

Z3str2: An Efficient Solver for Strings, Regular Expressions, and Length Constraints

**Yunhui Zheng · Vijay Ganesh ·
Sanu Subramanian · Omer Tripp ·
Murphy Berzish · Julian Dolby ·
Xiangyu Zhang**

Received: date / Accepted: date

Y. Zheng, J. Dolby
IBM T. J. Watson Research Center, Yorktown Heights, NY, USA
E-mail: {zhengyu, dolby}@us.ibm.com

O. Tripp
Google, USA
The research leading to this manuscript, and earlier versions of the manuscript, were completed while the author was at IBM Research.
E-mail: trippo@google.com

V. Ganesh, S. Subramanian, M. Berzish
University of Waterloo, Waterloo, ON, Canada
E-mail: {vganesh, ses, mtrberzi}@uwaterloo.ca

X. Zhang
Purdue University, West Lafayette, IN, USA
E-mail: xyzhang@cs.purdue.edu

Abstract In recent years, string solvers have become an essential component in many formal verification, security analysis, and bug-finding tools. Such solvers typically support a theory of string equations, the length function, and the regular-expression membership predicate. These enable considerable expressive power, which comes at the cost of slow solving time, and in some cases even non-termination. We present three techniques, designed for word-based SMT string solvers, to mitigate these problems: (i) detecting overlapping variables, which is essential to avoiding common cases of non-termination; (ii) pruning of the search space via bi-directional integration between the string and integer theories, enabling new cross-domain heuristics; and (iii) a binary search based heuristic, allowing the procedure to skip unnecessary string length queries and converge on consistent length assignments faster for large strings.

We have implemented above techniques atop the Z3-str solver, resulting in a significantly more robust and efficient solver, dubbed Z3str2, for the quantifier-free theory of string equations, the regular-expression membership predicate, and linear arithmetic over the length function. We report on a series of experiments over four sets of challenging real-world benchmarks, where we compare Z3str2 with five different string solvers: S3, CVC4, Kaluza, PISA and Stranger. Each of these tools utilizes a different solving strategy and/or string representation (based e.g. on words, bit vectors or automata). The results point to the efficacy of our proposed techniques, which yield dramatic performance improvement. We also demonstrate performance improvements enabled by Z3str2 in the context of symbolic execution for string-manipulating programs. We observe that the techniques presented here are of broad applicability, and can be integrated into other string solvers to improve their performance.

1 Introduction

Reasoning over strings is becoming increasingly important due to security threats imposed by improper handling of untrusted string values [8, 23, 32]. In response, different powerful string solvers have been developed, including HAMPI [23], Kaluza [33], PISA [35], Stranger [38], CVC4 [26], S3 [36], Norn [6] and Z3-str [41]. These tools primarily solve the satisfiability problem over string (aka word) equations, with some of them also providing support for regular expression (RE) membership predicates and linear arithmetic over the length function (e.g. [26, 36]). While these tools have improved dramatically in recent years, the demand for even more efficient solvers continues to grow unabated as software continues to grow in size and complexity.

Motivated by this need for efficient string solvers, we present two new techniques to solve combined string, regular expression, and integer constraints. These techniques are applicable primarily to SMT solvers that treat strings without abstractions or representation conversions, which we refer to collectively as *word-based string solvers*. Examples of such solvers include the Z3-str, CVC4 and S3 string solvers.

For the sake of completeness, we compare and contrast our techniques against solvers that use automata (e.g., PISA and Stranger) and bit-vector (e.g., Kaluza) string representations. Word-based string solvers have several important advantages: First, unlike bit-vector-based solvers, they can precisely model unbounded strings and string equalities without over-approximation, a feature that is crucial to string analysis of web

applications. Second, by modeling strings and length in native domains, word-based string solvers can leverage the state of the art in integer constraint solving, and further enable hybrid techniques via powerful SMT engines. Finally, such solvers can take advantage of well-developed application-specific rewrite rules.

At the same time, a fundamental problem of word-based string solvers (unlike those based on bit vectors, which impose a finite domain) is that it is unclear, at present, whether the satisfiability problem for the quantifier-free theory of word equations, regular-expression membership predicate, and length function is decidable [28]. All current practical string solvers suffer from incompleteness and nontermination. Addressing these problems is of primary importance, calling for new techniques to effectively explore the solution space. In light of this motivation, we have developed three techniques that address the problems of nontermination and search-space explosion.

First, a well-known reason for nontermination is overlapping variables [14, 41, 6], as we illustrate with the equation $a \cdot X = X \cdot b$, where a, b are constant strings and X is a string variable. Stated intuitively, the solution for X has to be in the form of $a \cdot X_1 \cdot b$, where X_1 is a string variable. The reduction step results in $a \cdot X_1 = X_1 \cdot b$, which is in the same form as the original equation, and thus leads to nontermination. However, this equation is obviously unsatisfiable. We revisit this example in section 4, which highlights the need for a robust procedure to detect overlapping variables.

The second technique, given the tight interplay between string and integer values (in index-sensitive string operations), is bi-directional solver-level integration between the string and integer theories. This can be leveraged to drastically reduce the search space for typical constraints obtained from practical applications.

The last technique is a binary search based heuristic that helps to achieve faster convergences in selecting consistent length assignment for large strings. As we observe, a huge number of length assignment polling iterations are redundant. We adapt a binary search algorithm and search the unbounded length space more efficiently.

We have implemented all of these techniques atop the Z3-str solver as the Z3str2 solver for the satisfiability problem over a quantifier-free theory of word equations, regular-expression membership predicate as well as the length function. We report on a comprehensive set of experiments that validate the efficacy of our proposed techniques by comparing Z3str2 with Kaluza, PISA, Stranger, S3 and CVC4 over four sets of benchmarks derived from the real world.¹ Besides, we illustrate performance improvements in the context of symbolic execution of string manipulating programs. We show that we are able to avoid exploring a large number of execution paths while precisely capturing the semantics of common string utility functions. We emphasize that our techniques are applicable also to other word-based string solvers such as S3 and CVC4.

Contributions To summarize, this paper makes the following principal contributions:

1. Formal representation of solving word equations: We use a boundary label based representation to formally model the key steps of solving word equations.

2. Guided search: We present three techniques designed for string solvers that treat strings as primitive types. The first is a sound and complete method to detect overlapping variables, which improves performance and avoids exploration of certain

¹ The Z3str2 code, as well as the data pertaining to our experiments, are all available at [1].

paths that may lead to nontermination. The second technique is a two-way integration between the string and integer theories, which enables effective pruning based on cross-domain heuristics. The last one is a binary search based pruning heuristic that allows the procedure to skip redundant string length assignment queries.

3. Z3-str integration: We have integrated all of the aforementioned techniques into the core solving algorithm of Z3-str. We describe the architecture of the resulting tool, Z3str2, and prove its soundness.

4. Symbolic execution performance improvements: We demonstrate that we can significantly simplify the constraint encoding/solving procedure and improve the performance of symbolic execution on string manipulating programs.

5. Experimental study: To validate the efficacy of our techniques, we have conducted a comprehensive set of experiments where we compare Z3str2 against five solvers — namely, S3, CVC4, Kaluza, PISA and Stranger — on four benchmark suites. The results show that Z3str2 is significantly faster than competing solvers (often by orders of magnitude) in all but few cases.

2 Preliminaries

In this section, we briefly introduce the syntax and semantics of the quantifier-free theory of word equations, regular expression membership, and length over integers. This theory will be the basis of other formalizations presented in this paper.

Syntax of Word Equations, RE membership, and Length. We fix a disjoint two-sorted set of variables $var = var_{str} \cup var_{int}$; var_{str} consists of string variables, denoted X, Y, S, \dots and var_{int} consists of integer variables, denoted m, n, \dots . We also fix a two-sorted set of constants $Con = Con_{str} \cup Con_{int}$. Moreover, $Con_{str} \subset \Sigma^*$ for some finite alphabet, Σ , whose elements are denoted f, g, \dots . Elements of Con_{str} will be referred to as *string constants* or *strings*. Elements of Con_{int} are nonnegative integers². The empty string is represented by ϵ , and has length 0. Terms may be string terms or integer terms. A string term is either an element of var_{str} , an element of Con_{str} , or a concatenation of string terms (denoted by the function *concat* or interchangeably by \cdot). An integer term is an element of var_{int} , an element of Con_{int} , the length function applied to a string term, or an addition, subtraction, multiplication, or division of two integer terms. The theory contains three types of atomic formulas, namely, word equations, length constraints, and regular expression (RE) membership predicates. REs are defined inductively, where constants and the empty string form the base case, and the operations of concatenation, alternation, and Kleene star are used to build up more complicated expressions (see details in [17]). REs may not contain variables. Z3str2 supports a list of common string-related operators such as *charAt*, *contains*, *startswith*, *endswith*, *indexOf*, *lastIndexOf*, *substring* and etc. They are desugared to word equations with length functions. Formulas are defined induc-

² This restriction is made only for consistency with the SMT-LIB standard [7] in a typical implementation. The only integer literals allowed in the SMT-LIB syntax are 0 and positive integers. Note that both the SMT-LIB standard and the syntax we give here still allow negative integer constants to be expressed as, for example, $(- 0 3)$ for -3 .

<i>Term:bool</i>	::=	Var:bool true false $(\text{Contains } \text{Term:string } \text{Term:string})$ $(\text{StartsWith } \text{Term:string } \text{Term:string})$ $(\text{EndsWith } \text{Term:string } \text{Term:string})$ $(\text{RegexIn } \text{Term:string } \text{Term:regex})$
<i>Term:int</i>	::=	Var:int Number $(\{+, -, \times, \div\} \text{Term:int } \text{Term:int})$ $(\text{Length } \text{Term:string})$ $(\text{IndexOf } \text{Term:string } \text{Term:string})$ $(\text{IndexOf2 } \text{Term:string } \text{Term:string } \text{Term:int})$ $(\text{LastIndexOf } \text{Term:string } \text{Term:string})$
<i>Term:string</i>	::=	Var:string ConstStr $(\text{Concat } \text{Term:string } \text{Term:string})$ $(\text{Substring } \text{Term:string } \text{Term:int } \text{Term:int})$ $(\text{Replace } \text{Term:string } \text{Term:string } \text{Term:string})$ $(\text{CharAt } \text{Term:string } \text{Term:int})$
<i>Term:regex</i>	::=	$(\text{Str2Regex } \text{ConstStr:string})$ $(\text{RegexStar } \text{Term:regex})$ $(\text{Term:regex})^+$ $(\text{Term:regex})^?$ $(\text{RegexConcat } \text{Term:regex } \text{Term:regex})$ $(\text{RegexUnion } \text{Term:regex } \text{Term:regex})$
<i>Expr:bool</i>	::=	Term:bool $(= \text{Term:bool } \text{Term:bool})$ $(\text{not } \text{Expr:bool})$ $(\text{and } \text{Expr:bool } \text{Expr:bool})$ $(\text{or } \text{Expr:bool } \text{Expr:bool})$ $(\text{ite } \text{Expr:bool } \text{Expr:bool } \text{Expr:bool})$ $(\text{implies } \text{Expr:bool } \text{Expr:bool})$ $(\{<, \leq, =, \geq, >\} \text{Term:int } \text{Term:int})$ $(= \text{Term:string } \text{Term:string})$
<i>Assertion</i>	::=	$(\text{assert } \text{Expr:bool})$

Fig. 1: Constraint syntax

tively over atomic formulas and are quantifier-free. Figure 1 summarizes the syntax of the constraint language.

Semantics of Word Equations, RE membership, and Length. For a word w , $\text{len}(w)$ denotes the length of w . The universe of discourse for the str sort is the set of strings Σ^* , and for the int sort it is the set of natural numbers. For a word equation $t_1 = t_2$, we refer to t_1 as the left hand side (LHS), and t_2 as the right hand side (RHS). We fix a finite alphabet Σ of characters over which strings are defined. Given a formula θ , an *assignment* for θ (with respect to Σ) is a map from the set of variables in θ to $\Sigma^* \cup \mathbb{N}$ (where string variables are mapped to strings and integer variables are mapped to numbers). Given such an assignment, θ can be interpreted as an assertion about Σ^* and \mathbb{N} . If this assertion is true, then θ is *satisfiable* or SAT. A formula with no satisfying assignment is *unsatisfiable* or UNSAT. Two formulas θ, ϕ are *equisatisfiable*

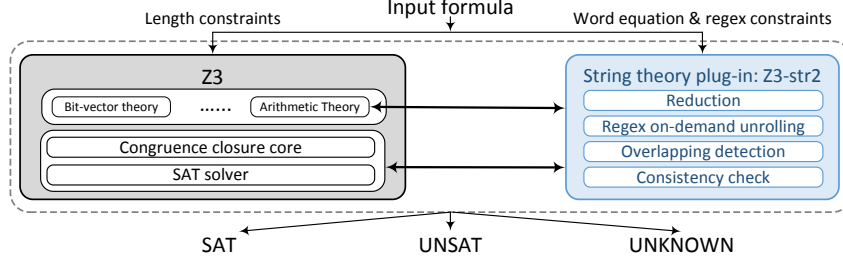


Fig. 2: Architecture of Z3str2.

if the assertion “ θ is SAT iff ϕ is SAT” holds. The *satisfiability problem* for a set S of formulas is to decide whether any given formula in S is SAT or not.

Soundness, Completeness, Termination, Decidability

A satisfiability procedure is *sound* if, whenever the procedure answers UNSAT, then the input is indeed unsatisfiable. Conversely, a satisfiability procedure is *complete* if, whenever the input is satisfiable, the procedure will answer SAT.

An algorithm is *terminating* if it returns an answer in finite time on every input.

A satisfiability procedure is a *decision procedure* if and only if it is sound, complete, and terminating. The satisfiability problem for a set of formulas is decidable if and only if there exists a decision procedure that solves its satisfiability problem.

3 Overview of the Z3str2 String Solver

The Z3str2 solver is a string plug-in built into the Z3 SMT Solver [12], with an efficient integration between the string plug-in and Z3’s integer solver. As can be seen from the architectural schematic of the Z3str2 string solver given in Figure 2 (the algorithm is explained in detail in Section 6.2.3), the first step is to purify the input into two parts, namely string constraints (word equations and RE membership) on the one hand, and integer linear arithmetic constraints over the length function on the other hand. Next, the word equations and RE constraints are given as input to the string plug-in. The plug-in may consult the Z3 core to detect equivalent terms. The word equations are solved using an algorithm described in detail in Section 4 below. The RE constraints are solved by unrolling as described also in Section 4. The length constraints are converted into a system of pure integer linear arithmetic inequations and solved using Z3’s integer solver. During the solving process, the string plug-in may generate length constraints that are incrementally added on demand to Z3’s integer solver, which are regularly checked for consistency with both the input length constraints and any previously generated ones.

On any well-formed input, as described in Section 2, Z3str2 may return SAT, UNSAT or UNKNOWN. Note that while Z3str2 can handle a boolean combination of atomic formulas, we refer only to conjunctions of literals in the rest of paper without loss of generality. Z3str2, like other competing solvers such as CVC4, is sound but not complete. If either Z3’s integer solver or our string plug-in determines that their

respective purified inputs are UNSAT, Z3str2 reports UNSAT. If the string plug-in detects that the input equations have complicated overlaps that its heuristics cannot handle, it reports UNKNOWN. This is a source of incompleteness in Z3str2’s implementation.

The third and only remaining possibility is that both Z3’s integer solver and the string plug-in determine that their respective purified inputs are SAT. However, this does not necessarily mean that the input is SAT. It could be that the solution produced by the integer solver is inconsistent with the solution produced by the string solver, or vice versa. For example, the integer solver may say that a particular string variable, say X , has length equal to 1, while the string plug-in might produce a specific assignment for X that is of length equal to 2. One way to overcome this problem is to iterate through all possible solutions until a consistent one is found, assuming one exists. However, given that the domain of strings and natural numbers is infinite, it is possible that such an iterative procedure may loop forever in the event there are no consistent solutions to be found. In other words, if the input is indeed SAT, the procedure discussed here will correctly establish consistency and determine that the input is SAT. Unfortunately, it is possible that, when the input is in fact UNSAT, both the integer solver and string plug-in may determine that their respective purified inputs are SAT and may then loop forever searching for a combined consistent solution. They are obviously not going to find a combined consistent solution in such cases given that the input is in fact UNSAT, and hence the iterative procedure may not terminate.

The above-described problem of non-termination due to the interaction between the integer and string parts of the theory is not specific to Z3str2. In fact, the problem of deciding the satisfiability problem for the quantifier-free theory of word equations and length function remains open after decades of research and is a major open problem in mathematical logic [28]. In conclusion, if Z3str2 reports that the input is UNSAT, then indeed the input is UNSAT (soundness). However, the converse is not necessarily true, i.e., just like all other competing practical solvers, Z3str2 is not complete.

