

Designing a Fast and Trustworthy String Solver

Andrew Reynolds & Cesare Tinelli

MOSCA 2023



Satisfiability Modulo Theories (SMT) Solvers

Many applications:

- Software verification
- Symbolic execution
- Security analysis
- Theorem proving

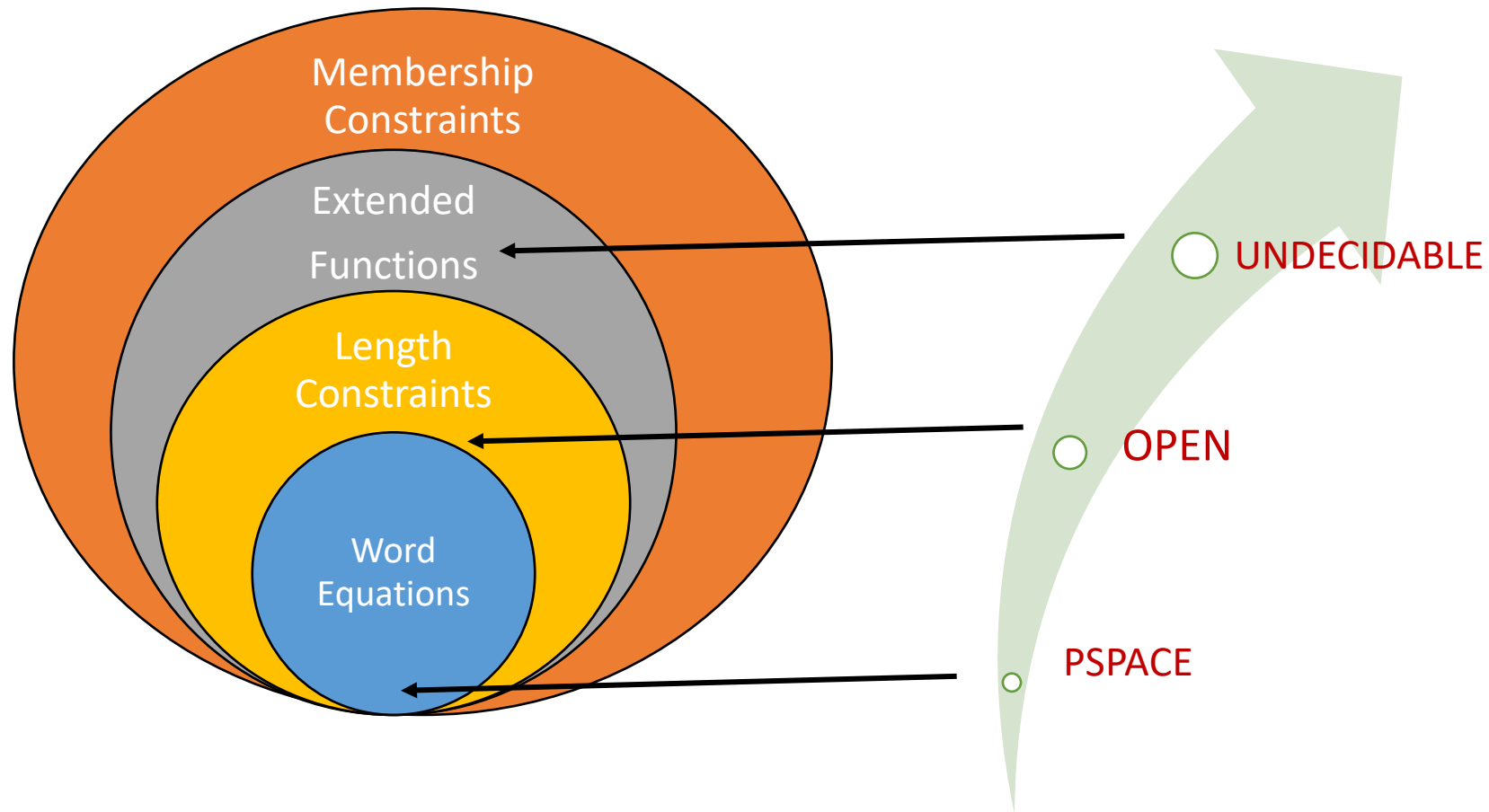
Traditionally:

- Efficient solvers for quantifier-free constraints over (combinations of) theories
 - Arithmetic, Arrays, Bit vectors

In this talk:

- SMT techniques **for string and RE constraints**

Strings and RE: Theoretical Challenges



Many applications require *extended string functions* and *RE memberships*

Ex.: $\text{toInt}(x) \neq 44$, $\text{toLowerCase}(x) = \mathbf{abc}$, $x \in \text{range}(\mathbf{A}, \mathbf{Z})$

The CVC4 and cvc5 SMT Solvers

Support for many theories and features

- UF, (non)linear arithmetic, arrays
- Bit-vectors, floating points
- (Multi)sets, relations, datatypes
- **Strings and regular expressions**

Co-developed at Stanford and Iowa

Project Leaders: Clark Barrett, Cesare Tinelli

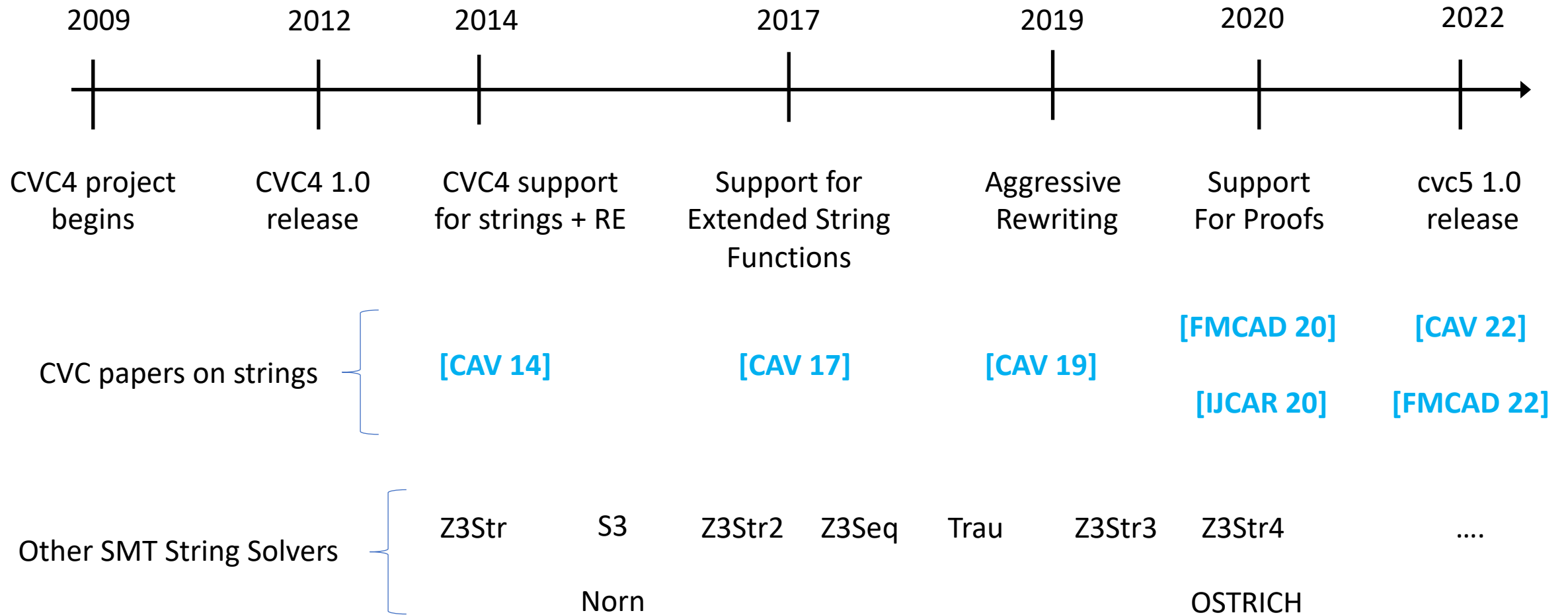
String solver developers: **Andrew Reynolds, Andres Noetzli**

Stanford
University



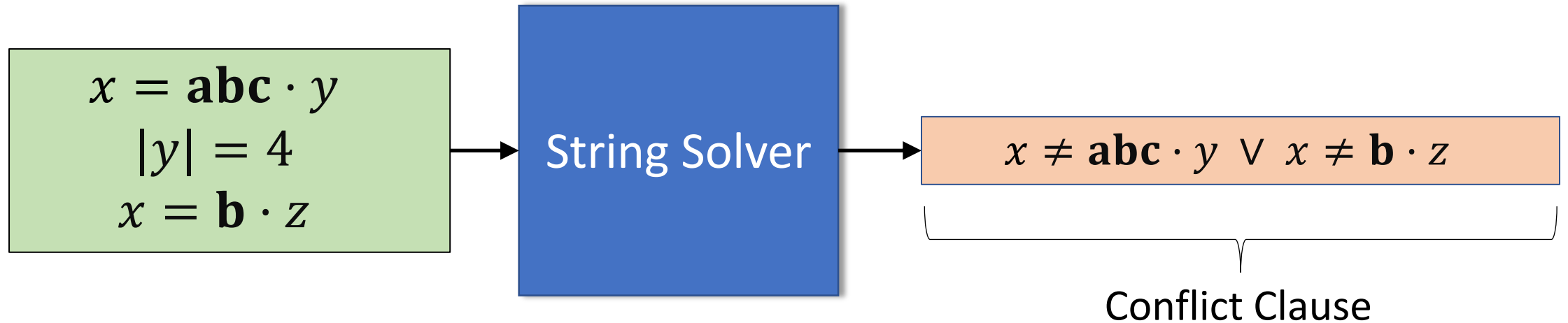
THE UNIVERSITY
OF IOWA

SMT Solvers for Strings: Timeline



String Reasoning in cvc5 in a Nutshell

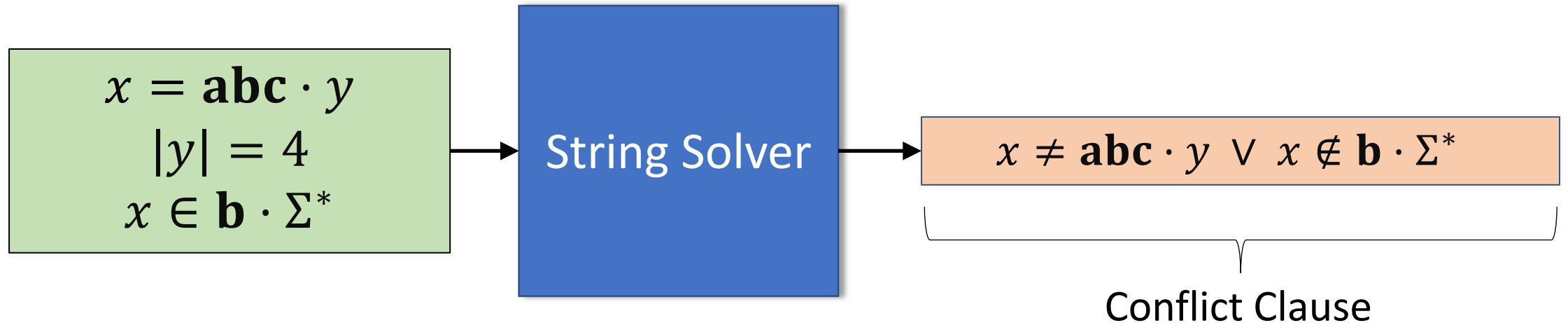
A Theory Solver for Strings [Liang et al., CAV'14]



