

A DPLL(T) Theory Solver for Strings and Regular Expressions

Tianyi Liang	}	University of Iowa
Andrew Reynolds		
Cesare Tinelli	}	New York University
Morgan Deters		
Clark Barrett		

Motivation : Security Applications

```
char buff[15];
char pass;
cout << "Enter the password :";
gets(buff);
if (regex_match(buff, std::regex("[A-Z]+")) ) {
    if(strcmp(buff, "PASSWORD")) {
        cout << "Wrong Password";
    } else {
        cout << "Correct Password";
        pass = 'Y';
    }
}
if(pass == 'Y') {
    /* Grant the root permission*/
}
}
```

Encode

```
(set-logic QF_S)
(declare-const input String)
(declare-const buff String)
(declare-const pass0 String)
(declare-const rest String)
(declare-const pass1 String)
(assert (= (str.len buff) 15))
(assert (= (str.len pass1) 1))
(assert (or (< (str.len input) 15)
           (= input (str.++ buff pass0 rest))))
(assert (str.in.re buff
                  (re.+ (re.range "A" "Z"))))
(assert (ite (= buff "PASSWORD")
            (= pass1 "Y")
            (= pass1 pass0)))
(assert (not (= buff "PASSWORD")))
(assert (= pass1 "Y"))
```

Extract

Solve

```
tiliang@milner:~/workspace/security/benchmarks/homemade$ ~/CVC4/bin/pt-cvc4 propsalex.smt2
sat
(model
  (define-fun input () String "AAAAAAAAAAAAAAAAAY")
  (define-fun buff () String "AAAAAAAAAAAAAAAA")
  (define-fun pass0 () String "Y")
  (define-fun rest () String "")
  (define-fun pass1 () String "Y")
)
```

Objectives

- Want solver to handle:
 - (Unbounded) string constraints
 - Length constraints
 - Regular language memberships, ...
 - Theoretical complexity of:
 - Word equation problem: **PSPACE**
 - ...with length constraints: **OPEN**
 - ...with other functions (e.g. `replace`): **UNDECIDABLE**
- } Focus of this talk

Objectives

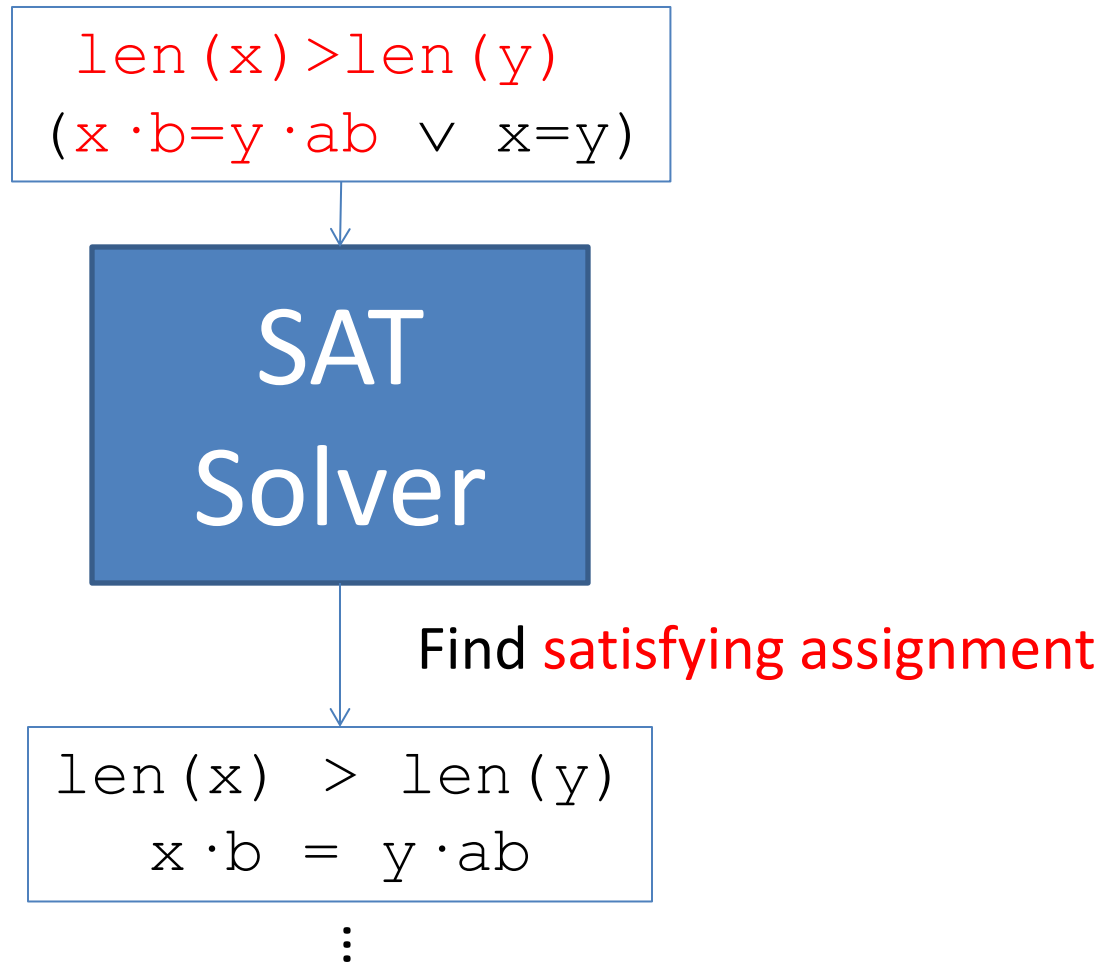
- Instead, focus on solver that is:
 - **Efficient** in practice
 - Tightly **integrated** into SMT architecture
 - Conflict analysis, T-propagation, lemma learning, combination of theories, ...
 - **Robust**