4 Word Equation Sub-solver in Z3str2

In this subsection, we focus on the word equation solving component of Z3str2. Starting with the work of Makanin [27], many decision procedures for word equations have been proposed [30, 34, 19]. While most procedures are not accompanied by practical implementations, they are a rich source of ideas for all the solvers that have recently been implemented. For example, the Z3str2 solver is based on the concepts of *boundary labels*, *generalized word equations and arrangements* (discussed in greater detail in subsection 6), that have their roots in the very first decision procedure for word equations by Makanin.

The key technique used by Z3str2 to solve a word equation W is to recursively convert W equisatisfiably into a disjunction of conjunctions of simpler equations, which we call “arrangements”. These arrangements are computed by aligning the concatenation function on the LHS and RHS of a given equation such that an occurrence of concatenation function in the LHS (resp. RHS) may “split” or “cut” variables on the RHS (resp. LHS).

As an illustration, consider the following formula composed of three equations: $Z = X \cdot Y \wedge Z = W \cdot c \wedge c \cdot Y = c \cdot b \cdot c$, where X, Y, Z and W are string variables, and b and c are characters. A simple rewriting is the following:

$$\begin{aligned} Z = X \cdot Y &\implies Z_1 = X \wedge Z_2 = Y \quad [1.1] \\ Z = W \cdot c &\implies Z_1 = W \wedge Z_2 = c \quad [1.2] \\ c \cdot Y = c \cdot b \cdot c &\implies Y = b \cdot c \quad [1.3] \end{aligned}$$

Observe that Z is split into Z_1 and Z_2 , which are constrained differently. However, this rewriting is not satisfiable because $Y = c$ from equations [1.1] and [1.2], and $Y = b \cdot c$ from equation [1.3]. Now observe that the alignment described above is not the only one we can consider. Below is a different alignment that leads to a new splitting and in fact yields a satisfying assignment:

$$Z = W \cdot c \implies Z_1 = W_1^{[1.4]} \wedge Z_2 = W_2 \cdot c^{[1.5]}$$

The difference now is that W is also split (into W_1 and W_2), and hence this splitting yields a satisfying assignment. In particular, from [1.1], [1.3] and [1.5], we have $W_2 = b$. Also note that X, Z_1 and W_1 become free variables as they are all equivalent but not constrained by any other variable.

What the above example highlights is that there are many different alignments of variable boundaries in the LHS (resp. RHS) that can split variables in the RHS (resp. LHS). We call every such alignment an *arrangement*. Here is the crucial fact about word equations: every equation can be equisatisfiably rewritten into a finite set of arrangements, where each arrangement is a finite set of word equations obtained from the splitting procedure described above. The Z3str2 solver exploits this fact, and solves word equations by converting them into finite sets of arrangements and inspecting each one individually to see if they are satisfiable. The input word equation is SAT if and only if at least one arrangement is SAT. This in a nutshell is how the word equations are solved by the Z3str2 solver, i.e., by recursively converting equations into a disjunction of arrangements (where each arrangement is a simpler set of equations) until a set of arrangements is derived where the satisfiability can be determined purely via inspection.

5 Supporting Regular Expression Membership Predicates

A RE membership predicate $X \in \mathcal{R}$ is reduced to word equations by a transformation function $\rho(X, \mathcal{R})$, where X is a string variable and \mathcal{R} is a regular expression. The function is defined as follows:

$$\begin{aligned} \rho(X, s) &::= X = s, \text{ where } s \text{ is a constant string} \\ \rho(X, \mathcal{R}_1 | \mathcal{R}_2) &::= \rho(X, \mathcal{R}_1) \vee \rho(X, \mathcal{R}_2) \\ \rho(X, \mathcal{R}_1 \cdot \mathcal{R}_2) &::= X = T_1 \cdot T_2 \wedge \rho(T_1, \mathcal{R}_1) \wedge \rho(T_2, \mathcal{R}_2) \\ \rho(X, \mathcal{R}^*) &::= X = \text{unroll}(\mathcal{R}, n) \wedge n \geq 0 \end{aligned}$$

where n is a fresh integer variable for each Kleene star operation; T_1 and T_2 are fresh string variables; $\text{unroll}(\mathcal{R}, n)$ represents the expression obtained by unrolling \mathcal{R} n times. After the RE membership predicates are replaced by word equations, the string solver proceeds as usual. When the solver explores various arrangements, the $\text{unroll}()$ functions are further simplified by the following rules.

$$\begin{aligned}
X = \text{unroll}(\mathcal{R}, n_1) &::= \text{if } (n_1 = 0) \text{ then } \{X = \epsilon\} \text{ else } \{X = T_3 \cdot \text{unroll}(\mathcal{R}, n_1 - 1) \wedge \rho(T_3, \mathcal{R})\} \\
X \cdot Y = \text{unroll}(\mathcal{R}, n_1) &::= \text{if } (n_1 = 0) \text{ then } \{X = \epsilon \wedge Y = \epsilon\} \text{ else } \{X = \text{unroll}(\mathcal{R}, n_2) \cdot T_3 \\
&\quad \wedge Y = T_4 \cdot \text{unroll}(\mathcal{R}, n_3) \wedge n_1 = n_2 + n_3 + 1 \wedge \rho(T_3 \cdot T_4, \mathcal{R})\}
\end{aligned}$$

Note that \mathcal{R} is essentially unrolled once in the else branch of both rules. Just like in other existing solvers that support RE, the unrolling process may not terminate, especially when there are no length constraints associated with the involved variables. We hence rely on a timeout mechanism.

6 Solving Word Equations and Overlapping Variable Detection

In this section, we formally explain the key steps of solving word equations. While simple and elegant for typical equations obtained from program analysis, the word equation solver described here may fall into infinite loops under certain circumstances.

We first use two examples to show the challenges introduced by overlapping variables. In fact, the problem of “overlapping” variables described below is recognized by logicians as the crucial source of complexity in solving word equations. Then, we describe how we solve word equations and avoid the nontermination brought by overlapping variables.

6.1 Examples that Highlight the Crucial Overlap Detection Problem

Example 1. Formula reduction may fall into infinite loops. Consider the following formula which was introduced in Section 1.

$$a \cdot X = X \cdot b$$

where X is a string variable.

In this example, X appears both as the LHS suffix and as the RHS prefix. This equation is not satisfiable. However, if we solve it using the word equation procedure described in Section 4, then the algorithm will not terminate. In particular, the equation has three arrangements. The first arrangement is where $X = \epsilon$, resulting in the equation $a = b$ which is unsatisfiable. The second arrangement is where the concatenation function in the LHS and RHS align exactly such that we get $X = a \wedge X = b$. The third arrangement is where the LHS occurrence of X cuts the RHS occurrence in the RHS illustrated in Fig. 3.

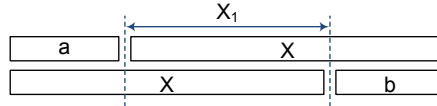


Fig. 3: Graphical representation of the input formula. The top row is the LHS and the bottom row is the RHS. Each box represents a substring. The diagram means that the two occurrences of X overlap.

Note that the suffix of X that occurs in the RHS (bottom part of Fig. 3) of the equation overlaps with the prefix of X in the LHS (top part of Fig. 3). We represent this overlapping part with a new variable X_1 . By applying some simple rewrites we derive the following:

$$a \cdot X = X \cdot b \implies X = a \cdot X_1^{[2.1]} \wedge X = X_1 \cdot b^{[2.2]}$$

From [2.1] and [2.2], we can infer $a \cdot X_1 = X_1 \cdot b$. Note this derived equation has the same form as the input formula. As a result, the above-mentioned procedure will not terminate, unless some steps are taken to detect such “overlaps” and determine their satisfiability without computing arrangements ad infinitum. \square

One could imagine heuristics to detect and handle relatively simple overlaps described above. However, in general, when many equations are involved with variables overlapping indirectly, the problem is not easy to detect or decide.

Example 2. The previous example shows a simple non-terminating case. In practice, non-terminating cases can be very complex, involving multiple equations. *Consider the following formula.*

$$X \cdot M = N \cdot Y \quad \wedge \quad Y \cdot K = L \cdot X$$

where X, Y, M, N, K , and L are all string variables.

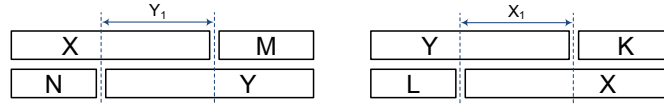


Fig. 4: Graphical representation of one arrangement of the input formula in Example 2. Both X and Y occur twice and they overlap in different ways.

The following rewriting shows why a naive solving procedure will not terminate for such an input.

Similarly to the arrangements presented in Example 1, Fig. 4 shows one possibility of both equations. The following formulas are the derived equations corresponding to the graphical representations in Fig. 4.

$$\begin{aligned} X \cdot M = N \cdot Y &\implies X = N \cdot Y_1^{[3.1]} \wedge Y = Y_1 \cdot M^{[3.2]} \\ Y \cdot K = L \cdot X &\implies Y = L \cdot X_1^{[3.3]} \wedge X = X_1 \cdot K^{[3.4]} \end{aligned}$$

The smaller equations can be rewritten as follows.

$$\begin{aligned} [3.1] \wedge [3.4] &\implies N \cdot Y_1 = X_1 \cdot K^{[3.5]} \\ [3.2] \wedge [3.3] &\implies L \cdot X_1 = Y_1 \cdot M^{[3.6]} \end{aligned}$$

Formula [3.7] and [3.8] can be derived from an arrangement of [3.5]. Similarly, [3.9] and [3.10] represent one possible arrangement of [3.6].

$$\begin{aligned}
[3.5] &\implies X_1 = N \cdot Y_2^{[3.7]} \wedge Y_1 = Y_2 \cdot K^{[3.8]} \\
[3.6] &\implies Y_1 = L \cdot X_2^{[3.9]} \wedge X_1 = X_2 \cdot M^{[3.10]}
\end{aligned}$$

The following new equivalences are derived.

$$\begin{aligned}
[3.7] \wedge [3.10] &\implies X_2 \cdot M = N \cdot Y_2^{[3.11]} \\
[3.8] \wedge [3.9] &\implies Y_2 \cdot K = L \cdot X_2^{[3.12]}
\end{aligned}$$

Observe that formula $([3.11] \wedge [3.12])$ is in the exact same form as the input. Since the loop appears after more than one round of reduction, we use the term “*multiple layer overlaps*” to denote such cases. A key contribution of our procedure is to detect and handle such overlaps. \square

In fact, overlapping variables get to the crux of the difficulty of solving word equations, for otherwise simple rewrites can solve such equations. Hence, any solution to detecting overlaps and deciding such equations is of universal value, and can be used as a subroutine by different types of string solvers. In the following subsections, we discuss a solution to the overlap detection problem.

6.2 Formal Description

Here we formally describe the procedure for solving word equations.

<i>Formula</i> \mathbb{F}	$:= \mathbb{E}_1 \wedge \mathbb{E}_2 \wedge \dots \mathbb{E}_n$
<i>WordEquation</i> \mathbb{E}	$:= \mathcal{W}_1 = \mathcal{W}_2$
<i>Character</i> c	$:= \{a, b, \dots\}$
<i>String Variable</i> X	$:= \{X, Y, \dots\}$
<i>Word</i> \mathcal{W}	$:= {}^{L_\varepsilon} {}^{L_c} \mathcal{W} {}^{L_X} \mathcal{W}$
<i>Label</i> l	$:= \{\triangleright_1^X, \triangleleft_1^X, \triangleright_2^c, \triangleleft_2^c, \dots\}$
<i>LabelSet</i> L	$:= \mathcal{P}(\text{Label})$

Fig. 5: Definitions.

Fig. 5 recursively defines the formulas in the quantifier-free theory of word equations that are input to the procedure we describe later. These definitions describe the syntax of equations. A formula \mathbb{F} is a conjunction of a set of word equations \mathbb{E} ³. An equation in \mathbb{E} dictates the equivalence of two words. A word \mathcal{W} is a sequence of string variables (e.g. X) or constant characters (e.g. c) delimited by label sets. Labels uniquely identify the boundaries of an occurrence of a character or a variable. Note that a character (variable) may occur multiple times in a formula. We use \triangleright_1^X and \triangleleft_1^X

³ While we do support AND-OR combination of word equations, without loss of generality it is sufficient to describe the procedure only for a conjunction of word equations.

to describe the left and right labels of the first occurrence of X respectively, following a pre-defined initial order. Below we show a sample equation.

$$\{\triangleright_1^a\}_a \{\triangleleft_1^a, \triangleright_1^X\} X \{\triangleleft_1^X\} = \{\triangleright_2^X\} X \{\triangleleft_2^X, \triangleright_1^b\} b \{\triangleleft_1^b\} \quad [\mathbb{E}_1]$$

Note that labels are introduced to facilitate formal definition of the procedure. The input equations provided by users do not contain any labels. Henceforth, we represent concatenation explicitly through the ‘ \cdot ’ operator or implicitly as in the definition of *Word*, if doing so does not cause any misunderstanding.

$LabelArrangement \mathcal{A} := \overline{\mathcal{P}(Label)}$
ARRGMTPRODUCT $\otimes : LabelArrangement \times LabelArrangement \rightarrow \mathcal{P}(LabelArrangement)$
$L_1 \cdot L_2 \cdot \dots \cdot L_n \otimes R_1 \cdot R_2 \cdot \dots \cdot R_m := L_1 \cup R_1 \cdot (L_2 \cdot \dots \cdot L_{n-1} \tilde{\otimes} R_2 \cdot \dots \cdot R_{m-1}) \cdot L_n \cup R_m$ $L_1 \cdot L_2 \cdot \dots \cdot L_n \tilde{\otimes} R_1 \cdot R_2 \cdot \dots \cdot R_m := L_1 \cdot (L_2 \cdot \dots \cdot L_n \tilde{\otimes} R_1 \cdot \dots \cdot R_m) \cup (L_1 \cup R_1) \cdot (L_2 \cdot \dots \cdot L_n \tilde{\otimes} R_2 \cdot \dots \cdot R_m) \cup R_1 \cdot (L_1 \cdot \dots \cdot L_n \tilde{\otimes} R_2 \cdot \dots \cdot R_m)$
ARRGMTMERGE $\Downarrow : LabelArrangement \times LabelArrangement \rightarrow \mathcal{P}(LabelArrangement)$
$\mathcal{A}_1 \Downarrow \mathcal{A}_2 := \{\mathcal{A}_x = L_1 \cdot \dots \cdot L_r \in \mathcal{A}_1 \otimes \mathcal{A}_2 \mid \mathcal{A}_x \text{ satisfies the following conditions :}$ (i) [Neighboring] $\forall \triangleright_i^c, \triangleleft_i^c \in \mathbb{F}, \triangleright_i^c \in L_t \rightarrow \triangleleft_i^c \in L_{t+1}$ (ii) [Consistency] $\neg \exists \triangleright_i^c, \triangleleft_i^c, \triangleright_j^d, \triangleleft_j^d \in \mathbb{F} \text{ and } c \neq d, \text{ s.t. } \{\triangleright_i^c, \triangleright_j^d\} \subset L_t \vee \{\triangleleft_i^c, \triangleleft_j^d\} \subset L_t$ (iii) [Uniqueness] $\neg \exists l \in \mathbb{F}, \text{ s.t. } l \in L_i \wedge l \in L_j \wedge i \neq j$ (iv) [Non-overlapping] $\neg \exists \triangleright_n^X, \triangleleft_n^X, \triangleright_m^X, \triangleleft_m^X \in \mathbb{F} \text{ and } m \neq n, \text{ s.t. } \triangleright_n^X \in L_i \wedge \triangleleft_n^X \in L_k \wedge \triangleright_m^X \in L_j \wedge i < j < k$
GETARRGMTFROMWORD $: Word \rightarrow LabelArrangement$
$GETARRGMTFROMWORD({}^L c \mathcal{W}_t) := L \cdot GETARRGMTFROMWORD(\mathcal{W}_t)$ $GETARRGMTFROMWORD({}^L X \mathcal{W}_t) := L \cdot GETARRGMTFROMWORD(\mathcal{W}_t)$ $GETARRGMTFROMWORD({}^L \varepsilon) := L$
PROJECTARRGMTTOVAR $\Downarrow : LabelArrangement \times StringVariable \rightarrow \mathcal{P}(LabelArrangement)$
$\Downarrow_X (L_1 \cdot \dots \cdot L_n) := \{L_t \cdot \dots \cdot L_{t+m} \mid \exists R_1 X R_2 \in \mathbb{F}, \text{ s.t. } R_1 \cap L_t \neq \emptyset \wedge R_2 \cap L_{t+m} \neq \emptyset\}$

Fig. 6: Arrangement Primitives.

