# SMT String Solving in CVC4

Andrew Reynolds & Cesare Tinelli

The University of Iowa

MOSCA 2019

May 6, 2019

# Satisfiability Modulo Theories (SMT) Solvers

Many applications:

- Software verification
- Automated theorem proving
- Symbolic execution
- Security analysis

In this talk:

- How SMT Solvers (CVC4) handle string constraints

# The CVC4 SMT Solver

Support for many theories and features

- UF, (non)linear arithmetic, arrays
- Bit-vectors, floating point
- Finite sets and relations, (co)datatypes

$\Rightarrow$ Strings and regular expressions

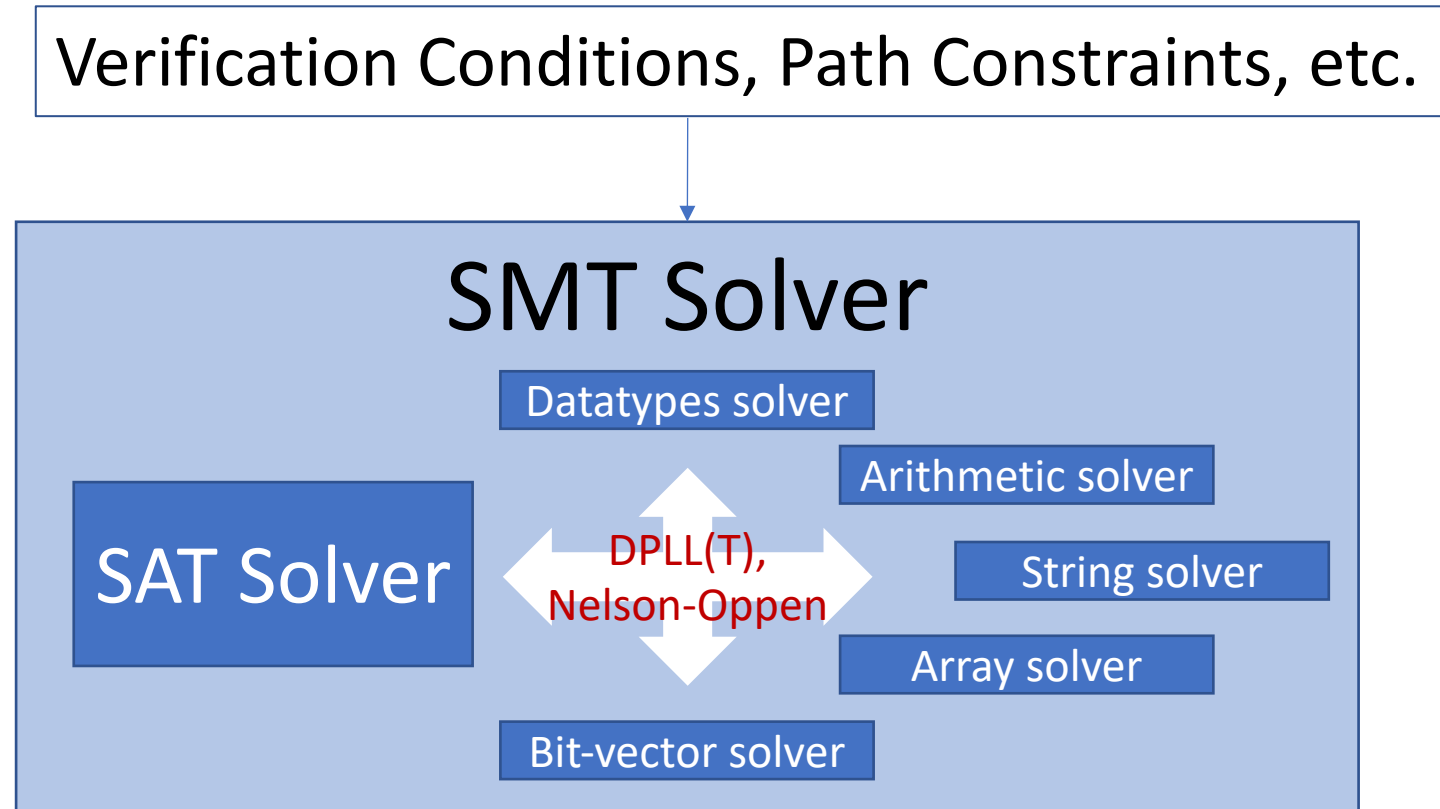Co-developed at Stanford and University of Iowa

- Project Leaders:

    Clark Barrett and Cesare Tinelli

- String solver developers:

    Andrew Reynolds, Tianyi Liang, Nestan Tsiskaridze, Andres Noetzli

# Overview

- How SMT string solvers work:
  - Basic architecture (DPLL(T))
  - Core Theory Solver for Word Equations with Length Constraints
  - Advanced Features
    - Finite model finding
    - Context-dependent simplification for extended string constraints
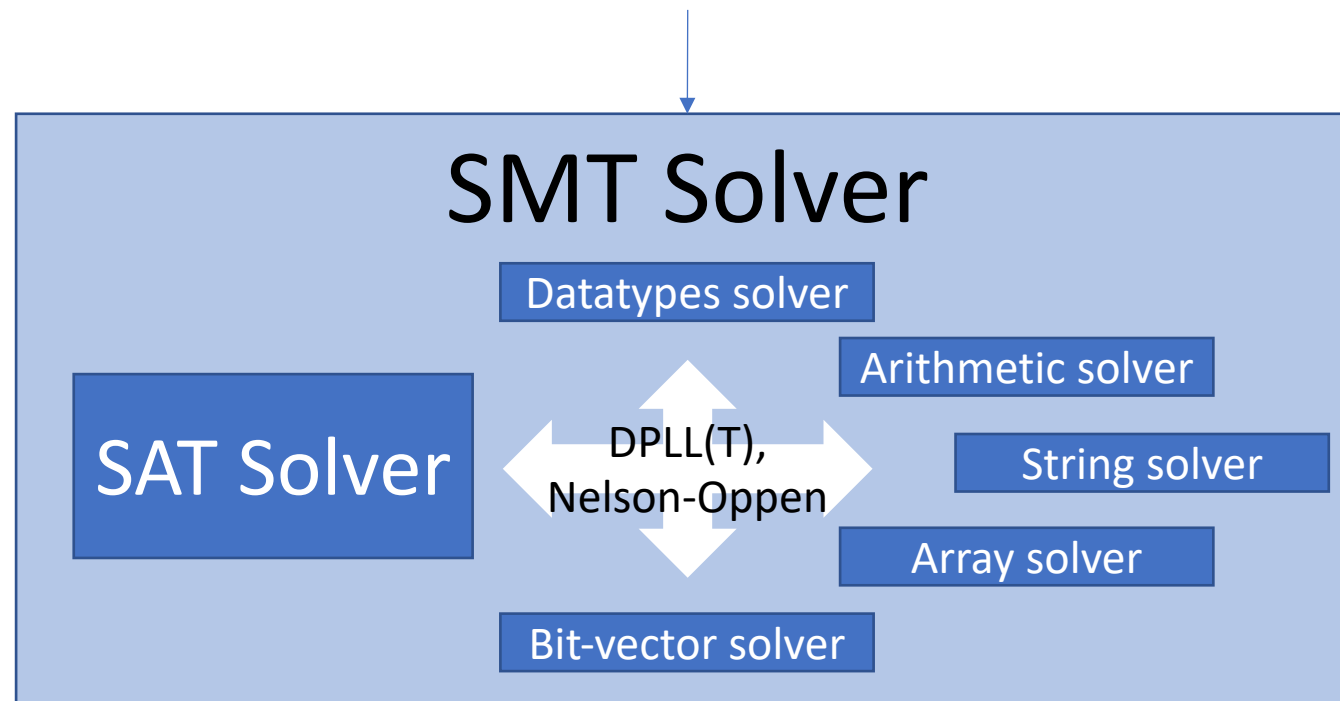    - Regular expression elimination

# SMT solvers

Efficient tools for satisfiability *modulo theories*

# SMT solvers

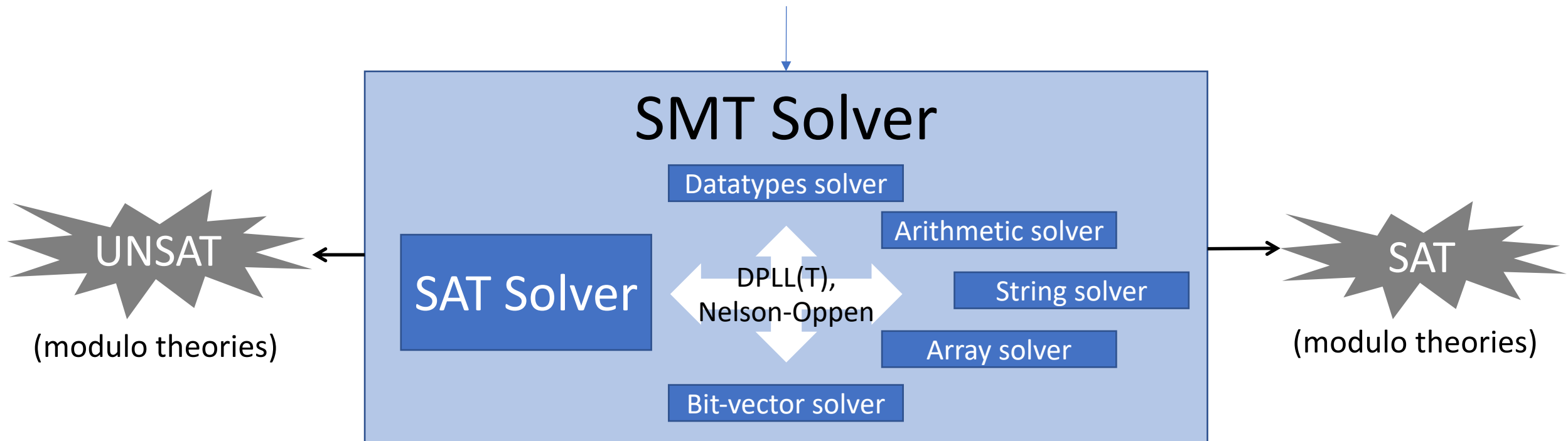Efficient tools for satisfiability *modulo theories*

$$( A[x] + B[x] > 0 \vee x + y > 0 ) \wedge ( \text{cons}(\text{``abc''}, d_1) \neq d_2 \vee x < 0)$$

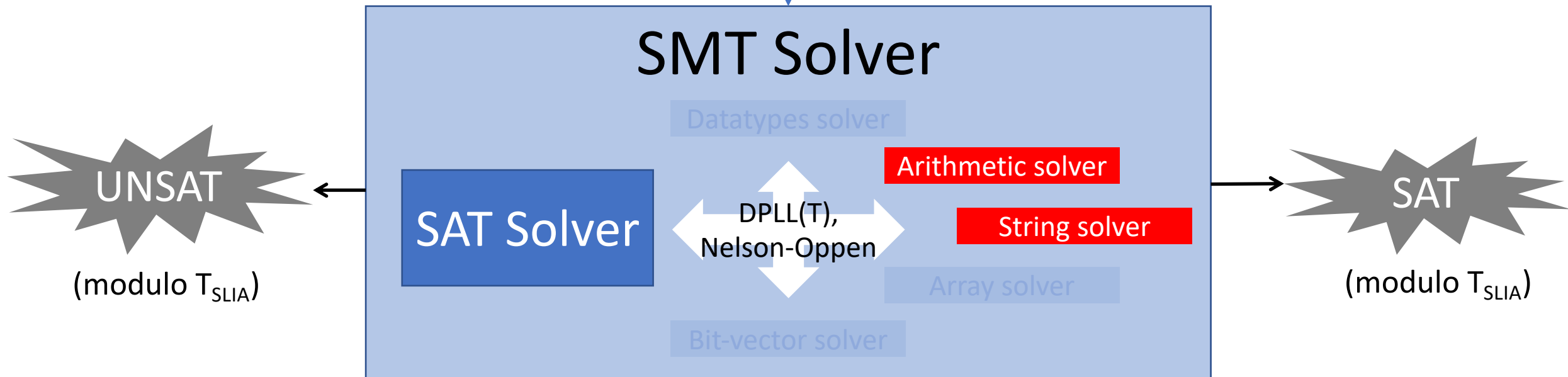# SMT solvers

Efficient tools for satisfiability *modulo theories*

$$( A[x] + B[x] > 0 \lor x + y > 0 ) \land ( \text{cons("abc", } d_1) \neq d_2 \lor x < 0)$$

# SMT solvers

Our focus: the theory of strings and linear arithmetic $T_{SLIA}$

$$x = \text{``ab''} \cdot z \wedge |x| + |y| \leq 5 \wedge (\text{``abcd''} \cdot x = y \vee |x| > 5)$$

# Theory of Strings + Linear Arithmetic ($T_{SLIA}$)

**Sorts:**

- Integers `Int`
- Strings `String`, interpreted as $A*$ for finite alphabet $A$

**Terms:**

- String Variables: $x$, $y$, $z$
- Integer Variables: $i$, $j$, $k$
- String Constants: "", "abc", "AAAAA", "http"
- String Concatenation: $x \cdot$"abc", $x \cdot y \cdot z \cdot w$
- String Length: $|x|$

**Formulas are:**

- Equalities and disequalities between string terms
- *Linear* arithmetic constraints: $|x| + 4 > |y|$

**Example:**   $x \cdot$"a" $= y$, $y \neq$ "b"$\cdot z$, $|y| > |x| + 2$

**Decidability:** unknown, regardless, many problems can be solved efficiently in practice

# T$_{SLIA}$ String Solver for DPLL(T)

Achieved as a Cooperation between:

SAT Solver

Arithmetic Solver

String Solver

# $T_{SLIA}$ String Solver for DPLL(T)

$$x = \text{"ab"} \cdot z$$
$$|x| + |y| \leq 5$$
$$\text{"abcd"} \cdot x = y \lor |x| > 5$$

Set of $T_{SLIA}$-formulas in clausal normal form (CNF)