Core Language for Theory of Strings

- Terms are:
 - Constants from a fixed finite alphabet Σ^* (a, ab, cbc...)
 - Free constants or “variables” (x, y, z, w...)
 - String concatenation
 $_ \cdot _ : \text{String} \times \text{String} \rightarrow \text{String}$
 - Length terms
 $\text{len}(_) : \text{String} \rightarrow \text{Int}$
- Example input:

$$\text{len}(x) > \text{len}(y) \quad \wedge \quad (x \cdot b = y \cdot ab \quad \vee \quad x = y)$$

DPLL(T): Find Satisfying Assignment



DPLL(T): Cooperating *Theory Solvers*

⋮

$$\begin{aligned} \text{len}(x) &> \text{len}(y) \\ x \cdot b &= y \cdot ab \end{aligned}$$

Theory
LIA

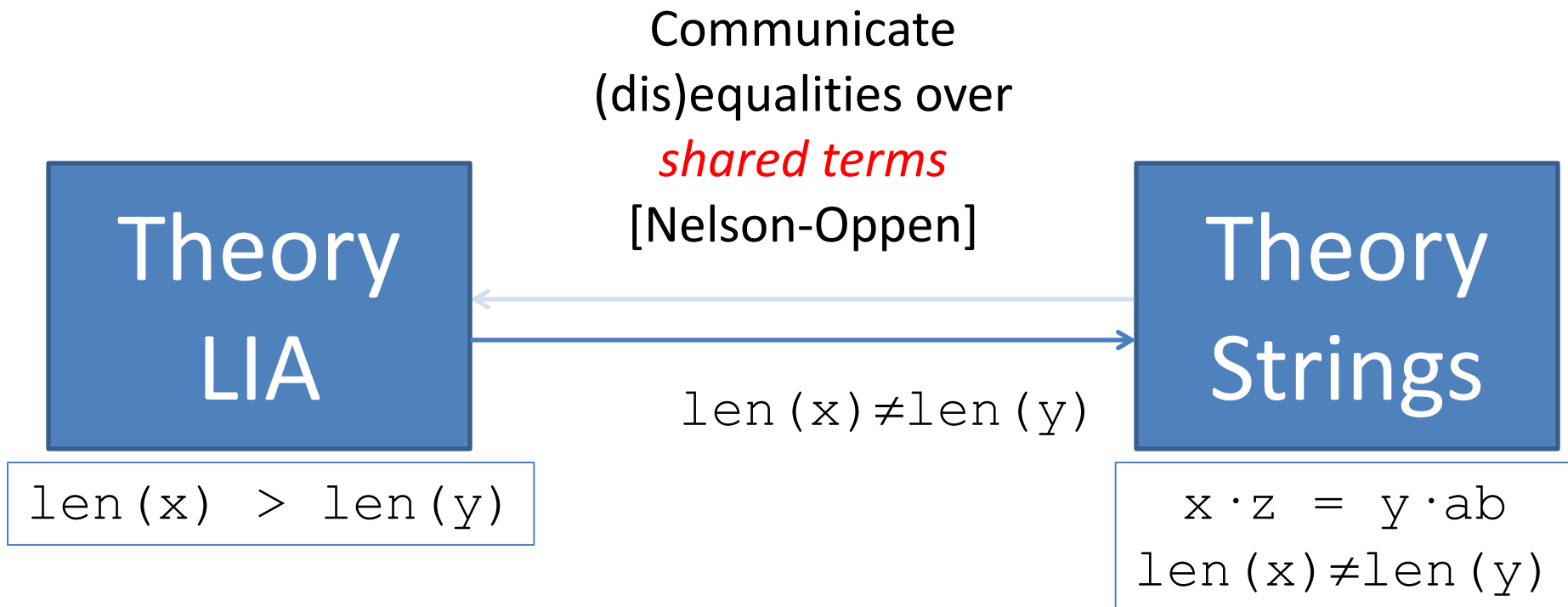
Purify and distribute
constraints to
corresponding **theory**
solvers