6.2.1 Label Arrangements

Boundary labels are particularly important for our technique. Intuitively, we leverage them to reason about the relative positions of the sub-parts in words, such that we can reduce the original equations to a set of smaller equations for the corresponding sub-parts until the equations become so fine-grained that the solution can be directly inferred. The input system of equations is UNSAT if none of the possible break-downs leads to valid solutions. In this subsection, we explain how our technique manipulates labels. We discuss how to split equations to smaller ones based on the arrangements and how to determine if equations are in *solvable form* in the next subsection.

We now formalize the concept of *label arrangement* to facilitate our discussion. A *label arrangement*, or *arrangement in short*, is a sequence of label sets as defined in Fig. 6. Each label set may contain labels from multiple words, such as the LHS and RHS words of an equation. Two labels in the same set imply the corresponding boundaries *align*. For example, assume the following equation $[\mathbb{E}_2]$.

$$\{\triangleright_1^a\} a \{\triangleleft_1^a, \triangleright_1^Y\} Y \{\triangleleft_1^Y\} = \{\triangleright_1^X\} X \{\triangleleft_1^X, \triangleright_1^b\} b \{\triangleleft_1^b\} \quad [\mathbb{E}_2]$$

One arrangement of its LHS and RHS words is the following.

$$\{\triangleright_1^a, \triangleright_1^X\} \cdot \{\triangleleft_1^a, \triangleleft_1^X, \triangleright_1^Y, \triangleright_1^b\} \cdot \{\triangleleft_1^Y, \triangleleft_1^b\} \quad [\mathcal{A}_1]$$

The corresponding alignment (with the labels) can be graphically illustrated in Fig. 7(a).

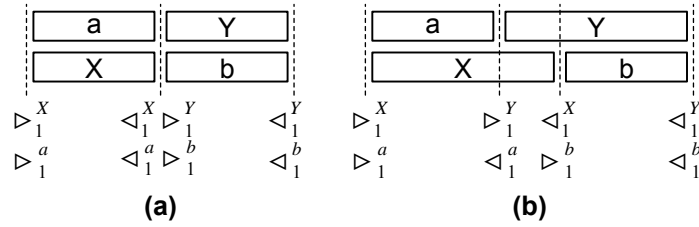


Fig. 7: Sample arrangements of $[\mathbb{E}_2]$

From this arrangement, we can easily derive two smaller equations $X = a$ and $Y = b$, which directly give a solution.

Arrangement Production. Given the LHS and RHS words of an equation, there are multiple (finite) ways to arrange their labels. For example, the following represents another arrangement of equation $[\mathbb{E}_2]$.

$$\{\triangleright_1^a, \triangleright_1^X\} \cdot \{\triangleleft_1^a, \triangleright_1^Y\} \cdot \{\triangleleft_1^X, \triangleright_1^b\} \cdot \{\triangleleft_1^Y, \triangleleft_1^b\} \quad [\mathcal{A}_2]$$

It corresponds to the following graphical alignment shown in Fig. 7(b). In the graphical representation, we can observe that the (left) boundary of Y “cut” the variable X into two parts which are hence constrained differently. We refer to this by saying that “ Y cuts X ”. In fact, we can consider the labels of a word in the initial formula to denote an initial arrangement.

Operation **GETARRGMTFROMWORD**(\mathcal{W}) in Fig. 6 describes how the initial arrangement of \mathcal{W} is recursively computed. It acquires the boundary label sets from left to right until the terminal symbol ε is encountered.

Next, we define the product operation \otimes that produces a set of possible arrangements from two arrangements. Intuitively, one can understand that given the initial arrangements of the LHS and RHS words of an equation, the operator produces the possible label set arrangements, each denoting a possible alignment of the boundaries of the two words. According to the definition in Fig. 6, the operator first unions the two left-most label sets and the two right-most label sets, corresponding to aligning the word boundaries. Note that there is only one way to align word boundaries. It then uses

a recursive operator $\tilde{\otimes}$ to arrange the internal label sets. The definition of $\tilde{\otimes}$ means that given the head label sets L_1 and R_1 in the two input arrangements, respectively, the head label set of the resulting arrangement can be L_1 , R_1 , or $L_1 \cup R_1$.

For our aforementioned equation $[\mathbb{E}_2]$, the product of the LHS and RHS word arrangements

$$\{\triangleright_1^a\} \cdot \{\triangleright_1^Y, \triangleleft_1^a\} \cdot \{\triangleleft_1^Y\} \otimes \{\triangleright_1^X\} \cdot \{\triangleright_1^b, \triangleleft_1^X\} \cdot \{\triangleleft_1^b\} = \{\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3\}$$

with $[\mathcal{A}_1]$ and $[\mathcal{A}_2]$ the arrangements defined earlier and $[\mathcal{A}_3]$ the following.

$$\{\triangleright_1^a, \triangleright_1^X\} \cdot \{\triangleleft_1^X, \triangleright_1^b\} \cdot \{\triangleleft_1^a, \triangleright_1^Y\} \cdot \{\triangleleft_1^Y, \triangleleft_1^b\} \quad [\mathcal{A}_3]$$

For ease of presentation, we assume *none of the variables has the empty solution ε* . All the algorithms and theoretical results can be easily extended to handle ε solutions. Note that if variables X and Y in $[\mathbb{E}_2]$ could be ε , the \otimes operator would produce more arrangements.

Valid Arrangements . Some of the arrangements produced by the \otimes operator may not be valid, violating intrinsic properties in the string domain. We hence define a new operator \Downarrow , called *arrangement merging*, that filters out the invalid arrangements. According to the definition in Fig. 6, the merge operator produces a subset of the arrangements generated by the \otimes operator, pruning those that do not satisfy the four conditions listed:

- **[Neighboring]** condition states that the left boundary label and the right boundary label of a character need to be in the neighboring label sets. Intuitively, it states that a character cannot be cut/split by any boundary. As such, the arrangement $[\mathcal{A}_3]$ for our previous example is not valid as it suggests character a is cut and split by the right boundary of X .
- **[Consistency]** condition states that the left/right boundaries of different characters cannot be in the same label set. Intuitively, different characters cannot align.
- **[Uniqueness]** condition specifies that a label can only appear once in an arrangement. Recall that multiple occurrences of the same variable/character have different labels.
- **[Non-Overlapping]** condition is more complex and is discussed below.

An arrangement satisfying the first three conditions is also called a *valid* arrangement.

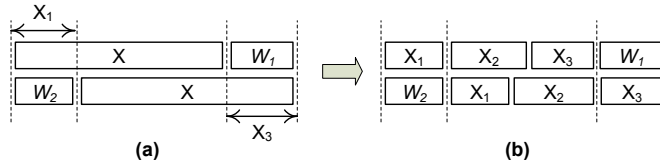


Fig. 8: The essence of overlapping arrangement. We omit the labels for readability.

Non-overlapping Arrangements. We have discussed in Section 6.1, the formula reduction process may get into an infinite loop. That is, the reduced equations (for the

sub-parts of the original equations) may have the same form of the original equations such that the reduction never ends. While we will formally discuss the theorem in Section 6, the intuition can be explained by Fig. 8. Figure (a) shows a boundary alignment of equation $X \cdot \mathcal{W}_1 = \mathcal{W}_2 \cdot X$. The two instances of X overlap and hence variable X is cut by two boundaries, one is \triangleleft_1^X (i.e. the RHS label of the X on the left) and the other is \triangleright_2^X . The two cuts may be at different positions inside X . As such, in the next step, X is split to $X_1 \cdot X_2 \cdot X_3$ with $X_1 = \mathcal{W}_2$ and $X_3 = \mathcal{W}_1$. Figure (b) shows the new alignment after the split. Note that the reduction yields an equation $X_2 \cdot X_3 = X_1 \cdot X_2$, which is essentially an re-occurrence of the form $X \cdot \mathcal{W}_1 = \mathcal{W}_2 \cdot X$, leading to infinite reduction.

As discussed in [41], overlapping solutions rarely happen in practice, especially in the context of program related string analysis. One key design choice is hence to detect and prune overlapping arrangements such that the procedure can avoid such non-terminations. Condition **[Non-overlapping]** in the definition of \Downarrow in Fig. 6 defines the non-overlapping requirement.

Observe that the merge operation \Downarrow generates a set of arrangements from two input arrangements. Our procedure selects one to proceed. If the selection does not lead to a SAT result, it is rolled back and a different selection will be taken. If none of the selections leads to SAT, the procedure returns UNSAT or indicates the presence of overlapping arrangements that have been pruned by the procedure. More details can be found in Section 6.2.3.

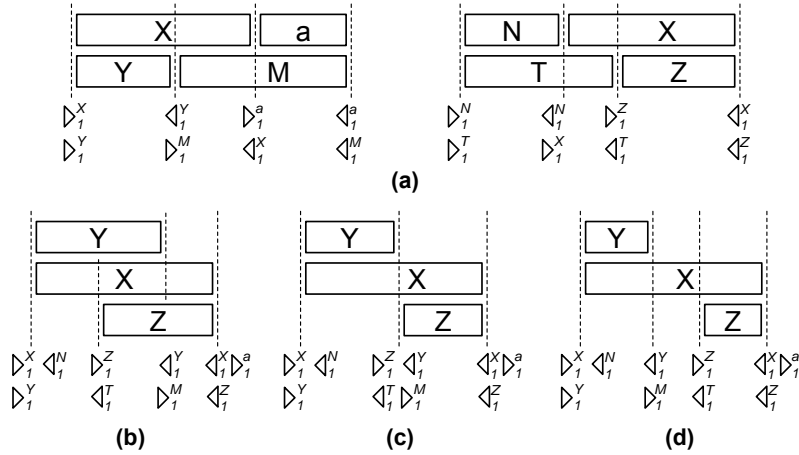


Fig. 9: Example for variable arrangement.

Per Variable Arrangement. A variable may appear multiple times in multiple equations. The label arrangements of multiple occurrences of the variable may suggest the variable being cut differently and divided into different sub-parts. Such information essentially denotes the constraints on the variable from multiple equations. For exam-

ple, if a variable X aligns with a character ‘ a ’ on its left boundary in one equation and aligns with a character ‘ b ’ on its right boundary in another equation, we know that X must start with ‘ a ’ and end with ‘ b ’.

In general, the constraints for the same variable from multiple equations can be very complex, requiring proper handling. For example, consider the following formula (the initial labels can be inferred and hence omitted for readability).

$$X \cdot a = Y \cdot M \wedge N \cdot X = T \cdot Z \quad [\mathbb{E}_3]$$

Note that X aligns with Y on their left boundaries in the first equation and aligns with variable Z on their right boundaries in the second equation. We need to consider the relative positions of Y and Z inside X because they lead to different constraints of the sub-parts of X .

Fig. 9 (a) shows a valid arrangement for $[\mathbb{E}_3]$, which is generated by the aforementioned step of merging the LHS and RHS words of the equations. Observe that in this arrangement, X is cut by both Y and Z . Figures (b), (c), and (d) depict the three different relative positions of Y and Z regarding X . Note that they represent different constraints such that we need to generate different (smaller) equations in formula reduction. Particularly, in figure (b), Y and Z overlap such that they constrain each other while both constrain part of X . In (c), Y and Z covers X but they don’t overlap. In (d), Y and Z constrain the head and the tail of X , respectively. But the middle part of X is not constrained by either.

Another novel aspect of our procedure is the fusion of the different constraints for the same variable from various equations. The observation is that *if we denote the constraint for a given variable from an equation as an arrangement, the variable arrangements from different equations can be merged using the merge operation \Downarrow* .

In Fig. 6, we define a variable arrangement projection operator: $\downarrow_X (\mathcal{A})$ acquires the parts of the arrangement \mathcal{A} related to variable X , that is, the set of subsequences of label sets that start with the LHS label of X and terminate with the corresponding RHS label of X . They contain all the boundaries that cut X in the given arrangement.

For example, let the arrangement of the equation on the left of Fig. 9 (a) be \mathcal{A}_{a1} , $\downarrow_X (\mathcal{A}_{a1}) = \{\{\triangleright_1^X, \triangleright_1^Y\} \cdot \{\triangleleft_1^Y, \triangleright_1^M\} \cdot \{\triangleleft_1^X, \triangleright_1^a\}\}$. In the definition of \downarrow in Fig. 6, given an arrangement $L_1 \cdot \dots \cdot L_n$ and a variable X , it identifies all the subsequences $L_t \cdot \dots \cdot L_{t+m}$, in which L_t share some common labels with the LHS label set of X in the formula and L_{t+m} share some common labels with the corresponding RHS label set of X .

Once we get the arrangements for all the occurrences of a variable, we can use the arrangement merge operator \Downarrow to merge them to generate all the possible (joint) arrangements inside the variable. In the example in Fig 9 (a), the projection of the arrangement of the equation on the right, \mathcal{A}_{a2} , yields $\downarrow_X (\mathcal{A}_{a2}) = \{\{\triangleright_1^X, \triangleleft_1^N\} \cdot \{\triangleleft_1^T, \triangleright_1^Z\} \cdot \{\triangleleft_1^X, \triangleleft_1^Z\}\}$. The operation $(\downarrow_X (\mathcal{A}_{a1})) \Downarrow (\downarrow_X (\mathcal{A}_{a2}))$ hence yields exactly the three arrangements in figures (b)-(d).

Again, the procedure selects one variable arrangement from the generated set to proceed. Backtracking is allowed.

REWRITEWD : $Word \times (StringVariable \rightarrow LabelArrangement) \rightarrow Word$	
REWRITEWD ($L_1 c^{L_2} \cdot W, \mathcal{A}_{var}[]$)	$:= L_1 c^{L_2} \cdot \mathbf{REWRITEWD}(L_2 W, \mathcal{A}_{var}[])$
REWRITEWD ($L_l X^{L_r} \cdot W, \mathcal{A}_{var}[]$)	$:= L_l \cup L_1 X_1^{L_2} X_2^{L_3} \dots X_{n-1}^{L_n \cup L_r} \cdot \mathbf{REWRITEWD}(L_n \cup L_r W, \mathcal{A}_{var}[]),$ in which $\mathcal{A}_{var}[X] = L_1 \dots L_n$
REWRITEWD ($L_\epsilon, \mathcal{A}_{var}[]$)	$:= L_\epsilon$
SPLITEQ : $Equation \times LabelArrangement \times (StringVariable \rightarrow LabelArrangement) \rightarrow Formula$	
SPLITEQ ($W_1 = W_2, \mathcal{A}_{eq}, \mathcal{A}_{var}[]$)	$:= \mathbf{RSPLITEQ}(\mathbf{REWRITEWD}(W_1, \mathcal{A}_{var}) = \mathbf{REWRITEWD}(W_2, \mathcal{A}_{var}), \mathcal{A}_{eq})$
RSPLITEQ ($L_i W_1^{L_j} \cdot W_2 = L_k W_3^{L_l} \cdot W_4, L_1 \cdot L_2 \cdot \mathcal{A}$)	$:= (L_i W_1^{L_j} = L_k W_3^{L_l}) \wedge \mathbf{RSPLITEQ}(L_j W_2 = L_l W_4, L_2 \cdot \mathcal{A}),$ with $(L_i \cap L_k \cap L_1 \neq \emptyset) \wedge (L_j \cap L_l \cap L_2 \neq \emptyset)$
SOLVABLE : $Formula \rightarrow Boolean$	
SOLVABLE (\mathbb{F})	$:= \begin{cases} true & \text{iff } \forall \mathbb{E} \in \mathbb{F}, \mathbb{E} \text{ must have the form of } "X = Y" \text{ or } "X = c" \\ false & \text{otherwise} \end{cases}$

Fig. 10: Formula Transformation Primitives.

6.2.2 Formula Transformation

In the previous subsection, we discuss how to generate arrangements from equations and generate arrangements for variables from equation arrangements. Once a variable arrangement is selected, the constraints for its sub-parts (i.e. their alignments with other variables and characters) are determined. Therefore, in this subsection, we discuss how to split equations to represent the constraints on subparts.

The process consists of two steps.

- In the first step, a variable is split to a set of new variables according to the variable arrangement. Each equation is rewritten by replacing each variable with the split variables.
- In the second step, each equation is split to a set of new equations, each constraining a sub-part of the original equation. This process is guided by determining the common labels between the label sets of the LHS and RHS words of the equations generated in the first step.

The **REWRITEWD** operator in Fig. 10 is for the first step. It takes a word and a mapping from variables to their corresponding arrangements, and generates a new word with the split variables. The operator is defined recursively. If the head of the word is a character, the resulting word inherits the same character and its labels. The operator continues to rewrite the remainder of the word. If the head of the word is a variable X , it splits the variable to a set of $n-1$ variables X_1, \dots, X_{n-1} , where n is the number of label sets in the variable arrangement of X . This implies that X has $n-1$ sub-parts, delimited by the label sets. Also, the new variables inherit both the label sets L_l and L_r from the original word and the label sets in the variable arrangement. These labels denote the boundaries that each new variable aligns with.