SAT Solver

Arithmetic Solver

String Solver

# T$_{SLIA}$ String Solver for DPLL(T)

x = "ab"·z
$|x| + |y| \leq 5$
"abcd"·x = y $\lor$ $|x| > 5$

SAT
Solver

Arithmetic
Solver

String
Solver

UNSAT

$\Rightarrow$ Either determines no satisfying assignments for input exist

# T$_{SLIA}$ String Solver for DPLL(T)



$x = \text{"ab"} \cdot z$

$|x| + |y| \leq 5$
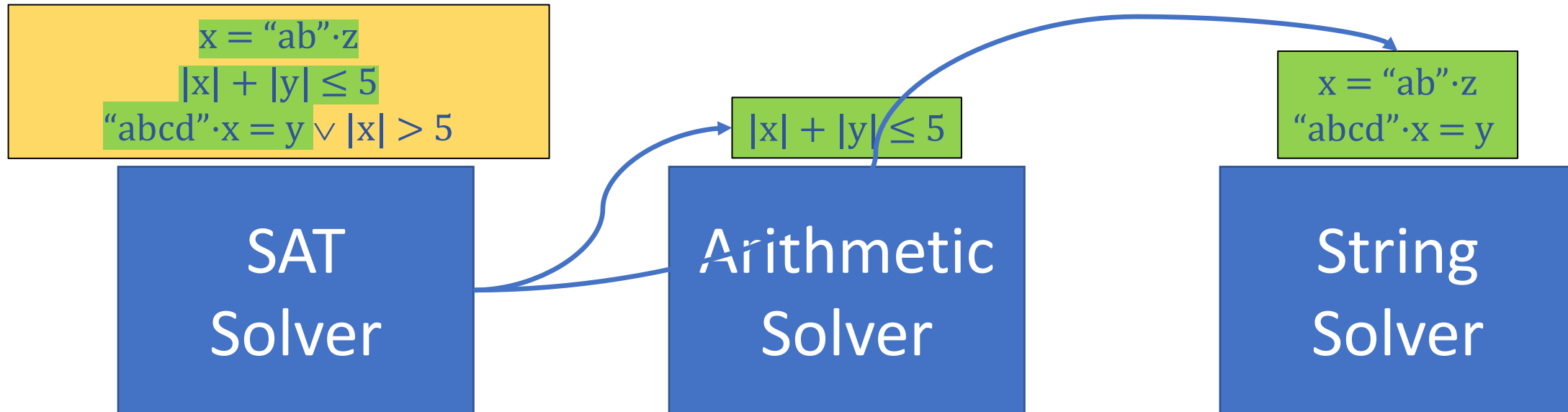
$\text{"abcd"} \cdot x = y \vee |x| > 5$
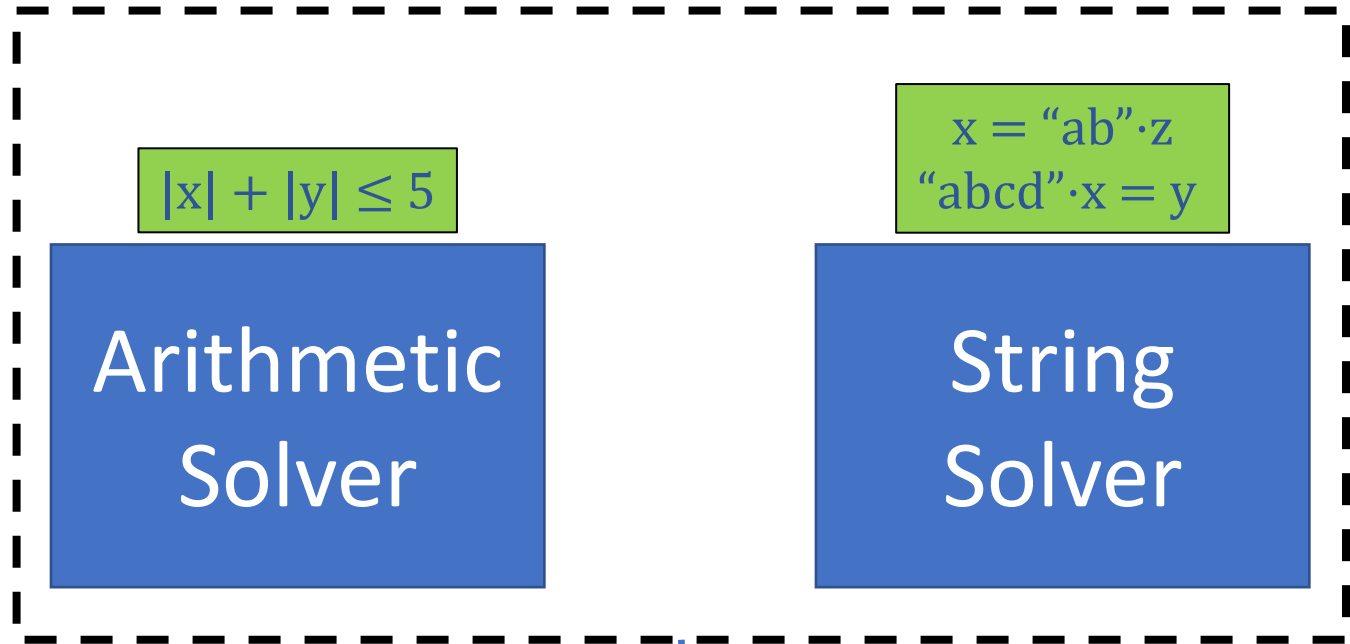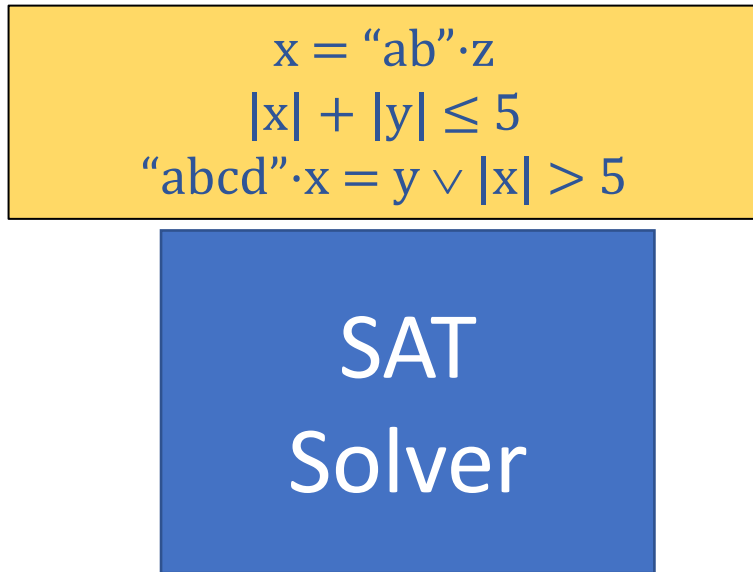
**SAT Solver**

**Arithmetic Solver**

**String Solver**

$\Rightarrow$ ... or returns a propositionally satisfying assignment

# T$_{SLIA}$ String Solver for DPLL(T)



$x = \text{"ab"} \cdot z$
$|x| + |y| \leq 5$
$\text{"abcd"} \cdot x = y \vee |x| > 5$

$|x| + |y| \leq 5$

$x = \text{"ab"} \cdot z$
$\text{"abcd"} \cdot x = y$

**SAT Solver**

**Arithmetic Solver**

**String Solver**

$\Rightarrow$ Constraints distributed to arithmetic and string solvers

# T$_{SLIA}$ String Solver for DPLL(T)

$x = \text{"ab"} \cdot z$
$|x| + |y| \leq 5$
$\text{"abcd"} \cdot x = y \lor |x| > 5$

**SAT Solver**

$|x| + |y| \leq 5$

**Arithmetic Solver**

$x = \text{"ab"} \cdot z$
$\text{"abcd"} \cdot x = y$

**String Solver**

SAT

$\Rightarrow$ Either find constraints are T$_{SLIA}$-satisfiable

# T$_{SLIA}$ String Solver for DPLL(T)



$x = \text{"ab"} \cdot z$
$|x| + |y| \leq 5$
$\text{"abcd"} \cdot x = y \vee |x| > 5$
$\neg(x = \text{"ab"} \cdot z) \vee |x| = |z| + 2$

$|x| + |y| \leq 5$

$x = \text{"ab"} \cdot z$
$\text{"abcd"} \cdot x = y$

**SAT Solver**

**Arithmetic Solver**

**String Solver**

$\neg(x = \text{"ab"} \cdot z) \vee |x| = |z| + 2$

$\Rightarrow$ or return *theory lemmas* (valid T$_{LIA}$/T$_S$-formulas) to SAT solver

# $T_{SLIA}$ String Solver for DPLL(T)



SAT Solver

$x = \text{``ab''} \cdot z$

$|x| + |y| \leq 5$

$\text{``abcd''} \cdot x = y \lor |x| > 5$

$\neg(x = \text{``ab''} \cdot z) \lor |x| = |z| + 2$

Arithmetic Solver

$|x| + |y| \leq 5$

$|x| = |z| + 2$

String Solver

$x = \text{``ab''} \cdot z$

$\text{``abcd''} \cdot x = y$

$\Rightarrow$ and repeat

# Inside a DPLL(T) Theory Solver

Given a set of T-literals $M_T$,
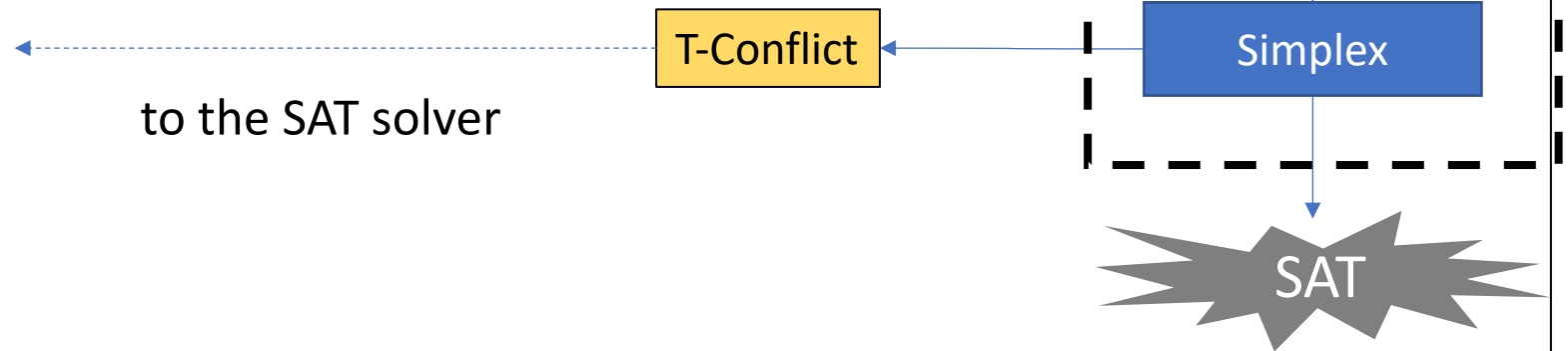
$$x = \text{"ab"} \cdot z$$
$$\text{"abcd"} \cdot x = y$$

Should the solver send a theory lemma to the SAT solver?

- no => return unknown, or

    return a *model* (a satisfying assignment)

- yes => which lemma?

    - In typical DPLL(T) theory solvers (e.g. LIA) theory lemmas $\Leftrightarrow$ *T-conflicts*
      $\neg(L_1 \wedge \ldots \wedge L_n)$ for some T-unsatisfiable $\{L_1, \ldots, L_n\} \subseteq M_T$
    - In string solver, theory lemmas may introduce new literals
    - Will describe a *strategy* for strings

# Arithmetic Theory Solver

$$M_{LIA}$$

Decision procedure:
    T-conflicts based on a standard procedure, e.g. Simplex

$$2*|x| + |y| \leq 5$$
$$|x| - |z| > 2$$

T-Conflict

Simplex

to the SAT solver

SAT

Properties:
- Sound, lemmas it generates are LIA-valid
- Model-sound, "SAT" can be trusted
- Terminating, in the context of DPLL(T)
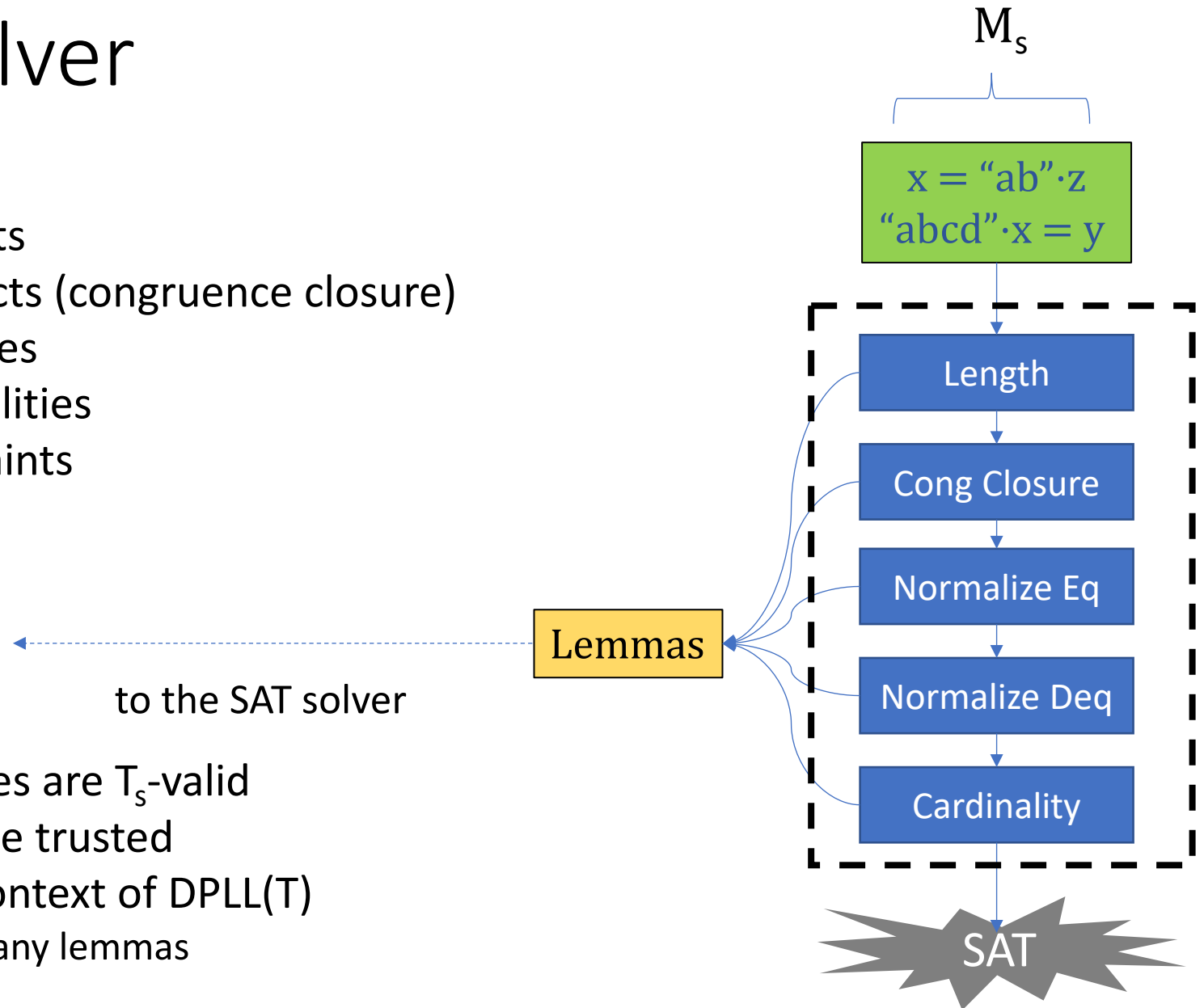      Only generates finitely many lemmas
- ∴ Complete

