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

Tianyi Liang
Andrew Reynolds
Cesare Tinelli
Morgan Deters
Clark Barrett

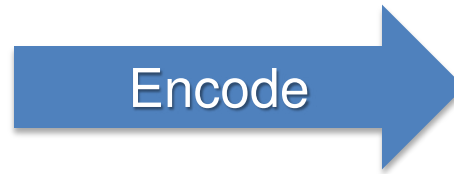# Motivation : Security Applications

```cpp
char buff[15];
char pass;

std::cout << "Enter the password :";
gets(buff);

if (std::regex_match(
        buff,
        std::regex("([A-Z]+)") )) {
    if(strcmp(buff, "PASSWORD")) {
        std::cout << "Wrong Password";
    }
    else {
        std::cout << "Correct Password";
        pass = 'Y';
    }
}

if(pass == 'Y') {
    /* Grant the root permission*/
}
```
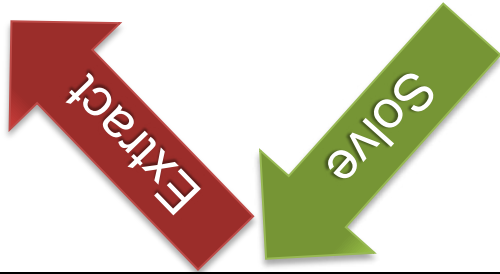
**Encode**

```
(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 (= 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"))
```

**Solve**

**Extract**

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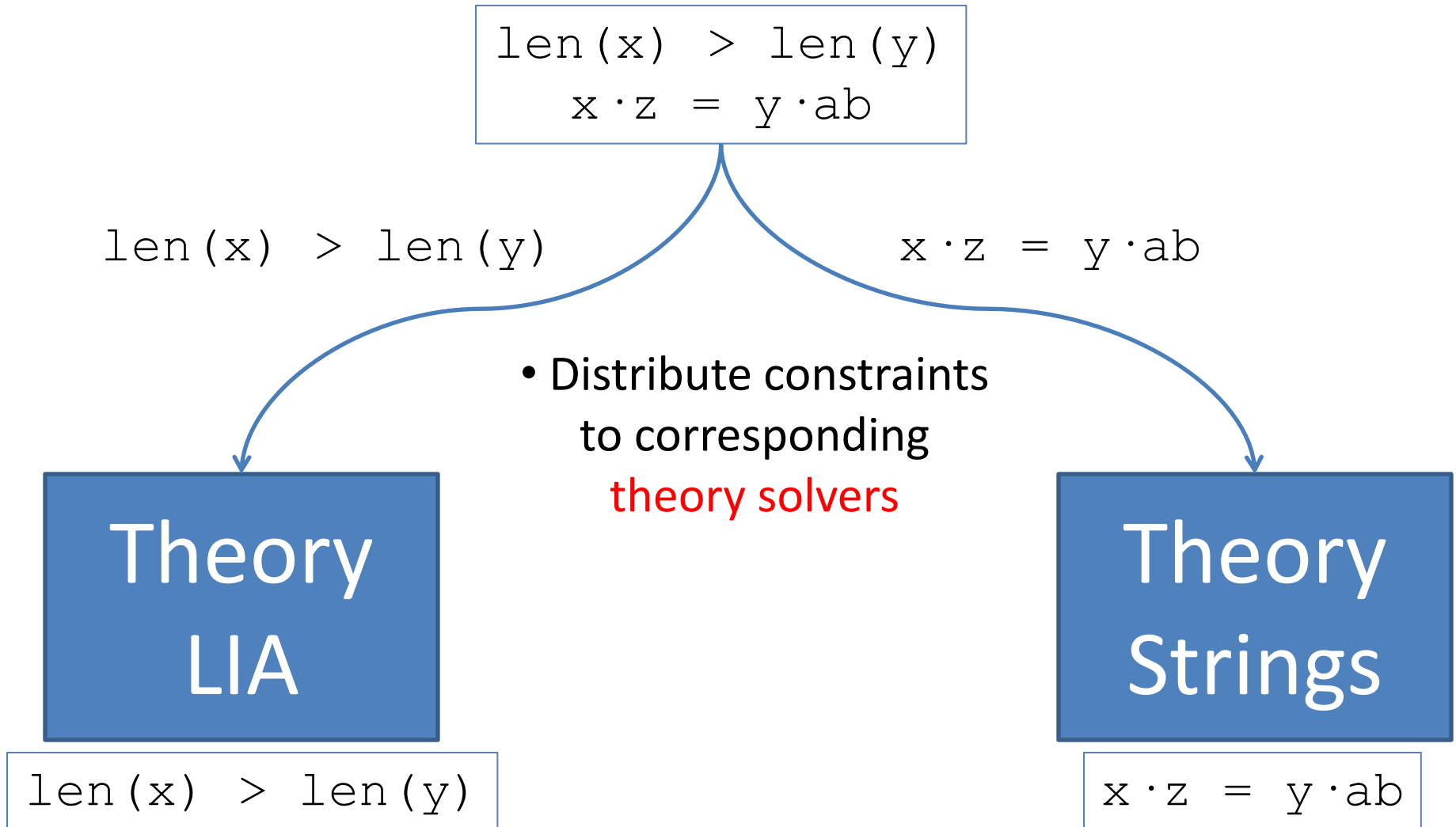
# Objectives