The **SPLITEQ** operator denotes the second step. It takes an equation, the corresponding equation arrangement, and the mapping for variable arrangements, and then generates a formula, which is the conjunction of a set of equations. It first rewrites the

LHS and RHS words of the equation and then uses the **RSPLITEQ** operator to split the rewritten equation. The **RSPLITEQ** operator takes the rewritten equation and the arrangement of the *original* equation, and generates a conjunction of new equations. The splitting is directed by the label sets of the original equation arrangement and the rewritten equation. In particular, given the first two labels $L_1 \cdot L_2$ of the original equation arrangement, it looks for the label sets in the rewritten LHS and RHS words that share some common labels with L_1 and L_2 . A new equation is generated to denote the equivalence of sub-parts delimited by the aforementioned label sets. The new equation is conjoined with the other new equations generated by applying the operator to the remainder of the equation.

Example. Consider the example in Fig. 9 (a). Assume the variable arrangement for X is the one in Fig. 9 (b). In other words,

$$\mathcal{A}_{var}[X] = \{\triangleright_1^X, \triangleright_1^Y, \triangleleft_1^N\} \cdot \{\triangleright_1^Z, \triangleleft_1^T\} \cdot \{\triangleleft_1^Y, \triangleright_1^M\} \cdot \{\triangleleft_1^X, \triangleleft_1^Z, \triangleright_1^a\}$$

The LHS word $X \cdot a$ of the first equation is rewritten as follows:

$$\begin{aligned} \mathcal{W}'_l &= \mathbf{REWRITEWD}(\{\triangleright_1^X\} X \{\triangleleft_1^X, \triangleright_1^a\} a \{\triangleleft_1^a\}, \mathcal{A}_{var}) \\ &= \{\triangleright_1^X, \triangleright_1^Y, \triangleleft_1^N\} X_1^{\{\triangleright_1^Z, \triangleleft_1^T\}} X_2^{\{\triangleleft_1^Y, \triangleright_1^M\}} X_3^{\{\triangleleft_1^X, \triangleleft_1^Z, \triangleright_1^a\}} a \{\triangleleft_1^a\} \end{aligned}$$

Note that X is split to X_1 , X_2 , and X_3 due to the cuts from Z and Y . Similarly, the RHS word $Y \cdot M$ is rewritten as follows:

$$\begin{aligned} \mathcal{W}'_r &= \mathbf{REWRITEWD}(\{\triangleright_1^Y\} Y \{\triangleleft_1^Y, \triangleright_1^M\} M \{\triangleleft_1^M\}, \mathcal{A}_{var}) \\ &= \{\triangleright_1^X, \triangleright_1^Y, \triangleleft_1^N\} Y \{\triangleleft_1^Y, \triangleright_1^M\} M_1^{\{\triangleleft_1^X, \triangleright_1^a\}} M_2^{\{\triangleleft_1^a, \triangleleft_1^M\}} \end{aligned}$$

M is split to M_1 and M_2 due to the cut by X .

Therefore, we split the first equation as follows. Note that $\mathcal{A} = \{\triangleright_1^X, \triangleright_1^Y\} \cdot \{\triangleleft_1^Y, \triangleright_1^M\} \cdot \{\triangleleft_1^X, \triangleright_1^a\} \cdot \{\triangleleft_1^a, \triangleleft_1^M\}$ is the equation arrangement.

$$\begin{aligned} &\mathbf{SPLITEQ}(Xa = YM, \mathcal{A}, \mathcal{A}_{var}) \\ &= \mathbf{RSPLITEQ}(\mathcal{W}'_l = \mathcal{W}'_r, \mathcal{A}) \\ &= (\{\triangleright_1^X, \triangleright_1^Y, \triangleleft_1^N\} X_1^{\{\triangleright_1^Z, \triangleleft_1^T\}} X_2^{\{\triangleleft_1^Y, \triangleright_1^M\}} = \{\triangleright_1^X, \triangleright_1^Y, \triangleleft_1^N\} Y \{\triangleleft_1^Y, \triangleright_1^M\}) \\ &\quad \wedge \mathbf{RSPLITEQ}(\cdots X_3 \cdots a \cdots = \cdots M_1 \cdots M_2 \cdots, \dots) \\ &= \dots \\ &= (X_1 \cdot X_2 = Y) \wedge (X_3 = M_1) \wedge (a = M_2) \tag{E_4} \end{aligned}$$

Note that we omit the labels in the last step for readability. The second equation is similarly split to the following.

$$(N = T_1) \wedge (T_2 = X_1) \wedge (X_2 \cdot X_3 = Z) \tag{E_5}$$

□

With the new set of equations, the procedure repeats the steps of (1) computing equation arrangement; (2) computing variable arrangement; (3) splitting formula. As such, the equations represent equivalence between smaller and smaller portions of the original words. The process terminates if the formula is in solvable form.

Solvable Form. As defined in Fig. 10, a formula (i.e. a conjunction of equations) is in solvable form if each equation is either an equivalence between a variable and a character, or an equivalence between two variables. For variables that are directly or indirectly (i.e. through other variables) equivalent to a character, their solution is the character. Our solving processing, namely the *[Consistency]* condition in arrangement production (Section 6.2.1), ensures that the same variable is not equivalent to different characters. Variables that are not equivalent to any character, directly or indirectly, are free variables such that we can assign any characters to them.

6.2.3 Algorithm

The procedure is summarized in Algorithm 1. It takes a formula over word equations \mathcal{Q}_w and the corresponding integer linear arithmetic constraints \mathcal{Q}_l over the length function, produces SAT, UNSAT or UNKNOWN results. UNKNOWN means that the algorithm has encountered overlapping arrangements and pruned those arrangements, even though it did not find any SAT solution.

The algorithm consists of three steps. In step one (lines 9-15), it generates the set of possible arrangements for each equation and selects one from the set and stores it to \mathcal{A}_{eq} to proceed. Each invocation to the **SELECT()** function is similar to a choice point in model checking, allowing back-tracking and selecting another option.

In step two (lines 16-17), variable arrangements for each variable are computed based on the previously selected equation arrangements. Again, each variable may have multiple arrangements.

Lines 18-25 denote the third step, in which the formula is split based on the selected variable arrangements.

In lines 26-29, the algorithm checks an overlapping arrangement has ever detected and pruned. If so, the procedure returns UNKNOWN, otherwise, it returns UNSAT.

Essence of the Procedure. The essence of the algorithm is to discover all the boundaries that are correlated and search for a total order of them. We say boundaries are correlated if they cut the same variable. Step one in the algorithm orders the boundaries within the LHS and RHS words using the merge operator \Downarrow . However, this does not discover all the correlated boundaries because a variable may be constrained differently in other equations. Therefore, in step two, it collects the ordered boundaries within a variable from all its occurrences and generates all the possible orders of the boundaries. For example in Fig. 9 (a), from the equation arrangement (step one), we only know Y cuts X . Through step two, we further know Z also cuts X and hence need to order the two cuts.

However, this does not expose all related boundaries. For example, suppose an additional equation $Z = J \cdot K$ is added to the formula. After the first round of the algorithm, $[\mathbb{E}_5]$ becomes

$$(N = T_1) \wedge (T_2 = X_1) \wedge (\dots X_2^{\{\triangleleft_1^Y, \triangleright_1^M\}} X_3 \dots = \dots Z_1^{\{\triangleleft_1^J, \triangleright_1^K\}} Z_2 \dots) \quad [\mathbb{E}_6]$$

In other words, the boundary between J and K cuts Z and hence cuts X . But it is not identified as correlated to the other boundaries that cut X in the first round because the cut is indirect (through Z). Fortunately, the correlation will be exposed and the boundaries will be ordered in the second round of the algorithm when $[\mathbb{E}_6]$

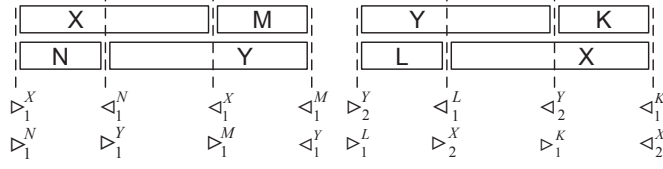
Algorithm 1 A high-level description of Z3str2's Algorithm

Input: Word equations Q_w , and the corresponding integer linear arithmetic constraints Q_l over the length function
Output: SAT / UNSAT / UNKNOWN
Definition: $\mathcal{A}_{eq}/\mathcal{A}_{var}$: mappings from equation/variable to set of label arrangements;
 $\mathcal{A}_{eq}/\mathcal{A}_{var}$: mappings from equation/variable to arrangement.

- 1: **procedure** SOLVESTRINGCONSTRAINT(Q_w, Q_l)
- 2: **if** equations in Q_w are all in solvable form **then**
- 3: **if** Q_w is UNSAT or Q_l is UNSAT **then**
- 4: **return** UNSAT
- 5: **if** Q_w and Q_l can be consistently determined as SAT together **then**
- 6: **return** SAT
- 7: Convert Q_w equisatisfiably in disjunctive normal form (DNF) formula Q_a
- 8: **for each** disjunct D in Q_a **do**
- 9: */* Convert each equation in D to arrangements consistent with the length constraints */*
- 10: **for each** equation $\mathbb{E} : \mathcal{W}_1 = \mathcal{W}_2 \in D$ **do**
- 11: $\mathcal{A}_{eq}[\mathbb{E}] \leftarrow \text{GETARRGMTFROMWORD}(\mathcal{W}_1) \Downarrow \text{GETARRGMTFROMWORD}(\mathcal{W}_2)$
- 12: **for** \mathcal{A} in $\mathcal{A}_{eq}[\mathbb{E}]$ **do**
- 13: Extract the length constraint implied by \mathcal{A} .
- 14: **if** \mathcal{A} has inconsistent length constraint **then**
- 15: Remove \mathcal{A} from $\mathcal{A}_{eq}[\mathbb{E}]$
- 16: */* SELECT(): Similar to a choice point in model checking, allowing backtracking and selecting another option.*/*
- 17: $\mathcal{A}_{eq}[\mathbb{E}] \leftarrow \text{SELECT}(\mathcal{A}_{eq}[\mathbb{E}])$
- 18: **for each** string variable X **do**
- 19: */* Merge per-equation arrangements to a set of possible global arrangements $\mathcal{A}_{var}[X]$ */*
- 20: */* Operator \Downarrow detects and prunes arrangements with overlaps from $\mathcal{A}_{var}[X]$ */*
- 21: $\mathcal{A}_{var}[X] \leftarrow (\Downarrow_X (\mathcal{A}_{eq}[\mathbb{E}_1])) \Downarrow \dots \Downarrow (\Downarrow_X (\mathcal{A}_{eq}[\mathbb{E}_n]))$
- 22: **for each** global arrangement combination selected from $\mathcal{A}_{var}[X], \dots$, for all variables X, \dots **do**
- 23: **Let** $D ::= E_1 \wedge \dots \wedge E_m$
- 24: */* Split each variable to sub-variables according to the selected global arrangement */*
- 25: $D' \leftarrow \text{SPLITEQ}(E_1, \mathcal{A}_{eq}[E_1], \mathcal{A}_{var}) \wedge \dots \wedge \text{SPLITEQ}(E_m, \mathcal{A}_{eq}[E_m], \mathcal{A}_{var})$
- 26: Convert Q_w equisatisfiably to a system Q'_w of simpler equation based on D'
- 27: Q'_l is the corresponding new set of length constraints
- 28: $r = \text{SOLVESTRINGCONSTRAINT}(Q'_w, Q'_l)$
- 29: **if** $r \equiv \text{SAT}$ **then**
- 30: **return** SAT
- 31: **if** overlapping variables have ever been detected **then**
- 32: **return** UNKNOWN
- 33: **else**
- 34: **return** UNSAT

is reduced. The algorithm repeats until all correlations are exposed and correlated boundaries are ordered.

Multiple Layer Overlaps. In Section 6.1, we present an example of infinite reduction caused by a multiple layer loop, that is, the loop contains multiple rounds of reduction. This is much more difficult to detect and prevent compared to direct looping such as $aX = Xb$. However, we observe the essence of such looping is still overlapping arrangement (Section 6.2.1) compared to direct looping. The difference is that correlated boundaries are gradually exposed such that the boundary labels that overlap may only manifest themselves after a few rounds of reduction. Our algorithm can detect them through the [Non-overlapping] condition in the merge operator \Downarrow .

Fig. 11: The graphic representations of \mathcal{A}_4 and \mathcal{A}_5

Example. We revisit the third example in Sec. 6.1.

$$X \cdot M = N \cdot Y \quad \wedge \quad Y \cdot K = L \cdot X$$

The equations with labels are

$$\{\triangleright_1^X\} X \{\triangleleft_1^X, \triangleright_1^M\} M \{\triangleleft_1^M\} = \{\triangleright_1^N\} N \{\triangleleft_1^N, \triangleright_1^Y\} Y \{\triangleleft_1^Y\} \quad [\mathbb{E}_7]$$

$$\{\triangleright_2^Y\} Y \{\triangleleft_2^Y, \triangleright_1^K\} K \{\triangleleft_1^K\} = \{\triangleright_1^L\} L \{\triangleleft_1^L, \triangleright_2^X\} X \{\triangleleft_2^X\} \quad [\mathbb{E}_8]$$

One arrangement of the LHS and RHS words of $[\mathbb{E}_7]$ and $[\mathbb{E}_8]$ are $[\mathcal{A}_4]$ and $[\mathcal{A}_5]$ respectively. Their graphic representations are shown in Fig. 11.

$$\{\triangleright_1^X, \triangleright_1^N\} \cdot \{\triangleleft_1^N, \triangleright_1^Y\} \cdot \{\triangleleft_1^X, \triangleright_1^M\} \cdot \{\triangleleft_1^M, \triangleleft_1^Y\} \quad [\mathcal{A}_4]$$

$$\{\triangleright_2^Y, \triangleright_1^L\} \cdot \{\triangleleft_1^L, \triangleright_2^X\} \cdot \{\triangleleft_2^Y, \triangleright_1^K\} \cdot \{\triangleleft_1^K, \triangleleft_2^X\} \quad [\mathcal{A}_5]$$

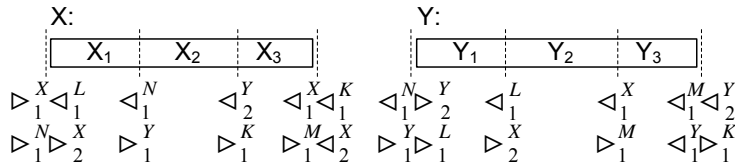
Now we apply the variable arrangement projection operator \downarrow and get the arrangements related to variable X and Y .

$$\downarrow_X (\mathcal{A}_4) = \{\mathcal{A}_7\}, \quad \mathcal{A}_7 = \{\triangleright_1^X, \triangleright_1^N\} \cdot \{\triangleleft_1^N, \triangleright_1^Y\} \cdot \{\triangleleft_1^X, \triangleright_1^M\}$$

$$\downarrow_X (\mathcal{A}_5) = \{\mathcal{A}_8\}, \quad \mathcal{A}_8 = \{\triangleleft_1^L, \triangleright_2^X\} \cdot \{\triangleleft_2^Y, \triangleright_1^K\} \cdot \{\triangleleft_1^K, \triangleleft_2^X\}$$

$$\downarrow_Y (\mathcal{A}_4) = \{\mathcal{A}_9\}, \quad \mathcal{A}_9 = \{\triangleleft_1^N, \triangleright_1^Y\} \cdot \{\triangleleft_1^X, \triangleright_1^M\} \cdot \{\triangleleft_1^M, \triangleleft_1^Y\}$$

$$\downarrow_Y (\mathcal{A}_5) = \{\mathcal{A}_{10}\}, \quad \mathcal{A}_{10} = \{\triangleright_2^Y, \triangleright_1^L\} \cdot \{\triangleleft_1^L, \triangleright_2^X\} \cdot \{\triangleleft_2^Y, \triangleright_1^K\}$$

Fig. 12: The graphic representations of $\mathcal{A}_{var}[X]$ and $\mathcal{A}_{var}[Y]$

Then we can use the merge operator (\Downarrow) to generate the arrangements in variables. One of the arrangement combination generated by $\mathcal{A}_7 \Downarrow \mathcal{A}_8$ and $\mathcal{A}_9 \Downarrow \mathcal{A}_{10}$ is as follows. Fig. 12 shows the corresponding graphic representations.