# String Theory Solver

Inference strategy:

1. Process length constraints
2. Check for equality conflicts (congruence closure)
3. Normalize string equalities
4. Normalize string disequalities
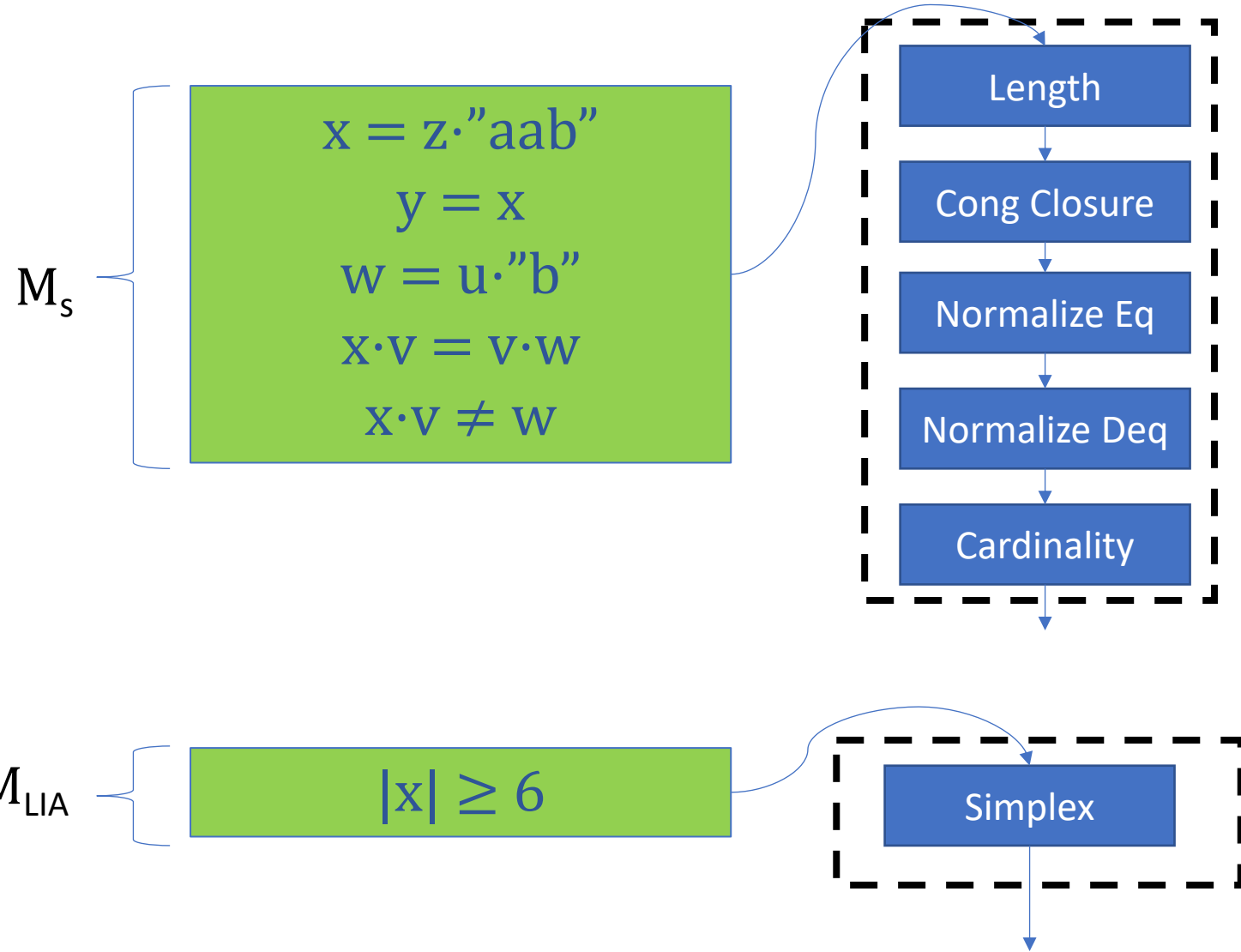5. Check cardinality constraints

$M_s$

$x = \text{"ab"} \cdot z$
$\text{"abcd"} \cdot x = y$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

Lemmas

to the SAT solver

SAT

Properties:

- Sound, lemmas it generates are $T_s$-valid
- Model-sound, "SAT" can be trusted
- Non-terminating, in the context of DPLL(T)
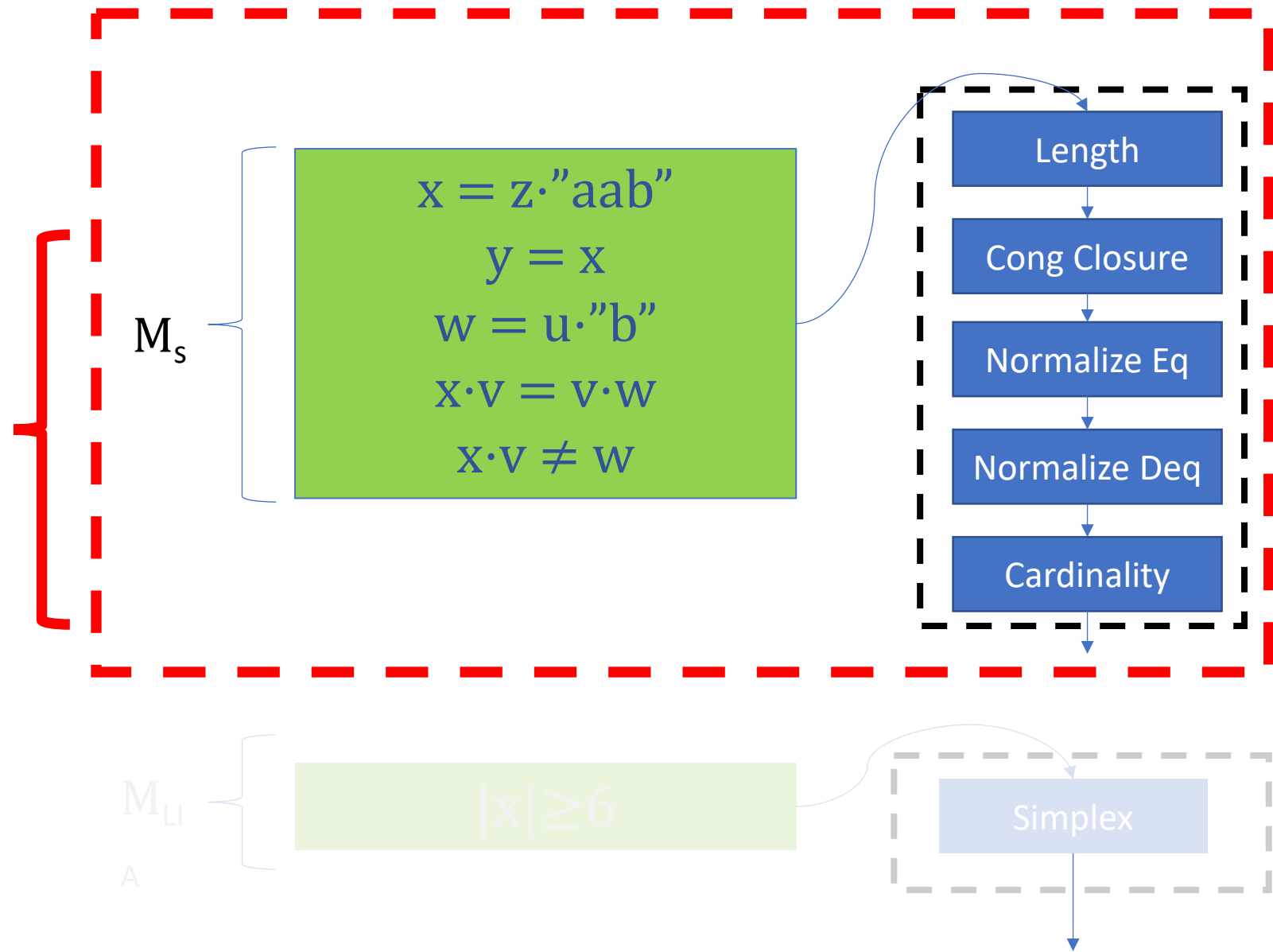  - May generate infinitely many lemmas

# String Solver

Running example:

$M_s$

$$x = z\cdot"aab"$$
$$y = x$$
$$w = u\cdot"b"$$
$$x\cdot v = v\cdot w$$
$$x\cdot v \neq w$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

$M_{LIA}$

$$|x| \geq 6$$

Simplex

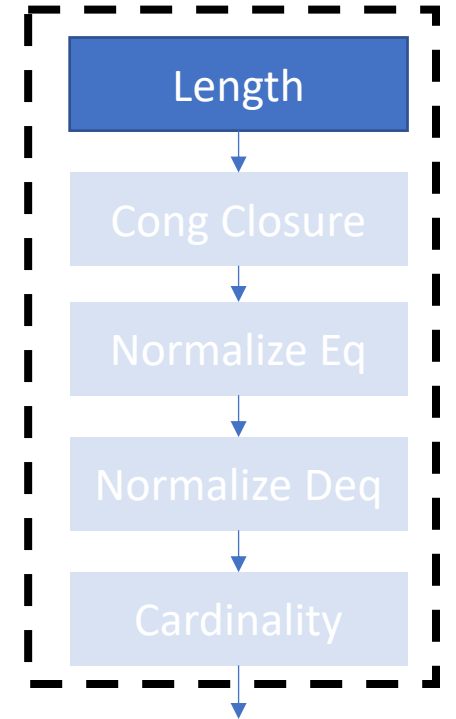# String Solver

Running example:

Will focus on string solver

[Liang et al. CAV2014]

# String Solver: Process Length

$$M_s \begin{cases} x = z \cdot \text{"aab"} \\ y = x \\ w = u \cdot \text{"b"} \\ x \cdot v = v \cdot w \\ x \cdot v \neq w \end{cases}$$

Length
↓
Cong Closure
↓
Normalize Eq
↓
Normalize Deq
↓
Cardinality

# String Solver: Process Length

$M_s$ 
$$x = z \cdot ``aab"$$
$$y = x$$
$$w = u \cdot ``b"$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

| Length |
| Cong Closure |
| Normalize Eq |
| Normalize Deq |
| Cardinality |

- For each term of type string in $M_s$:

  returns a lemma giving the definition of its length:

  | | | |
  |---|---|---|
  | $\|``b"\| = 1$ | $\|``aab"\| = 3$ | $\|x \cdot v\| = \|x\| + \|v\|$ |
  | $\|z \cdot ``aab"\| = \|z\| + 3$ | $\|u \cdot ``b"\| = \|u\| + 1$ | $\|v \cdot w\| = \|v\| + \|w\|$ |

- For each variable of type string in $M_s$:
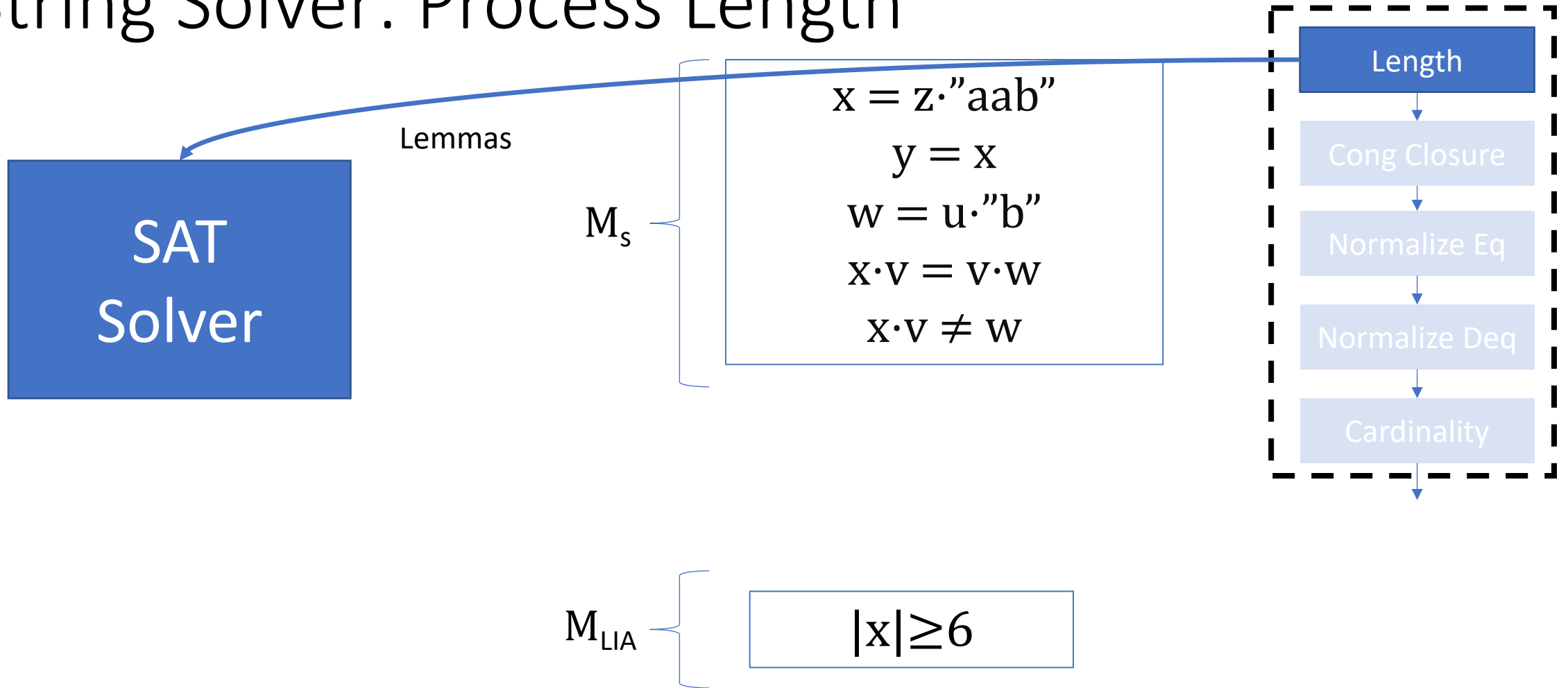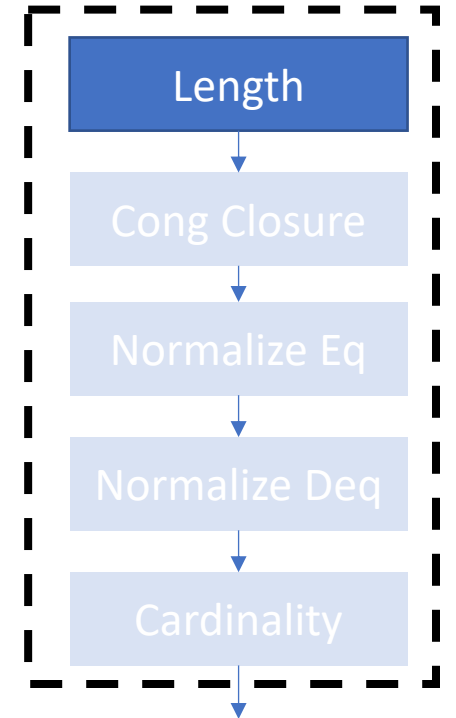
  returns an emptiness splitting lemma:

  $$x = ``" \vee \|x\| \geq 1 \qquad y = ``" \vee \|y\| \geq 1 \qquad \ldots$$

# String Solver: Process Length



SAT
Solver

Lemmas

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

$M_s$
$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

$M_{LIA}$
$$|x| \geq 6$$

# String Solver: Process Length

**SAT Solver**

$M_s$
$$x = z \cdot "aab"$$
$$y = x$$
$$w = u \cdot "b"$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

new propositional assignment

adds **new constraints** in arithmetic solver

$M_{LIA}$
$$|x| \geq 6$$
$$|"b"| = 1$$
$$|"aab"| = 3$$
$$|x \cdot v| = |x| + |v|$$
$$|z \cdot "aab"| = |z| + 3$$
$$|u \cdot "b"| = |u| + 3$$
$$|v \cdot w| = |v| + |w|$$
$$|x| \geq 1$$
$$...$$

UNSAT?