Designed a string solver for **concat + length + RE constraints** that is:

- **refutation and model sound** ("unsat" and "sat" can be trusted)
- **not terminating** in general
- **efficient** in practice

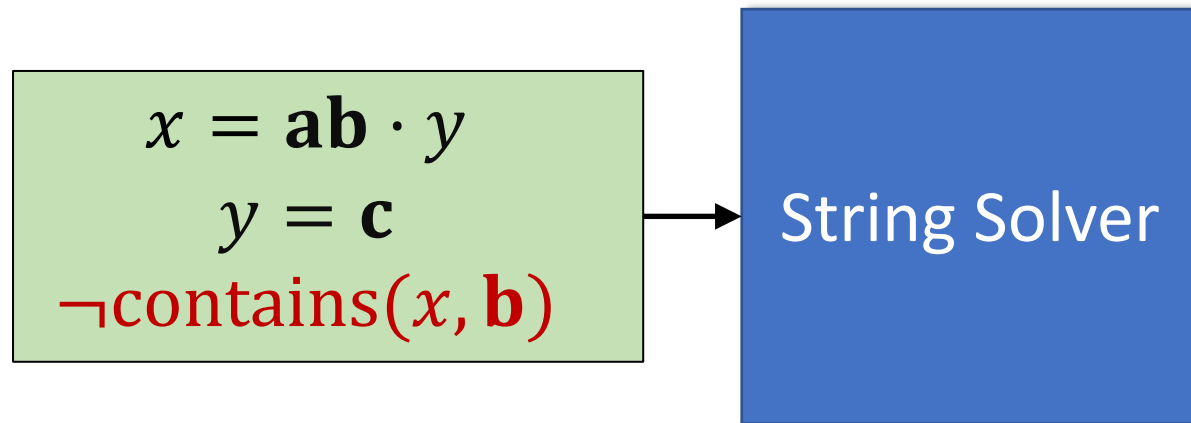
Perfecting support for REs [Liang et al., FroCoS'15]



Symbolic approach to RE constraint solving

- Yields a **decision procedure** over a reasonable fragment
- Gives rise to an **incremental** RE subsolver

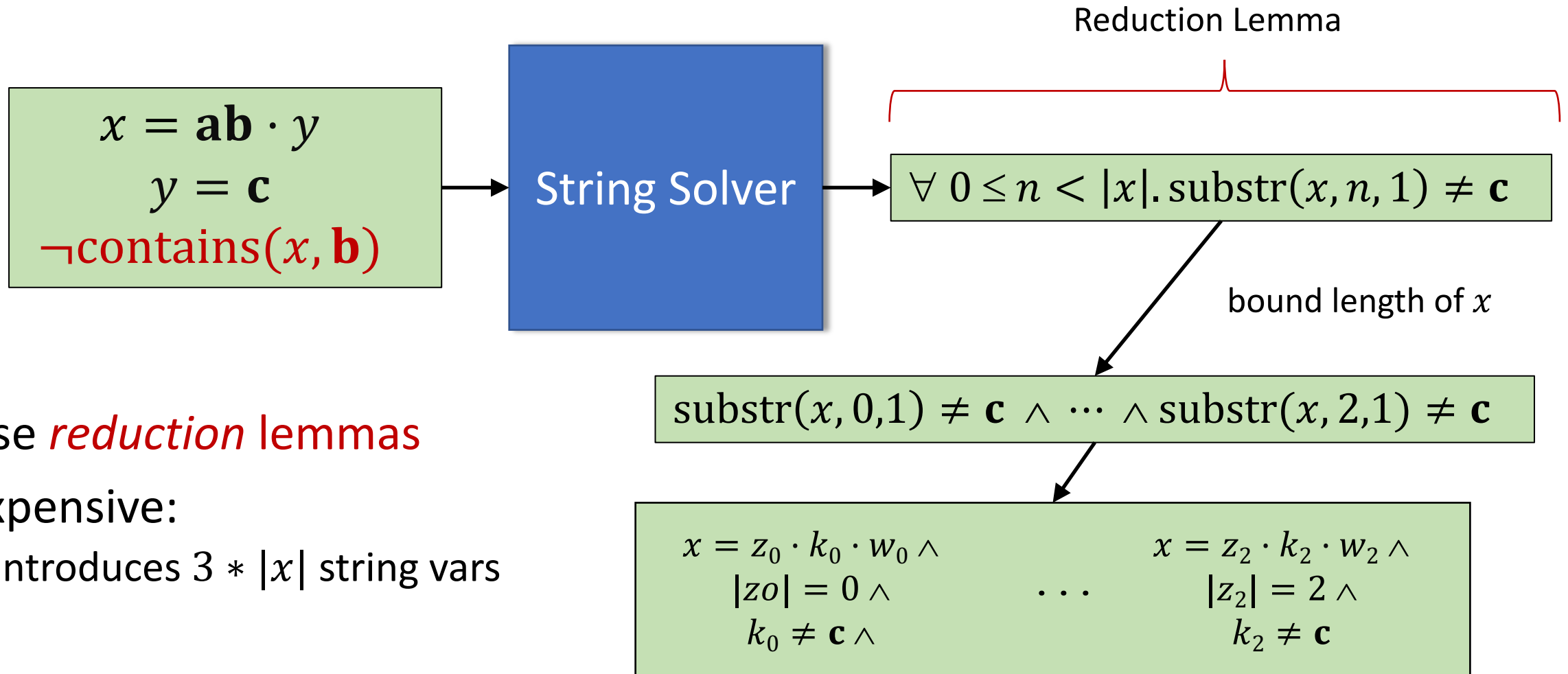
Extended Theory of Strings [Reynolds et al., CAV'17]



Support for *extended* string functions commonly used in applications

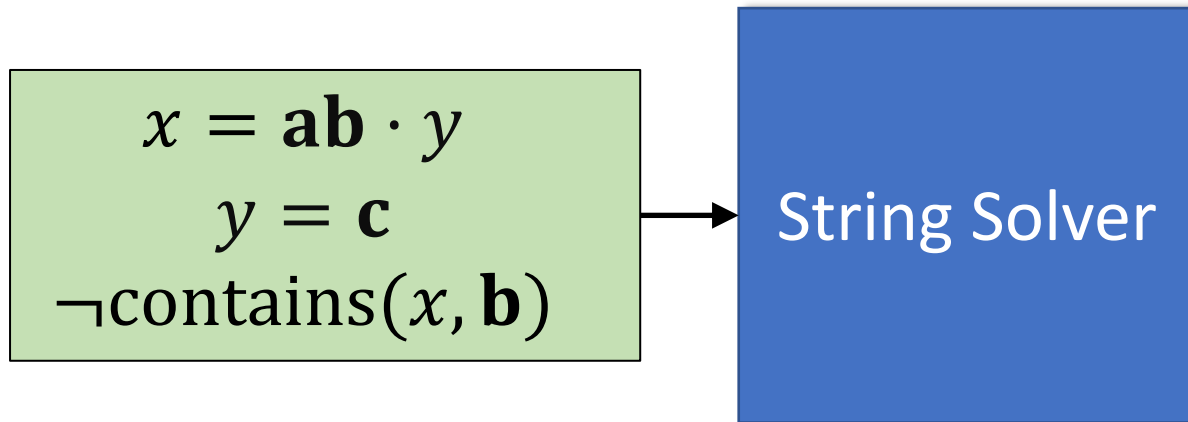
- $\text{substr}(x, n, m)$ substring of x at position n of length at most m
- $\text{contains}(x, y)$ true if string x contains substring y
- $\text{indexof}(x, y, n)$ position of string y in string x , starting from position n
- $\text{replace}(x, y, z)$ result of replacing first occurrence of y in x by z

Extended Theory of Strings [Reynolds et al., CAV'17]



- Use *reduction lemmas*
- Expensive:
Introduces $3 * |x|$ string vars

Extended Theory of Strings [Reynolds et al., CAV'17]



Context-dependent simplification crucial for performance

- $x = \mathbf{ab} \cdot y, y = \mathbf{c} \models x = \mathbf{abc}$
- $\neg \text{contains}(x, \mathbf{b}) \rightarrow \neg \text{contains}(\mathbf{abc}, \mathbf{b}) \rightarrow \neg \top \rightarrow \perp$

