcal{C}}$ results in the upper bounds of the shortest solution h_C for $\psi^{\mathcal{C}}$. If $|h_C(\mathbf{x})| \leq b(\mathbf{x})$ holds for all $\mathbf{x} \in \Gamma$, then $\psi^{\mathcal{C}}$ has no solution and unsatisfiability of ψ' follows from Proposition 6.3. Because ψ and ψ' are equisatisfiable, we can conclude that ψ is unsatisfiable.

Otherwise, we increase the bounds on the variables that occur in $\psi^{\mathcal{C}}$ while keeping bounds on the other variables unchanged: We construct b_{k+1} with $b_k(\mathbf{x}) \leq b_{k+1}(\mathbf{x}) \leq |h_C(\mathbf{x})|$ for all $\mathbf{x} \in \Gamma$, such that $b_k(\mathbf{y}) < b_{k+1}(\mathbf{y})$ holds for at least one $\mathbf{y} \in V(\psi^{\mathcal{C}})$. By strictly increasing at least one variable's bound, we eventually either reach the upper bounds of $\psi^{\mathcal{C}}$ and return unsatisfiability, or we eliminate it as an unsatisfiable implication of ψ . As there are only finitely many possibilities for \mathcal{C} and thus for $\psi^{\mathcal{C}}$, our procedure is guaranteed to terminate.

We do not explicitly construct formula ψ^C to compute bounds on h_C as we know the set $\text{lits}(\psi^C) = \mathcal{C}^+ \cup \mathcal{C}^-$. Finding upper bounds still requires enumerating all subsets of $\text{lits}(\psi^C)$, but we have $|\text{lits}(\psi^C)| \leq |\text{lits}(\psi)|$ and usually $\text{lits}(\psi^C)$ is much smaller than $\text{lits}(\psi)$. For example, consider the formula

$$\psi = z \neq abd \wedge (x \doteq a \vee x \doteq ab^*) \wedge x \doteq y \wedge (y \doteq bbc \vee z \in a(b|c)^*d) \wedge y \doteq ab.^?*$$

which is unsatisfiable for the bounds $b(x) = b(y) = 1$ and $b(z) = 4$. The unsatisfiable core C returned after solving $\llbracket \psi \rrbracket^b$ results in the formula $\psi^C = (x \doteq a \vee x \in ab^*) \wedge x \doteq y \wedge y \in ab.^?*$ containing four literals. Finding upper bounds for ψ^C thus amounts to enumerating just 2^4 subsets, which is substantially less than considering all 2^7 subsets of $\text{lits}(\psi)$ upfront. The conjunction of a subset of $\text{lits}(\psi^C)$ yielding the largest upper bounds is $x \in ab^* \wedge x \doteq y \wedge y \in ab.^?*$, which simplifies to $x \in ab^* \cap ab.^?*$ and has a solution of length at most 2 for x and y . With bounds $b(x) = b(y) = 2$ and $b(z) = 4$, the formula is satisfiable.

7 Implementation

We have implemented our approach in a solver called NFA2SAT. NFA2SAT is written in RUST and uses CADICAL [9] as the backend SAT solver. We use the incremental API provided by CADICAL to solve problems under assumptions. Soundness of NFA2SAT follows from Theorem 5.1. For completeness, we rely on CADICAL's *failed* function to efficiently determine *failed assumptions*, i.e., assumption literals that were used to conclude unsatisfiability.