- Want solver to handle:
  - (Unbounded) string constraints
  - Length constraints
  - Regular language memberships, …
- Theoretical complexity of:
  - Word equation problem is in <span style="color:red">PSPACE</span>
  - …with length constraints is <span style="color:red">OPEN</span>
  - …with extended functions (e.g. `replace`) is <span style="color:red">UNDECIDABLE</span>
- Instead, focus on:
  - Solver that is efficient in practice
  - Tightly integrated into SMT solver architecture
    - Conflict analysis, T-propagation, lemma learning, …
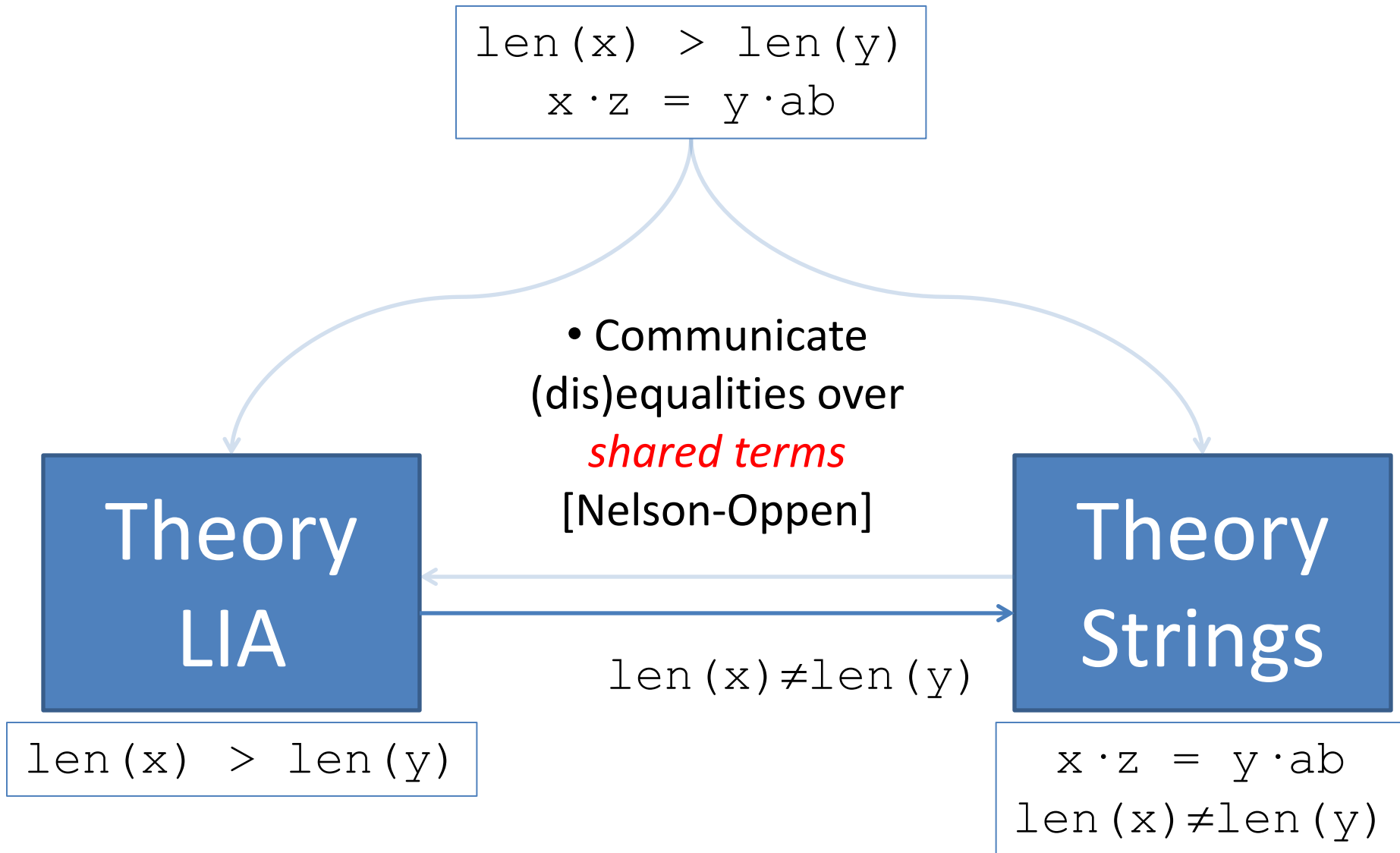
# Core Language for Theory of Strings