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

# String Solver: Congruence Closure

$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

$M_s$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

# String Solver: Congruence Closure

$$M_s \begin{cases} x = z \cdot \text{"aab"} \\ y = x \\ w = u \cdot \text{"b"} \\ x \cdot v = v \cdot w \\ x \cdot v \neq w \end{cases}$$
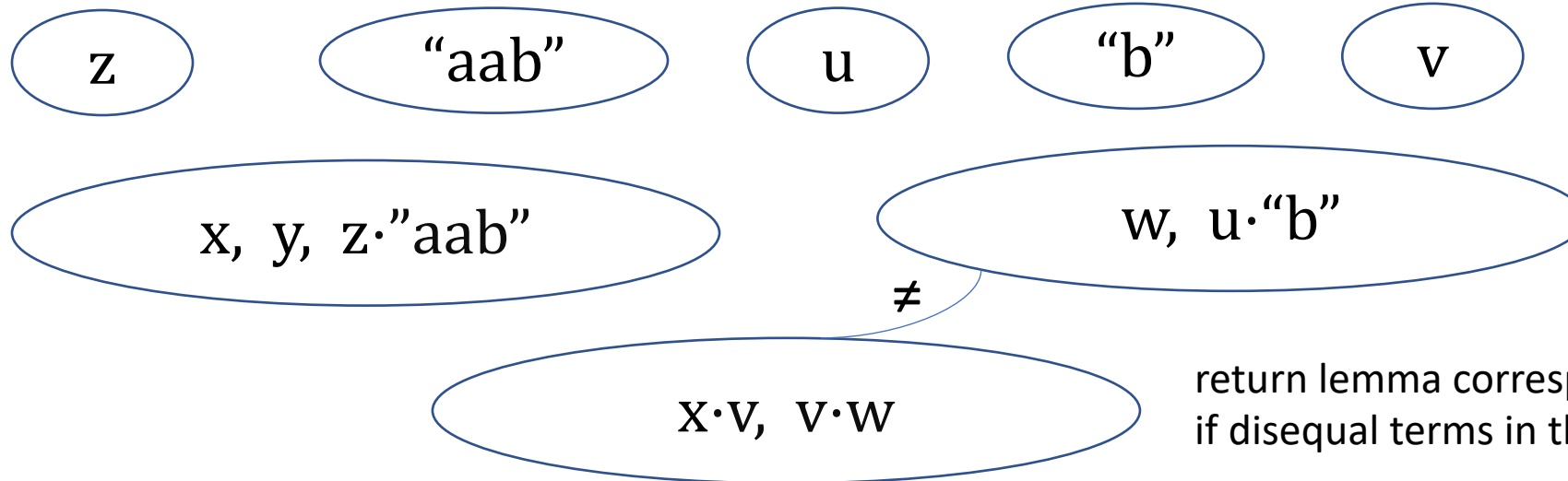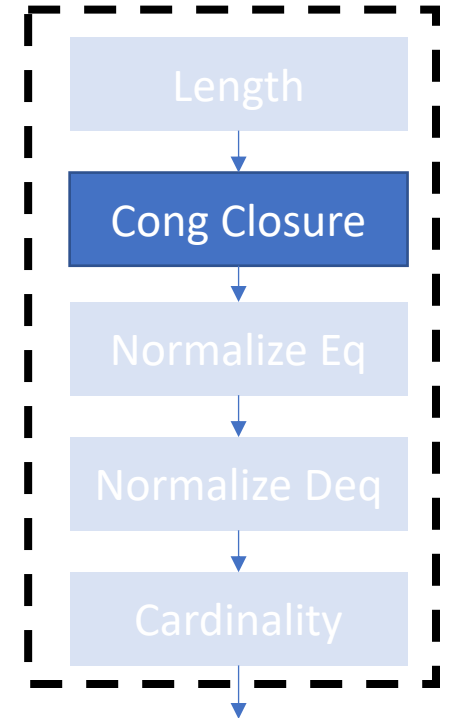
- Group terms by *equivalence classes*:

z

"aab"

u

"b"

v

x,  y,  z·"aab"

w,  u·"b"

$\neq$

x·v,  v·w

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

# String Solver: Congruence Closure

$$M_s \begin{cases} x = z \cdot \text{"aab"} \\ y = x \\ w = u \cdot \text{"b"} \\ x \cdot v = v \cdot w \\ x \cdot v \neq w \end{cases}$$
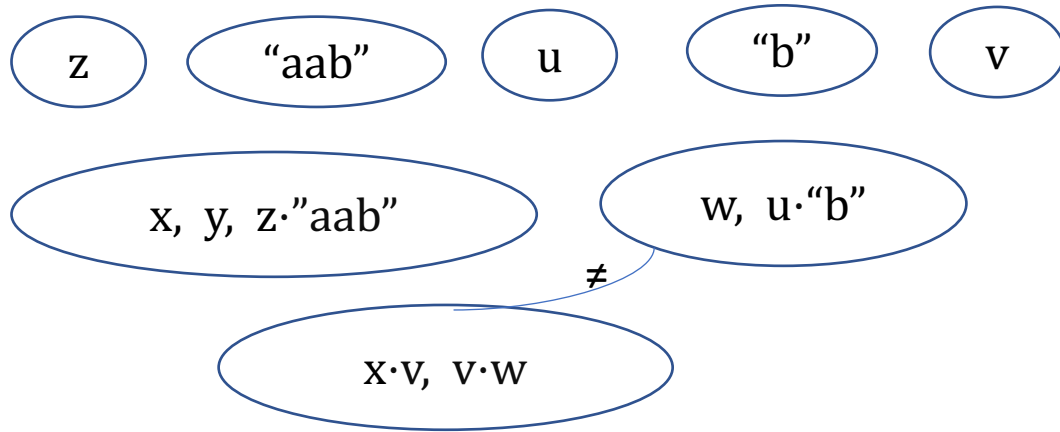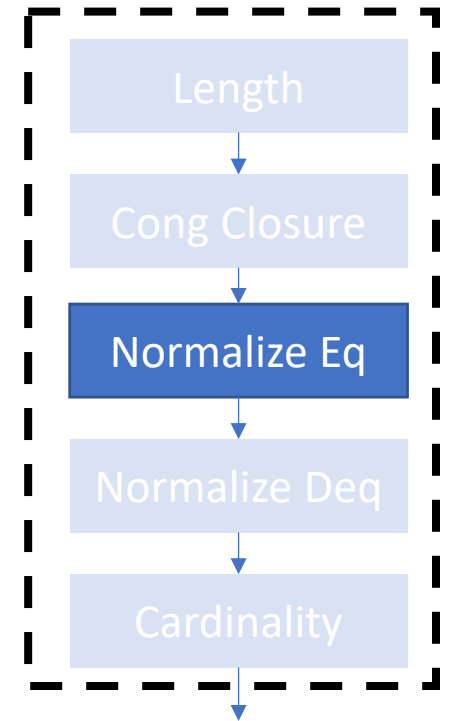
Length

**Cong Closure**

Normalize Eq

Normalize Deq

Cardinality

- Group terms by *equivalence classes*:

$z$   "aab"   $u$   "b"   $v$

$x, \ y, \ z \cdot \text{"aab"}$      $w, \ u \cdot \text{"b"}$

$\neq$

$x \cdot v, \ v \cdot w$

return lemma corresponding to $T_s$-conflict
if disequal terms in the same equivalence class

# String Solver: Normalize Equality

$z$  "aab"  $u$  "b"  $v$

$x, y, z\cdot$"aab"   $w, u\cdot$"b"

$\neq$

$x\cdot v, v\cdot w$

$$x = z\cdot\text{"aab"}$$
$$y = x$$
$$w = u\cdot\text{"b"}$$
$$x\cdot v = v\cdot w$$
$$x\cdot v \neq w$$

Length

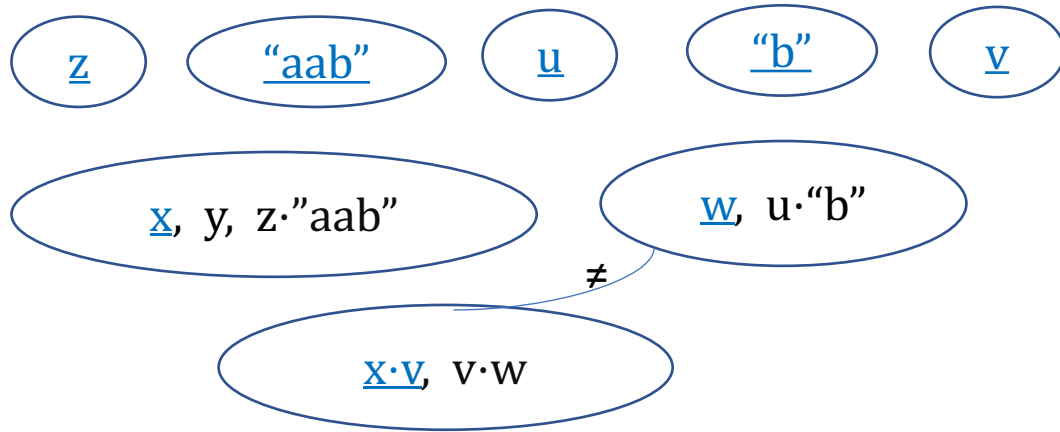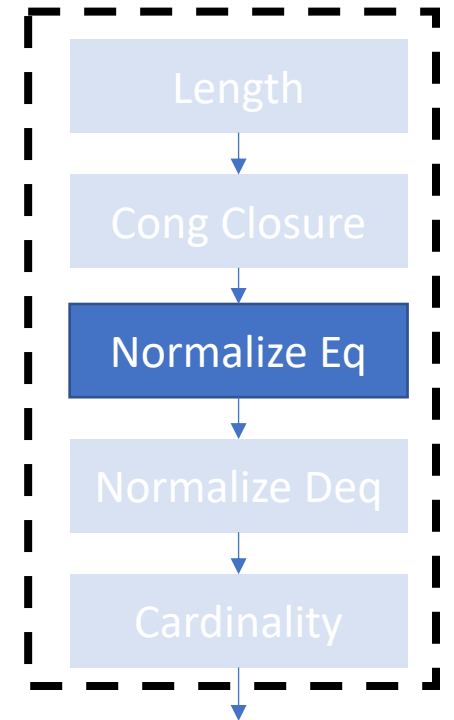Cong Closure

Normalize Eq

Normalize Deq

Cardinality

# String Solver: Normalize Equality



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

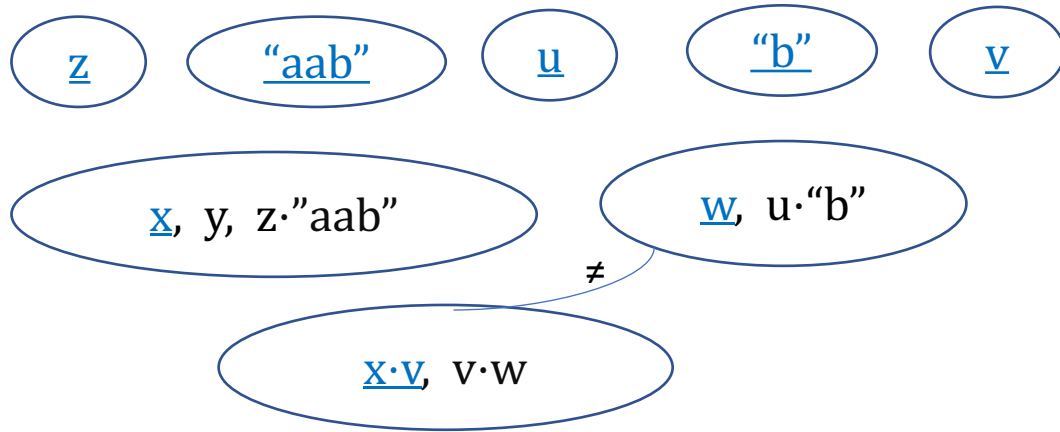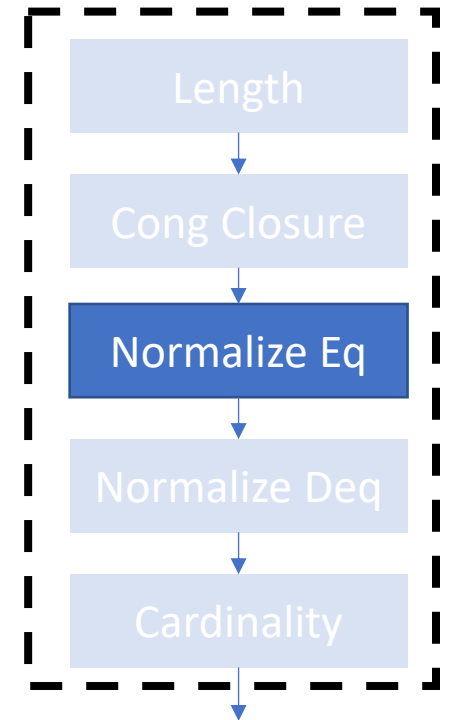Length → Cong Closure → **Normalize Eq** → Normalize Deq → Cardinality