Theory
Strings

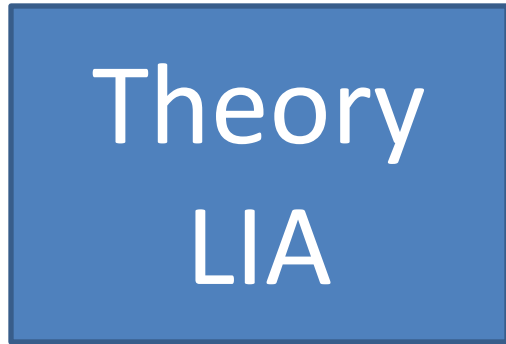
$$\text{len}(x) > \text{len}(y)$$

$$x \cdot z = y \cdot ab$$

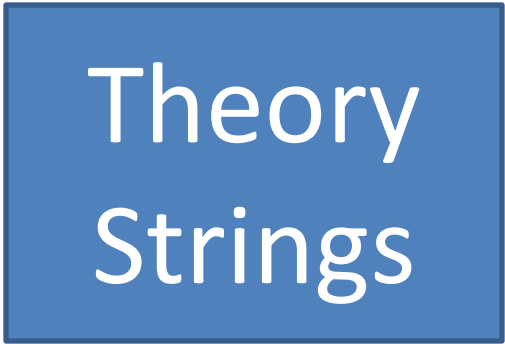
DPLL(T): Cooperating *Theory Solvers*



DPLL(T): Cooperating *Theory Solvers*



A { $\text{len}(x) > \text{len}(y)$



S { $x \cdot z = y \cdot ab$

L { $\text{len}(x) \neq \text{len}(y)$

Summary of Approach

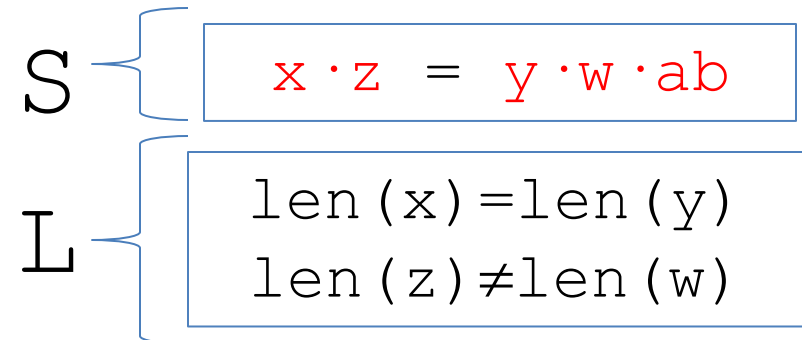
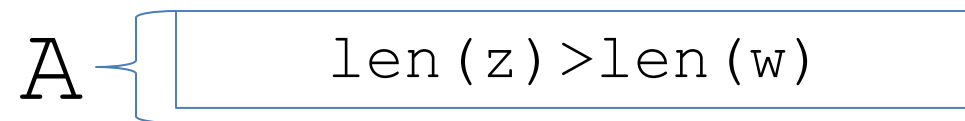
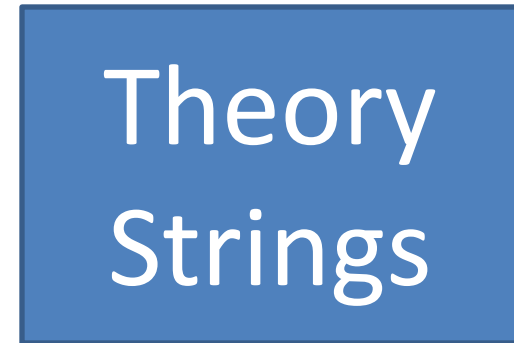
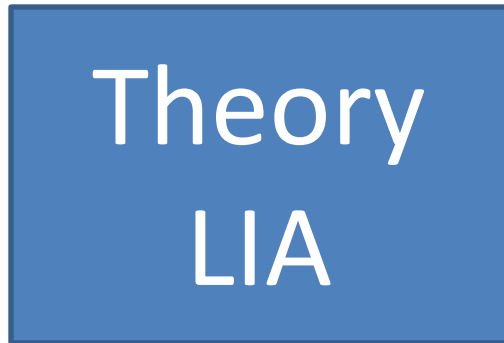
- Approach for $\bar{A} \cup S \cup \bar{L}$ in **four steps**:

1. Check arithmetic constraints \bar{A}
2. Normalize equalities in S
3. Normalize disequalities in S
4. Check cardinality of Σ

Check Length Constraints

1. Check length constraints
2. Normalize equalities
3. Normalize disequalities
4. Check cardinality of Σ

- Add **equalities** to \mathbb{A} based on terms from S



Check Length Constraints

1. Check length constraints
2. Normalize equalities
3. Normalize disequalities
4. Check cardinality of Σ

Theory
LIA

Theory
Strings

$$\begin{aligned} \text{len}(z) &> \text{len}(w) \\ \text{len}(x) + \text{len}(z) &= \\ \text{len}(y) + \text{len}(w) + 2 \end{aligned}$$

$$S \quad x \cdot z = y \cdot w \cdot ab$$

$$L \quad \begin{aligned} \text{len}(x) &= \text{len}(y) \\ \text{len}(z) &\neq \text{len}(w) \end{aligned}$$

\Rightarrow Check if A is satisfiable

Normalize Equalities

1. Check length constraints
2. **Normalize equalities**
3. Normalize disequalities
4. Check cardinality of Σ

$$S \left\{ \begin{array}{l} x \cdot z = y \cdot w \cdot ab \end{array} \right.$$

$$L \left\{ \begin{array}{l} \text{len}(x) = \text{len}(y) \\ \text{len}(z) \neq \text{len}(w) \end{array} \right.$$

- To check satisfiability of equalities in S ,
 - Add additional equalities to S
 - Until pairs of equiv. terms have same **normal form**

Normalize Equalities

1. Check length constraints
2. **Normalize equalities**
3. Normalize disequalities
4. Check cardinality of Σ

$$S \quad \left\{ \begin{array}{l} x \cdot z = y \cdot w \cdot ab \end{array} \right.$$

$$L \quad \left\{ \begin{array}{l} \text{len}(x) = \text{len}(y) \\ \text{len}(z) \neq \text{len}(w) \end{array} \right.$$