- Terms are:
  - Constants from a fixed finite alphabet $\Sigma^*$ (a, ab, cbc…)
  - Free constants or "variables" (x, y, z…)
  - String concatenation

    `_`$\cdot$`_` : String $\times$ String $\rightarrow$ String
  - Length terms

    `len(_)` : String $\rightarrow$ Int
- Example input:

$$\texttt{len(x) > len(y)}$$
$$\texttt{x}\cdot\texttt{z = y}\cdot\texttt{ab}$$

# Cooperating *Theory Solvers*

len(x) > len(y)
x·z = y·ab

len(x) > len(y)                                    x·z = y·ab

- Distribute constraints
  to corresponding
  theory solvers

**Theory LIA**

**Theory Strings**

len(x) > len(y)

x·z = y·ab

# Cooperating *Theory Solvers*

```
len(x) > len(y)
x·z = y·ab
```

• Communicate (dis)equalities over *shared terms* [Nelson-Oppen]

**Theory LIA**

**Theory Strings**

$len(x) \neq len(y)$

```
len(x) > len(y)
```

```
x·z = y·ab
len(x)≠len(y)
```

# Summary of Approach

- Determines satisfiability of $A \cup S$, where
  - $A$ is a set of linear <span style="color:red">arithmetic constraints</span>
  - $S$ is a set of <span style="color:red">(dis)equalities</span> over:
    - String terms
    - Length terms