- Compute *normal forms* for equivalence classes
  - A normal form is a concatenation of string terms $r_1 \cdot ... \cdot r_n$
    where each $ri_i$ is the representative of its equivalence class
    **Restriction:** string constants must be chosen as representatives
  - An equivalence class can be assigned a normal form $r_1 \cdot ... \cdot r_n$ if:
    Each non-variable term in it can be expanded (modulo equality and rewriting) to $r_1 \cdot ... \cdot r_n$

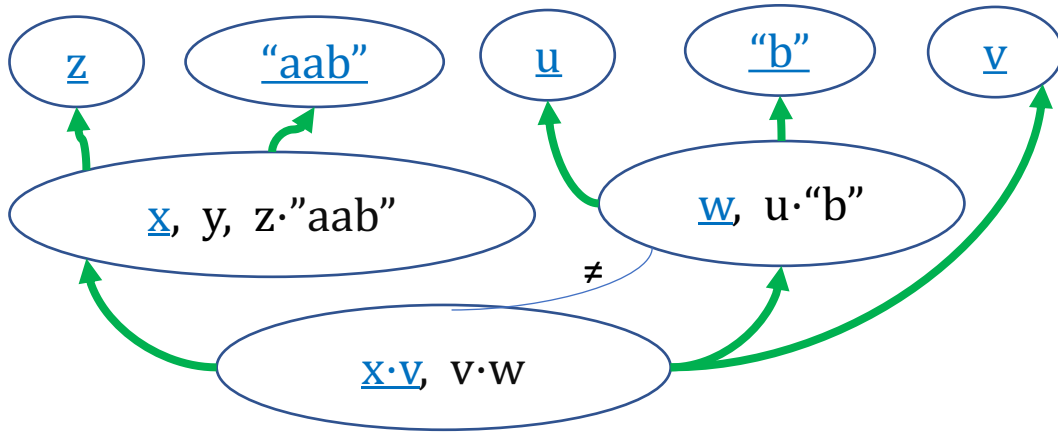# String Solver: Normalize Equality



$$x = z\cdot\text{"aab"}$$
$$y = x$$
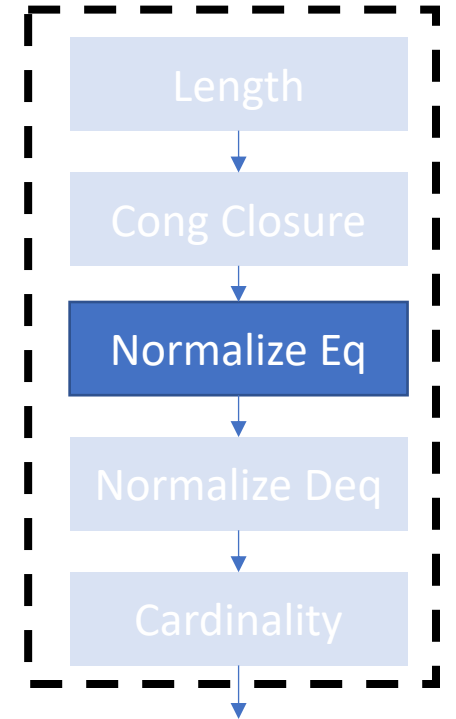$$w = u\cdot\text{"b"}$$
$$x\cdot v = v\cdot w$$
$$x\cdot v \neq w$$

Normal forms computed by a bottom-up procedure

# String Solver: Normalize Equality



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

Length → Cong Closure → **Normalize Eq** → Normalize Deq → Cardinality

## Normal forms computed by a bottom-up procedure

- First, compute containment relation induced by concatenation terms
  - To compute a n.f. for eq-class of $x \cdot v$, we must first compute a n.f. for eq-class of $x$ and $v$
  - This relation is guaranteed to be acyclic due to length elaboration step (cycle $\Rightarrow$ LIA-conflict)

# String Solver: Normalize Equality



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
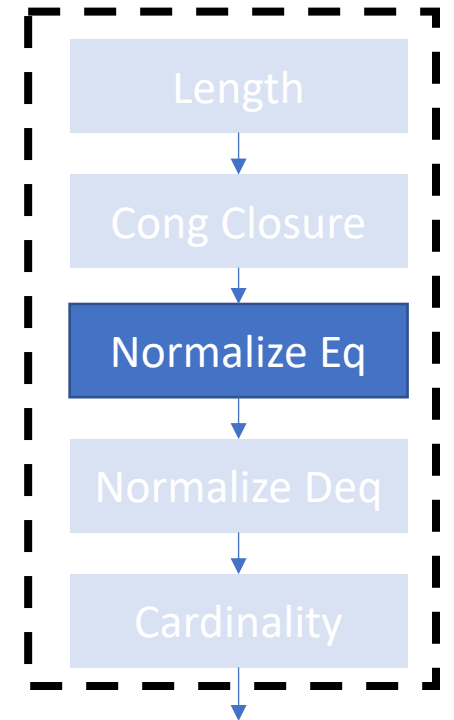$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

## Normal forms computed by a bottom-up procedure

- First, compute containment relation induced by concatenation terms
  - To compute a n.f. for eq-class of $x \cdot v$, we must first compute a n.f. for eq-class of $x$ and $v$
  - This relation is guaranteed to be acyclic due to length processing step (cycle $\Rightarrow$ LIA-conflict)
- Base case: eqc containing only variables can be assigned representative as a normal form
- Inductive case: compare the expanded form $t_1, \ldots, t_n$ of each non-variable term $t$
  - If $t_1 \cong \ldots \cong t_n$, assign to $t$. If there exists distinct $t_i$, $t_j$, then propagate or split

# String Solver: Normalize Equality



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

Single non-variable string term $\Rightarrow$ assign

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

# String Solver: Normalize Equality



z | "aab" | u | "b" | v

x, y, z·"aab"

w, u·"b"

x·v, v·w

≠

Single non-variable string term ⇒ assign

$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

# String Solver: Normalize Equality



z | "aab" | u | "b" | v

$x$, y, z·"aab"

$w$, u·"b"

$x·v$, v·w

$\neq$

$$x = z·\text{"aab"}$$
$$y = x$$
$$w = u·\text{"b"}$$
$$x·v = v·w$$
$$x·v \neq w$$

Length
↓
Cong Closure
↓
Normalize Eq
↓
Normalize Deq
↓
Cardinality

# String Solver: Normalize Equality



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

Length

Cong Closure

**Normalize Eq**

Normalize Deq

Cardinality

# String Solver: Normalize Equality



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
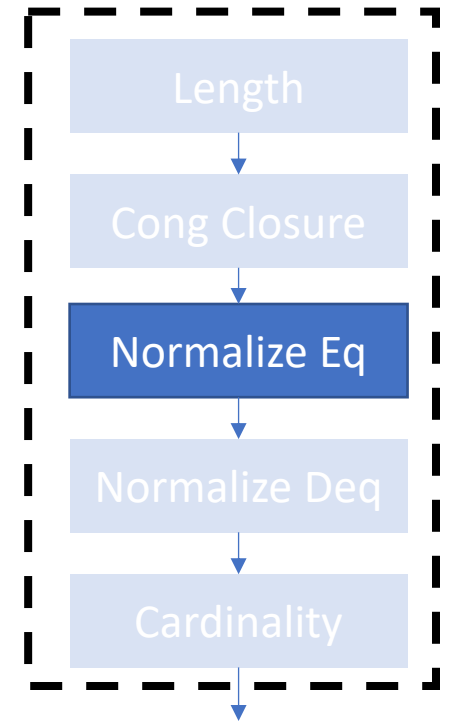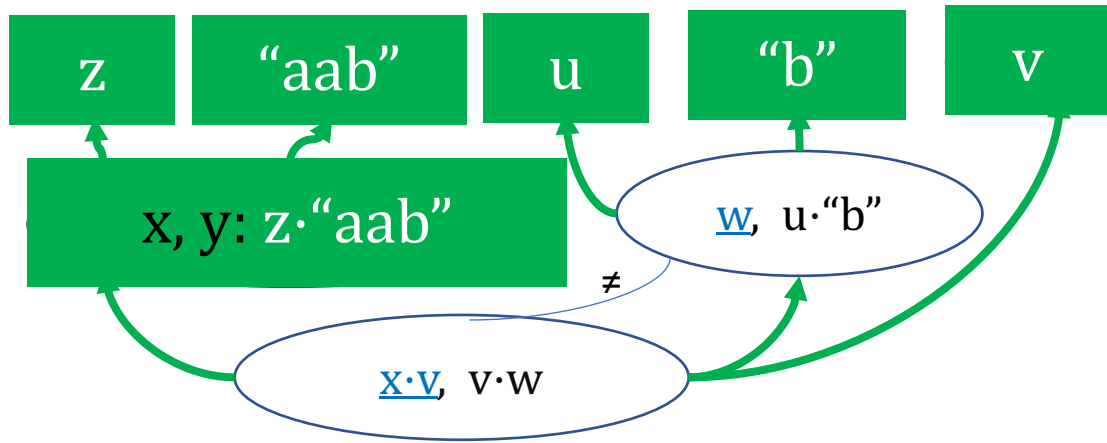$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

# String Solver: Normalize Equality



z   "aab"   u   "b"   v

x, y: z·"aab"     w: u·"b"

≠

$\underline{x \cdot v}$, v·w

$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

Equivalence class with two non-variable terms with distinct expanded forms:

- x·v = (z·"aab")·v     = z·"aab"·v
- v·w = v·(u·"b")       = v·u·"b"

# String Solver: Normalize Equality



z    "aab"    u    "b"    v

x, y: z·"aab"     w: u·"b"

≠

$z·"aab"·v \overset{?}{=} v·u·"b"$

$x = z·"aab"$
$y = x$
$w = u·"b"$
$x·v = v·w$
$x·v \neq w$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

# String Solver: Normalize Equality



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
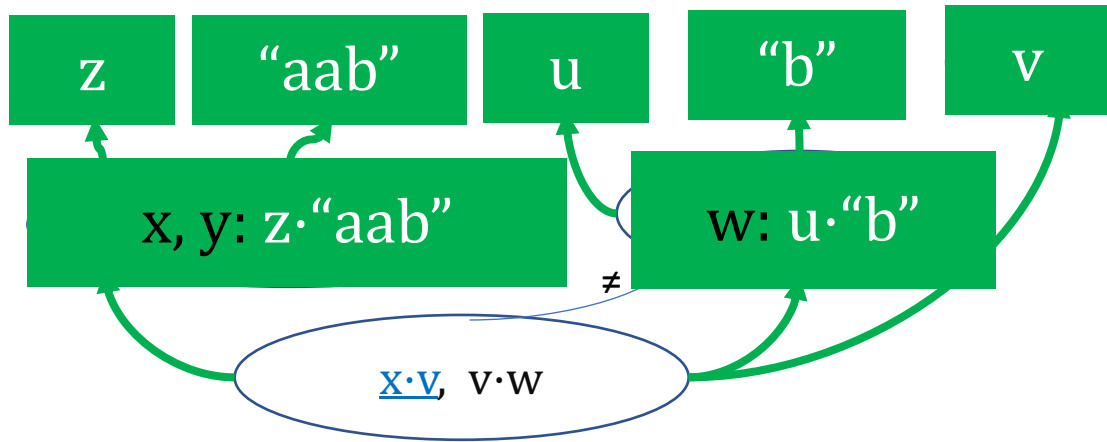$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

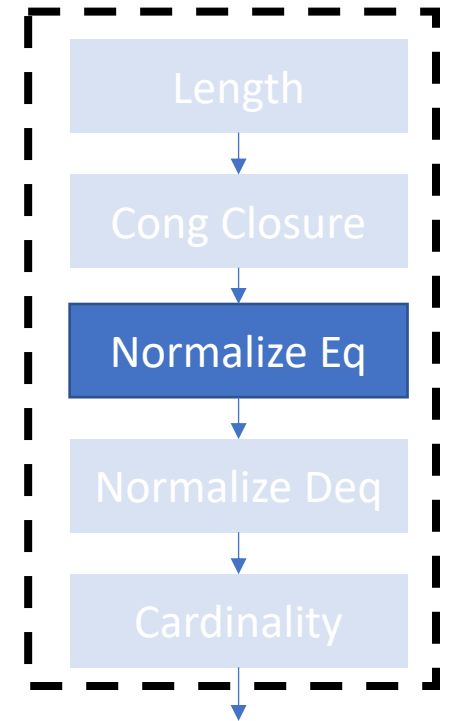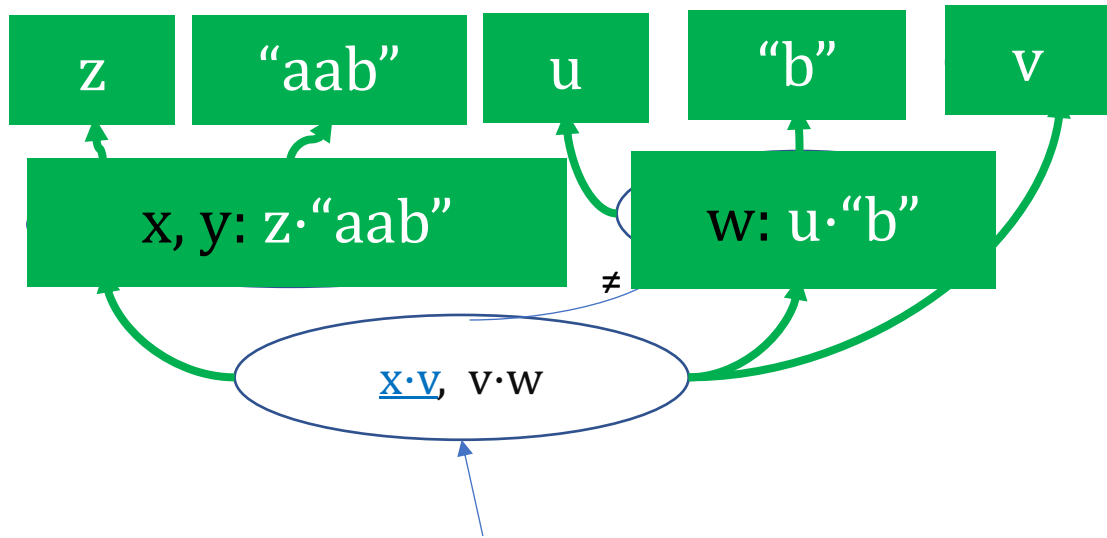z · "aab" · v $\overset{?}{=}$ v · u · "b"