Proof Certificates for Unsat String Constraints

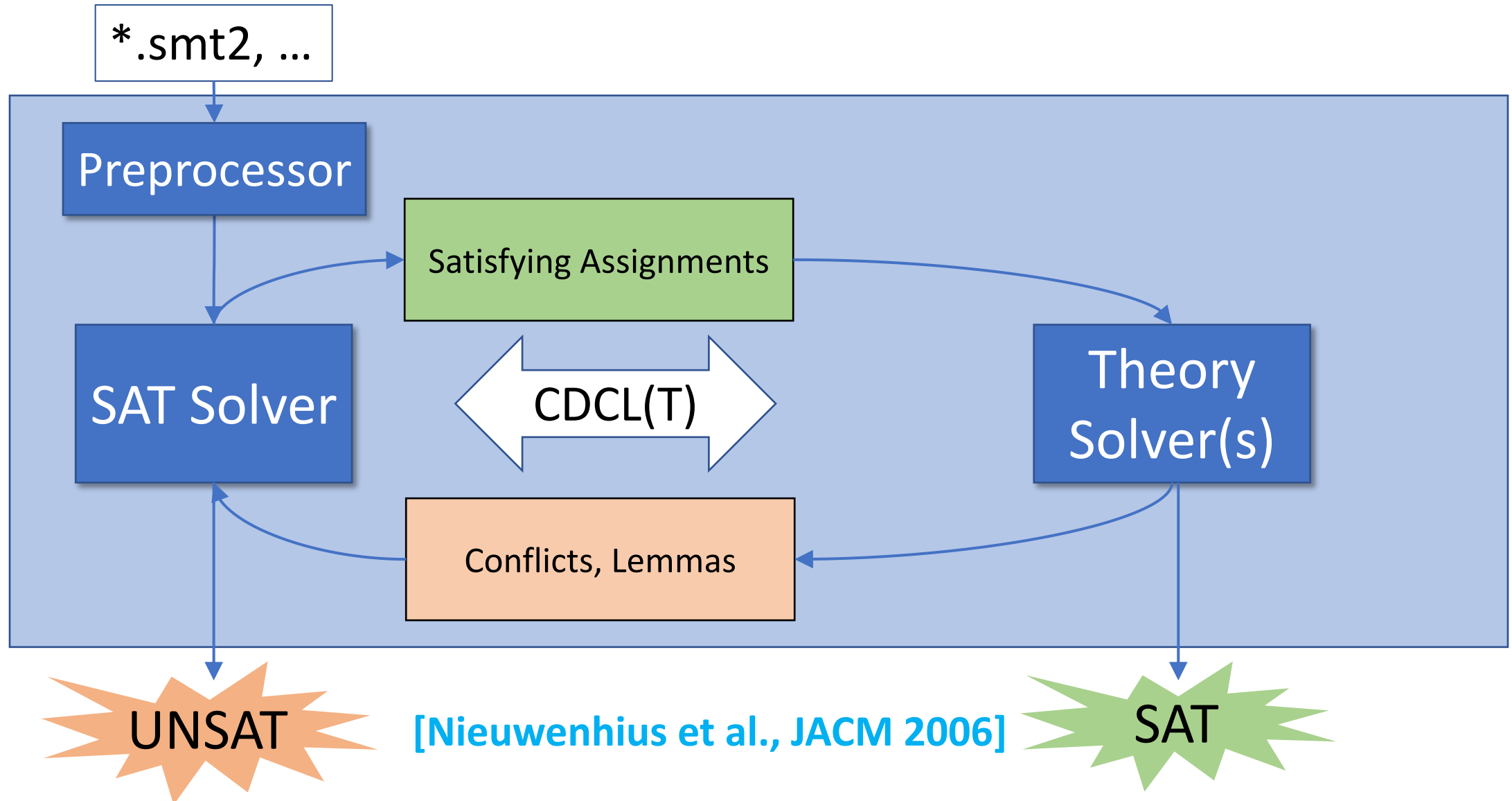
- Part of general effort to make cvc5 fully proof producing
[\[Barbosa et al., CACM'23\]](#)
- Covers great majority of the system
- Several proof granularity levels
- Evaluated on many SMT-LIB theories, including strings
[\[Barbosa et al., IJCAR'22\]](#)
- Fine-grained proofs for rewrites, for strings
[\[Noetzli et al., FMCAD'22\]](#)

Recent Developments for Theory of Strings

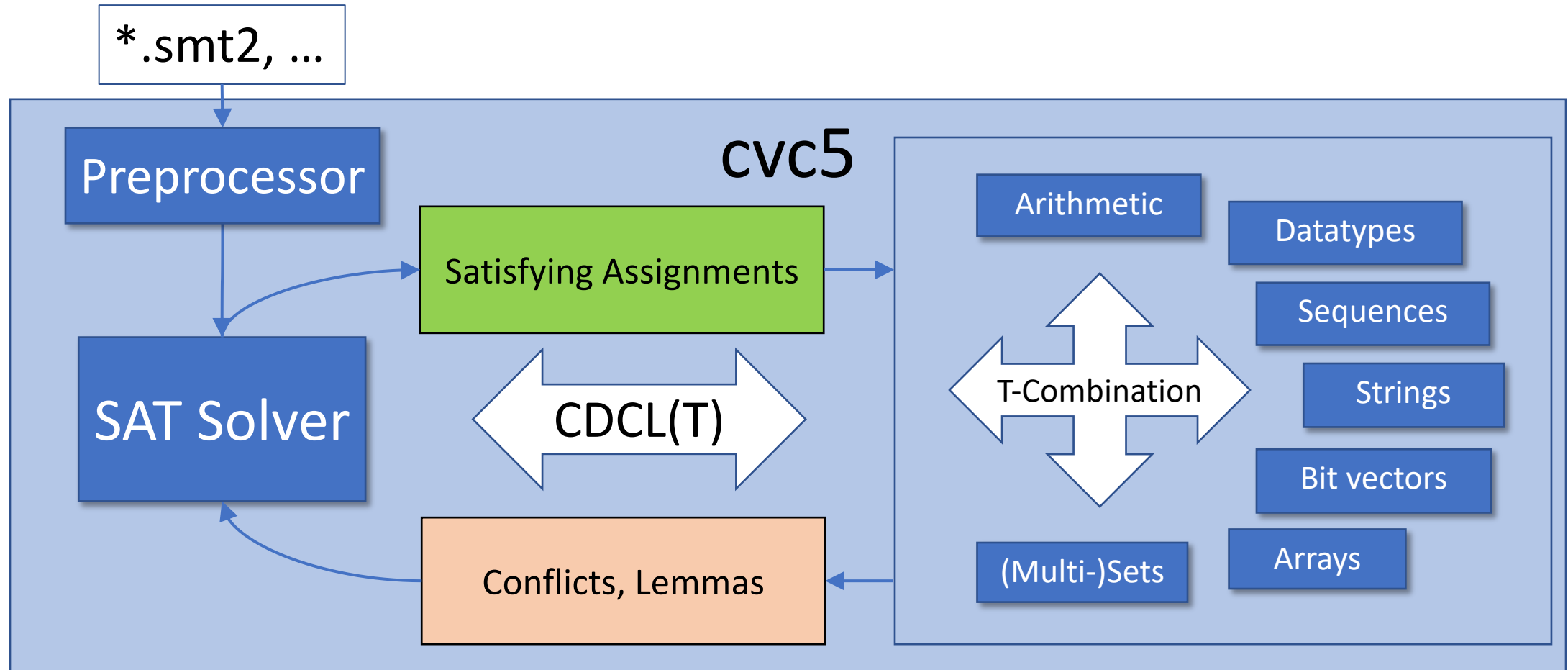
- Context-dependent simplifications
 - Use aggressive rewriting [Reynolds et al., CAV 2019]
 - Applied eagerly [Noetzli et al., CAV 2022]
- Reduction lemmas
 - Leverage string-to-code point (`code`) conversion [Reynolds et al., IJCAR 2020]
 - Improved encodings [Reynolds et al., FMCAD 2020]
 - Applied lazily based on model [Noetzli et al., CAV 2022]

SMT Solvers Architecture

Architecture of most SMT solvers

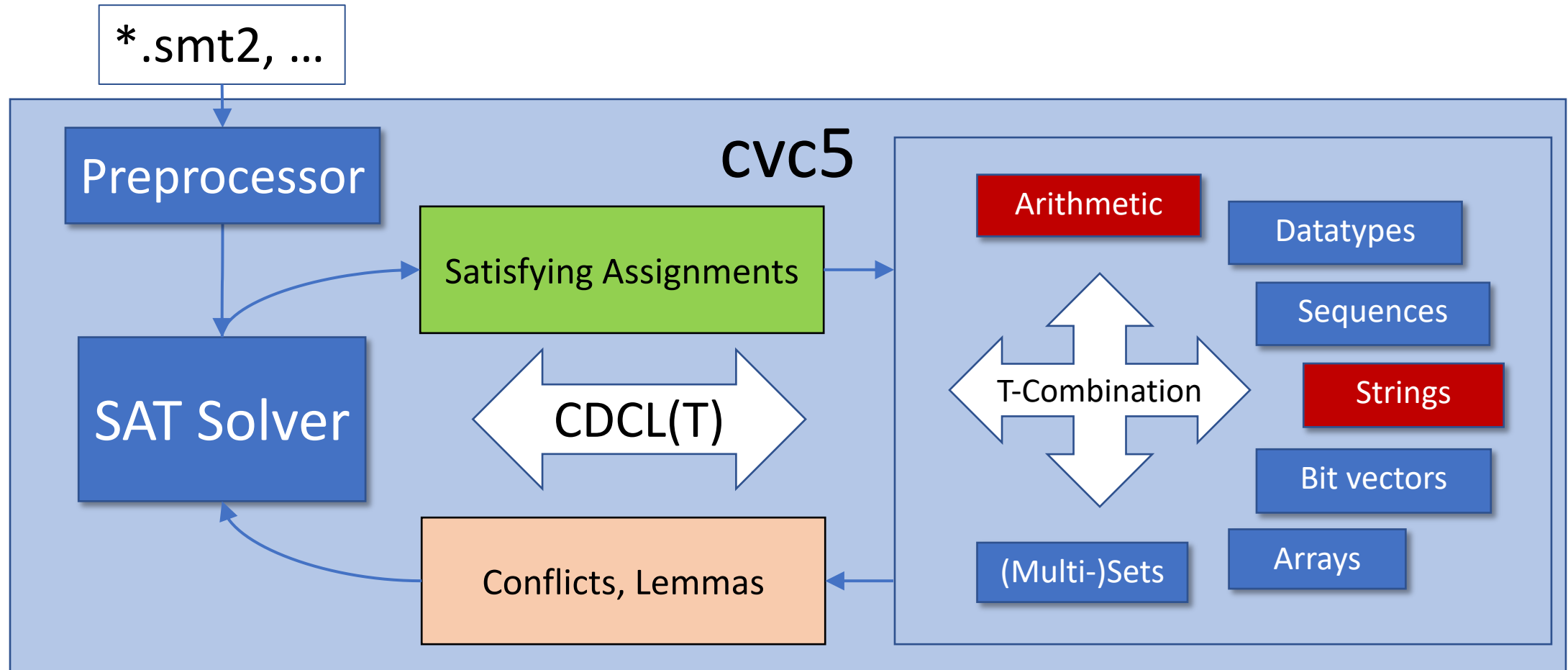


Architecture of cvc5



Centralized methods (Nelson-Oppen, polite) for combining theories

Architecture of cvc5



Focus of this talk: **theory of strings and regular expressions**

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

Sorts:

- Integers Int
- Strings String , interpreted as Σ^* for finite alphabet Σ

Terms:

String variables: x, y, z, u, w

Integer variables: i, j, k

String constants: $\varepsilon, \mathbf{abc}, \mathbf{AcBAA}, \mathbf{http}$

String concatenation: $x \cdot \mathbf{abc}, x \cdot y \cdot z \cdot w$

String length: $|x|$

Formulas:

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

Example:

$$x \cdot \mathbf{a} = y, \quad y \neq \mathbf{b} \cdot z, \quad |y| > |x| + 2$$

Although decidability is **unknown**, many problems can be solved **efficiently in practice**

CDCL(T) String Solvers

Cooperation between:



SAT
Solver

Arithmetic
Solver

String
Solver

CDCL(T) String Solvers

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

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

SAT
Solver

Arithmetic
Solver

String
Solver

CDCL(T) String Solvers

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

SAT
Solver

Arithmetic
Solver

String
Solver

UNSAT

Either determines no satisfying assignments for input exist ...