Normalize Equalities

1. Check length constraints
2. **Normalize equalities**
3. Normalize disequalities
4. Check cardinality of Σ

$$S \left\{ \begin{array}{l} x \cdot z = y \cdot w \cdot ab \\ x = y \end{array} \right.$$

$$L \left\{ \begin{array}{l} \text{len}(x) = \text{len}(y) \\ \text{len}(z) \neq \text{len}(w) \end{array} \right.$$



II **Propagate**, since $\text{len}(x) = \text{len}(y)$



Normalize Equalities

1. Check length constraints
2. **Normalize equalities**
3. Normalize disequalities
4. Check cardinality of Σ

$$S \left\{ \begin{array}{l} x \cdot z = y \cdot w \cdot ab \\ x = y \end{array} \right.$$

$$L \left\{ \begin{array}{l} \text{len}(x) = \text{len}(y) \\ \text{len}(z) \neq \text{len}(w) \end{array} \right.$$



|| : We have that $\text{len}(z) \neq \text{len}(w)$



Normalize Equalities

1. Check length constraints
2. **Normalize equalities**
3. Normalize disequalities
4. Check cardinality of Σ

$$S \left\{ \begin{array}{l} x \cdot z = y \cdot w \cdot ab \\ x = y \\ z = w \cdot z' \end{array} \right.$$

$$L \left\{ \begin{array}{l} \text{len}(x) = \text{len}(y) \\ \text{len}(z) \neq \text{len}(w) \end{array} \right.$$



\parallel **Decide** $z = w \cdot z'$



Normalize Equalities

1. Check length constraints
2. **Normalize equalities**
3. Normalize disequalities
4. Check cardinality of Σ

$$S \left\{ \begin{array}{l} x \cdot z = y \cdot w \cdot ab \\ x = y \\ z = w \cdot z' \end{array} \right.$$

$$L \left\{ \begin{array}{l} \text{len}(x) = \text{len}(y) \\ \text{len}(z) \neq \text{len}(w) \end{array} \right.$$