$$\mathcal{A}_{var}[X] = \{\triangleright_1^X, \triangleright_1^N, \triangleleft_1^L, \triangleright_2^X\} \cdot \{\triangleleft_1^N, \triangleright_1^Y\} \cdot \{\triangleleft_2^Y, \triangleright_1^K\} \cdot \{\triangleleft_1^X, \triangleright_1^M, \triangleleft_1^K, \triangleleft_2^X\}$$

$$\mathcal{A}_{var}[Y] = \{\triangleleft_1^N, \triangleright_1^Y, \triangleright_2^Y, \triangleright_1^L\} \cdot \{\triangleleft_1^L, \triangleright_2^X\} \cdot \{\triangleleft_1^X, \triangleright_1^M\} \cdot \{\triangleleft_1^M, \triangleleft_1^Y, \triangleleft_2^Y, \triangleright_1^K\}$$

The LHS word $X \cdot M$ and RHS word $N \cdot Y$ can be rewritten as follows.

$$\begin{aligned}
\mathcal{W}_l^{\mathbb{E}_7} &= \mathbf{REWRITEWD}(\{\triangleright_1^X\} X^{\{\triangleleft_1^X, \triangleright_1^M\}} M^{\{\triangleleft_1^M\}}, \mathcal{A}_{var}) \\
&= \{\triangleright_1^X, \triangleright_1^N, \triangleleft_1^L, \triangleright_2^X\} X_1^{\{\triangleleft_1^N, \triangleright_1^Y\}} X_2^{\{\triangleleft_2^Y, \triangleright_1^K\}} X_3^{\{\triangleleft_1^X, \triangleright_1^M, \triangleleft_1^K, \triangleleft_2^X\}} M^{\{\triangleleft_1^M\}} \\
\mathcal{W}_r^{\mathbb{E}_7} &= \mathbf{REWRITEWD}(\{\triangleright_1^N\} N^{\{\triangleleft_1^N, \triangleright_1^Y\}} Y^{\{\triangleleft_1^Y\}}, \mathcal{A}_{var}) \\
&= \{\triangleright_1^N\} N^{\{\triangleleft_1^N, \triangleright_1^Y, \triangleright_2^Y, \triangleright_1^L\}} Y_1^{\{\triangleleft_1^L, \triangleright_2^X\}} Y_2^{\{\triangleleft_1^X, \triangleright_1^M\}} Y_3^{\{\triangleleft_1^M, \triangleleft_1^Y, \triangleleft_2^Y, \triangleright_1^K\}}
\end{aligned}$$

The equation $[\mathbb{E}_7]$ can be split as what is showed in Fig. 13.

$$\begin{aligned}
&\mathbf{SPLITQ}(XM = NY, \mathcal{A}_4, \mathcal{A}_{var}) \\
&= \mathbf{RSPLITQ}(\mathcal{W}_l^{\mathbb{E}_7} = \mathcal{W}_r^{\mathbb{E}_7}, \mathcal{A}_4) \\
&= (\{\triangleright_1^X, \triangleright_1^N, \triangleleft_1^L, \triangleright_2^X\} X_1^{\{\triangleleft_1^N, \triangleright_1^Y\}} = \{\triangleright_1^N\} N^{\{\triangleleft_1^N, \triangleright_1^Y, \triangleright_2^Y, \triangleright_1^L\}} \\
&\quad \wedge \mathbf{RSPLITQ}(\cdots X_2 \cdots X_3 \cdots M \cdots = \cdots Y_1 \cdots Y_2 \cdots Y_3 \cdots, \dots)) \\
&= \dots \\
&= (\cdots X_1 \cdots = \cdots N \cdots) \wedge (\cdots M \cdots = \cdots Y_3 \cdots) \wedge \\
&\quad \boxed{(\{\triangleleft_1^N, \triangleright_1^Y\} X_2^{\{\triangleleft_2^Y, \triangleright_1^K\}} X_3^{\{\triangleleft_1^X, \triangleright_1^M, \triangleleft_1^K, \triangleleft_2^X\}} = \{\triangleleft_1^N, \triangleright_1^Y, \triangleright_2^Y, \triangleright_1^L\} Y_1^{\{\triangleleft_1^L, \triangleright_2^X\}} Y_2^{\{\triangleleft_1^X, \triangleright_1^M\}})} \quad [\mathbb{E}_9]
\end{aligned}$$

Fig. 13: Splitting Equation $[\mathbb{E}_6]$

For equation $[\mathbb{E}_9]$, we can get its arrangements if we repeat the same procedure for $[\mathbb{E}_7]$ and $[\mathbb{E}_8]$. Among them, an interesting one is the following:

$$\{\triangleright_1^Y, \triangleleft_1^N, \triangleright_2^Y, \triangleright_1^L\} \cdot \{\triangleleft_1^L, \triangleright_2^X\} \cdot \{\triangleleft_2^Y, \triangleright_1^K\} \cdot \{\triangleleft_1^X, \triangleright_1^M, \triangleleft_2^X, \triangleleft_1^K\} \quad [\mathcal{A}_{11}]$$

We further compute the arrangement projection to variable X . As shown in Fig. 14, intuitively, the variable project is to determine the order between the boundaries in red, which are created in the previous cycle.

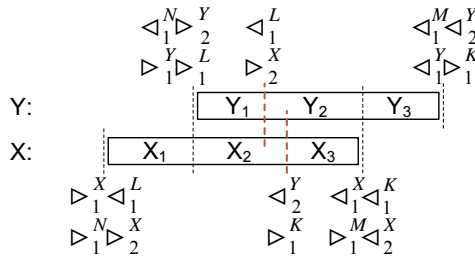


Fig. 14: The graphic representations of $\mathcal{A}_{var}[X]$

If we project these arrangements back to variable X , we get $\mathcal{A}_{var}[X]$:

$$\{\boxed{\triangleright_1^X}, \triangleright_1^N, \triangleleft_1^L, \triangleright_2^X\} \cdot \{\triangleright_1^Y, \triangleleft_1^N, \triangleright_2^Y, \triangleright_1^L\} \cdot \{\triangleleft_1^L, \boxed{\triangleright_2^X}\} \cdot \{\triangleleft_2^Y, \triangleright_1^K\} \cdot \{\boxed{\triangleleft_1^X}, \triangleright_1^M, \triangleleft_2^X, \triangleleft_1^K\}$$

The boxed labels violate the **[Non-overlapping]** condition. A loop is detected.

6.2.4 Soundness of the Z3str2 Algorithm

In this section we sketch the soundness proof of Z3str2's algorithm given in Algorithm 1. For a detailed formal analysis we refer the reader to the associated tech-report. The soundness property of any decision procedure in an SMT solving context can be stated as "If the procedure returns UNSAT, then input formula is indeed UNSAT".

Theorem 1 *Algorithm 1 is sound, i.e., when Algorithm 1 reports UNSAT, the input constraint is indeed UNSAT.*

Proof. To see that Z3str2 is sound, we show that the UNSAT returned at line 4 and line 29 are both sound. First observe that line 4 returns an UNSAT if either string or integer constraints are determined to be UNSAT. For string constraints, we use the algorithms described in [14] to decide the satisfiability of word (dis)equations in the solved form. The soundness of line 4 relies on the soundness of the procedure [14] and the integer solver (here Z3).

For the UNSAT returned at line 29, we show transformations impacting it are all satisfiability-preserving. If a transformation is satisfiability-preserving, it means its output formula is satisfiable if and only if its input formula is satisfiable. In particular, transformations at (i) line 7 (ii) line 10 (iii) line 17 and (iv) lines 20-21 are satisfiability-preserving: (i) The disjunctive normal form conversion at line 7 is obviously satisfiability-preserving. (ii) The conversion in line 10 is probably the most involved in terms of establishing soundness. This step is a variant of the idea of sound transformation of word equations to arrangements mentioned in Makanin's paper [27]. We can show that arrangement generation is satisfiability-preserving because each arrangement is a finite set of equations implied by the input system of equations. In addition, we extract length constraints from arrangements and they may conflict with the existing integer constraints. If so, we drop inconsistent arrangements based on the UNSAT results determined by the integer theory. Similarly, since we assume the integer theory is sound, this step is also satisfiability-preserving. (iii) At line 17, we systematically enumerate all feasible orders among boundary labels according to the definition of the arrangement merge (Fig. 6). This step is satisfiability-preserving. (iv) In lines 20 and 21, this step derives simpler equations by a satisfiability-preserving rewriting. Please see the technical report for proof details [40]. Note the REs are reduced to word equations so that they can be handled by this same procedure. In addition, although we prune arrangements at line 17, the answer can only be SAT or UNKNOWN once this happens. The algorithm is still sound. Therefore, we return UNSAT exactly when we can prove this to be the case.

7 String and Integer Theory Integration

Basic Length Rules. For strings X and Y , we assert the following: (1) $|X| \geq 0$ (2) $|X| = 0 \leftrightarrow X = \epsilon$ (3) $X = Y \rightarrow |X| = |Y|$ (4) $|X \cdot \dots \cdot Y| = |X| + \dots + |Y|$.

String and Integer Theory Integration: As discussed in Section 3, finding a consistent solution for both strings and integers can be expensive due to the infinite search space. The goal of string and integer theory integration is to achieve synergy from the two such that the procedure can converge faster. In particular, one theory will generate new assertions in the domain of the other theory, and vice versa. Inside the string theory, the set of arrangements that is explored is constrained by the assertions on string lengths, which are provided by the integer theory. On the other hand, the string theory will derive new length assertions when it makes progress in exploring new arrangements. These assertions are provided to the integer theory so that the search space is pruned.

Consider $X \cdot Y = M \cdot N$, where X, Y, M and N are nonempty string variables. It has three possible arrangements: $[a1] X = M \cdot T_1 \wedge N = T_1 \cdot Y$; $[a2] X = M \wedge N = Y$; $[a3] M = X \cdot T_2 \wedge Y = T_2 \cdot N$. Assume the integer theory infers that $|X| > |M|$ or $|Y| < |N|$. Thus, only $[a1]$ is consistent with the length conditions. The string solver only needs to explore one arrangement instead of three. On the other hand, assume the string solver is exploring arrangement $[a1]$. It generates a new assertion $[a1] \rightarrow |X| = |M \cdot T_1| \wedge |N| = |T_1 \cdot Y| \wedge |X| > |M| \wedge |N| > |Y|$, which in turn triggers the Z3 core to add an integer assertion.

Note that different string solvers implement string and integer integrations in vastly different ways [8, 39, 26, 36]. Bjørner et al. [8] focus on integration in a staged manner. Yu et al. [39] focus on integration via automata manipulations. CVC4 [26], S3 [36], and Z3str2 are integrations within the DPLL(T) architecture, where the algorithm only solves parts of the formula on demand and learns new constraints as it solves such that these implied constraints often cut the search dramatically. Compared to CVC4 and S3, our integration is tighter, powered by the bi-directional heuristics.

7.1 Search Space Pruning via Binary Search

In this subsection, we present a relatively new heuristic that is enabled by integration between the string and integer theories, and the ability for the string theory solver to use information directly from the integer solver, combine it with facts about strings, and generate new terms that guide the search. Although this work is preliminary and requires further evaluation, we have found it to be very promising and to demonstrate noticeable performance improvements.

We observed that Z3str2 runs slowly when the input constraints can only be satisfied by very long strings. As mentioned in Algorithm 1, when a consistent global arrangement is obtained, for each string variable that cannot be split further, Z3str2 tries to assign concrete values to it. This is a two-step procedure: Z3str2 first tries to fix its length and makes sure that no conflict will be triggered by this assignment. Next, Z3str2 proposes concrete values from a pool of all possible constant strings whose lengths meet the requirement obtained in the previous step.

Z3str2 already takes advantage of information from the integer theory to some extent when performing this search. However, it is possible to refine this information and use it to prune the search space much more effectively. In the following subsections, we discuss the issues related to the previous search strategy, and then present a

new search heuristic that converges much more rapidly for cases where the previous strategy performs poorly.

7.1.1 Performance Challenges Imposed by Solving Long Strings

The integer theory represents a variable's possible value by a range, and gradually refines this range during the search. After the problem's satisfiability is established, it starts to construct the model. In some sense, the integer theory decides variables' assignments lazily. However, this will block the string theory from making further progress: even though the string theory can access the integer theory via the bi-directional integration, usually it cannot obtain an assignment until a concrete length (instead of a range) has been fixed by the integer theory.

Therefore, the string theory in Z3str2 and integer theory in the Z3 core should work together to find consistent string lengths. Z3-str does so using a naive linear search approach. In particular, the string theory component continuously and incrementally queries the integer theory to get a suitable length for each variable through Z3 core.

For a string variable X to be assigned, Z3str2 linearly samples the search space and adds assertions for each value in the sample space. For this purpose, the first length axiom added is of the form $(|X| \geq i * k) \rightarrow (|X| = i * k) \vee \dots \vee (|X| = (i + 1) * k - 1) \vee (|X| \geq (i + 1) * k)$, where k is a predefined non-negative integer representing the window size and i is a non-negative integer denoting the number of iteration. The Z3 core picks one assertion from the given set, verifies it against the existing context in the integer theory and backtracks if a conflict is detected. If $(|X| \geq (i + 1) * k)$ is selected and included in the context, the string part goes to the next iteration and repeats the same procedure.

However, for constraints asserting large length, we found that the linear search consumes a significant amount of solving time in proposing gradually increasing length values, and such iterations impose unnecessary performance overhead on the whole decision procedure.

Consider the constraint $|X| \geq 52000$ that is derived from the current search context. Assume k is set to 4. In the first iteration, the length query added is $(|X| \geq 0) \rightarrow (|X| = 0) \vee \dots \vee (|X| = 3) \vee (|X| \geq 4)$. Suppose initially the core picks the assertion $X = 2$ but it is not consistent with the length constraint. Therefore, the core backtracks and tries another assertion, for instance, $|X| \geq 4$. Then, the string theory asserts the next length axiom, $(|X| \geq 4) \rightarrow (|X| = 4) \vee \dots \vee (|X| = 7) \vee (|X| \geq 8)$ and repeats the same procedure. In essence, the string plug-in needs to add the assertion in a linear order till it finds a valid one. The same process has to be done for all string variables, which significantly degrades the performance.

7.1.2 Binary Search Based Heuristics

As discussed above, gradually increasing the length introduces redundant length query iterations as well as a large number of unnecessary length constraints for string variables. Therefore, we designed a new search space pruning technique based on the binary search with heuristic support for Z3 core's backtracking functionality.

However, as the strings are unbounded, we cannot fix an upper bound on their lengths. So, the traditional binary search based technique is not directly applicable. Instead, we design new heuristics, driven by binary search and backtracking, to guide the search. In particular, we chose a window with concrete lower and upper bound in which the size of the window varies dynamically in each iteration depending on the search criteria. We perform the binary search within the window while sliding the window from lower to higher values to find an upper bound. In other words, the lower and upper boundaries of the window are modified when the window slides. In addition, the window size expands or shrinks by the order of two based on the feedback provided by the integer theory. The backtracking ability of the heuristics helps to guide the search and refine our choice of window.

Algorithm 2 Binary search based heuristics for string-integer solver

Input: A string variable X .
Output: The consistent assignment of the length of X ($|X|$)

```

1: procedure GETLENGTHASSIGNMENT( $X, low, high$ )
2:   if  $|X|$  has a concrete assignment  $l$  in the integer theory then
3:     return  $l$ 
4:    $mid = low + (high - low)/2$ 
5:   assert an axiom:  $(|X| < mid) \vee (|X| = mid) \vee (|X| > mid)$ 
6:   query the integer theory and check the context
7:   if  $(|X| > mid)$  is selected by the Integer theory then
8:     return GETLENGTHASSIGNMENT( $X, mid + 1, high \times 2$ )
9:   else if  $|X| = mid$  is selected by the Integer theory then
10:    return  $mid$ 
11:   else
12:    return GETLENGTHASSIGNMENT( $X, low, mid - 1$ )

```

The binary search based heuristics for the theory of strings and integers is summarized in Algorithm 2. It is a recursive procedure driven by the feedback obtained from the integer theory. In turn, it also proposes possible ranges and advances the search inside the integer theory. In lines 2 - 3, it checks whether the length of X has already been assigned by the integer theory. If so, it directly returns the assignment found in the current context. Otherwise, it asserts an additional axiom and starts a modified version of binary search. It probes the context and adjusts the search window accordingly based on the feedback. In particular, it finds the middle element in the current window and asserts a length axiom at line 5. At line 6, the string theory checks the outcome of the axiom just asserted and prunes the search space in lines 7 - 12. Note that the higher bound is doubled if the lower half can be eliminated at line 8. This guarantees that the search can eventually cover a large string length in the unbounded search space.

8 Enhancing Symbolic Execution of C/C++ String Manipulating Programs

Improper string manipulations are an important cause of software defects, which make them a target for program analysis. A study [31] on concolic testing tools and their

limitations shows that string operations make up a significant proportion of system-level code. The difficulty of analyzing these operations adds additional overhead to current symbolic execution engines.