CDCL(T) String Solvers

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

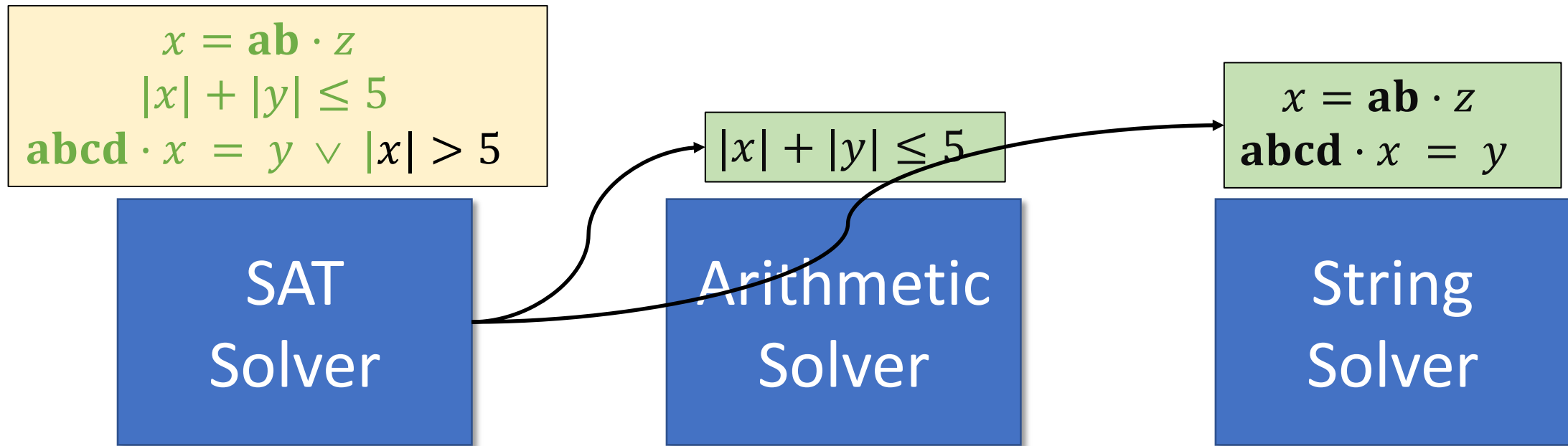
SAT
Solver

Arithmetic
Solver

String
Solver

... or returns a propositionally satisfying assignment

CDCL(T) String Solvers

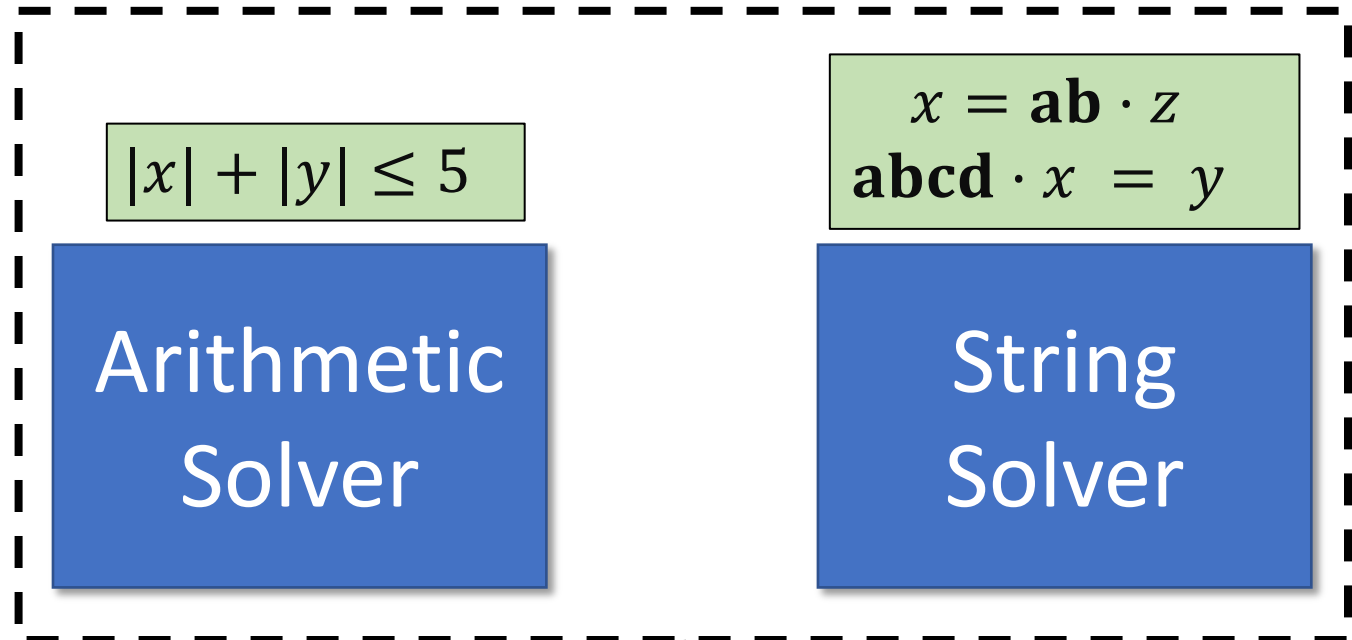


⇒ Constraints distributed to arithmetic and string solvers

CDCL(T) String Solvers

$$\begin{aligned} x &= \mathbf{ab} \cdot z \\ |x| + |y| &\leq 5 \\ \mathbf{abcd} \cdot x &= y \vee |x| > 5 \end{aligned}$$

SAT
Solver



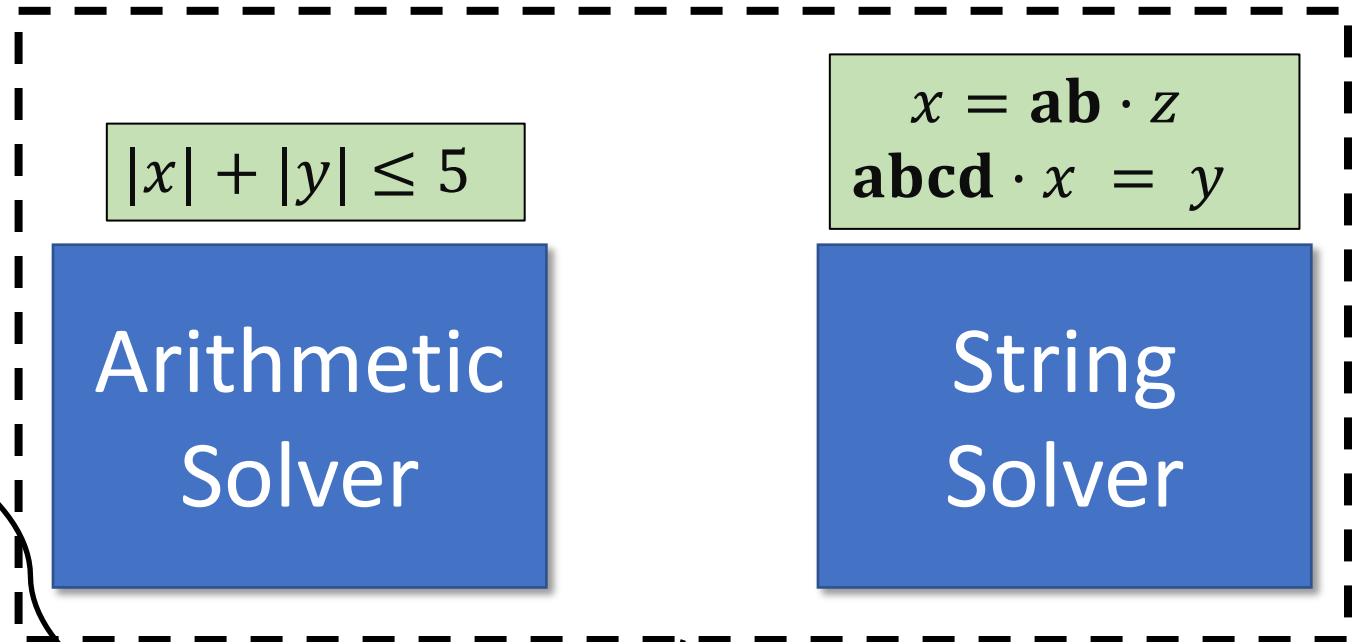
Either find constraints are T_{SLIA} -satisfiable ...