The procedure works as follows. Given a formula ψ , we first introduce one fresh Boolean selector variable s_l for each theory literal $l \in \text{lits}(\psi)$. Then, instead of adding the encoded definitions of the theory literals directly to the SAT solver, we precede them with their corresponding selector variables: for a positive literal a , we add $s_a \rightarrow (\mathbf{d}(a) \rightarrow \llbracket a \rrbracket)$, and for a negative literal $\neg a$, we add $s_{\neg a} \rightarrow (\neg \mathbf{d}(a) \rightarrow \llbracket \neg a \rrbracket)$ (considering assumptions introduced by $\llbracket a \rrbracket$ as unit clauses). In the resulting CNF formula, the new selector variables are present in all clauses that encode their corresponding definition, and we use them as assumptions for every incremental call to the SAT solver, which does not affect satisfiability. If such an assumption failed, then we know that at least one of the corresponding clauses in the propositional formula was part of an unsatisfiable core, which enables us to efficiently construct the sets \mathcal{C}^+ and \mathcal{C}^- of positive and negative atoms present in the unsatisfiable core. As noted previously, we have $\text{lits}(\psi^C) = \mathcal{C}^+ \cup \mathcal{C}^-$ and hence the sets are sufficient to find bounds on a shortest model for ψ^C .

This approach is efficient for obtaining $\text{lits}(\psi^C)$ but since CADICAL does not guarantee that the set of failed assumptions is minimal, $\text{lits}(\psi^C)$ is not minimal in general. Moreover, even a minimal $\text{lits}(\psi^C)$ can contain too many elements for processing all subsets. To address this issue, we enumerate the subsets only if $\text{lits}(\psi^C)$ is small (by default, we use a limit of ten literals). In this case, we construct the automata M_i used in Theorem 6.1 for each subset, facilitating the techniques described in [7] for quickly ruling out unsatisfiable ones. Otherwise,

instead of enumerating the subsets, we resort to sound approximations of upper bounds, which amounts to over-approximating the number of states without explicitly constructing the automata (c.f. [14]).

Once we have obtained upper bounds on the length of the solution of ψ^C , we increment bounds on all variables involved, except those that have reached their maximum. Our default heuristics computes a new bound that is either double the current bound of a variable or its maximum, whichever is smaller.

8 Experimental Evaluation

We have evaluated our solver on a large set of benchmarks from the ZALIGVINDER [22] repository². The repository contains 120,287 benchmarks stemming from both academic and industrial applications. In particular, all the string problems from the SMT-LIB repository,³ are included in the ZALIGVINDER repository. We converted the ZALIGVINDER problems to the SMT-LIB 2.6 syntax and removed duplicates. This resulted in 82,632 unique problems out of which 29,599 are in the logical fragment we support.

We compare NFA2SAT with the state-of-the-art solvers CVC5 (version 1.0.3) and Z3 (version 4.12.0). The comparison is limited to these two solvers because they are widely adopted and because they had the best performance in our evaluation. Other string solvers either don't support our logical fragment (CERTISTR, WOORPJE) or gave incorrect answers on the benchmark problems considered here. Older, no-longer maintained, solvers have known soundness problems, as reported in [7] and [27].

We ran our experiment on a Linux server, with a timeout of 1200 s seconds CPU time and a memory limit of 16 GB. Table 1 shows the results. As a single tool, NFA2SAT solves more problems than CVC5 but not as many as Z3. All three tools solve more than 98% of the problems.

The table also shows results of portfolios that combine two solvers. In a portfolio configuration, the best setting is to use both Z3 and NFA2SAT. This combination solves all but 20 problems within the timeout. It also reduces the total run-time from 283,942 s for Z3 (about 79 h) to 28,914 s for the portfolio (about 8 h), that is, a 90% reduction in total solve time. The other two portfolios—namely, Z3 with CVC5 and NFA2SAT with CVC5—also have better performance than a single solver, but the improvement in runtime and number of timeouts is not as large.