||

||



|| Reflexive



Normalize Equalities

1. Check length constraints
2. **Normalize equalities**
3. Normalize disequalities
4. Check cardinality of Σ

S {

$$x \cdot z = y \cdot w \cdot ab$$

$$x = y$$

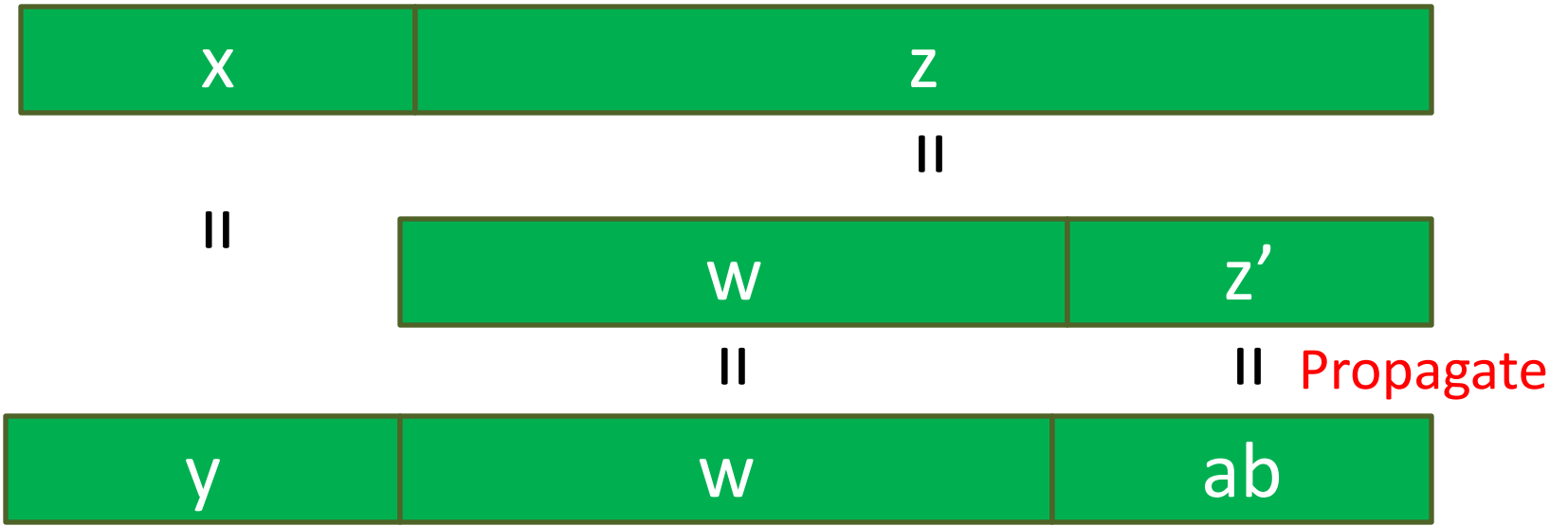
$$z = w \cdot z'$$

$$z' = ab$$

L {

$$\text{len}(x) = \text{len}(y)$$

$$\text{len}(z) \neq \text{len}(w)$$



\Rightarrow Normal form

Normalize Disequalities

1. Check length constraints
2. Normalize equalities
3. **Normalize disequalities**
4. Check cardinality of Σ

- Disequalities normalized analogously

$$S \left\{ \begin{array}{l} x \cdot z \cdot z' \neq y \cdot w \cdot ab \\ x = y \end{array} \right.$$

$$L \left\{ \begin{array}{l} \text{len}(x) = \text{len}(y) \\ \text{len}(z) = \text{len}(w) \end{array} \right.$$



Normalize Disequalities

1. Check length constraints
2. Normalize equalities
3. **Normalize disequalities**
4. Check cardinality of Σ

$$S \left\{ \begin{array}{l} x \cdot z \cdot z' \neq y \cdot w \cdot ab \\ x = y \end{array} \right.$$

$$L \left\{ \begin{array}{l} \text{len}(x) = \text{len}(y) \\ \text{len}(z) = \text{len}(w) \end{array} \right.$$



|| Given



Normalize Disequalities

1. Check length constraints
2. Normalize equalities
3. **Normalize disequalities**
4. Check cardinality of Σ

$$S \left\{ \begin{array}{l} x \cdot z \cdot z' \neq y \cdot w \cdot ab \\ x = y \\ z \neq w \end{array} \right.$$

$$L \left\{ \begin{array}{l} \text{len}(x) = \text{len}(y) \\ \text{len}(z) = \text{len}(w) \end{array} \right.$$