CDCL(T) String Solvers

$$\begin{aligned} x &= \mathbf{ab} \cdot z \\ |x| + |y| &\leq 5 \\ \mathbf{abcd} \cdot x &= y \vee |x| > 5 \\ x \neq \mathbf{ab} \cdot z &\vee |x| = |z| + 2 \end{aligned}$$

SAT Solver



... or return *theory lemmas*

(valid T_{LIA}/T_S -formulas) to SAT solver ...

$$x = \mathbf{ab} \cdot z \Rightarrow |x| = |z| + 2$$

CDCL(T) String Solvers

$x = \mathbf{ab} \cdot z$
 $|x| + |y| \leq 5$
 $\mathbf{abcd} \cdot x = y \vee |x| > 5$
 $x \neq \mathbf{ab} \cdot z \vee |x| = |z| + 2$

SAT
Solver

$|x| + |y| \leq 5$
 $|x| = |z| + 2$

Arithmetic
Solver

$x = \mathbf{ab} \cdot z$
 $\mathbf{abcd} \cdot x = y$

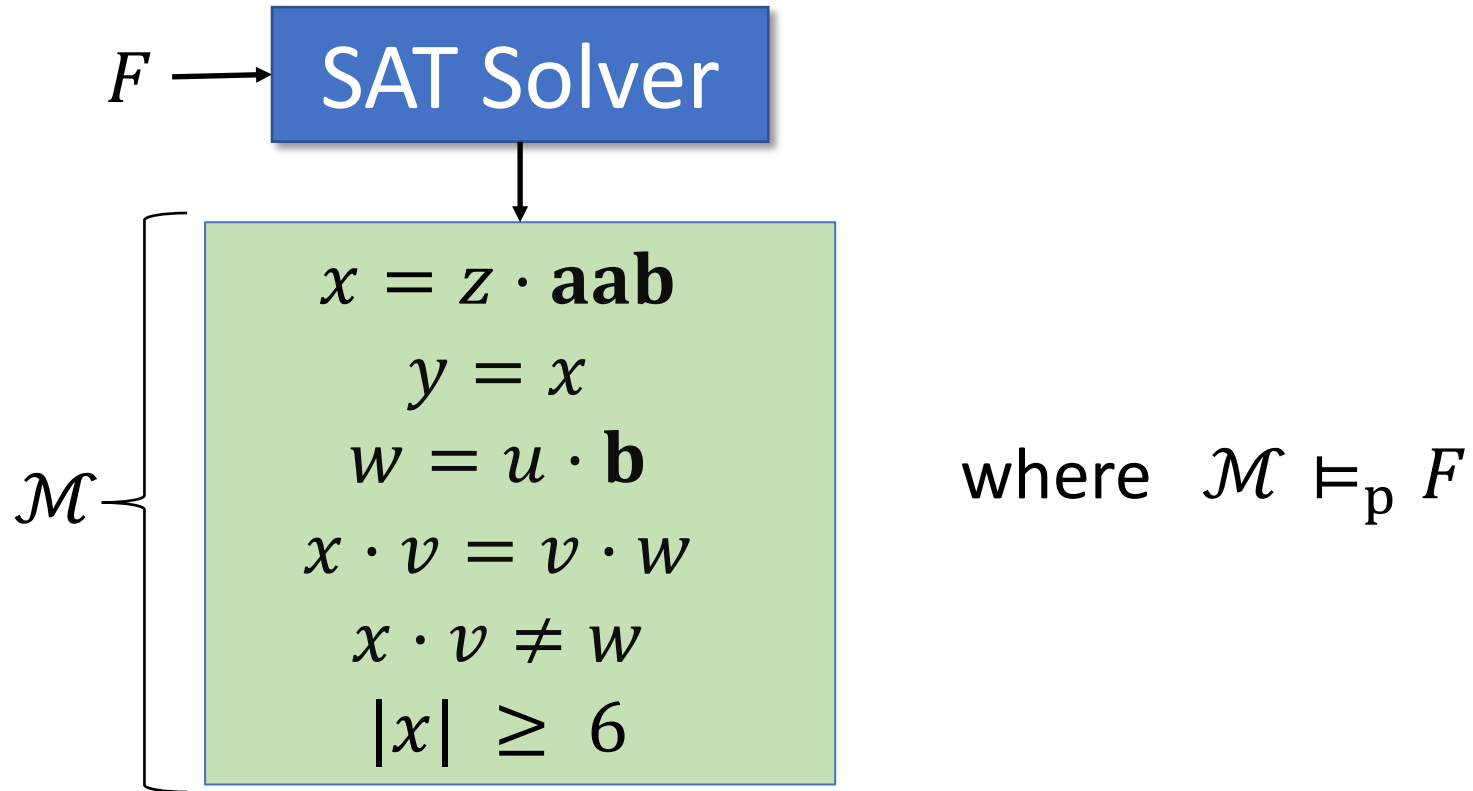
String
Solver

... and repeat

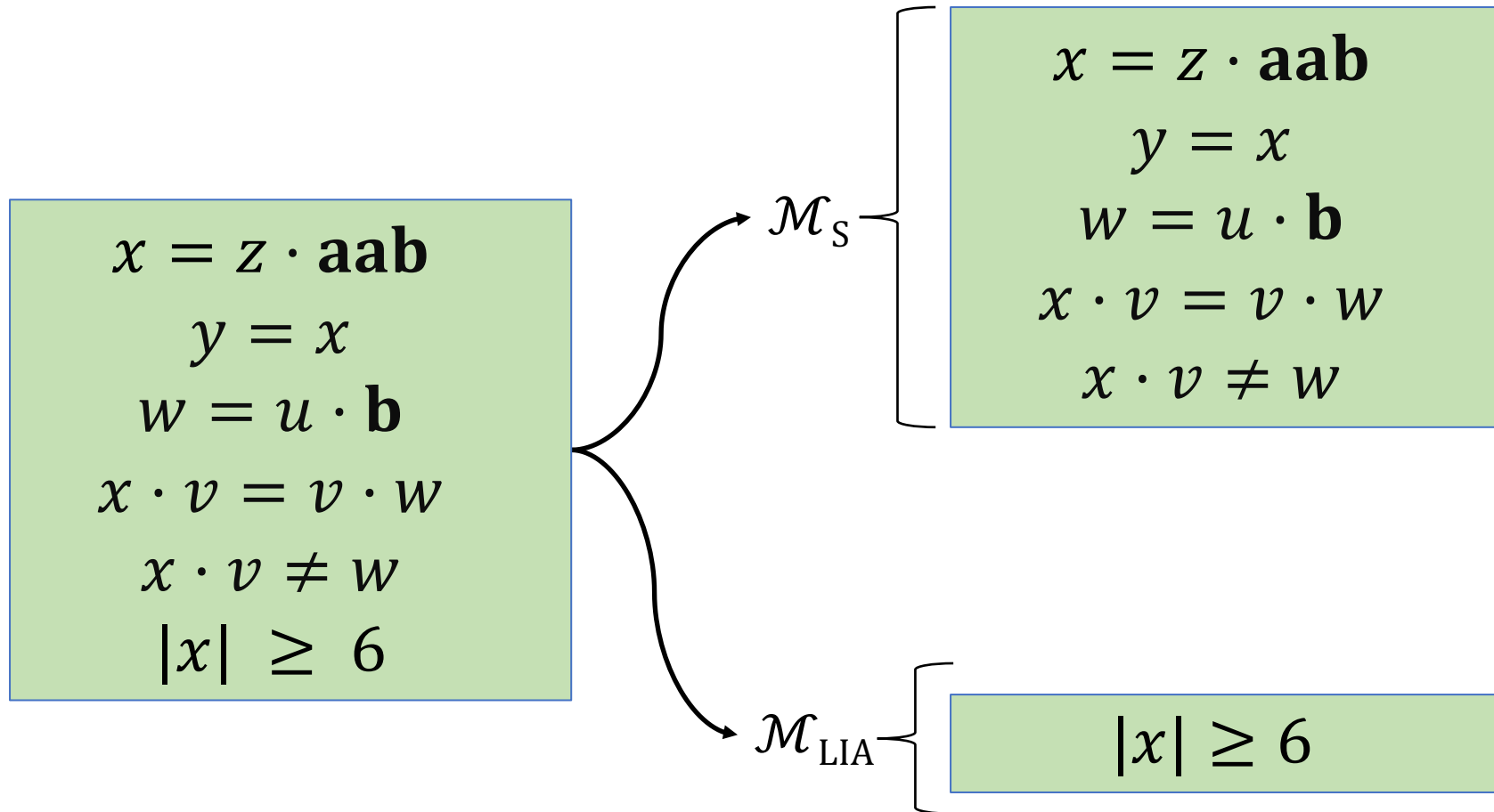
A Theory Solver for Strings

[Liang, Reynolds, Deters, Tinelli and Barrett, CAV 14]