Length
Cong Closure
Normalize Eq
Normalize Deq
Cardinality

| z | "aab" | v | $\overset{?}{=}$

Goal: split strings so that **all** aligning components are equal

| v | u | "b" |

# String Solver: Normalize Equality

z     "aab"     u     "b"     v

x: z·"aab"     w: u·"b"

≠

$$z \cdot \text{"aab"} \cdot v \overset{?}{=} v \cdot u \cdot \text{"b"}$$

$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$

Length
↓
Cong Closure
↓
Normalize Eq
↓
Normalize Deq
↓
Cardinality

- Consider three cases for making these two terms equal:

| z | "aab" | v |
|---|---|---|

‖      When $|z| = |v|$

| v | u | "b" |
|---|---|---|

# String Solver: Normalize Equality

z    "aab"    u    "b"    v

x: z·"aab"    w: u·"b"

≠

$$z·"aab"·v \stackrel{?}{=} v·u·"b"$$

$$x = z·"aab"$$
$$y = x$$
$$w = u·"b"$$
$$x·v = v·w$$
$$x·v \neq w$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

- Consider three cases for making these two terms equal:

z    "aab"    v

z    v'    When $|z| < |v|$

=

v    u    "b"

# String Solver: Normalize Equality



z    "aab"    u    "b"    v

x: z·"aab"    w: u·"b"

$\ne$

$$z·\text{"aab"}·v \overset{?}{=} v·u·\text{"b"}$$

$$x = z·\text{"aab"}$$
$$y = x$$
$$w = u·\text{"b"}$$
$$x·v = v·w$$
$$x·v \ne w$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

- Consider three cases for making these two terms equal:

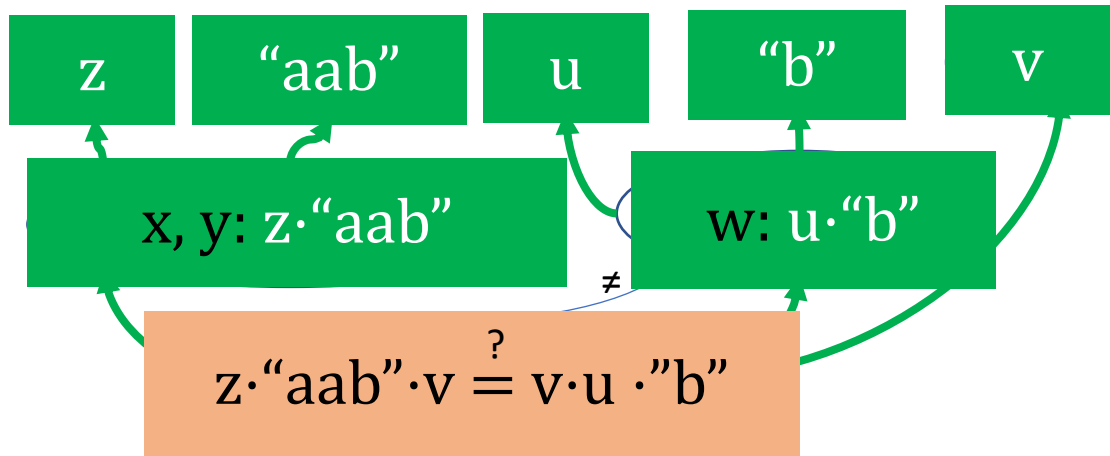| z | "aab" | v |

||

| v | z' | When $|z| > |v|$

| v | u | "b" |

# String Solver: Normalize Equality

$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$
$$\boxed{z = v}$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

- Consider:

| z | "aab" | v |
|---|-------|---|

=

| v | u | "b" |
|---|---|-----|

# String Solver: Normalize Equality

$v$,z     _"abb"_     $u$     _"b"_

$x$, y, z·"aab"     $w$, u·"b"

$x·v$, z·w    ≠

$$x = z·"aab"$$
$$y = x$$
$$w = u·"b"$$
$$x·v = v·w$$
$$x·v \neq w$$
$$z = v$$

| Length |
| --- |
| Cong Closure |
| **Normalize Eq** |
| Normalize Deq |
| Cardinality |

Recompute <span style="color:red">congruence closure</span>

# String Solver: Normalize Equality



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
$$x \cdot v \neq w$$
$$z = v$$

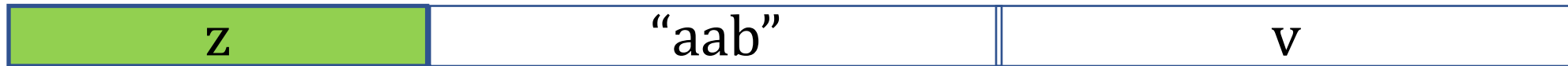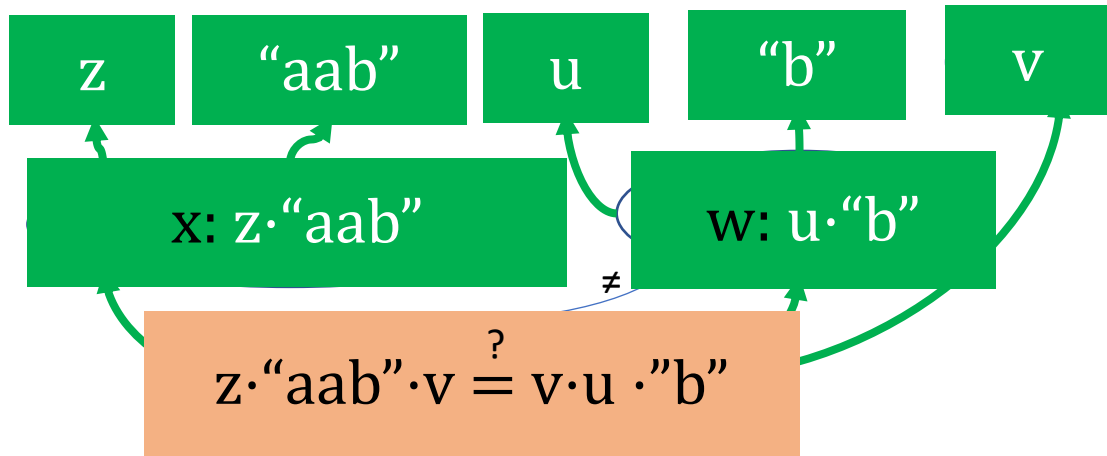Recompute congruence closure and normal forms

# String Solver: Normalize Equality



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v = v \cdot w$$
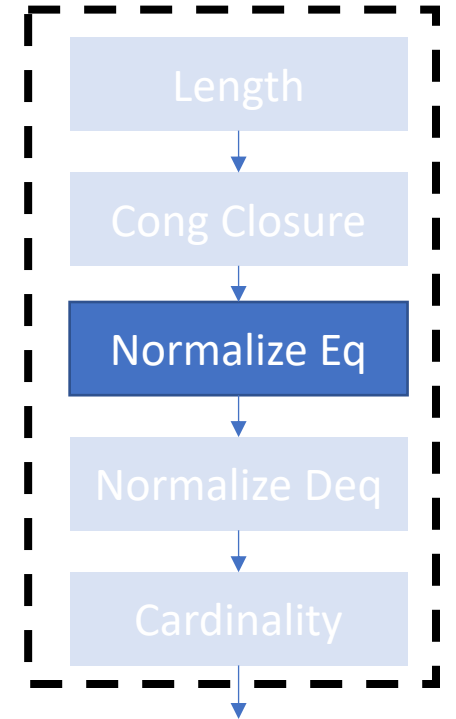$$x \cdot v \neq w$$
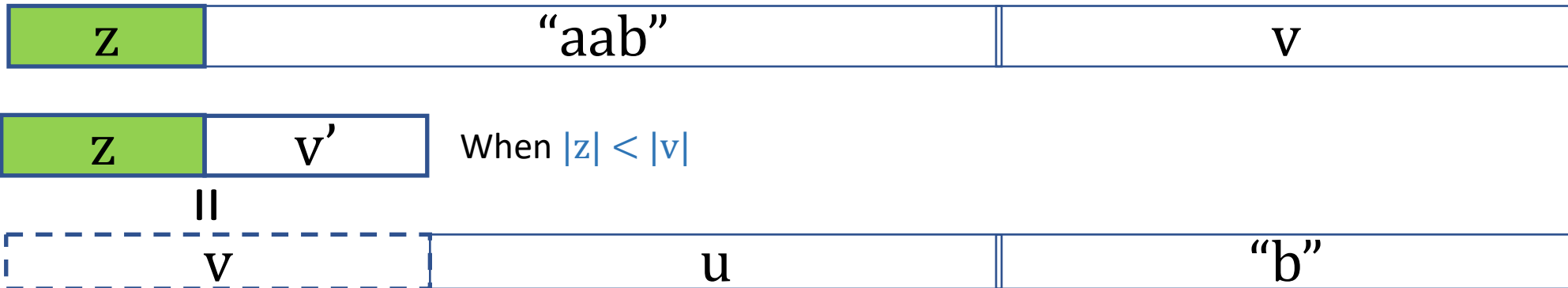$$z = v$$

Recompute congruence closure and *normal forms*

# String Solver: Normalize Equality



$x = z \cdot \text{"aab"}$
$y = x$
$w = u \cdot \text{"b"}$
$x \cdot v = v \cdot w$
$x \cdot v \neq w$
$z = v$

# String Solver: Normalize Equality

v   "aab"   u   "b"

x: v·"aab"   w: u·"b"

≠

$$v·"aab"·v \overset{?}{=} v·u·"b"$$

$$x = z·"aab"$$
$$y = x$$
$$w = u·"b"$$
$$x·v = v·w$$
$$x·v \neq w$$
$$z = v$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

v   "aab"   v

repeat the process on these components

v   u   "b"

# Splitting on String Equalities

Choosing how to process equalities is

highly non-trivial and critical to performance:

- Prefer propagations over splits

  Infer $x \cdot w = y \cdot w \Rightarrow x = y$  before  $x \cdot w = z \cdot v \Rightarrow (x = z \cdot x' \lor z = x \cdot z')$

- Can consider both the prefix and suffix of strings

  Infer $w \cdot x = w \cdot y \Rightarrow x = y$

- Use length entailment [Zheng et al 2015]

  If $|x| > |y|$ is entailed by the arith. solver, then $x \cdot w = y \cdot v \land |x| > |z| \Rightarrow x = y \cdot x'$

# Splitting on String Equalities

Choosing how to process equalities is

highly non-trivial and critical to performance:

- Propagation based on adjacent constants
  $x \cdot \text{"b"} = \text{"aab"} \cdot y \Rightarrow x = \text{"aa"} \cdot x'$, since "b" cannot overlap with prefix "aa"

- Special treatment for looping word equations [Liang et al 2014]
  - splitting leads to non-termination; reduce to RE membership instead
  - e.g. $x \cdot \text{"ba"} = \text{"ab"} \cdot x \Rightarrow x \in (\text{"ab"})^* \cdot \text{"a"}$

- Deduced string equalities are not sent as unit lemmas
  instead they are maintained internally

# String Solver: Normalize Disequalities

modified example

$$x = z \cdot "aab"$$
$$y = x$$
$$w = u \cdot "b"$$
$$x \cdot v \neq v \cdot w$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

# String Solver: Normalize Disequalities

$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v \neq v \cdot w$$

Length

Cong Closure

Normalize Eq

**Normalize Deq**

Cardinality

Disequalities are handled analogously to equalities
- If $|x \cdot v| \neq |v \cdot w|$, then trivially $x \cdot v \neq v \cdot w$
- Otherwise, consider the normal forms of $x \cdot v$ and $v \cdot w$ from previous step

# String Solver: Normalize Disequalities



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v \neq v \cdot w$$

Disequalities are handled analogously to equalities

# String Solver: Normalize Disequalities



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v \neq v \cdot w$$

Length → Cong Closure → Normalize Eq → **Normalize Deq** → Cardinality

Disequalities are handled analogously to equalities

Goal: find **any** aligning component that is disequal

# String Solver: Normalize Disequalities



$$x = z \cdot \text{"aab"}$$
$$y = x$$
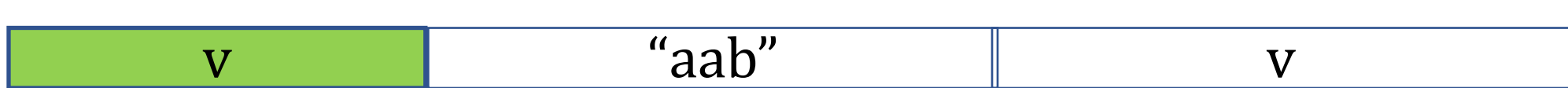$$w = u \cdot \text{"b"}$$
$$x \cdot v \neq v \cdot w$$
$$v \neq z$$

Length → Cong Closure → Normalize Eq → **Normalize Deq** → Cardinality

Disequalities are handled analogously to equalities

| z | "aab" | v |
|---|---|---|

$\neq$   $|z| = |v|$ and $z \neq v$

| v | u | "b" |
|---|---|---|

# String Solver: Cardinality

$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v \neq v \cdot w$$
$$v \neq z$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

# String Solver: Cardinality

$$x = z \cdot "aab"$$
$$y = x$$
$$w = u \cdot "b"$$
$$x \cdot v \neq v \cdot w$$
$$v \neq z$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

- $M_S$ may be unsatisfiable since alphabet A is <span style="color:red">finite</span>

- For instance, if:
  - A is a finite alphabet of 256 characters, and
  - $M_S$ entails the existence of 257 distinct strings of length 1
  $\Rightarrow$ Then $M_S$ is unsatisfiable

$\therefore \; (\text{distinct}(s_1, ..., s_{257}) \wedge |s_1| = ... = |s_{257}|) \Rightarrow |s_1| > 1$

# String Solver: Return SAT

$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v \neq v \cdot w$$
$$v \neq z$$

| Length |
|--------|
| Cong Closure |
| Normalize Eq |
| Normalize Deq |
| Cardinality |

SAT

If all steps finish with no new lemmas:

1. $M_s$ is $T_s$-satisfiable
2. Model can be computed based on normal forms
   - String constants assigned to eq classes whose normal form is a variable
     Length fixed by model from arithmetic solver
   - Each variable interpreted as the valuation of the normal form of their eq class