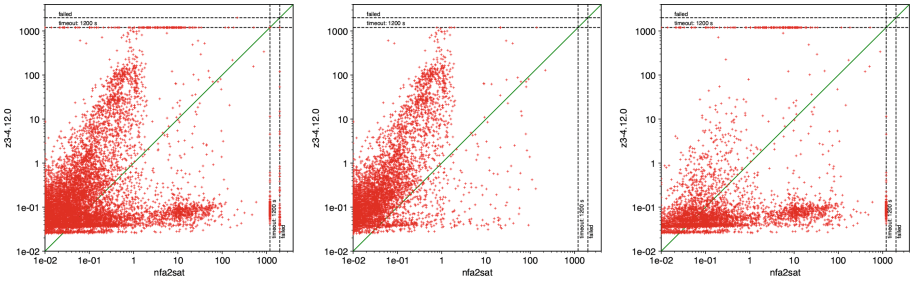
Figure 4a illustrates why NFA2SAT and Z3 complement each other well. The figure shows three scatter plots that compare the runtime of NFA2SAT and Z3 on our problems. The plot on the left compares the two solvers on *all* problems, the one in the middle compares them on *satisfiable* problems, and the one on the right compares them on *unsatisfiable* problems. Points in the left plot are concentrated close to the axes, with a smaller number of points near the diagonal, meaning that Z3 and NFA2SAT have different runtime on most problems. The other two

² <https://github.com/zaligvinder/zaligvinder>.

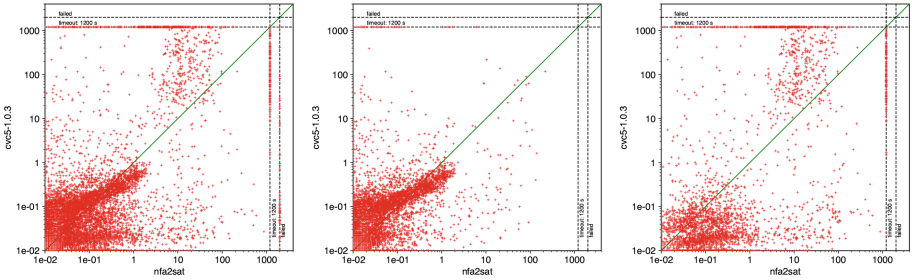
³ https://clc-gitlab.cs.uiowa.edu:2443/SMT-LIB-benchmarks/QF_S.

Table 1. Evaluation on ZALIGVINDER benchmarks. The three left columns show results of individual solvers. The other three columns show results of portfolios combining two solvers.

	CVC5	Z3	NFA2SAT	CVC5 Z3	NFA2SAT CVC5	NFA2SAT Z3
SAT	22895	22927	22922	22934	22934	22934
UNSAT	6259	6486	6405	6526	6598	6645
Timeout	445	185	206	139	67	20
Out-of-memory	0	1	66	n/a	n/a	n/a
Total Solved	29154	29413	29327	29460	29532	29579
Total Runtime (s)	655877	283942	275420	169553	126655	28914



(a) NFA2SAT vs. Z3



(b) NFA2SAT vs. CVC5

Fig. 4. Comparison of runtime (in seconds) with Z3 and CVC5. The left plots include all problems, the middle plots include only satisfiable problems, and the right plots include only unsatisfiable problems. The lines marked “failed” correspond to problems that are not solved because a solver ran out of memory. The lines marked “timeout” correspond to problems not solved because of a timeout (1200 s).

plots show this even more clearly: NFA2SAT is faster on satisfiable problems while Z3 is faster on unsatisfiable problems. Figure 4b shows analogous scatter plots comparing NFA2SAT and CVC5. The two solvers show similar performance on a large set of easy benchmarks although CVC5 is faster on problems that both

solvers can solve in less than 1 s. However, CVC5 times out on 38 problems that NFA2SAT solves in less than 2 s. On unsatisfiable problems, CVC5 tends to be faster than NFA2SAT, but there is a class of problems for which NFA2SAT takes between 10 and 100 s whereas CVC5 is slower.

Overall, the comparison shows that NFA2SAT is competitive with CVC5 and Z3 on these benchmarks. We also observe that NFA2SAT tends to work better on satisfiable problems. For best overall performance, our experiments show that a portfolio of Z3 and NFA2SAT would solve all but 20 problems within the timeout, and reduce the total solve time by 90%.