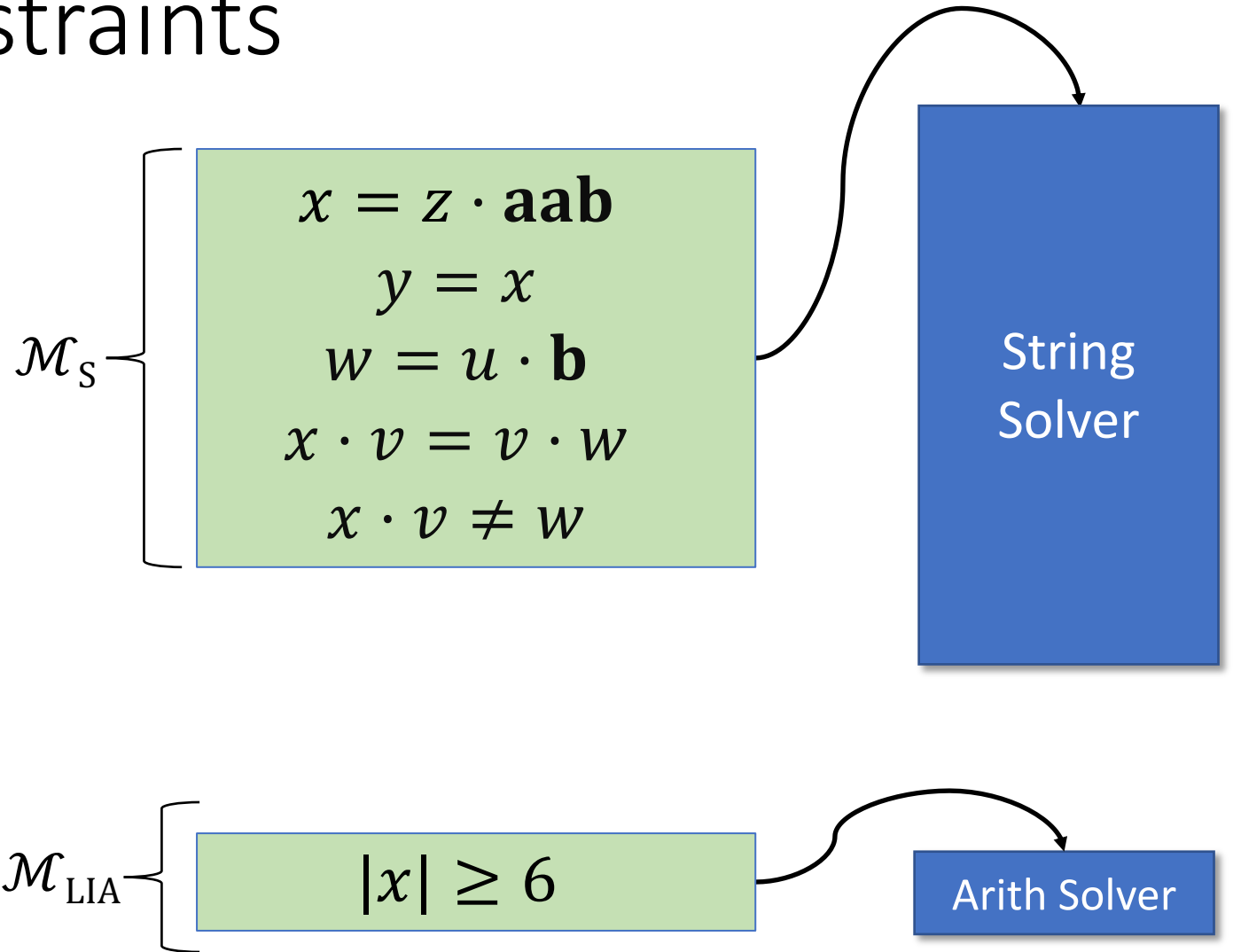
Solving String Constraints



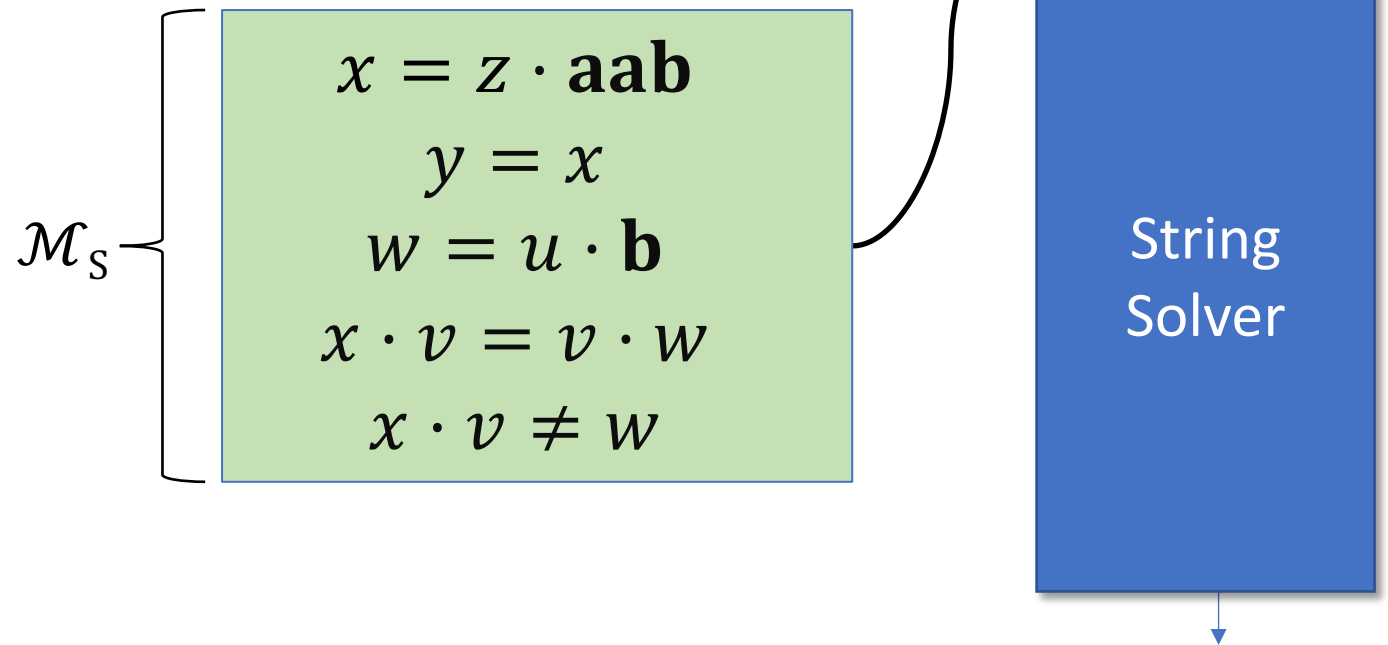
Solving String Constraints



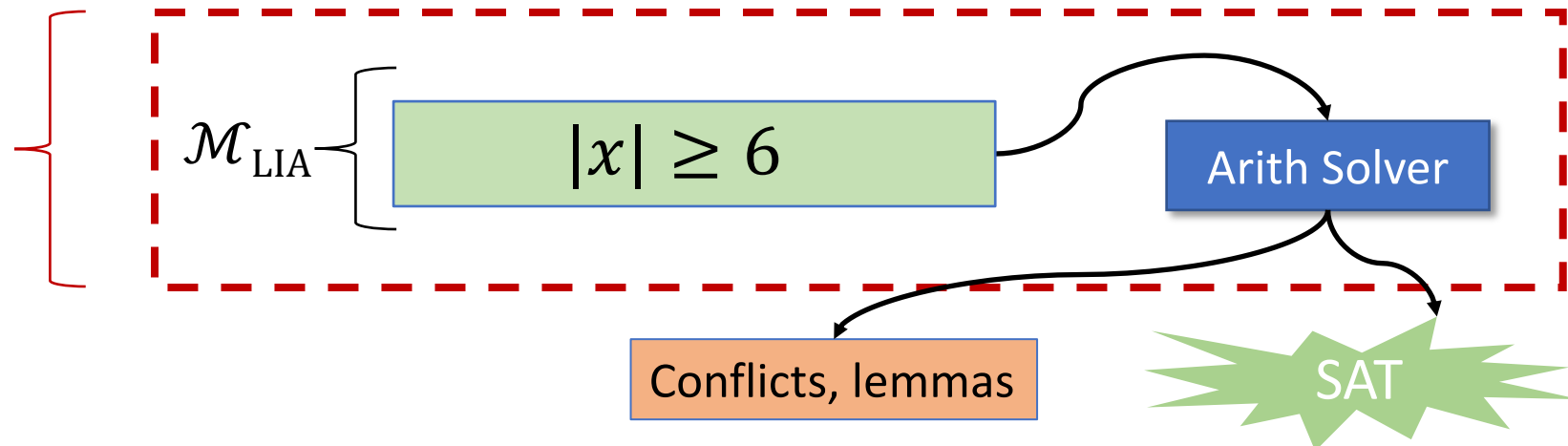
Solving String Constraints



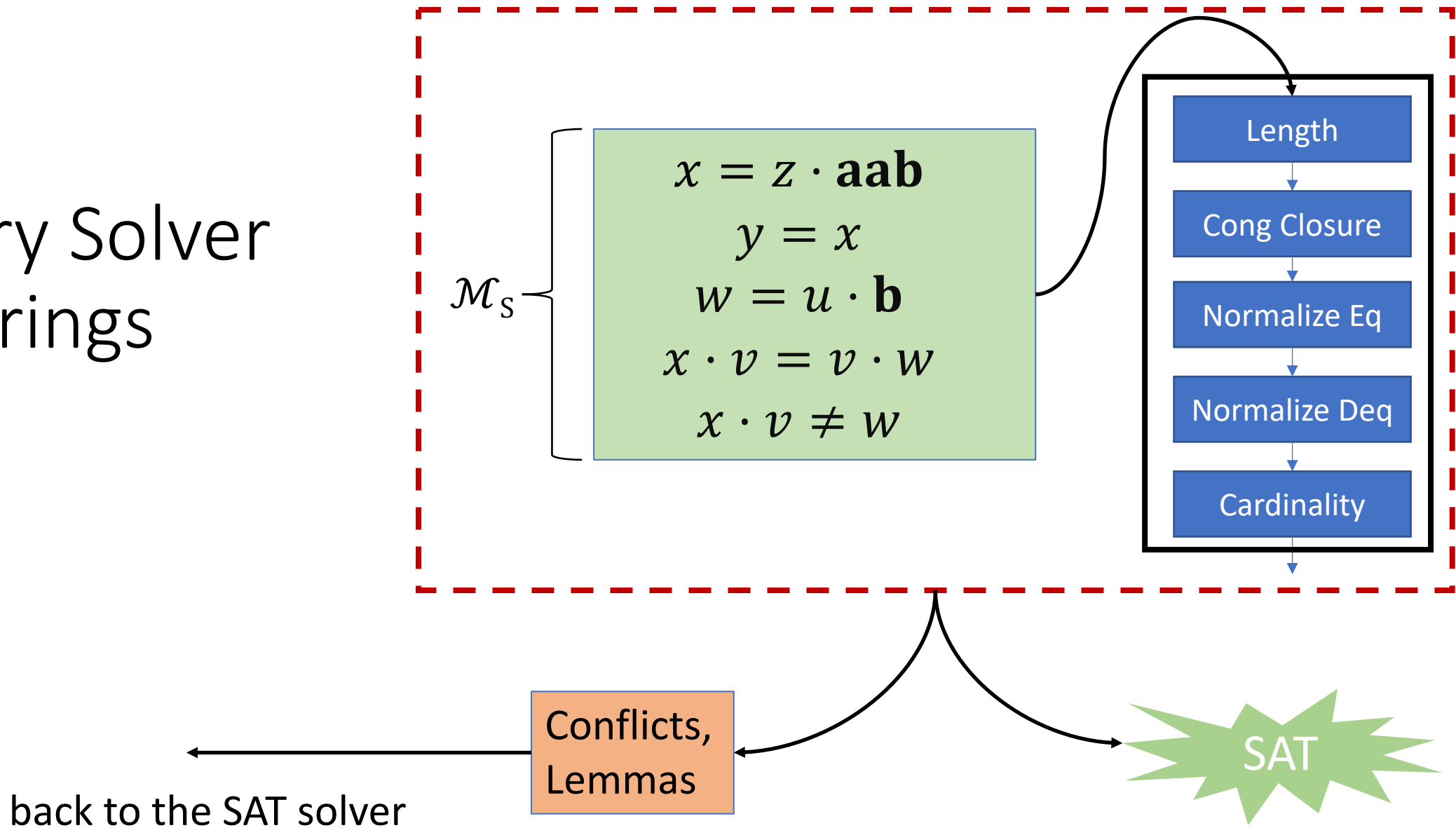
Solving String Constraints



Theory Solver for Linear Integer Arithmetic (Simplex)