In this section, we discuss how a constraint solver over string and integer theory can provide significant performance improvement in the symbolic execution of string manipulating programs written in C/C++. We compare the constraint expressiveness and the overall procedure among existing tools and our software to show the enhancement enabled by Z3str2.

8.1 Limitations Introduced by Modeling Strings as Bit-Vectors

In current symbolic execution techniques, there is a semantic gap between the high-level notion of strings and low-level representation of program states and memory. Using a bit-vector solver is useful for low-level program analysis as it captures constraints with bit-level precision and efficiently reasons about arithmetic overflows and bitwise operations. Bit-vector solvers also allow for reasoning about symbolic memory regions.

Symbolic execution engines (e.g., KLEE [9] and S2E [10]) collect constraints as bit-vectors by accumulating branch predicates and solving for the branch conditions using constraint solvers such as STP [13] and Z3 [12]. However, these engines perform poorly on programs containing string functions as they don't capture the high-level semantics of strings in accordance with the low-level bit-vector representation of all program data [10].

In particular, current symbolic execution engines that track constraints at the bit-vector level suffer from the path explosion problem. One of the main reason is that representing strings with bit-vectors implies that string operations have to be modeled as loops that check string values on character level. For example, commonly used string functions, such as *strlen* or *strcmp*, mainly consist of loops that iterate over characters until either a particular condition is met or the current character is the null terminator, marking the end of the string. To do so, the symbolic execution engine usually forks a fresh state per each possible length, leading to path explosion, especially when the string is unbounded.

For example, an invocation of the string library function *strlen* on a symbolic string S of size N will generate a total of $N + 1$ paths, one for each possible value of the length between 0 and N , regardless of how the length of the string is actually used throughout the rest of the program. In essence, uninteresting regions of code are explored by enumerating each branch without efficiently pruning the search space. Prior works from the S2E [10] group explained these performance bottlenecks with strings in their symbolic analysis. These restrictions could be lifted by applying string solvers that reason about strings as a primitive data type, similar to integers and bit-vectors, in the context of symbolic execution.

Moreover, it is difficult to detect certain classes of security vulnerabilities arising from certain overflow and underflow errors with state-of-the-art symbolic execution engines. For instance, if we analyze the history of integer-related overflows in CVE database, we observe that operations on extremely long strings are one of the primary

sources of integer overflow vulnerabilities. Since the strings have to be converted to bit-vector representations, it requires the symbolic execution engine to fix concrete string lengths well in advance of where the strings are actually used. Essentially, the engine has to enumerate a large number of possible string lengths without knowing if they are necessary from the point of view of the string theory. As a result, existing symbolic execution engines powered by bit-vector-based solvers are not able to reason efficiently about arithmetic overflow/underflow errors caused by improper handling of string values. This important class of security vulnerabilities inspired us to propose a new theory combination for the underlying solver which allows it to efficiently reason about arithmetic overflow/underflow errors and memory corruption caused by incorrect string operations.

```
bool check_login(char* username, char* password){
    unsigned short len = strlen(username)+1;
    if(len > 32){
        invalid_login_attempt();
        exit(-1);
    }
    char* _username= (char*) malloc(len);
    strcpy(_username, username);
    ...
}
```

Fig. 15: Login method demonstrating an overflow vulnerability

We use a motivating example to explain the limitations of existing string solvers in the context of low-level program analysis and the importance of the string-bitvector theory combination. Consider the *check_login()* function shown in Figure 15. Here, the program calculates the length of the user-controlled input value *username*, and adds 1 to accommodate the trailing *null* character. A new buffer is allocated for the resulting string and the program copies *username* into it using the *strcpy* function. This code behaves as intended for normal-sized input. However, an integer overflow occurs if the user submits a *username* that is 65,535 characters long. The variable *len* is declared as *unsigned short*, which is a 16-bit integer that can hold any value between 0 and 65,535. When a string of length 65,535 is submitted as *username*, the result of *strlen(username)+1* overflows to 0. This results in allocation of a buffer of size 0 on the heap with the call to *malloc()*. Attempting to copy *username* into this buffer results in memory corruption and the potential for an exploit.

8.2 Enhancements Enabled by String-Integer Solvers in Symbolic Execution

Compared to bit-vector-based string constraint solving, modeling strings as a primitive type and tight string-integer theory solver integration make constraint modeling more expressive, as well as eliminating a significant number of unnecessary path enumerations in the context of symbolic execution.

Language expressiveness. Strings modeled as bit-vector arrays are inherently bounded, whereas the length of strings modeled as a primitive type can be arbitrarily (finitely) large. As a result of this, a bit-vector-based string solver cannot precisely encode the overflows caused by a limited-precision type. Due to the upper bound on string length, bit-vector-based encoding has difficulty reasoning about the overflow that occurs in the previous example.

Eliminating loops over characters. Since the solver is able to handle strings of unbounded length and reason about string operations at a high level, it is possible to abstract away low-level loops in some string-handling methods. For example, we can completely abstract away loops over individual characters in the heavily-used *strcmp()* and *strlen()* functions, as the solver has high-level knowledge of the semantics of these methods. As a result, by eliminating the loops we also eliminate a large number of potential symbolic paths that would otherwise be explored.

9 Experimental Results

In this section, we describe the implementation of Z3str2, as well as experiments to validate the efficacy of the first two new techniques proposed in this paper, namely, overlapping variable detection and deeper string/integer theory integration in Section 9.1. However, typical benchmarks used to evaluate search space pruning techniques rarely involve large strings. To demonstrate the effectiveness of the optimized search heuristics, in Section 9.2, we compare the performance of the binary search and the linear search heuristics using a set of benchmarks that systematically cover different string length settings. We also compare with other solvers on a set of constraints requiring large strings. Finally, we evaluate the enhancements in the context of symbolic execution enabled by a constraint solver over string and integer theories in Section 9.3.

9.1 Evaluating Search Space Pruning Techniques

Both techniques improve solver efficiency in isolation, as well as when used simultaneously. We report both combined and isolated contributions.

1. Detection of Overlapping Variables. During solving, Z3str2 prunes away arrangements with overlapping variables, leading to a smaller search space. Thus, if the technique is effective, we would be able to observe that other solvers time out on the cases reported as UNKNOWN by Z3str2. In Z3str2, an UNKNOWN result is returned when at least one arrangement with overlapping variables is found (and pruned), and no satisfying solution can be found over all arrangements with non-overlapping variables.

2. Evaluating String and Integer Theory Integration. The contribution of string-integer theory integration will be illustrated by demonstrating improvement in the performance of solving both SAT and UNSAT instances, in comparison with other solvers.

Table 1: Comparison study settings.

	Z3str2	Z3str2 w/o integration	CVC4	S3	Kaluza	PISA-MONA	Stranger
Kaluza suite	✓	✓	✓	✓	✓		
PISA suite	✓		✓			✓	
AppScan suite	✓		✓				
Kausler suite	✓		✓				✓

9.1.1 Experimental Setting

We compare Z3str2 against five state-of-the-art string solvers, namely, CVC4[26], S3[36], Kaluza [33], PISA-MONA [35], and Stranger [38] across four different suites of benchmarks obtained from Kudzu/Kaluza[33], PISA[35], AppScan Source[2] and Kausler’s [22] projects. Given the rich and diversified landscape of string problems, we chose to validate our approach using benchmarks from real-world applications with different characteristics. Additionally, the total number of tests on which we compared Z3str2 with other solvers is approximately 69,000.

Table 1 summarizes the experimental settings. Z3str2 and CVC4 both accept constraints in similar formats, and so we run Z3str2 and CVC4 on all suites. PISA-MONA and Stranger can only be invoked by APIs. Although S3 and Kaluza accept input from a file, the input language they understand has limited expressiveness. Therefore, we compare Z3str2 to them using the benchmarks over which those solvers were originally evaluated.

Kaluza Benchmark suite. The Kaluza constraints were generated by a JavaScript symbolic execution engine [33], where length, concatenation and (finite) RE membership queries occur frequently. Despite the redundancy across many of the problems, it is still a “classical” test suite over which string solvers are commonly benchmarked. Both CVC4 and S3 were originally evaluated only on this suite, which consists of approximately 50,000 problems in the Kaluza format. The CVC4 team selected 47,284 of them, and translated them into the CVC4 format. The S3 team did the translation to S3 format. We wrote translators from CVC4 to Z3str2, and from Z3str2 to CVC4. The timeout threshold for comparison over this suite was set at 20 seconds per problem, which was the threshold used in CVC4 [26].

PISA Benchmark suite. While the Kaluza suite is large and includes string problems of varying sizes, it only contains a small subset of string operations. To make the comparison more comprehensive, we included constraints from real-world Java sanitizer methods that were used in the evaluation of the PISA system [35]. Sanitizers cleanse user input to remove the threat of an injection attack. They are usually complex and utilize various primitive string operations. We generated two groups of constraints: First, as in the PISA paper, we encode the semantics of the sanitizers and check the return value(s) against predefined attack patterns (such as cross-site scripting (XSS)). In the second group, we also encode input constraints per the application defining the sanitizer. For the PISA suite, we set a timeout of 200 seconds due to its higher complexity.

Table 2: Results on Kaluza suite[33].

	Z3str2		Z3str2 w/o integration		CVC4		S3		Kaluza	
	✓	×	✓	×	✓	×	✓	×	✓	×
sat	34859	0	32752	0	33190	0	32503	488	21651	n/a
unsat	11799	0	11313	0	11625	0	11351	412	12099	10909
unknown	626†		395†		0		0		0	
timeout	0		2824		2469		989		340	
tool reports error	0		0		0		2		2285	
crash	0		0		0		1539		0	
Total time (sec)	4288.8(1x)		61232.8(14.3x)		52478.8 (12.2x)		22543.4 (5.3x)		46753.9 (10.9x)	

† ‘unknown’ indicates Z3str2 detected and avoided overlapping arrangements.

AppScan Benchmark suite. The third suite of benchmarks is derived from security warnings output by IBM Security AppScan Source Edition [2], an application sold commercially by IBM. These reflect potentially vulnerable information flows, represented as traces of program statements, which yield more representative real-world constraints than focusing on sanitizers only. We ran AppScan on popular websites to obtain traces. Similar to the PISA benchmarks, the AppScan constraints also utilize a rich set of string operators. As with PISA, the timeout here was also set at 200 seconds per benchmark.

Kausler Benchmark suite. The final suite is extracted from eight Java programs by Scott Kausler [22]. They represent path conditions obtained from dynamic symbolic execution, and are pure string constraints [21]. Unlike other benchmarks, Kausler’s suite does not dump string constraints to file but instead calls the solvers via an API. The suite contains 174 path condition encoding files, and the resulting constraints are input to the solvers in-memory via their APIs. The comparison [22] was originally done using a driver interface [3]. However, we observed bugs ranging from JNI issues for Stranger to generation of invalid constraints for Z3str2. We made our best attempt to compare Z3str2 with Stranger using modified interfaces [1] patched by both the Stranger team and us.

We performed all the experiments on a workstation running Ubuntu 12.04 with an Intel i7-3770 CPU and 8GB of RAM memory. For reproducibility, we have made the Z3str2 source code publicly available⁴. For the other solvers, we used V1.5-prerelease of CVC4; the version of S3 from the original paper [36]; the Kaluza version from the CVC4 paper with “var” as the query string; and Stranger from [4].

9.1.2 Performance results

Kaluza suite. In Table 2, “tool reports error” counts the number of inputs on which the solver reports an error. “crash”, instead, refers to runtime errors such as segfaults. For “sat” and “unsat”, × denotes the number of provably incorrect results (either an “unsat” response where the problem has a verified solution or a “sat” response with

⁴ <https://sites.google.com/site/z3strsolver/>

an infeasible solution, as defined in [26]), and $\sqrt{}$ the rest. The comparison involves Z3str2 without bidirectional string-integer theory integration, CVC4 and S3, but not PISA, as PISA cannot model string lengths or symbolic arithmetic operations which are intensive in the suite.

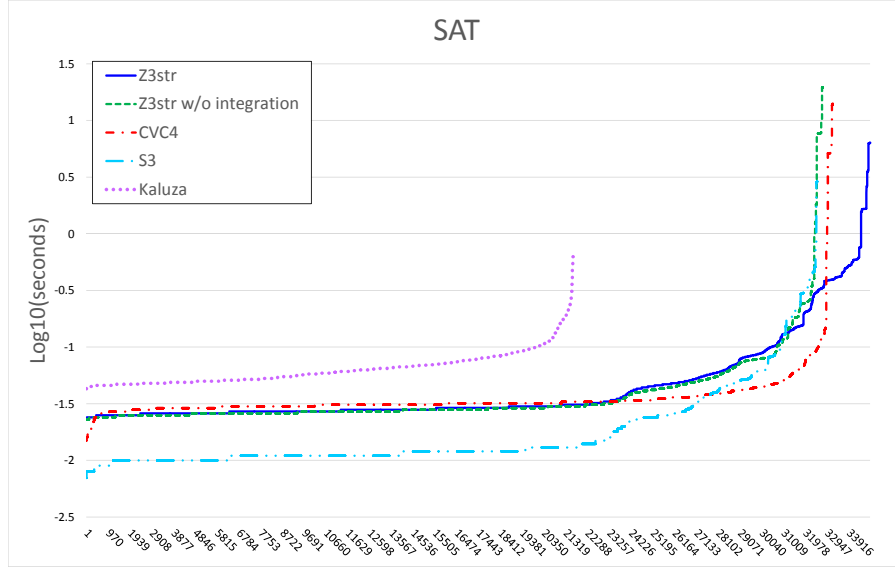


Fig. 16: Cactus plots for the Kaluza benchmark (w/o incorrect results): SAT instances

According to Table 2, neither Z3str2 nor CVC4 report any provably incorrect result, though Z3str2 is more effective and can solve more cases (46658 compared to 44815) without timeouts. Though Z3str2 additionally has 626 unknown cases, CVC4 times out on all these cases. Z3str2 without bidirectional integration solves 2593 fewer cases and times out more often. S3 has errors in both directions, as well as a total of 989 timeouts, while Kaluza suffers from fewer timeouts (340) but has many sat-as-unsat errors (10909). Kaluza therefore is unsound. Since Kaluza only provides assignments for variables matching the query, SAT answers are not verifiable. Both S3 and Kaluza also have tool errors (2 and 2285, respectively). In addition, S3 crashes on 1539 cases. To compare performance on the sat and unsat Kaluza cases across the different solvers, we created the cactus plots in Figures 16a (sat) and 17a (unsat). The vertical axis measures solving time on a logarithmic scale; the horizontal axis counts the total number of instances that (individually) took at most the indicated time to be solved. Incorrect results are excluded. In both categories, Z3str2 and CVC4 have comparable performance, while Z3str2 solves more cases and is faster on complex cases. S3 and Z3str2 without string-integer integration are slower. Kaluza has the worst performance.

PISA suite. Table 3 presents the results on the PISA benchmarks. The “string operators stats” column lists the involved operations and their number of occurrences. In

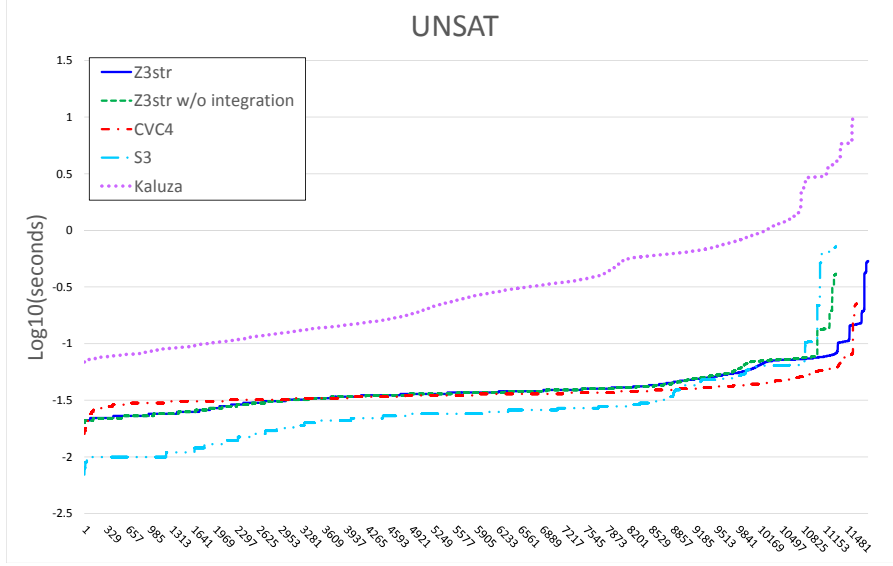


Fig. 17: Cactus plots for the Kaluza benchmark (w/o incorrect results): UNSAT instances

addition, we also count the number of variables and predicates for each format. In this comparison, we included CVC4 and PISA, but not S3 and Kaluza, as we were not able to model popular string operations such as `indexof` using their language. Besides, for PISA, while one group of constraints is equivalent to the MONA program generated by PISA, enabling proper comparison, the other group requires changes to the PISA translation algorithm (to fix the input constraints as well as the negative output constraints), and thus the respective comparisons were not possible. From Table 3, we have the following observations. First, Z3str2 reports 8 sat cases compared to 6 by CVC4 and 2 timeouts. For the 6 sats in common, Z3str2 solves them in 1.069s while CVC4 requires 51.394s. Second, MONA and Z3str2 are in agreement. MONA runs faster on the sat cases, though it cannot generate satisfying string assignments, and has comparable performance to Z3str2 on the unsat cases.