9 Conclusion

We have presented the first eager SAT-based approach to string solving that is both sound and complete for a reasonably expressive fragment of string theory. Our experimental evaluation shows that our approach is competitive with the state-of-the-art lazy SMT solvers Z3 and CVC5, outperforming them on satisfiable problems but falling behind on unsatisfiable ones. A portfolio that combines our approach with these solvers—particularly with Z3—would thus yield strong performance across both types of problems.

In future work, we plan to extend our approach to a more expressive logical fragment, including more general word equations. Other avenues of research include the adaption of model checking techniques such as IC3 [10] to string problems, which we hope would lead to better performance on unsatisfiable instances. A particular benefit of the eager approach is that it enables the use of mature techniques from the SAT world, especially for proof generation and parallel solving. Producing proofs of unsatisfiability is complex for traditional CDCL(T) solvers because of the complex rewriting and deduction rules they employ. In contrast, efficiently generating and checking proofs produced by SAT solvers (using the DRAT format [32]) is well-established and practicable. A challenge in this respect would be to combine unsatisfiability proofs from a SAT solver with proof that our reduction to SAT is sound. For parallel solving, we plan to explore the use of a parallel incremental solver (such as iLINGELING [9]) as well as other possible ways to solve multiple bounds in parallel.

References

1. Abdulla, P.A., et al.: Trau: SMT solver for string constraints. In: 2018 Formal Methods in Computer Aided Design (FMCAD), pp. 1–5 (2018). <https://doi.org/10.23919/FMCAD.2018.8602997>
2. Backes, J., et al.: Semantic-based automated reasoning for AWS access policies using SMT. In: 2018 Formal Methods in Computer Aided Design (FMCAD), pp. 1–9 (2018). <https://doi.org/10.23919/FMCAD.2018.8602994>

3. Barbosa, H., et al.: CVC5: A versatile and industrial-strength SMT solver. In: Fisman, D., Rosu, G. (eds.) Tools and Algorithms for the Construction and Analysis of Systems - 28th International Conference, TACAS 2022, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2022, Munich, Germany, April 2-7, 2022, Proceedings, Part I. Lecture Notes in Computer Science, vol. 13243, pp. 415–442. Springer (2022). https://doi.org/10.1007/978-3-030-99524-9_24
4. Barrett, C., Fontaine, P., Tinelli, C.: The SMT-LIB Standard: Version 2.6. Tech. rep., Department of Computer Science, The University of Iowa (2017). www.smt-lib.org
5. Berzish, M., et al.: String theories involving regular membership predicates: From practice to theory and back. In: Lecroq, T., Puzynina, S. (eds.) Combinatorics on Words, pp. 50–64. Springer International Publishing, Cham (2021)
6. Berzish, M., et al.: Towards more efficient methods for solving regular-expression heavy string constraints. *Theoretical Computer Science* **943**, 50–72 (2023). <https://doi.org/10.1016/j.tcs.2022.12.009>, <https://www.sciencedirect.com/science/article/pii/S030439752200723X>
7. Berzish, M., et al.: An SMT solver for regular expressions and linear arithmetic over string length. In: Silva, A., Leino, K.R.M. (eds.) Computer Aided Verification, pp. 289–312. Springer International Publishing, Cham (2021)
8. Biere, A.: Bounded model checking. In: Biere, A., Heule, M., van Maaren, H., Walsh, T. (eds.) Handbook of Satisfiability, Frontiers in Artificial Intelligence and Applications, vol. 185, pp. 457–481. IOS Press (2009). <https://doi.org/10.3233/978-1-58603-929-5-457>
9. Biere, A., Fazekas, K., Fleury, M., Heisinger, M.: CaDiCaL, Kissat, Paracooba, Plingeling and Treengeling entering the SAT Competition 2020. In: Balyo, T., Froyls, N., Heule, M., Iser, M., Järvisalo, M., Suda, M. (eds.) Proc. of SAT Competition 2020 - Solver and Benchmark Descriptions. Department of Computer Science Report Series B, vol. B-2020-1, pp. 51–53. University of Helsinki (2020)
10. Bradley, A.R.: SAT-based model checking without unrolling. In: Jhala, R., Schmidt, D. (eds.) VMCAI 2011. LNCS, vol. 6538, pp. 70–87. Springer, Heidelberg (2011). https://doi.org/10.1007/978-3-642-18275-4_7
11. Chen, T., Hague, M., Lin, A.W., Rümmer, P., Wu, Z.: Decision procedures for path feasibility of string-manipulating programs with complex operations. *Proc. ACM Program. Lang.* **3**(POPL) (jan 2019). <https://doi.org/10.1145/3290362>
12. Day, J.D., Ehlers, T., Kulczynski, M., Manea, F., Nowotka, D., Poulsen, D.B.: On solving word equations using SAT. In: Filot, E., Jungers, R., Potapov, I. (eds.) Reachability Problems, pp. 93–106. Springer International Publishing, Cham (2019)
13. Eén, N., Sörensson, N.: Temporal induction by incremental SAT solving. *Electronic Notes in Theoretical Computer Science* **89**(4), 543–560 (2003). [https://doi.org/10.1016/S1571-0661\(05\)82542-3](https://doi.org/10.1016/S1571-0661(05)82542-3), <https://www.sciencedirect.com/science/article/pii/S1571066105825423>, bMC’2003, First International Workshop on Bounded Model Checking
14. Gao, Y., Moreira, N., Reis, R., Yu, S.: A survey on operational state complexity. *CoRR* **abs/1509.03254** (2015), <http://arxiv.org/abs/1509.03254>
15. Hojjat, H., Rümmer, P., Shamakhi, A.: On strings in software model checking. In: Lin, A.W. (ed.) Programming Languages and Systems, pp. 19–30. Springer International Publishing, Cham (2019)