\parallel

\nparallel **Decide** $z \neq w$



\Rightarrow *Normal form*

Check Cardinality of Σ

1. Check length constraints
2. Normalize equalities
3. Normalize disequalities
4. Check cardinality of Σ

- S may be unsatisfiable since Σ is **finite**
 - For instance, if:
 - Σ is a finite alphabet of 256 characters, and
 - S entails that 257 distinct strings of length 1 exist
- Then:
- S is unsatisfiable
- Performed as a last step of our procedure

Rule-Based Procedure

$$\begin{array}{c}
 \dots \\
 \text{F-Unify} \frac{\text{F } s = (w, u, u_1) \quad \text{F } t = (w, v, v_1) \quad s \approx t \in \mathcal{C}(S) \quad S \models \text{len } u \approx \text{len } v}{S := S, u \approx v} \\
 \\
 \text{F-Split} \frac{\text{F } s = (w, u, u_1) \quad \text{F } t = (w, v, v_1) \quad s \approx t \in \mathcal{C}(S) \quad S \models \text{len } u \neq \text{len } v}{u \notin \mathcal{V}(v_1) \quad v \notin \mathcal{V}(u_1)} \\
 \quad \quad \quad S := S, u \approx \text{con}(v, z) \quad || \quad S := S, v \approx \text{con}(u, z) \\
 \\
 \text{F-Loop} \frac{\text{F } s = (w, x, u_1) \quad \text{F } t = (w, v, v_1, x, v_2) \quad s \approx t \in \mathcal{C}(S) \quad x \notin \mathcal{V}((v, v_1))}{S := S, x \approx \text{con}(z_2, z), \text{con}(v, v_1) \approx \text{con}(z_2, z_1), \text{con}(u_1) \approx \text{con}(z_1, z_2, v_2)} \\
 \quad \quad \quad R := R, z \text{ in star}(\text{set } \text{con}(z_1, z_2)) \quad C := C, t \\
 \dots
 \end{array}$$