# String Solver: Return SAT

z    "aab"    v    "b"    u

x: z·"aab"          w: u·"b"

x·v: z·"aab" ·v          v·w: v·u·"b"

$$x = z\cdot\text{"aab"}$$
$$y = x$$
$$w = u\cdot\text{"b"}$$
$$x\cdot v \neq v\cdot w$$
$$v \neq z$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

SAT

If all steps finish with no new lemmas:

1. $M_s$ is $T_s$-satisfiable
2. Model can be computed based on normal forms
   - String constants assigned to eq classes whose normal form is a variable
     - Length fixed by model from arithmetic solver
   - Each variable interpreted as the valuation of the normal form of their eq class

# String Solver: Return SAT

$|z| = 1$   $|v| = 1$   $|u| = 3$

z   "aab"   v   "b"   u

x: z·"aab"   w: u·"b"

x·v: z·"aab" ·v   v·w: v·u·"b"

$x = z·"aab"$
$y = x$
$w = u·"b"$
$x·v \neq v·w$
$v \neq z$

**Example:**

model

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

SAT

Simplex

# String Solver: Return SAT

|z| = 1  |v| = 1  |u| = 3

"c"   "aab"   v   "b"   u

x: "c"·"aab"   w: u·"b"

x·v: "c"·"aab"·v   v·w: v·u·"b"

$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v \neq v \cdot w$$
$$v \neq z$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

SAT

**Example:**
- z assigned to "c"

# String Solver: Return SAT

$|v| = 1$        $|u| = 3$



| "c" | "aab" | "d" | "b" | u |

x: "c"·"aab"

w: u·"b"

x·v: "c"·"aab" ·"d"

v·w: "d"·u·"b"

$$x = z·"aab"$$
$$y = x$$
$$w = u·"b"$$
$$x·v \neq v·w$$
$$v \neq z$$

Length
Cong Closure
Normalize Eq
Normalize Deq
Cardinality

SAT

**Example:**
- z assigned to "c"
- v assigned to "d"

# String Solver: Return SAT

$|u| = 3$

"c"    "aab"    "d"    "b"    "aaa"

x: "c"·"aab"    w: "aaa"·"b"

x·v: "c"·"aab"·"d"    v·w: "d"·"aaa"·"b"

$x = z \cdot \text{"aab"}$
$y = x$
$w = u \cdot \text{"b"}$
$x \cdot v \neq v \cdot w$
$v \neq z$

**Example:**
- z assigned to "c"
- v assigned to "d"
- u assigned to "aaa"

Length
↓
Cong Closure
↓
Normalize Eq
↓
Normalize Deq
↓
Cardinality
↓

SAT

Cardinality step ensures enough enough constants exist

# String Solver: Return SAT

"c"   "aab"   "d"   "b"   "aaa"

x: "c"·"aab"    w: "aaa"·"b"

x·v: "c"·"aab"·"d"    v·w: "d"·"aaa"·"b"

$$x = z·"aab"$$
$$y = x$$
$$w = u·"b"$$
$$x·v \neq v·w$$
$$v \neq z$$

Length

Cong Closure

Normalize Eq

Normalize Deq

Cardinality

SAT

**Example:**

- z assigned to "c"
- v assigned to "d"
- u assigned to "aaa"
- Variables assigned to value of the normal form of their eq classes:
  - x,y assigned to "caab", w assigned to "aaab"

# String Solver: Return SAT



$$x = z \cdot \text{"aab"}$$
$$y = x$$
$$w = u \cdot \text{"b"}$$
$$x \cdot v \neq v \cdot w$$
$$v \neq z$$

Length → Cong Closure → Normalize Eq → Normalize Deq → Cardinality

SAT

**Example:**

- z assigned to "c"
- v assigned to "d"
- u assigned to "aaa"
- Variables assigned to value of the normal form of their eq classes:
  - x,y assigned to "caab", w assigned to "aaab"

Saturation criteria of procedure ensures this model satisfies $M_s$

# Advanced Topics

- Finite model finding for strings
- Context-dependent simplification for extended string constraints
- Regular expression elimination

# Finite Model Finding for Strings

# Finite Model Finding for Strings

**Idea:** Incrementally bound the lengths of input string variables $x_1, \ldots, x_n$

$\Rightarrow$ Improved solver's ability to answer "SAT" for problems with small models

Search for models where sum of lengths is 0

$\Sigma_{i=1\ldots n}|x_i| \leq 0$     $\neg\Sigma_{i=1\ldots n}|x_i| \leq 0$

$\Sigma_{i=1\ldots n}|x_i| \leq 1$     $\neg\,\Sigma_{i=1\ldots n}|x_i| \leq 1$

Search for models where sum of lengths is 1

$\Sigma_{i=1\ldots n}|x_i| \leq 2$     $\neg\,\Sigma_{i=1\ldots n}|x_i| \leq 2$

etc.

# Finite Model Finding

- Minimize sum of lengths $\Sigma_{i=1...n} |x_i| \leq 0$
- Which variables have unbounded length?

$$x = \text{“ab”} \cdot z$$
$$x = y \cdot u \cdot v \lor u \neq \text{“abc”}$$
$$w = x \cdot \text{”ab”} \lor w = y \cdot \text{”cde”}$$

# Finite Model Finding

- Minimize sum of lengths $\Sigma_{i=1\ldots n} |x_i| \leq 0$
- Which variables have unbounded length?

$$x = \text{``ab''} \cdot z$$
$$x = y \cdot u \cdot v \vee u \neq \text{``abc''}$$
$$w = x \cdot \text{''ab''} \vee w = y \cdot \text{''cde''}$$

- Can include a subset of the overall input variables in this sum
  Above, upper bound on $|x + u|$ implies upper bounds on the length of $z, y, w, v$
- Reduces the overall sum of lengths

# Context-Dependent Simplification for Extended String Constraints

# Extended String Constraints

- *Basic* terms
  - String and integer variables, constants, concatenation, length, and LIA-terms
- *Extended* string terms:
  - Substring: $\text{substr}(x, 1, 3)$

    (the substring of $x$ starting at pos. $1$ of length at most $3$)
  - String contains: $\text{contains}(x, "abc")$

    ($\text{true}$ iff $x$ contains the substring "abc")
  - Find "index of": $\text{indexof}(x, "d", 5)$

    (the pos. of the first occurrence of "d" in $x$, starting from position $5$, or $-1$ if it does not exist)
  - String replace: $\text{replace}(x, "a", "b")$

    (the result of replacing the first occurrence of "a" in $x$, if any, with "b")

**Example:** $\neg\text{contains}(\text{substr}(x, 0, 3), "a") \wedge 0 \le \text{indexof}(x, "ab", 0) < 4$

# Processing Extended String Constraints

$\neg\text{contains}(x, \text{"a"})$

# Processing Extended String Constraints

- Naively, by reduction to basic constraints + bounded $\forall$

$$\neg contains(x, \text{``a"})$$

# Processing Extended String Constraints

- Naively, by reduction to basic constraints + bounded $\forall$

$\neg\text{contains}(x, \text{"a"})$

$\forall 0 \leq n < |x|.\ \text{substr}(x, n, 1) \neq \text{"a"}$

Expand contains

# Processing Extended String Constraints

- Naively, by reduction to basic constraints + bounded $\forall$

$\neg\text{contains}(x, \text{"a"})$

$\forall 0 \leq n < |x|.\ \text{substr}(x, n, 1) \neq \text{"a"}$

Expand contains

$\text{substr}(x, 0, 1) \neq \text{"a"} \wedge \ldots \wedge \text{substr}(x, 4, 1) \neq \text{"a"}$

Assuming bound $|x| \leq 5$

# Processing Extended String Constraints

- Naively, by reduction to basic constraints + bounded $\forall$

$\neg\text{contains}(x, \text{"a"})$

$\forall 0 \leq n < |x|.\ \text{substr}(x, n, 1) \neq \text{"a"}$

Expand contains

$\text{substr}(x, 0, 1) \neq \text{"a"} \wedge \ldots \wedge \text{substr}(x, 4, 1) \neq \text{"a"}$

Assuming bound $|x| \leq 5$

$$x = z_{11} \cdot k_1 \cdot z_{21} \wedge \qquad\qquad x = z_{14} \cdot k_4 \cdot z_{24} \wedge$$
$$|z_{11}| = 0 \wedge \qquad \ldots \qquad |z_{14}| = 4 \wedge$$
$$k_1 \neq \text{"a"} \wedge \qquad\qquad k_4 \neq \text{"a"}$$

Expand substr

# Processing Extended String Constraints

- Naively, by reduction to basic constraints + bounded $\forall$

$\neg\text{contains}(x, \text{"a"})$

$\forall 0 \leq n < |x|.\ \text{substr}(x, n, 1) \neq \text{"a"}$

Expand contains

$\text{substr}(x, 0, 1) \neq \text{"a"} \wedge \ldots \wedge \text{substr}(x, 4, 1) \neq \text{"a"}$

Assuming bound $|x| \leq 5$

$x = z_{11} \cdot k_1 \cdot z_{21} \wedge$
$|z_{11}| = 0 \wedge$
$k_1 \neq \text{"a"} \wedge$

$\ldots$

$x = z_{14} \cdot k_4 \cdot z_{24} \wedge$
$|z_{14}| = 4 \wedge$
$k_4 \neq \text{"a"}$

Expand substr

- Approach used by many current solvers
  [Bjorner et al. 2009, Zheng et al. 2013, Li et al. 2013, Trinh et al. 2014]

# (Eager) Expansion of Extended Constraints

$\neg$contains(x,"a")
x = y·"d"
y = "ab" $\vee$ y = "ac"

| SAT Solver | Arithmetic Solver | String Solver |
|---|---|---|

# (Eager) Expansion of Extended Constraints

$\neg$contains(x,"a")
x = y·"d"
y = "ab" $\lor$ y = "ac"

$x = z_{11} \cdot k_1 \cdot z_{21}$
$|z_{11}| = 0$
$k_1 \neq$ "a"

...

$x = z_{14} \cdot k_4 \cdot z_{24}$
$|z_{14}| = 4$
$k_4 \neq$ "a"

x = y·"d"
y = "ab" $\lor$ y = "ac"

Expand and eliminate extended symbols

**SAT Solver**

**Arithmetic Solver**

**String Solver**

# (Eager) Expansion of Extended Constraints

$\neg$contains(x,"a")

$x = y\cdot$"d"

$y = $"ab" $\vee$ $y = $"ac"

---

$x = z_{11}\cdot k_1\cdot z_{21}$ ... $x = z_{14}\cdot k_4\cdot z_{24}$

$|z_{11}| = 0$ $|z_{14}| = 4$

$k_1 \neq$ "a" $k_4 \neq$ "a"

$x = y\cdot$"d"

$y = $"ab" $\vee$ $y = $"ac"

Expand and eliminate extended symbols

**SAT Solver**

**Arithmetic Solver**

**String Solver**

UNSAT

or

SAT

# (Eager) Expansion of Extended Constraints

$\neg contains(x, \text{``a''})$
$x = y \cdot \text{``d''}$
$y = \text{``ab''} \vee y = \text{``ac''}$

$x = z_{11} \cdot k_1 \cdot z_{21}$           $x = z_{14} \cdot k_4 \cdot z_{24}$
$|z_{11}| = 0$      ...      $|z_{14}| = 4$
$k_1 \neq \text{``a''}$            $k_4 \neq \text{``a''}$

$x = y \cdot \text{``d''}$
$y = \text{``ab''} \vee y = \text{``ac''}$

Must deal with a large constraint set

## SAT Solver

## Arithmetic Solver

## String Solver

UNSAT

or

SAT

85

# (Eager) Expansion of Extended Constraints

$\neg$contains(x,"a")
x = y·"d"
y = "ab" $\lor$ y = "ac"

...what if we simplify the input?

$x = z11 \cdot k1 \cdot z21$         $x = z14 \cdot k4 \cdot z24$
$|z11| = 0$          ...          $|z14| = 4$
$k1 \neq$" a"                    $k4 \neq$"a"

$x = y \cdot$"d"
$y =$ "ab" $\lor$ y = "ac"

**SAT Solver**

**Arithmetic Solver**

**String Solver**

UNSAT          or          SAT

# SMT Solvers + Simplification

All SMT solvers implement *simplification* techniques

(also called *normalization* or *rewrite* rules)

$\neg$contains(x, "a")

x = y·"d"

y = "ab" $\vee$ y = "ac"

# SMT Solvers + Simplification

All SMT solvers implement *simplification* techniques

(also called *normalization* or *rewrite* rules)

$$\neg contains(x, \text{``a''})$$
$$x = y \cdot \text{``d''}$$
$$y = \text{``ab''} \lor y = \text{``ac''}$$

$$\neg contains(y \cdot \text{``d''}, \text{``a''})$$
$$y = \text{``ab''} \lor y = \text{``ac''}$$

since $x = y \cdot \text{``d''}$

# SMT Solvers + Simplification

All SMT solvers implement *simplification* techniques

(also called *normalization* or *rewrite* rules)

$\neg$contains(x, "a")
x = y·"d"
y = "ab" $\vee$ y = "ac"

$\neg$contains(y·"d","a")$\wedge$
y="ab" $\vee$ y="ac"

since x = y·"d"

$\neg$contains(y,"a")$\wedge$
y="ab" $\vee$ y="ac"

since contains(y·"d", "a") $\Leftrightarrow$ contains(y, "a")

# SMT Solvers + Simplification

All SMT solvers implement *simplification* techniques

(also called *normalization* or *rewrite* rules)

$\neg \text{contains}(x, \text{"a"})$
$x = y \cdot \text{"d"}$
$y = \text{"ab"} \lor y = \text{"ac"}$

$\neg \text{contains}(y \cdot \text{"d"}, \text{"a"}) \land$
$y = \text{"ab"} \lor y = \text{"ac"}$

since $x = y \cdot \text{"d"}$

$\neg \text{contains}(y, \text{"a"}) \land$
$y = \text{"ab"} \lor y = \text{"ac"}$

since $\text{contains}(y \cdot \text{"d"}, \text{"a"}) \Leftrightarrow \text{contains}(y, \text{"a"})$

## Leads to smaller inputs

Some problems can be solved by simplification alone

# (Lazy) Expansion + Simplification

¬contains(x, "a")
x = y·"d"
y = "ab" ∨ y = "ac"

SAT
Solver

Arithmetic
Solver

String
Solver

# (Lazy) Expansion + Simplification

$\neg$contains(x, "a")
x = y·"d"
y = "ab" $\vee$ y = "ac"

$\neg$contains(y, "a")
y = "ab" $\vee$ y = "ac"

Simplify the input

## SAT Solver

## Arithmetic Solver

## String Solver

# (Lazy) Expansion + Simplification

¬contains(y,"a")
y="ab" ∨ y="ac"

SAT
Solver

Arithmetic
Solver

¬contains(y, "a")
y = "ab"

String
Solver

# (Lazy) Expansion + Simplification



¬contains(y, "a")
y = "ab" ∨ y = "ac"

¬contains(y, "a")
y = "ab"

SAT Solver

Arithmetic Solver

String Solver

$\text{contains}(y, \text{"a"}) \Leftrightarrow$

$(y = z_{11} \cdot k_1 \cdot z_{21} \wedge$
$|z_{11}| = 0 \wedge$
$k_1 \neq \text{"a"} \wedge$

$\dots$

$y = z_{14} \cdot k_4 \cdot z_{24} \wedge$
$|z_{14}| = 4 \wedge$
$k_4 \neq \text{"a"})$

# (Lazy) Expansion + Simplification

¬contains(y, "a")
y = "ab" ∨ y = "ac"

¬contains(y, "a")
y = "ab"

## SAT Solver

## Arithmetic Solver

## String Solver

Still have a large constraint

$contains(y, "a") \Leftrightarrow$

$(y = z_{11} \cdot k_1 \cdot z_{21} \land$
$|z_{11}| = 0 \land$
$k_1 \neq "a" \land$

...