AppScan suite. The results of the third comparison over the AppScan benchmarks appear in Table 4. Correlation between variables is relatively high, as demonstrated by the “var” and “pred” columns. As shown in the table, Z3str2 reports sat on 8 cases while CVC4 agrees on 4 and times out on the rest. The performance gap between the solvers on sat cases in agreement is significant: Z3str2 completes in 0.707s, whereas the CVC4 solving time is 154.852s.

Kausler suite. Table 5 shows the average solving time of the instances generated by the constraint generation tool in the Kausler suite. We were able to run the Stranger tool on five sets in this suite. Both CVC4 and Z3str2 are faster than Stranger on two of these sets, while Stranger is faster than Z3str2 and CVC4 on the remaining three.

Table 3: Results on constraints generated from sanitizers detected by PISA [35].

input	string operators stats (omitting eq and dis-eq)	Z3str2				CVC4				PISA-MONA			
		var	pred	result	time (s)	var	pred	result	time (s)	var	pred	result	time (s)
pisa-000.smt2	contains (3), indexof (1), substring (1)	4	12	sat (✓)	0.164	4	12	sat (✓)	0.264	9	301	sat (?+)	0.029
pisa-001.smt2	contains (1), indexof (1), substring (1)	4	9	sat (✓)	0.114	4	9	sat (✓)	0.032	— ⁺	— ⁺	— ⁺	— ⁺
pisa-002.smt2	contains (4)	2	10	sat (✓)	0.114	2	10	sat (✓)	50.871	—	—	—	—
pisa-003.smt2	contains (3), concat (1)	3	11	unsat (✓)	0.064	3	11	timeout	200.00	—	—	—	—
pisa-004.smt2	contains (2), indexof (1), length (1), lastIndexof [†] (1), substring (2)	7	22	unsat (✓)	0.038	10	32	timeout	200.00	9	331	unsat (✓)	0.041
pisa-005.smt2	indexof (1), lastIndexof [†] (1), length (1), substring (2),	7	23	sat (✓)	0.115	10	33	sat (✓)	0.165	—	—	—	—
pisa-006.smt2	indexof (1), lastIndexof [†] (1), length (1), substring (2), contains (1)	7	24	unsat (✓)	0.039	11	36	timeout	200.00	9	331	unsat (✓)	0.038
pisa-007.smt2	indexof (2), lastIndexof [†] (1), length (1), substring (2), contains (1)	8	26	unsat (✓)	0.042	11	36	timeout	200.00	9	324	unsat (✓)	0.039
pisa-008.smt2	replace (5), contains (2)	6	13	sat (✓)	0.214	6	13	timeout	200.00	9	283	sat (?+)	0.031
pisa-009.smt2	replace (2), concat (1), contains (2)	3	8	sat (✓)	0.447	3	8	sat (✓)	0.046	9	292	sat (?+)	0.054
pisa-010.smt2	replace (2), concat (1)	3	6	sat (✓)	0.165	3	6	timeout	200.00	—	—	—	—
pisa-011.smt2	replace (1), concat (2)	3	6	sat (✓)	0.115	3	6	sat (✓)	0.016	—	—	—	—

⁺ We could not generate constraints without changing PISA. No string solutions are generated so it's not verifiable.

[†] CVC4 doesn't provide operator 'lastIndexof'. We encode it with operators "concat", "length" and "contains".

Table 4: Results on constraints derived from AppScan traces [2].

input	string operators stats (omitting eq and dis-eq)	Z3str2				CVC4			
		var	pred	result	time (s)	var	pred	result	time (s)
t01.smt2	indexof (4), substring (3)	7	37	sat (✓)	0.265	7	37	timeout	200.00
t02.smt2	concat (3), membership (1), regexConcat (2), regexUnion (14), str2Regex (17), length (1)	5	47	sat (✓)	0.215	5	33	sat (✓)	0.026
t03.smt2	concat (3), membership (1), regexConcat (2), regexUnion (14), str2Regex (17), length (1)	5	46	sat (✓)	2.519	5	32	timeout	200.00
t04.smt2	concat (5), membership (1), regexConcat (2), regexUnion (14), str2Regex (17), length (1)	6	50	sat (✓)	4.574	6	35	timeout	200.00
t05.smt2	concat (3), membership (1), regexConcat (2), indexof (1) regexUnion (14), str2Regex (17), length (2), substring (1)	8	56	sat (✓)	2.770	8	42	timeout	200.00
t06.smt2	concat (1), indexof (3), endsWith (5)	5	33	sat (✓)	0.214	5	33	sat (✓)	3.021
t07.smt2	concat (6), regexStar (2), str2Regex (4), endsWith (2) regexUnion (2), membership (2), startsWith (2)	8	32	sat (✓)	0.114	8	29	sat (✓)	0.115
t08.smt2	concat (2), regexStar (2), str2Regex (4), endsWith (2) regexUnion (2), membership (2), startsWith (2)	5	23	sat (✓)	0.164	5	22	sat (✓)	151.663

However, these findings in the comparisons with Stranger should be qualified. First, Stranger crashes or hangs on 98 files. Z3str2 neither crashes nor hangs nor times out on any of the generated instances. We have omitted these 98 files from our comparison. Additionally, Stranger over-approximates disequalities (\neq operator) among variables that can represent multiple strings [5]. We observe that such cases commonly exist in all sets (the percentages of instances with \neq operators in each set are 83.4%, 61.7%, 79.0%, 96.0% and 95.0% respectively, and many fall into this category of disequalities among variables that represent multiple strings). This implies that Stranger produces unsound results. We believe that some of these constraints are easy for Stranger thanks to this over-approximation. By contrast Z3str2 correctly implements all operators and predicates in its input language. Finally, Stranger requires

Table 5: Average solving time per constraint instance generated from Kausler suite.

	instance #	Stranger (ms)	CVC4 (ms)	Z3str2 (ms)
beasties	7230	51.8	6.4	6.4
jerichoHTMLParser	1275	5.9	9.5	10.7
mathParser	9138	1.4	12.8	39.9
mathQuizGame	21	9.4	8.5	7.1
naturalCLI	2367	3.0	10.0	23.4

that integers occurring as indices and length bounds be constant, whereas Z3str2 and most other competing solvers support integers symbolically thus providing expressive power that is essential in practice.

The results also show that CVC4 and Z3str2 have comparable performance on three sets but CVC4 is faster than Z3str2 on the other two. A closer look at the constraints in *mathParser* and *naturalCLI* suggests that CVC4 usually has better performance when the variables are loosely constrained. We plan to introduce heuristics for such cases to improve Z3str2’s performance in future versions.

The general trend, across all benchmark suites, is that CVC4 has comparable performance to Z3str2 although CVC4 times out far more often than Z3str2, whereas S3 is significantly slower. These results establish the efficacy of both techniques presented in this paper.

Detection of overlapping variables. Z3str2 can decide either sat or unsat on 98.7%, 100% and 100% of the instances in the Kaluza, PISA and AppScan suites, respectively. CVC4, in comparison, achieves 94.8%, 50% and 50%. For unknowns reported by Z3str2 on the Kaluza instances, which occur in merely 1.3% of the cases, CVC4 times out on all of them. This lends support to our design choice of purposely pruning away parts of the solution space (those with overlapping arrangements) to avoid nontermination.

String and Integer Theory Integration. As the comparisons between Z3str2 versions with and without the integration clearly demonstrate, there is significant gain thanks to tightening the integer and string theory integration, which enables generation of implied constraints in both domains for more aggressive elimination of assignments unsatisfying for combined string-integer constraints.

9.2 Evaluating Binary Search Based heuristics for Large Strings

The binary search based approach with backtracking functionality helps to prune the search space efficiently when the input contains operations on very long strings. We used the string-integer solver for the comparison and performance evaluation of the binary search approach. We performed two sets of experiments for this evaluation, which are discussed below.

9.2.1 Performance Evaluation of String Lengths in Z3str2

The first set of experiments demonstrates the inefficiency of the linear search heuristic in the previous version of Z3str2. For this purpose, we chose a simple constraint

Table 6: Performance of Z3str2 over string lengths

Benchmark	Length (l)	Binary Search (seconds)	Naive Search (seconds)
test1	10	0.042	0.040
test2	50	0.045	0.041
test3	100	0.039	0.051
test4	500	0.041	0.254
test5	1000	0.055	0.928
test6	5000	0.056	46.699
test7	10000	0.088	Timeout
test8	50000	0.323	Timeout
test9	100000	0.842	Timeout
test10	500000	21.084	Timeout
test11	1000000	105.636	Timeout

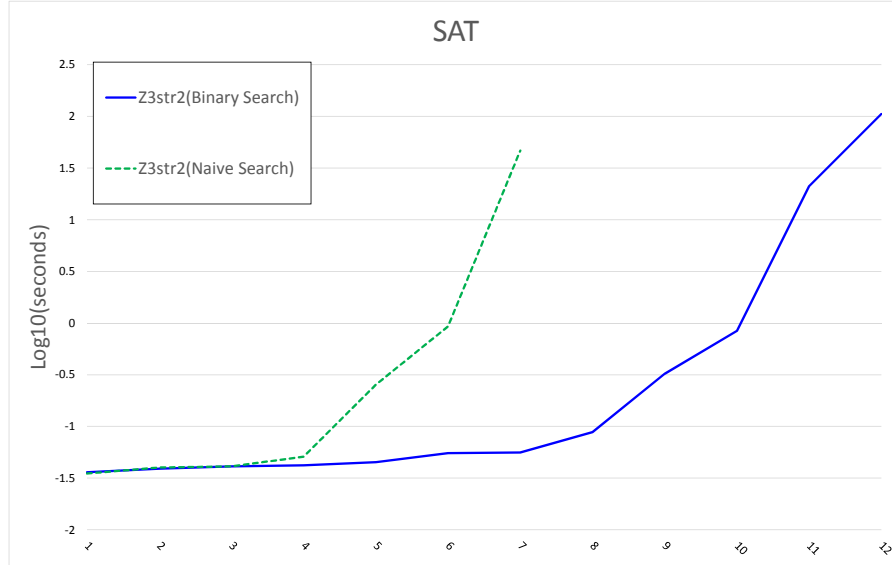


Fig. 18: Cactus plots for the length test

set containing only the length function. The constraints we chose are of the form $\text{length}_{str} = l$, where the choice of l is as described in Table 6. We performed a series of 11 tests and used 200 seconds as the timeout for each test case. The stable release of the Z3str2 solver in the linear search mode solved constraints of length up to 5000, but timed out on all larger instances. The solver took more than 12 hours to terminate when the length is on the order of 10^5 . However, the solver in the binary search mode performed very efficiently and solved all instances for l up to and including 10^6 . A comparison of solving time for the two search heuristics in Z3str2 is shown in Table 6 and the cactus plots in Figure 18.

9.2.2 Evaluation on benchmark suite

To measure the efficacy of our approach, we designed a benchmark suite containing 205 test cases, containing handcrafted constraints involving operations on large strings. The benchmark consists of constraints on various string operations supported by Z3str2 such as *Length*, *Concat*, *Indexof*, *Substring*, *EndsWith*, *StartsWith*, *Replace*, etc. Furthermore, we translated the constraints into the CVC4 input language and compared the performance of the benchmarks against CVC4 version 1.5. We used 200 seconds as the timeout for each constraint input.

Table 7: Performance on benchmark suite

	Z3str2 (Binary Search)	Z3str2 (Naive Search)	CVC4
SAT	169	138	126
UNSAT	34	34	19
UNKNOWN	2	2	0
Timeout	0	31	60
Tool reports error	0	0	0
Crash	0	0	0
Total time (factor)	41.697 (1x)	9569.639 (x229)	12014.893 (x264)

† 'unknown' indicates Z3str2 detected and avoided overlapping arrangements.

The results of the comparison are presented in Table 7. The *SAT* and *UNSAT* column denotes the number of SAT and UNSAT answers respectively. According to Table 7, Z3str2 with binary search solves all test instances very quickly compared to the others. Z3str2 with linear search timed out after 200 seconds on 31 test cases, whereas CVC4 timed out on 60 test cases. The Z3str2 solver detects overlapping arrangements in both search modes and reports “UNKNOWN” in several cases. However, neither of these solvers crashed or reported errors in any of the test instances.

The cactus plots for the sat and unsat results are presented in the Figure 19 and Figure 20 respectively. Comparing the time taken by each solver, Z3str2 with the binary search heuristic exhibits the best overall performance. The overall time of the binary search based approach is 41.7 seconds for these 205 test cases, whereas the linear search method used a total of 9569.6 seconds and CVC4 needed 12014.9 seconds. The binary search based approach used by the new version of Z3str2 is about 229 times faster than the linear search approach of Z3str2 and 288 times faster than CVC4.

9.3 Evaluating Improvements in the Context of Symbolic Execution

We demonstrate the performance improvements enabled by a string constraint solver for the combined string-integer theory with respect to the example shown in Fig. 15 by applying both our technique and KLEE, a state-of-the-art symbolic execution engine, to this code. The goal was to detect the heap corruption threat in that code.

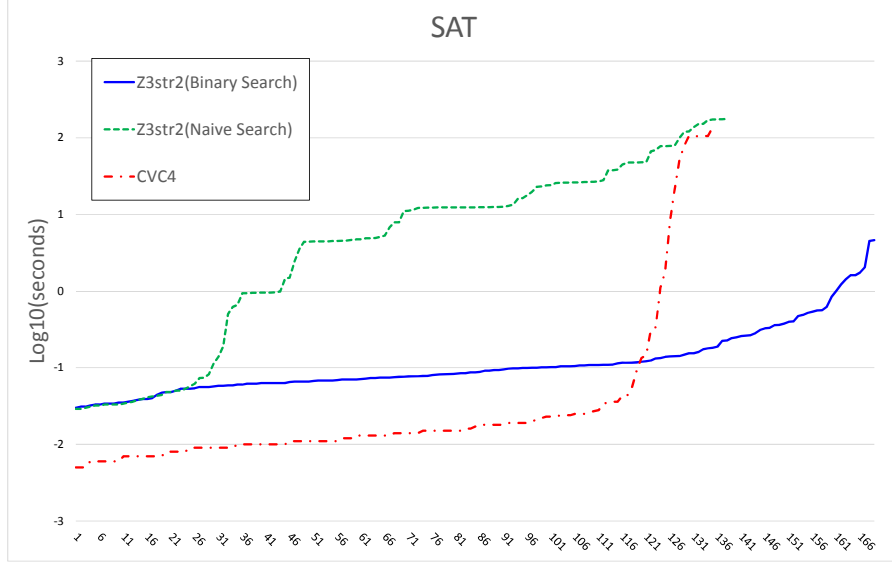


Fig. 19: Cactus plots: SAT instances

We faithfully encoded the program snippet in `check_login()` as string/bit-vector constraints. We then checked whether the buffer pointed to by `_username` is susceptible to overflow.

Notice that `len` is an unsigned short variable, and thus ranges from 0 to $2^{16} - 1$ (65,535). As it represents the buffer size, it determines the number of concrete execution paths KLEE has to enumerate, as well as the search space for library-aware solving. By contrast, the constraints generated by a library-aware SMT solver declaratively model `strlen` as part of the SMT solver logic.

To characterize performance trends, we consider two different precision settings for string length: 8-bit and 16-bit. We used 120 minutes as the timeout value. There was no need to go beyond 16 bits since KLEE was already significantly slower at 16 bits relative to the library-aware SMT solver. Note that KLEE performs poorly on this test because it has to explore a large number of paths, not because the individual path constraints are difficult to solve (as illustrated in Fig. 21).