$$x \cdot z = y \cdot ab$$
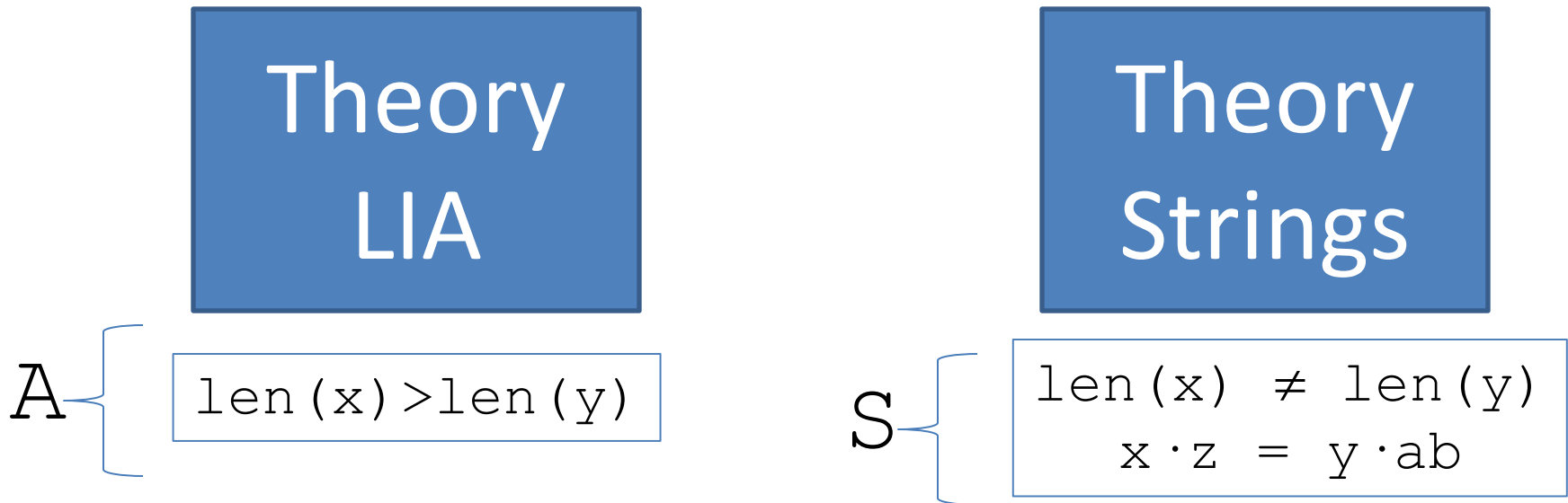$$len(x) \neq len(y)$$

- Uses procedure consisting of <span style="color:red">four steps</span>:

  1. Check length constraints $A$
  2. Normalize equalities in $S$
  3. Normalize disequalities in $S$
  4. Check cardinality of $\Sigma$

# Check Length Constraints

1. **Check length constraints**
2. **Normalize equalities**
3. **Normalize disequalities**
4. **Check cardinality of $\Sigma$**

- Add equalities to $\mathbb{A}$ regarding the length of (non-variable) terms from S

Theory LIA

Theory Strings

$\mathbb{A}$ {
```
len(x)>len(y)
```

$S$ {
```
len(x) ≠ len(y)
x·z = y·ab
```

# Check Length Constraints

Theory LIA

Theory Strings

$A$ 
```
len(x)>len(y)
len(x)+len(z)=len(y)+2
```

$S$ 
```
len(x) ≠ len(y)
x·z = y·ab
```

$\Rightarrow$ Check if $A$ is satisfiable

# Normalize Equalities

- To show: satisfiability of (dis)equalities $S$ between string terms

**Theory Strings**

$$len(x) \neq len(y)$$
$$\mathbf{x \cdot z = y \cdot ab}$$

- To ensure equality $t=s$ has model:
  - If $t$ and $s$ are non-variable,
    - Must be equivalent to flat forms $F[t], F[s]$
      - $F[t]$ and $F[s]$ are syntactically equivalent
- Flat form $F[t]$ computed by expanding/flattening $t$

# Normalize Equalities

- Modified example:

$$\texttt{len(x) = len(y)}$$
$$\texttt{z·w = y·ab}$$
$$\texttt{z = x·a}$$

- Flat form of terms from first equality are not the same:
  - $\texttt{F[z·w]}$ is: $\texttt{x·a·w}$
  - $\texttt{F[y·ab]}$ is: $\texttt{y·ab}$
- Procedure continues based on three cases:
  - We know the length of x and y are equal : *conclude $x=y$*
  - We know the length of x and y are disequal : conclude
    $$\texttt{∃k.(( x=y·k ∨ y=x·k ) ∧ len(k)>0 )}$$
  - We know neither : guess their lengths are equal, restart

# Normalize Equalities

- After concluding `x=y`,

$$\texttt{len(x) = len(y)}$$
$$\texttt{z·w = y·ab}$$
$$\texttt{z = x·a}$$
$$\texttt{x = y}$$

- Flat form of terms from first equality are now, e.g.:
  - `F[z·w]` is: `y·a·w`
  - `F[y·ab]` is: `y·ab`
- Will conclude `w=b`, after which `F[z·w]=F[y·ab]`

# Normalize Equalities

- For $t=s$, procedure makes progress* towards:
  - Towards forcing flat forms $F[t]$ and $F[s]$ equal, or
  - Discovering conflicts