Theory Solver for Strings

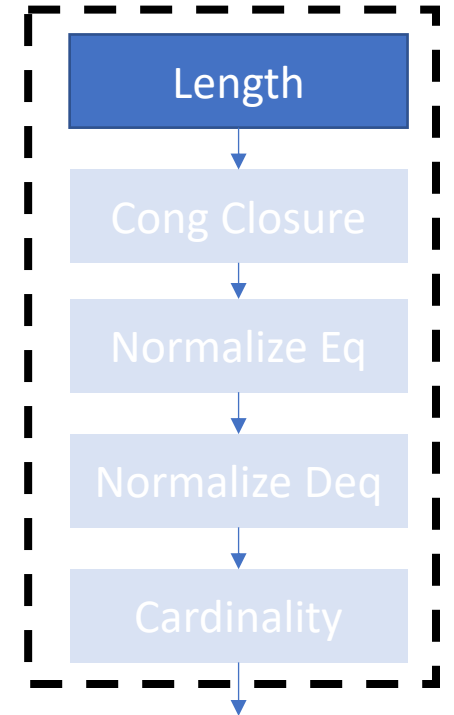


String Theory Solver Inference Strategy

1. Elaborate **length** constraints
 2. Check for **equality conflicts** (compute congruence closure)
 3. Normalize string **equalities**
 4. Normalize string **disequalities**
 5. Check **cardinality** constraints
- Each step may add lemma or a conflict
 - If no step adds a lemma or conflict, the current constraint set $(\mathcal{M}_S \cup \mathcal{M}_S)$ is sat

1. Elaborate Length Constraints

$$\mathcal{M}_S \left\{ \begin{array}{l} x = z \cdot \mathbf{aab} \\ y = x \\ w = u \cdot \mathbf{b} \\ x \cdot v = v \cdot w \\ x \cdot v \neq w \end{array} \right.$$



1. Elaborate Length Constraints

$$\mathcal{M}_s \left\{ \begin{array}{l} x = z \cdot \mathbf{aab} \\ y = x \\ w = u \cdot \mathbf{b} \\ x \cdot v = v \cdot w \\ x \cdot v \neq w \end{array} \right.$$

- For each term of type string in \mathcal{M}_s ,
add lemma providing the **definition of its length**:

$$|\mathbf{b}| = 1$$

$$|\mathbf{aab}| = 3$$

$$|x \cdot v| = |x| + |v|$$

$$|z \cdot \mathbf{aab}| = |z| + 3$$

$$|u \cdot \mathbf{b}| = |u| + 3$$

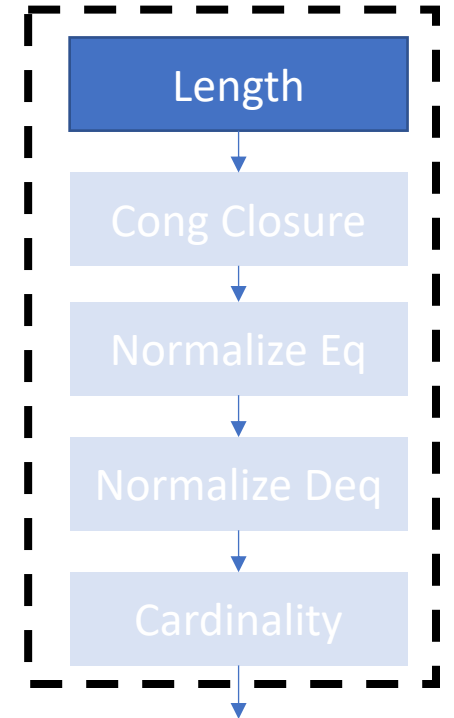
$$|v \cdot w| = |v| + |w|$$

- For each variable of type string in \mathcal{M}_s ,
add an **emptiness splitting lemma**:

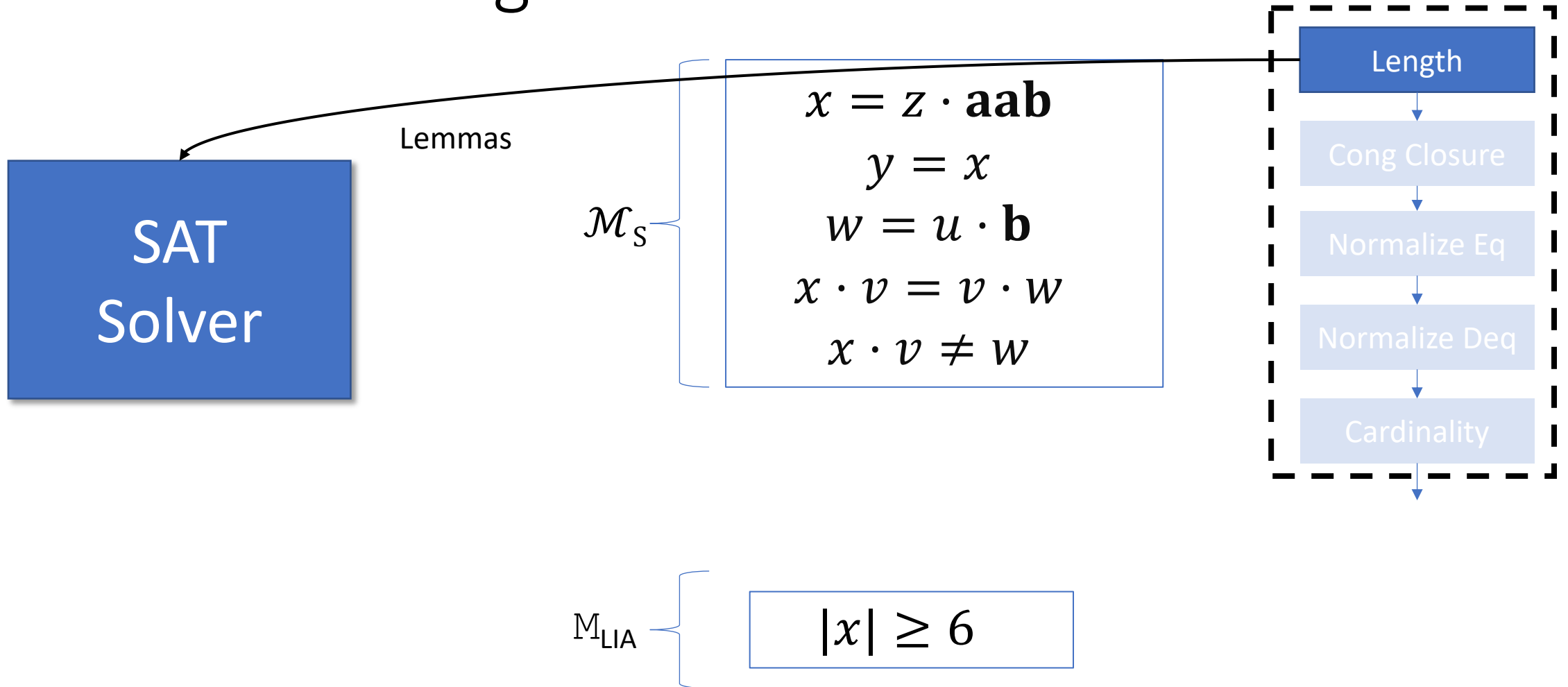
$$x = \epsilon \vee |x| \geq 1$$

$$y = \epsilon \vee |y| \geq 1$$

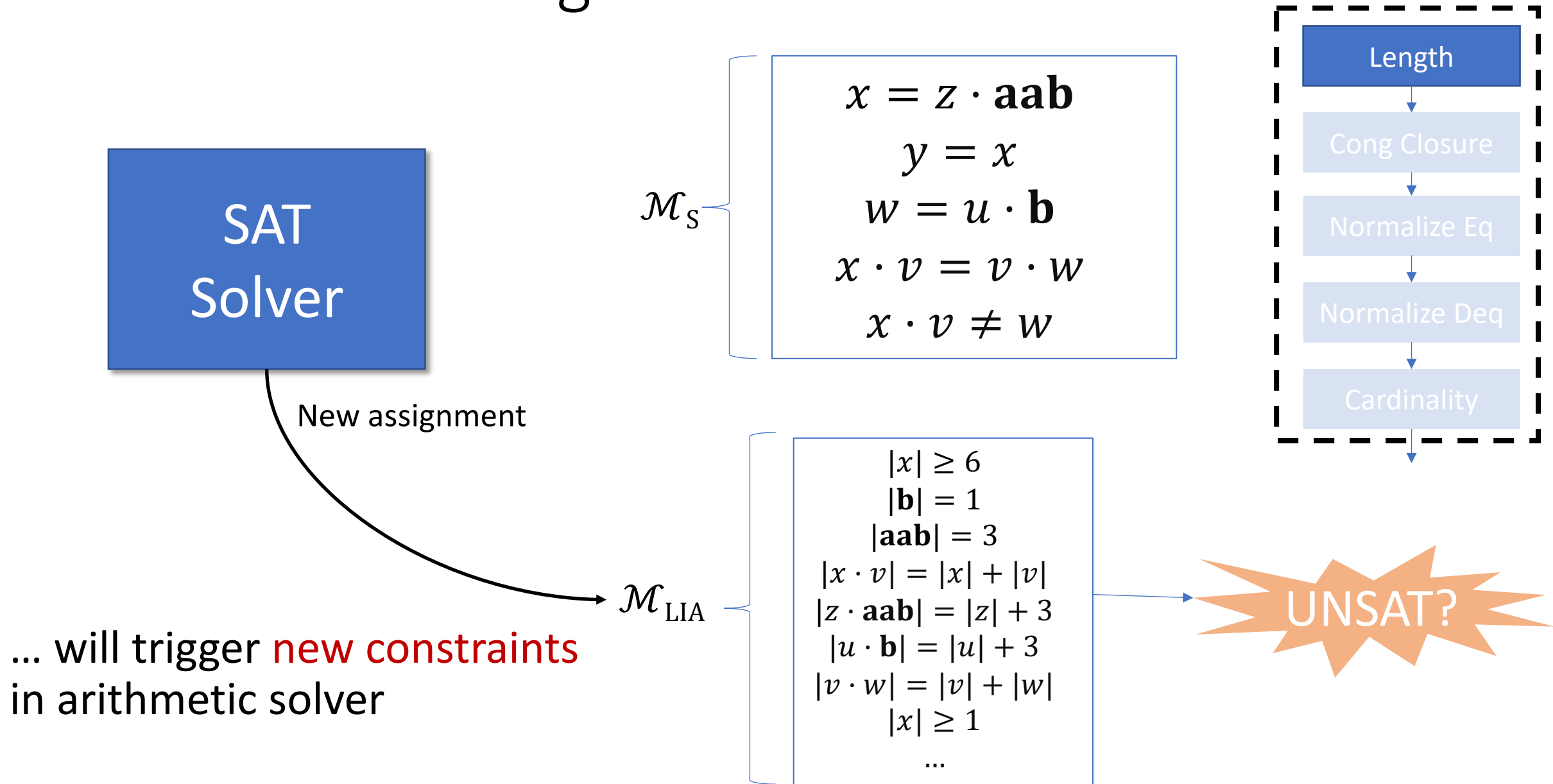
...