- Approach is **algebraic**
 - Rules model interaction of string + arithmetic solvers
 - A closed derivation tree \Rightarrow problem is **UNSAT**
 - A state where no rule applies \Rightarrow problem is **SAT**

Theoretical Results

- Our approach is:
 - **Refutation sound**
 - When it answers “UNSAT”, it can be trusted
 - Even for strings of unbounded length
 - **Solution sound**
 - When it answers “SAT”, it can be trusted
- (A version of) our approach is:
 - **Solution complete**
 - When problem is “SAT”, it will eventually find a model
 - Somewhat trivially, by finite model finding
- Our approach is:
 - **Refutation incomplete**
 - When problem is “UNSAT”, it is **not** guaranteed to derive refutation

Experimental Results

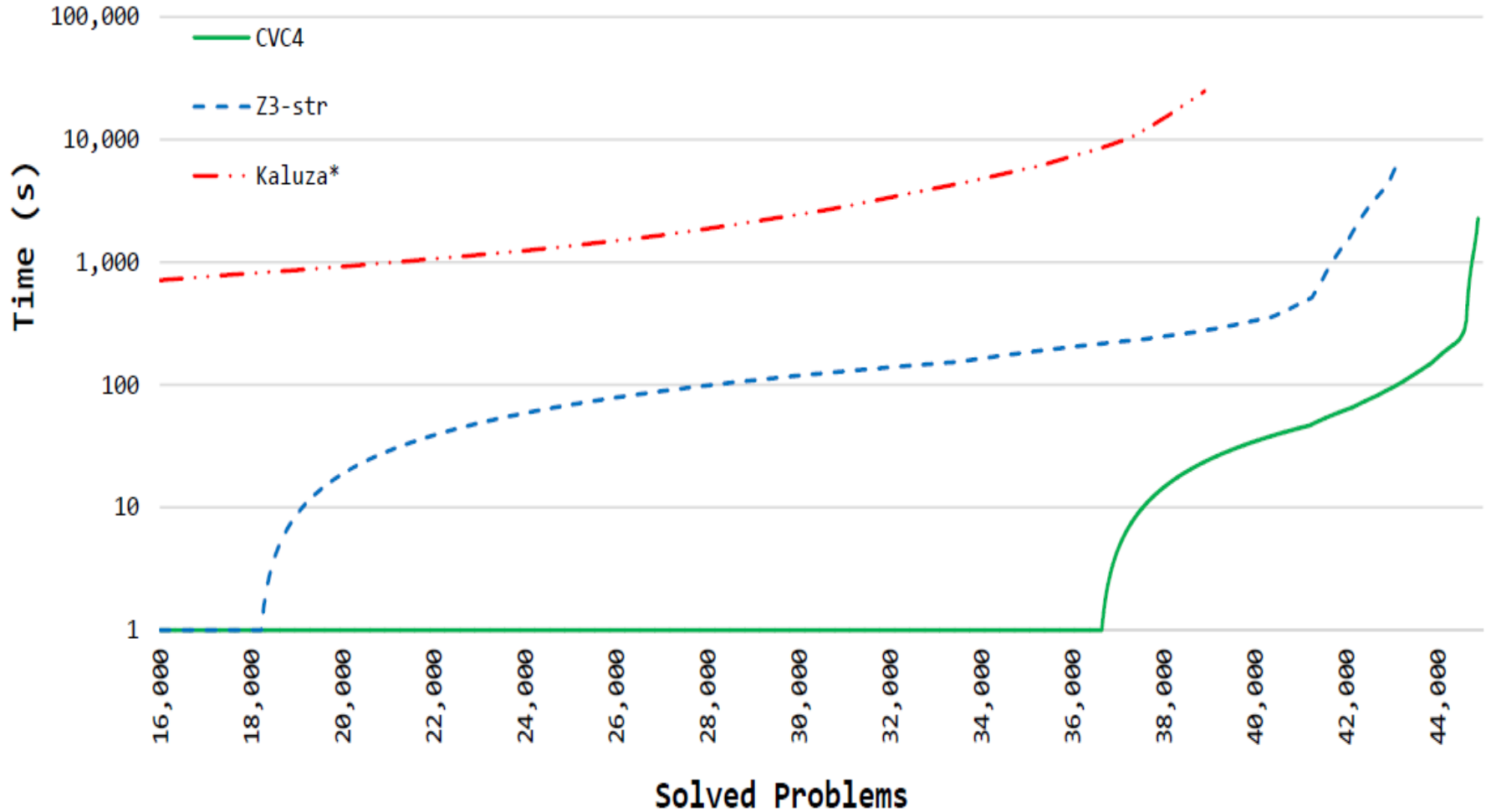
- Implemented in SMT solver CVC4
- Tested:
 - 50,000 benchmarks from Kudzu
 - Correspond to VCs in web security applications
- Compared against solvers:
 - Kaluza (UBerkeley)
 - Z3-STR (Purdue, Waterloo)