- If $F[t_1]=\ldots=F[t_n]$ for an eq class $E=\{t_1 \ldots t_n\}$:
  - We refer to $F[t_1]$ as the normal form $N[t_1]$ of $E$

- If normal form exists for each eq class,
  - Then a model exists for all equalities from $S$
    - Constructed trivially, given normal form

* exception: looping word equations (explained later)

# Normalize Disequalities

- For disequalities in `S`

  - A disequality `t≠s` is normalized if:

    - `len(t)≠len(s)`, or

    - `N[t]=`$t_1 \cdot u \cdot t_2$ and `N[s]=`$s_1 \cdot v \cdot s_2$, where:
      - `len(`$t_1$`)=len(`$t_2$`)`,
      - `len(u)=len(v)`, and
      - `u≠v`

- For example:

  | |
  |---|
  | `len(z)≠len(y)` ✔ |
  | `z≠y` ✔ |
  | `x·a·z≠x·b·z` ✔ |
  | `x·w≠y·b` ✖ |

# Normalize Disequalities

- To normalize disequalities,
  - Proceed by cases, similar to Step 2
    - In example, we would succeed, for example if:
      - `len(x·w)≠len(y·b),` or
      - `len(x)=len(y)` and x≠y,
      - …
  - Continue until all disequalities are normalized

```
len(z)≠len(y)      ✔
      z≠y          ✔
x·a·z≠x·b·z        ✔
      x·w≠y·b      ✘
```

# Check Cardinality of $\Sigma$

- $S$ may be unsatisfiable since $\Sigma$ is finite

- For instance,

  *If*

    - $\Sigma$ 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

# Challenge: Looping Word Equations

- Say we are given: $x \cdot a = b \cdot x$

# Challenge: Looping Word Equations

- Say we are given: $\boxed{\texttt{x·a = b·x}}$

- Flat forms are:

    $\texttt{F[x·a] = }{\color{red}\texttt{x}}\texttt{·a}$

    $\texttt{F[b·x] = }{\color{red}\texttt{b}}\texttt{·x}$

- Compare $\texttt{len(x)}$ and $\texttt{len(b)}$, i.e. $\texttt{1}$

    - If $\texttt{len(x)=1}$, then $\texttt{x=a}$ and $\texttt{x=b} \Rightarrow$ <span style="color:red">conflict</span>

    - If $\texttt{len(x)}\neq\texttt{1}$

        - If $x$ is a prefix of b (i.e. it is empty), then $\texttt{a=b} \Rightarrow$ <span style="color:red">conflict</span>

        - If b is a prefix of x, then $\texttt{x=b·k}$ for some $\texttt{k}$

# Challenge: Looping Word Equations

- Now we have:

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

- Flat forms of first equation are:

$$F[x \cdot a] = b \cdot k \cdot a$$
$$F[b \cdot x] = b \cdot b \cdot k \implies \text{Problem: looping!}$$

- Solution:

  - Recognize when these cases occur

  - Reduce to regular language membership:
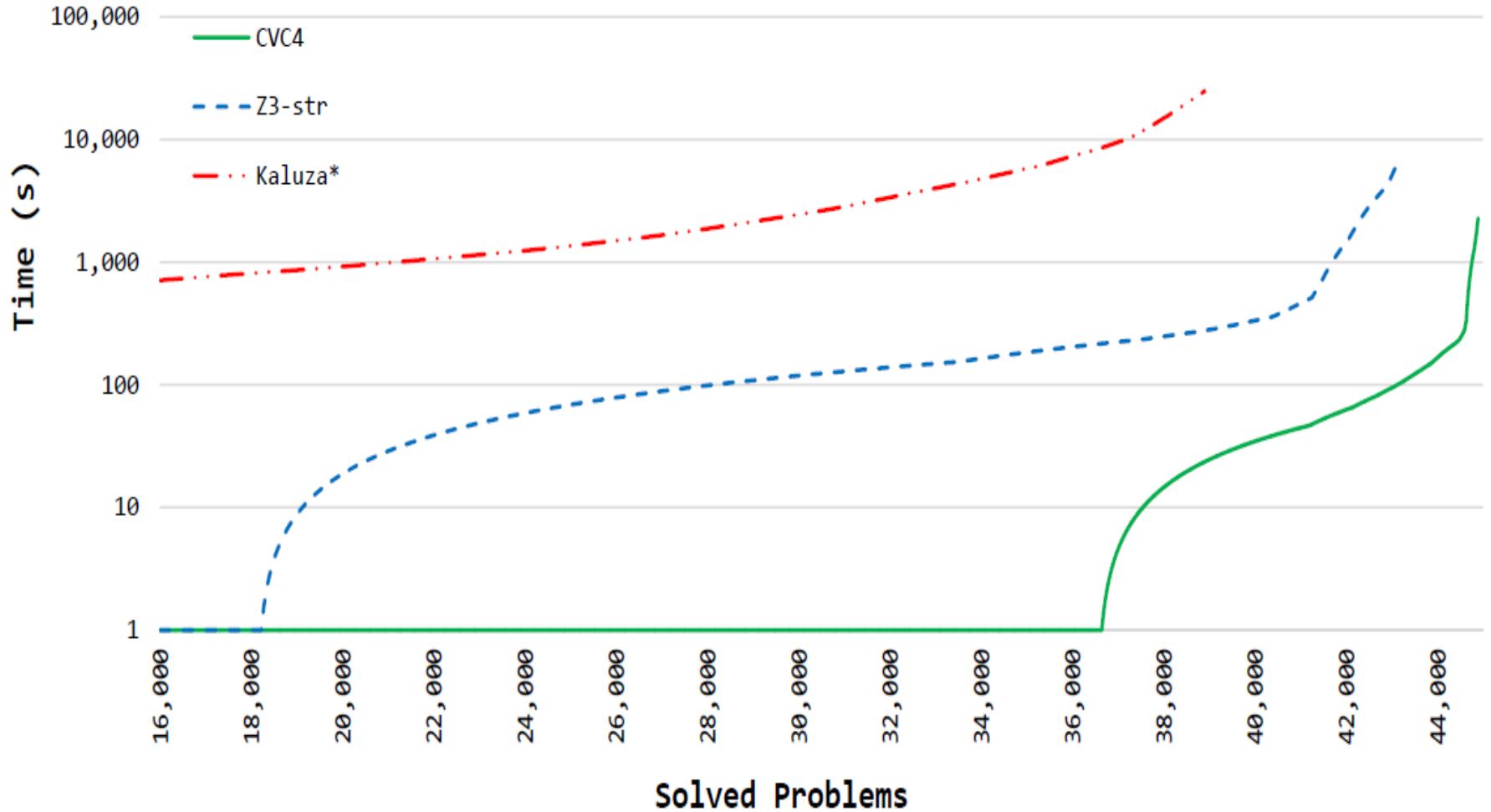
$$x \cdot a = b \cdot x \Leftrightarrow \exists yz.(a = y \cdot z \;\land\; b = z \cdot y \;\land\; x \in (z \cdot y) * z)$$

# 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

# Theoretical Results

- Our approach is:
  - <span style="color:green">Refutation sound</span>
    - When it answer "UNSAT", it can be trusted
      - Even for strings of unbounded length
  - <span style="color:green">Solution sound</span>
    - When it answers "SAT", it can be trusted
- (A version of) our approach is:
  - <span style="color:green">Solution complete</span>
    - When it is "SAT", it will eventually get a model
      - Somewhat trivially, by finite model finding
- Our approach is <span style="color:red">not</span>:
  - <span style="color:red">Refutation complete</span>
    - When it is "UNSAT", it is <span style="color:red">not</span> guaranteed to derive refutation
      - Would like to identify fragments (i.e. non-cyclical) where it is

# Further Work

- Handling regular language membership $t \in R*$
  - Currently handled, but naively (unrolling)
- Handling extended functions
  - `substr, contains, replace, prefixOf, suffixOf, str.indexOf, str.to.int, int.to.str`
  - Many are challenging, for instance:
    $$\neg\texttt{contains(x,y)}$$
    - Intuitively, requires (universal) quantification over the positions of `x`

# Questions?

- For more details, see CAV 2014 paper
- CVC4 is publicly available at:
  `http://cvc4.cs.nyu.edu/web/`