Table 8: Vulnerability Detection using KLEE (solving bit-vector constraints) and Z3str2

Precision	Upper Bound	String-Integer Theory		KLEE (bit-vector)
		Z3str2 (binary search)	Z3str2 (linear search)	
8-bit	255	1.1s	167s	300s
16-bit	65535	787.6s	timeout (7200s)	timeout (7200s)

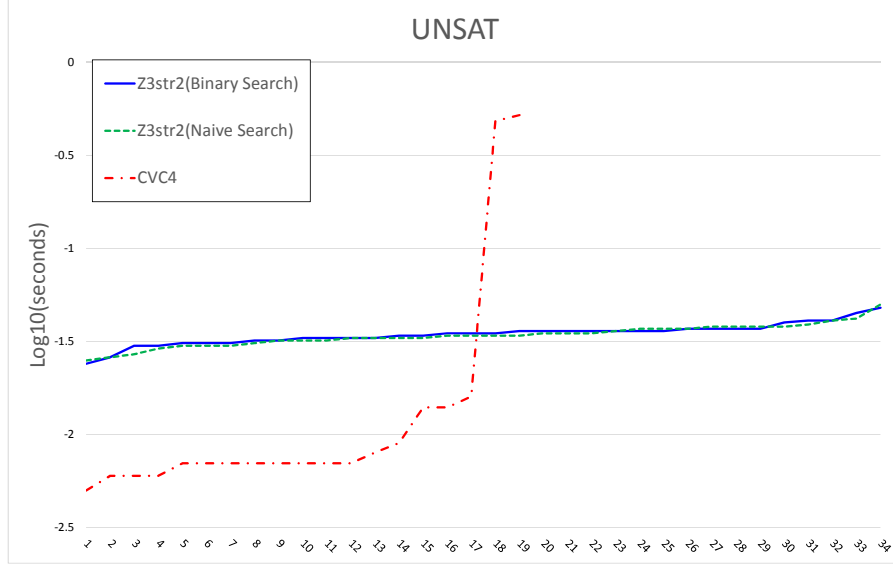


Fig. 20: Cactus plots: UNSAT instances

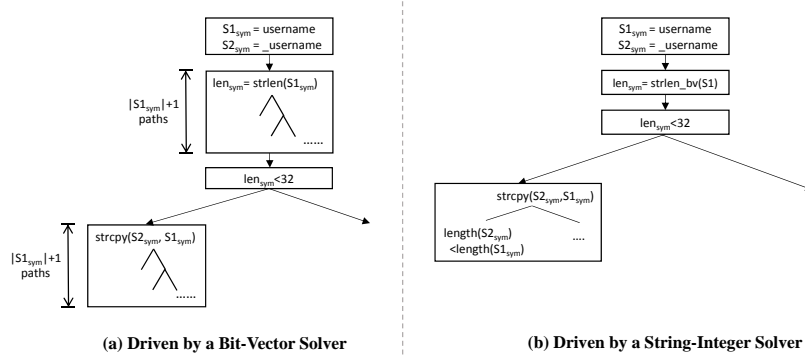


Fig. 21: Comparison between Symbolic Execution Driven by a Bit-Vector solver (KLEE) and a String Solver over the String and Integer Theories.

The results are provided in Table 8. Under both precision settings, KLEE is consistently and significantly slower than the library-aware solving technique. In particular, if we represent numeric values using 16 bits, then KLEE is not able to find a path that triggers the bug in 120 minutes, while Z3str2 with binary search heuristics can find a satisfying solution in 787.6 seconds.

The benefit thanks to native modeling of strings is clear. Suppose both `username` and `_username` are symbolic string variables. In Figure 21(a), as KLEE forks a new state for each character, an invocation of `strlen` on a symbolic string $S1_{sym}$ of size $|S1_{sym}|$ will generate and check $|S1_{sym}| + 1$ path constraints (one per each possible

length value between 0 and $|S1_{sym}|$). In Figure 21(b), in contrast, the constraint encoding essentially captures the semantics of the program without explicitly handling the loops in `strlen` and `strcmp`. Thus, only one query is needed to check whether the length $S2_{sym}$ can be smaller than the length $S1_{sym}$.

10 Related Work

The theory considered in this paper – namely the quantifier-free (QF) theory T_{wlr} over word equations, membership predicate over REs, and length function – is a multi-sorted theory with string (str) and numeric (num) sorts. Makanin was the first to show, in 1977, that the QF theory of word equations is decidable [27]. Since, many have improved upon this seminal result [34, 20, 29, 30, 19]. In particular, Plandowski proved that this problem is in PSPACE [30]. Despite decades of effort, the satisfiability problem for T_{wlr} remains open [27, 30, 14, 19]. Still, many practical solvers have been proposed.

Automata-based Solvers. Regular languages (or automata), as well as context-free grammars (CFGs), can be used to represent strings and handling regex-related operations. JSA [11] computes CFGs for string variables in Java programs. Hooimeijer et al. [16] suggest an optimization whereby automata are built lazily. A primary challenge faced by automata-based approaches, which we do not suffer from, is to capture the connections between strings and other domains, e.g., integers. To overcome this limitation, refinements have been proposed. JST [15] extends JSA. It asserts length constraints in each automaton, and handles numeric constraints after conversion. PISA [35] encodes Java programs into M2L formulas that it discharges to the MONA solver to obtain path- and index-sensitive string approximations. PASS [25, 24] combines automata and parametrized arrays for efficient treatment of unsat cases. Stranger is a powerful extension of string automata with arithmetic automata [38, 39].

Bit-vector based solvers. Another group of solvers converts string constraints to constraints into other domains such as integers or bit-vectors. HAMPI [23] is an efficient solver that represents strings as bit-vectors, though it requires the user to provide an upper bound on string lengths. Early versions of Kaluza [33] extended both STP [13] and HAMPI to support mixed string and numeric constraints represented as bit-vector. A similar approach powers Pex [8], though strings are reduced to integer abstractions.

Word-based string solvers. CVC4 [26] handles constraints over the theory of unbounded strings with length and RE membership. Solving is based on multi-theory reasoning backed by the DPLL(T) architecture combined with existing SMT theories. The Kleene star operator in RE formulas is dealt with via unrolling as in Z3str2. S3 [36] is another word-based solver, and it can be viewed as an extension of an early version of Z3-str. Roughly speaking, CVC4, S3 and Z3str2 embody similar approaches, and hence CVC4 and S3 can also benefit from the techniques proposed in this paper.

11 Conclusion and Future Work

We have described three techniques that dramatically improve the efficiency of word-based string solvers: (i) a sound and complete procedure to detect overlapping vari-

ables, thereby automatically identifying and avoiding sources of nontermination; (ii) tight bi-directional integer/string theory integration, which prunes a vast array of inconsistent search candidates; and (iii) further optimization of the string length search strategy based on a binary search heuristic.

We have implemented these techniques on top of Z3-str as Z3str2. We show the efficacy of these techniques through an extensive set of experiments, comparing Z3str2 with the CVC4, S3, Kaluza, PISA and Stranger solvers over four benchmark suites derived from real-world applications. In addition, we empirically show the performance improvements enabled by Z3str2 in the context of symbolic execution of string manipulation programs. We compare the performance of symbolic execution powered by a bit-vector solver and Z3str2. We also evaluate our ability to capture the semantics of string library functions without having to explore a huge number of program execution paths in order to enumerate all possible string lengths.

The technique we present for sound and complete detection of overlapping variables is applicable to other word-based string solvers in addition to Z3str2. As we discussed, the problem of overlapping variables is a key difficulty in reasoning about strings and a source of non-termination in existing tools. The solution we present can be used to avoid these potentially infinite search paths and allow the solve to continue; in future we plan to explore additional heuristics that give the solver the ability to reason about a fragment of cases with overlapping variables in an efficient manner.

Bi-directional integration between the string and integer theory solvers is an extension beyond the classical DPLL(T) model for SMT solvers, which limits the ability for independent theory solvers to communicate. Our experiments have shown that this technique is highly effective at pruning the search space for the solver and preliminary evidence suggests that the binary search heuristic enabled by this integration can improve performance further when dealing with very long strings. We plan to leverage this technique further, and investigate more ways to share information with other theory solvers and additional methods by which our theory solver can guide the search.

While summary-based symbolic execution has been studied previously (e.g. as part of the S-Looper tool [37]), we are not aware of existing work in which SMT solvers are extended with first-class support for programming-language library functions declaratively as part of their logic. One recent application of a similar concept is discussed in [18], in which models of design patterns are abstracted into a symbolic execution engine. The ability to perform a similar analysis at the level of detail of individual library methods can enhance library-aware symbolic execution as part of a library-aware SMT solver. We intend to explore this idea further in the future to broaden its applicability beyond the current context.

References

1. Z3str2 String Constraint Solver. <https://sites.google.com/site/z3strsolver/>.
2. IBM Security AppScan Source. <http://www-03.ibm.com/software/products/en/appscan-source>.
3. Kausler Suite. <https://github.com/BoiseState/string-constraint-solvers>.
4. LibStranger. <https://github.com/vlab-cs-ucsb/LibStranger>.
5. Personal communications with the Stranger team. 2015.

6. Parosh Aziz Abdulla, Mohamed Faouzi Atig, Yu-Fang Chen, Lukás Holík, Ahmed Rezine, Philipp Rümmer, and Jari Stenman. String constraints for verification. In *Proceedings of the 26th International Conference on Computer Aided Verification, CAV'14*, pages 150–166, 2014.
7. Clark Barrett, Pascal Fontaine, and Cesare Tinelli. The Satisfiability Modulo Theories Library (SMT-LIB). www.SMT-LIB.org, 2016.
8. Nikolaj Bjørner, Nikolai Tillmann, and Andrei Voronkov. Path feasibility analysis for string-manipulating programs. In *Proceedings of the 15th International Conference on Tools and Algorithms for the Construction and Analysis of Systems, TACAS '09*, pages 307–321, 2009.
9. Cristian Cadar, Daniel Dunbar, and Dawson Engler. Klee: Unassisted and automatic generation of high-coverage tests for complex systems programs. In *Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation, OSDI'08*, pages 209–224, Berkeley, CA, USA, 2008. USENIX Association.
10. Vitaly Chipounov, Volodymyr Kuznetsov, and George Candea. S2e: A platform for in-vivo multi-path analysis of software systems. In *Proceedings of the Sixteenth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS XVI*, pages 265–278, New York, NY, USA, 2011. ACM.
11. Aske Simon Christensen, Anders Møller, and Michael I. Schwartzbach. Precise analysis of string expressions. In *Proceedings of the 10th International Conference on Static Analysis, SAS'03*, pages 1–18, 2003.
12. Leonardo De Moura and Nikolaj Bjørner. 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, TACAS'08*, pages 337–340, 2008.
13. Vijay Ganesh and David L. Dill. A decision procedure for bit-vectors and arrays. In *Proceedings of the 19th International Conference on Computer Aided Verification, CAV'07*, pages 519–531, 2007.
14. Vijay Ganesh, Mia Minnes, Armando Solar-Lezama, and Martin Rinard. Word equations with length constraints: what's decidable? In *HVC'12*, 2012.
15. Indradeep Ghosh, Nastaran Shafiei, Guodong Li, and Wei-Fan Chiang. JST: An Automatic Test Generation Tool for Industrial Java Applications with Strings. In *Proceedings of the 2013 International Conference on Software Engineering, ICSE '13*, pages 992–1001, 2013.
16. Pieter Hooimeijer and Westley Weimer. Solving string constraints lazily. In *Proceedings of the IEEE/ACM international conference on Automated software engineering, ASE '10*, pages 377–386, 2010.
17. J.E. Hopcroft, R. Motwani, and J.D. Ullman. *Introduction to automata theory, languages, and computation*. Pearson/Addison Wesley, 2007.
18. Jinseong Jeon, Xiaokang Qiu, Jonathan Fetter-Degges, Jeffrey S. Foster, and Armando Solar-Lezama. Synthesizing framework models for symbolic execution. In *Proceedings of the 38th International Conference on Software Engineering, ICSE '16*, pages 156–167, New York, NY, USA, 2016. ACM.
19. Artur Jež. Recompression: Word equations and beyond. In *Developments in Language Theory, Lecture Notes in Computer Science*, pages 12–26. 2013.
20. Juhani Karhumäki, Filippo Mignosi, and Wojciech Plandowski. The expressibility of languages and relations by word equations. *J. ACM*, 47(3):483–505, May 2000.
21. Scott Kausler. Evaluation of string constraint solvers using dynamic symbolic execution. Master's thesis, Boise State University, 2014.
22. Scott Kausler and Elena Sherman. Evaluation of string constraint solvers in the context of symbolic execution. In *Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering, ASE '14*, pages 259–270, New York, NY, USA, 2014. ACM.
23. Adam Kiezun, Vijay Ganesh, Philip J. Guo, Pieter Hooimeijer, and Michael D. Ernst. Hampi: A solver for string constraints. In *Proceedings of the Eighteenth International Symposium on Software Testing and Analysis, ISSTA '09*, pages 105–116, 2009.
24. Guodong Li, Esben Andreasen, and Indradeep Ghosh. SymJS: Automatic symbolic testing of javascript web applications. In *Proceedings of the 22Nd ACM SIGSOFT International Symposium on Foundations of Software Engineering, FSE 2014*, pages 449–459, 2014.
25. Guodong Li and Indradeep Ghosh. PASS: String solving with parameterized array and interval automaton. In *9th International Haifa Verification Conference, HVC '13*, pages 15–31. 2013.
26. Tianyi Liang, Andrew Reynolds, Cesare Tinelli, Clark Barrett, and Morgan Deters. A dpll(t) theory solver for a theory of strings and regular expressions. In *Proceedings of the 26th International Conference on Computer Aided Verification, CAV'14*, pages 646–662. 2014.
27. G.S. Makanin. The problem of solvability of equations in a free semigroup. *Math. Sbornik*, 103:147–236, 1977. English transl. in *Math USSR Sbornik* 32 (1977).

28. Yu. Matiyasevich. Word equations, fibonacci numbers, and hilbert's tenth problem. In *Workshop on Fibonacci Words*, 2007.
29. Wojciech Plandowski. Satisfiability of word equations with constants is in pspace. *J. ACM*, 51(3):483–496, May 2004.
30. Wojciech Plandowski. An efficient algorithm for solving word equations. In *Proceedings of the 38th Annual ACM Symposium on Theory of Computing*, STOC '06, pages 467–476, 2006.
31. Xiao Qu and Brian Robinson. A case study of concolic testing tools and their limitations. In *Empirical Software Engineering and Measurement (ESEM), 2011 International Symposium on*, pages 117–126. IEEE Computer Society, 2011.
32. Gideon Redelinghuys, Willem Visser, and Jaco Geldenhuys. Symbolic execution of programs with strings. In *Proceedings of the South African Institute for Computer Scientists and Information Technologists Conference*, SAICSIT '12, pages 139–148, 2012.
33. Prateek Saxena, Devdatta Akhawe, Steve Hanna, Feng Mao, Stephen McCamant, and Dawn Song. A symbolic execution framework for javascript. In *Proceedings of the 2010 IEEE Symposium on Security and Privacy*, SP '10, pages 513–528, 2010.
34. K. Schulz. Makanin's algorithm for word equations-two improvements and a generalization. In K. Schulz, editor, *Word Equations and Related Topics*, volume 572 of *Lecture Notes in Computer Science*, pages 85–150. Springer Berlin / Heidelberg, 1992.
35. Takaaki Tateishi, Marco Pistoia, and Omer Tripp. Path- and index-sensitive string analysis based on monadic second-order logic. *ACM Trans. Softw. Eng. Methodol.*, 22(4):33:1–33:33, October 2013.
36. Minh-Thai Trinh, Duc-Hiep Chu, and Joxan Jaffar. S3: A symbolic string solver for vulnerability detection in web applications. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security*, CCS '14, pages 1232–1243, 2014.
37. Xiaofei Xie, Yang Liu, Wei Le, Xiaohong Li, and Hongxu Chen. S-looper: Automatic summarization for multipath string loops. In *Proceedings of the 2015 International Symposium on Software Testing and Analysis*, ISSTA 2015, pages 188–198, New York, NY, USA, 2015. ACM.
38. Fang Yu, Muath Alkhalaf, and Tevfik Bultan. Stranger: an automata-based string analysis tool for php. In *Proceedings of the 16th international conference on Tools and Algorithms for the Construction and Analysis of Systems*, TACAS'10, pages 154–157, 2010.
39. Fang Yu, Tevfik Bultan, and Oscar H. Ibarra. Symbolic string verification: Combining string analysis and size analysis. In *Proceedings of the 15th International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, TACAS '09, pages 322–336, 2009.
40. Yunhui Zheng, Vijay Ganesh, Sanu Subramanian, Omer Tripp, Julian Dolby, and Xiangyu Zhang. Effective Search-space Pruning for Solvers of String Equations, Regular Expressions and Length Constraints, <https://sites.google.com/site/z3strsolver/publications>. Technical report, 2015.
41. Yunhui Zheng, Xiangyu Zhang, and Vijay Ganesh. Z3-str: A z3-based string solver for web application analysis. In *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*, ESEC/FSE 2013, pages 114–124, 2013.