16. Jez, A.: Word Equations in Nondeterministic Linear Space. In: Chatzigiannakis, I., Indyk, P., Kuhn, F., Muscholl, A. (eds.) 44th International Colloquium on Automata, Languages, and Programming (ICALP 2017). Leibniz International Proceedings in Informatics (LIPIcs), vol. 80, pp. 95:1–95:13. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany (2017). <https://doi.org/10.4230/LIPIcs.ICALP.2017.95>, <http://drops.dagstuhl.de/opus/volltexte/2017/7408>
17. Kan, S., Lin, A.W., Rümmer, P., Schrader, M.: Certistr: A certified string solver. In: Proceedings of the 11th ACM SIGPLAN International Conference on Certified Programs and Proofs, pp. 210–224. CPP 2022, Association for Computing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3497775.3503691>
18. Karhumäki, J., Mignosi, F., Plandowski, W.: The expressibility of languages and relations by word equations. *J. ACM* **47**(3), 483–505 (may 2000). <https://doi.org/10.1145/337244.337255>
19. Kiezun, A., Ganesh, V., Guo, P.J., Hooimeijer, P., Ernst, M.D.: Hampi: A solver for string constraints. In: Proceedings of the Eighteenth International Symposium on Software Testing and Analysis, pp. 105–116. ISSTA '09, Association for Computing Machinery, New York, NY, USA (2009). <https://doi.org/10.1145/1572272.1572286>
20. Klieber, W., Kwon, G.: Efficient CNF encoding for selecting 1 from N objects. In: Fourth Workshop on Constraints in Formal Verification (CFV) (2007)
21. Kulczynski, M., Lotz, K., Nowotka, D., Poulsen, D.B.: Solving string theories involving regular membership predicates using SAT. In: Legunsen, O., Rosu, G. (eds.) *Model Checking Software*, pp. 134–151. Springer International Publishing, Cham (2022)
22. Kulczynski, M., Manea, F., Nowotka, D., Poulsen, D.B.: Zaligvinder: A generic test framework for string solvers. *J. Softw.: Evolution and Process* **n/a**(n/a), e2400. <https://doi.org/10.1002/smr.2400>, <https://onlinelibrary.wiley.com/doi/abs/10.1002/smr.2400>
23. Makanin, G.S.: The problem of solvability of equations in a free semi-group. *Math. USSR, Sb.* **32**, 129–198 (1977). <https://doi.org/10.1070/SM1977v032n02ABEH002376>
24. Mora, F., Berzish, M., Kulczynski, M., Nowotka, D., Ganesh, V.: Z3str4: A multi-armed string solver. In: Huisman, M., Păsăreanu, C., Zhan, N. (eds.) *Formal Methods*, pp. 389–406. Springer International Publishing, Cham (2021)
25. de Moura, L., Bjørner, N.: Z3: An efficient SMT solver. In: *Proceedings of the Theory and Practice of Software, 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pp. 337–340. TACAS'08/ETAPS'08, Springer-Verlag, Berlin, Heidelberg (2008)
26. Murray, N.V.: Completely non-clausal theorem proving. *Artificial Intelligence* **18**(1), 67–85 (1982). [https://doi.org/10.1016/0004-3702\(82\)90011-X](https://doi.org/10.1016/0004-3702(82)90011-X), <https://www.sciencedirect.com/science/article/pii/000437028290011X>
27. Nötzli, A., Reynolds, A., Barbosa, H., Barrett, C.W., Tinelli, C.: Even faster conflicts and lazier reductions for string solvers. In: Shoham, S., Vizel, Y. (eds.) *Computer Aided Verification - 34th International Conference, CAV 2022, Haifa, Israel, August 7-10, 2022, Proceedings, Part II. Lecture Notes in Computer Science*, vol. 13372, pp. 205–226. Springer (2022). https://doi.org/10.1007/978-3-031-13188-2_11, https://doi.org/10.1007/978-3-031-13188-2_11
28. Plaisted, D.A., Greenbaum, S.: A structure-preserving clause form translation. *Journal of Symbolic Computation* **2**(3), 293–304 (1986). [https://doi.org/10.1016/S0747-7171\(86\)80028-1](https://doi.org/10.1016/S0747-7171(86)80028-1), <https://www.sciencedirect.com/science/article/pii/S0747717186800281>