Experimental Results

	CVC4	Z3-STR		Kaluza	
Result		Incorrect ³		Incorrect ³	
unsat	11,625 ¹	317	11,769 ²	7,154	13,435 ²
sat	33,271	1,583	31,372	n/a ⁴	25,468 ⁴
unknown	0		0		3
timeout	2,388		2,123		84
error	0		120 ⁵		1,140

1. For the problems where CVC4 answers UNSAT, neither Z3-STR nor Kaluza answer SAT
2. We cannot verify the problems where CVC4 does not answer UNSAT
3. We verified these errors by asserting a model back as assertions to the tool
4. We cannot verify these answers due to bugs in Kaluza's model generation
5. One is because of non-trivial regular expression, and 119 are because of escaped characters

Experimental Results



Further Work

- Theoretical:
 - Identify fragments when approach is refutation complete
 - [Abdullah et al CAV14]
- Regular language membership $t \in R^*$
 - Currently handled, but naively (unrolling)
- More functions
 - `substr`, `contains`, `replace`, `prefixOf`,
`suffixOf`, `str.indexOf`, `str.to.int`, `int.to.str`
- Generalize to theory of sequences

Thank You!

- CVC4 is publicly available at:
<http://cvc4.cs.nyu.edu/>



Challenge: Looping Word Equations

- Say we are given:

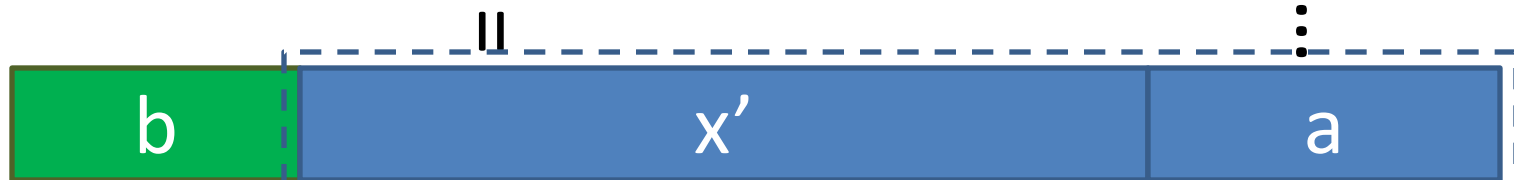
$$x \cdot a = b \cdot x$$



Challenge: Looping Word Equations

$$x \cdot a = b \cdot x$$

$$x = b \cdot x'$$



Variant of Original Equation!



Challenge: Looping Word Equations

$$x \cdot a = b \cdot x$$

- Solution:
 - Recognize when these cases occur
 - Reduce to regular language membership:

$$x \cdot a = b \cdot x \Leftrightarrow \exists yz. (a = y \cdot z \wedge b = z \cdot y \wedge x \in (z \cdot y)^* z)$$