$y = z_{14} \cdot k_4 \cdot z_{24} \land$
$|z_{14}| = 4 \land$
$k_4 \neq "a")$

# (Lazy) Expansion + Simplification

What if we simplify based on the context?

$\neg\text{contains}(y,\text{"a"})\wedge$
$y=\text{"ab"}\vee y=\text{"ac"}$

$\neg\text{contains}(y, \text{"a"})$
$y = \text{"ab"}$

| SAT Solver | Arithmetic Solver | String Solver |
|---|---|---|

$\text{contains}(y, \text{"a"}) \Leftrightarrow$

$(y = z_{11}{\cdot}k_1{\cdot}z_{21} \wedge$
$|z_{11}| = 0 \wedge$
$k_1 \neq \text{"a"} \wedge$

...

$y = z_{14}{\cdot}k_4{\cdot}z_{24} \wedge$
$|z_{14}| = 4 \wedge$
$k_4 \neq \text{"a"})$

# (Lazy) Expansion + Context-Dependent Simplification

¬contains(y, "a")
y = "ab" ∨ y = "ac"

¬contains(y, "a")
y = "ab"

**SAT Solver**

**Arithmetic Solver**

**String Solver**

Since contains(y, "a") is true when y = "ab" ...

# (Lazy) Expansion + Context-Dependent Simplification



¬contains(y, "a")
y = "ab" ∨ y = "ac"
¬y = "ab" ∨ contains(y, "a")

¬contains(y, "a")
y = "ab"

SAT Solver

Arithmetic Solver

String Solver

y= "ab" ⇒ contains(y, "a")

# (Lazy) Expansion + Context-Dependent Simplification

¬contains(y, "a")
y = "ab" ∨ y = "ac"
¬y = "ab" ∨ contains(y, "a")

¬contains(y, "a")
y = "ab"

**SAT Solver**

**Arithmetic Solver**

**String Solver**

# (Lazy) Expansion + Context-Dependent Simplification

¬contains(y, "a")
y = "ab" ∨ y = "ac"
¬y = "ab" ∨ contains(y, "a")

**SAT Solver**

**Arithmetic Solver**

¬contains(y, "a")
y = "ac"
¬y = "ab"

**String Solver**

# (Lazy) Expansion + Context-Dependent Simplification

¬contains(y, "a")
y = "ab" ∨ y = "ac"
¬y = "ab" ∨ contains(y, "a")

¬contains(y, "a")
y = "ac"
¬y = "ab"

**SAT Solver**

**Arithmetic Solver**

**String Solver**

contains(y, "a") is also true when y = "ac" …

# (Lazy) Expansion + Context-Dependent Simplification



¬contains(y, "a")
y = "ab" ∨ y = "ac"
¬y = "ab" ∨ contains(y, "a")
¬y = "ac" ∨ contains(y, "a")

¬contains(y, "a")
y = "ac"
¬y = "ab"

SAT Solver

Arithmetic Solver

String Solver

$y = \text{``ac''} \Rightarrow \text{contains}(y, \text{''a''})$

# (Lazy) Expansion + Context-Dependent Simplification

¬contains(y, "a")
y = "ab" ∨ y = "ac"
¬y = "ab" ∨ contains(y, "a")
¬y = "ac" ∨ contains(y, "a")

¬contains(y, "a")
y = "ac"
¬y = "ab"

SAT
Solver

Arithmetic
Solver

String
Solver

UNSAT

# (Lazy) Expansion + Context-Dependent Simplification

¬contains(y, "a")

y = "ab" ∨ y = "ac"

¬y = "ab" ∨ contains(y, "a")

¬y = "ac" ∨ contains(y, "a")

Did not need to fully expand  contains!

¬contains(y, "a")

y = "ac"

¬y = "ab"

**SAT Solver**

**Arithmetic Solver**

**String Solver**

UNSAT

**context-dependent simplification**

[Reynolds et al CAV 2017]

104

# Results on Symbolic Execution [Reynolds et al. CAV 17]



- cvc4+fs (finite model finding + context-dependent simpl.) solves    23,802 benchmarks in 5h8m
- Without finite model finding, solves    23,266 benchmarks in 8h46m
- Without either finite model finding or cd-simplification, solves    22,607 benchmarks in 6h38m

# Many Simplification Rules for Strings

Unlike arithmetic:

| x + x + 7*y = y - 4 | ┄┄┄┄┄┄┄┄┄► | 2*x + 6*y + 4 = 0 |

... simplification rules for strings are highly non-trivial:

| substr(x· "abcd", 1 + len(x),2) | ┄┄┄┄► | "bc" |

| contains("abcde", "b"· x· "a") | ┄┄┄┄► | ⊥ |

| contains(x·"ac"·y, "b") | ┄┄┄┄► | contains(x, "b") ∨ contains(y, "b") |

| indexof("abc"·x, "a"·x,1) | ┄┄┄┄► | -1 |

| replace("a"·x, "b",y) | ┄┄┄┄► | con("a", replace(x, "b", y)) |

# Simplification based on High-Level Abstractions

Rules based on high-level abstractions

- When viewing strings as #characters (e.g. reasoning about their length):

$$contains(substr(x, i, j), x \cdot "a") \dashrightarrow ""$$

since the second argument is longer than the first

- When considering the containment relationship between strings:

$$contains(replace(x, y, z), z) \dashrightarrow contains(x, y) \vee contains(x, z)$$

- When viewing strings as multisets of characters:

$$x \cdot x \cdot y \cdot "ab" = x \cdot "bbbbbb" \cdot y \dashrightarrow \perp$$

since LHS contains at least 1 more occurrences of "a"

# Impact of Aggressive Simplification

| Set | | all | -arith | -contain | -msets | z3 | OSTRICH |
|---|---|---|---|---|---|---|---|
| CMU | sat | 7947 | 7746 | **7948** | 7946 | 4585 | |
| | unsat | **66** | 31 | **66** | **66** | 52 | |
| | × | 173 | 409 | 172 | 174 | 3549 | |
| TERMEQ | sat | **10** | **10** | **10** | **10** | 1 | |
| | unsat | **49** | 36 | 27 | **49** | 36 | |
| | × | 22 | 35 | 44 | 22 | 44 | |
| SLOG | sat | **1302** | **1302** | **1302** | **1302** | 1100 | 1289 |
| | unsat | **2082** | **2082** | **2082** | **2082** | 2075 | **2082** |
| | × | 7 | 7 | 7 | 7 | 216 | 20 |
| APLAS | sat | **132** | **132** | **132** | **132** | 10 | |
| | unsat | **292** | 291 | 171 | 171 | 94 | |
| | × | 159 | 160 | 280 | 280 | 479 | |
| Total | sat | 9391 | 9190 | **9392** | 9390 | 5696 | 1289 |
| | unsat | **2489** | 2440 | 2346 | 2368 | 2257 | 2082 |
| | × | 361 | 611 | 503 | 483 | 4288 | 8870 |

[Reynolds et al. CAV 19]

**-arith:** w/o arithmetic simplifications
**-contain:** w/o contain-based simplifications
**-mset:** w/o multiset-based simplifications

CVC4 implements >3000 lines of C++ for simplification rules (and growing)
Important aspect of modern string solving

# Regular Expression Elimination

# Regular Expression Elimination

CVC4 supports regular expressions, via:

- Decomposing memberships
  E.g. $x \in R_1 \cup R_2 \Rightarrow x \in R_1 \lor x \in R_2$, $x \in R_1 \cap R_2 \Rightarrow x \in R_1 \land x \in R_2$

- Intersection (modulo equality):
  E.g. $(x \in R_1 \land y \in R_2 \land x = y) \Rightarrow x \in \text{compute\_intersection}(R_1, R_2)$

- Unfolding
  E.g. $x \in R^* \Rightarrow x = \text{""} \lor (x = x_1 \cdot x_2 \land x_1 \in R \land x_2 \in R^*)$ for fresh $x_1, x_2$

- Elimination based on reduction to extended string constraints

# Regular Expression Elimination

**Idea:** reduce RE to extended string constraints

Possible for many regular expression memberships that occur in practice:

$$x \in A \cdot A^* \cdot A \qquad \Longleftrightarrow \qquad |x| \geq 2$$

$$x \in A^* \cdot "abc" \cdot A^* \qquad \Longleftrightarrow \qquad \text{contains}(x, "abc")$$

$$x \in A^* \cdot "a" \cdot A^* \cdot "bcd" \cdot A^* \qquad \Longleftrightarrow \qquad \begin{array}{l} \text{contains}(x, "a") \wedge \\ \text{contains}(\text{substr}(x, \text{indexof}(x, "a", 1) + 1, |x|), "bcd") \end{array}$$
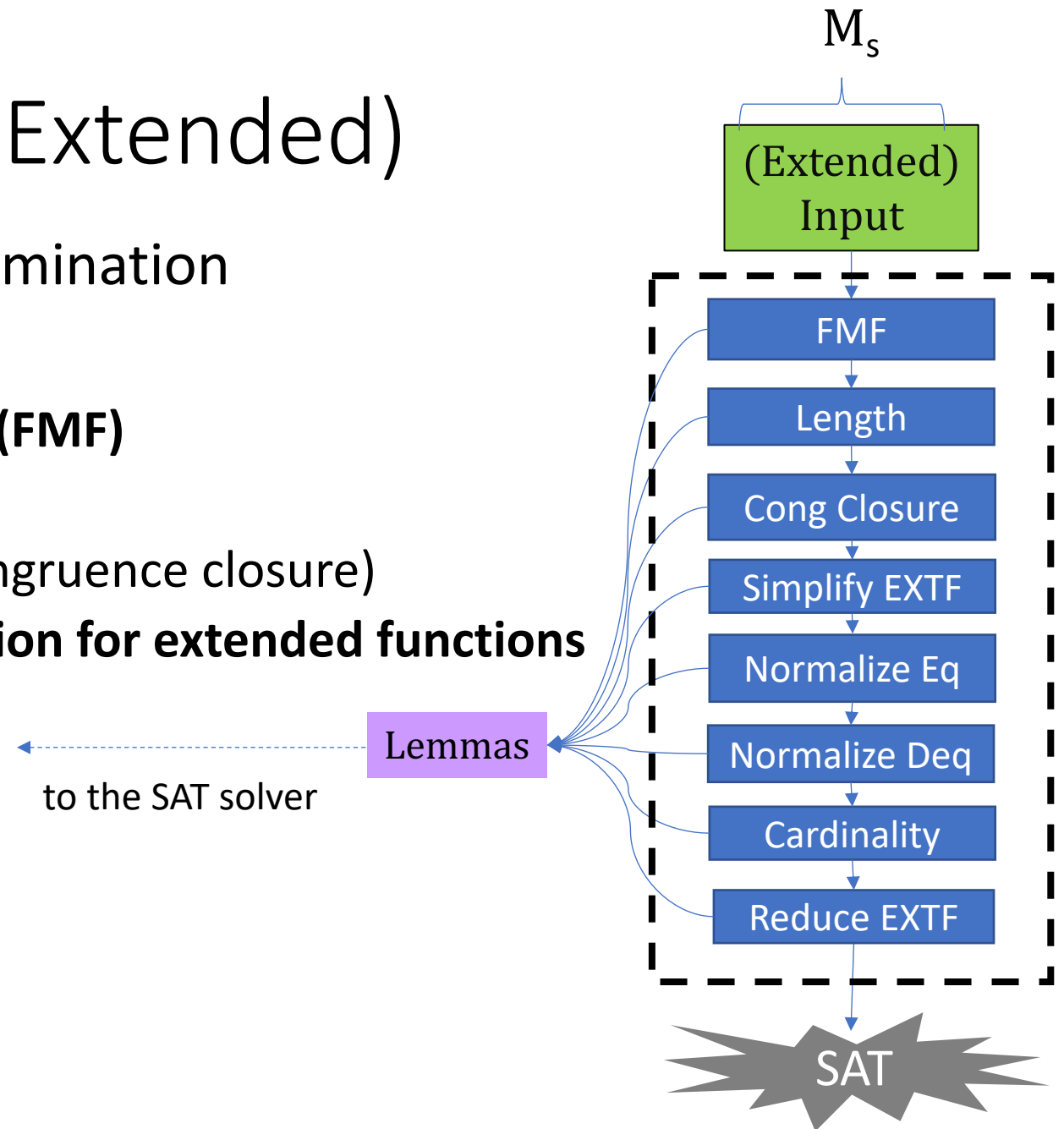
CVC4 supports (aggressive) elimination techniques for RE like those above

Utilizes existing support for extended functions

# String Theory Solver (Extended)

$M_s$

- Preprocess based on reg-exp elimination

- Then, run inference strategy:

  **1. Split on sum of lengths bound (FMF)**

  2. Process length constraints

  3. Check for equality conflicts (congruence closure)

  **4. Context-dependent simplification for extended functions**

  5. Normalize string equalities

  6. Normalize string disequalities

  7. Check cardinality constraints

  **8. Reduce extended functions**

(Extended) Input

FMF

Length

Cong Closure

Simplify EXTF

Normalize Eq

Normalize Deq

Cardinality

Reduce EXTF

Lemmas

to the SAT solver

SAT

# Conclusions

- CVC4 supports DPLL(T) theory solver for strings and regular expressions
  - Efficient in practice (incomplete) procedure for word equations with length
  - More advanced features like FMF, context-dependent simplification, RE elimination
  - Also supports: str.code, str.<=, str.to-int, str.from-int, str.replaceall

- Open-source, available at [https://cvc4.github.io/](https://cvc4.github.io/)

Thanks for listening