29. Plandowski, W.: Satisfiability of word equations with constants is in PSPACE. In: 40th Annual Symposium on Foundations of Computer Science (Cat. No.99CB37039), pp. 495–500 (1999). <https://doi.org/10.1109/SFFCS.1999.814622>
30. Rungta, N.: A billion SMT queries a day (invited paper). In: Shoham, S., Vizel, Y. (eds.) Computer Aided Verification. pp. 3–18. Springer International Publishing, Cham (2022). https://doi.org/10.1007/978-3-031-13185-1_1
31. Saxena, P., Akhawe, D., Hanna, S., Mao, F., McCamant, S., Song, D.: A symbolic execution framework for JavaScript. In: 2010 IEEE Symposium on Security and Privacy, pp. 513–528 (2010). <https://doi.org/10.1109/SP.2010.38>
32. Wetzler, N., Heule, M., Jr., W.A.H.: Drat-trim: Efficient checking and trimming using expressive clausal proofs. In: Sinz, C., Egly, U. (eds.) Theory and Applications of Satisfiability Testing - SAT 2014 - 17th International Conference, Held as Part of the Vienna Summer of Logic, VSL 2014, Vienna, Austria, July 14–17, 2014. Proceedings. Lecture Notes in Computer Science, vol. 8561, pp. 422–429. Springer (2014). https://doi.org/10.1007/978-3-319-09284-3_31, https://doi.org/10.1007/978-3-319-09284-3_31
33. Yu, F., Alkhalaf, M., Bultan, T., Ibarra, O.H.: Automata-based symbolic string analysis for vulnerability detection. *Formal Methods Syst. Design* **44**(1), 44–70 (2014). <https://doi.org/10.1007/s10703-013-0189-1>, <https://doi.org/10.1007/s10703-013-0189-1>

Open Access This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