1. Elaborate Length Constraints



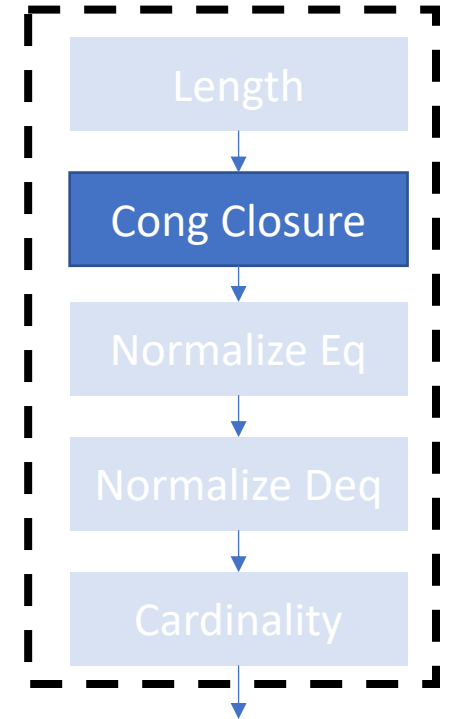
1. Elaborate Length Constraints



... will trigger **new constraints** in arithmetic solver

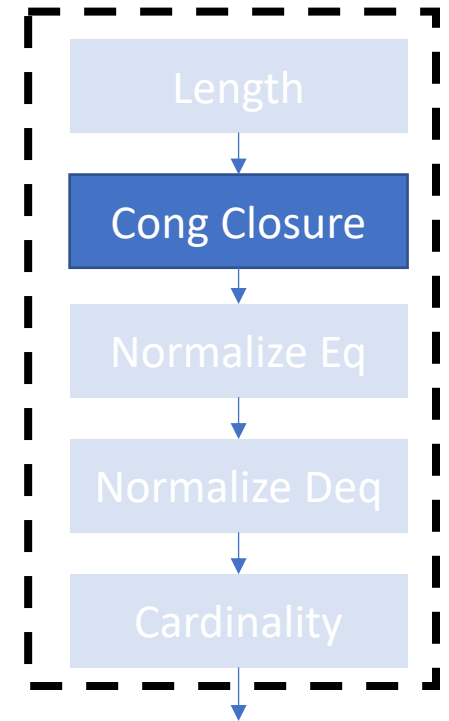
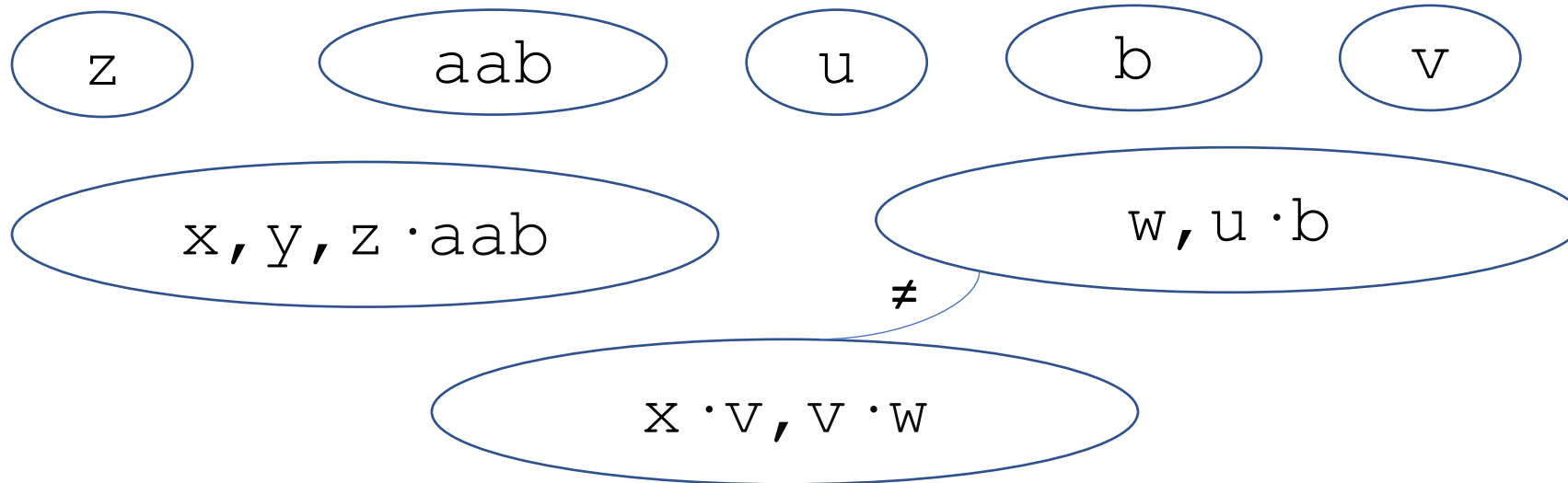
2. Compute Congruence Closure

$$\mathcal{M}_S \left\{ \begin{array}{l} x = z \cdot \mathbf{aab} \\ y = x \\ w = u \cdot \mathbf{b} \\ x \cdot v = v \cdot w \\ x \cdot v \neq w \end{array} \right.$$

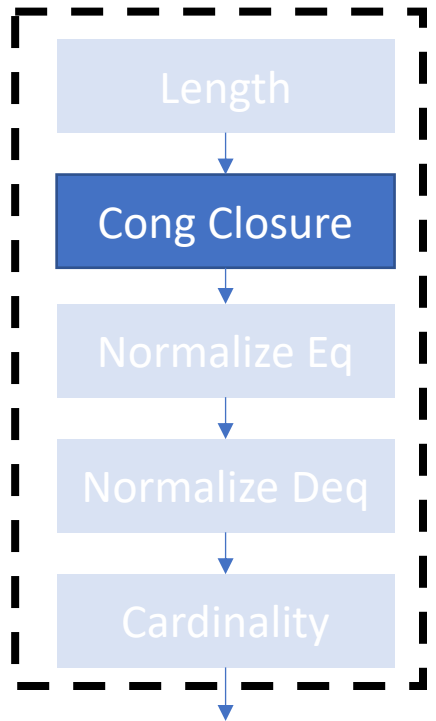


$$\mathcal{M}_S \left\{ \begin{array}{l} x = z \cdot \mathbf{aab} \\ y = x \\ w = u \cdot \mathbf{b} \\ x \cdot v = v \cdot w \\ x \cdot v \neq w \end{array} \right.$$

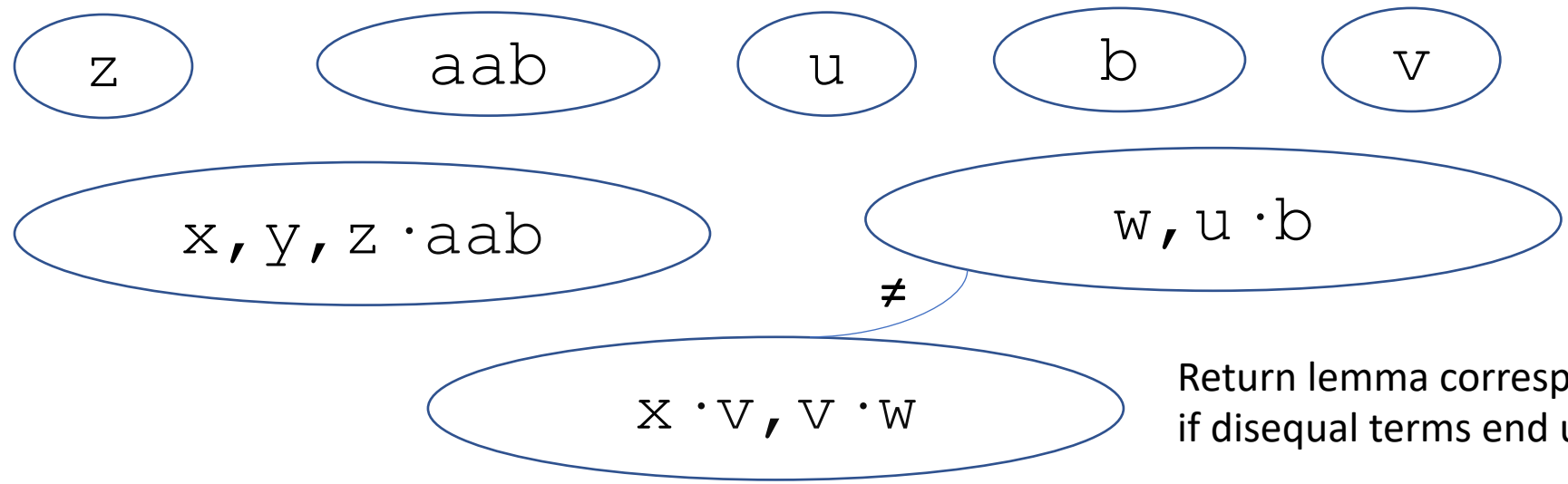
Group terms by **equivalence classes**:



$$\mathcal{M}_S \left\{ \begin{array}{l} x = z \cdot \mathbf{aab} \\ y = x \\ w = u \cdot \mathbf{b} \\ x \cdot v = v \cdot w \\ x \cdot v \neq w \end{array} \right.$$

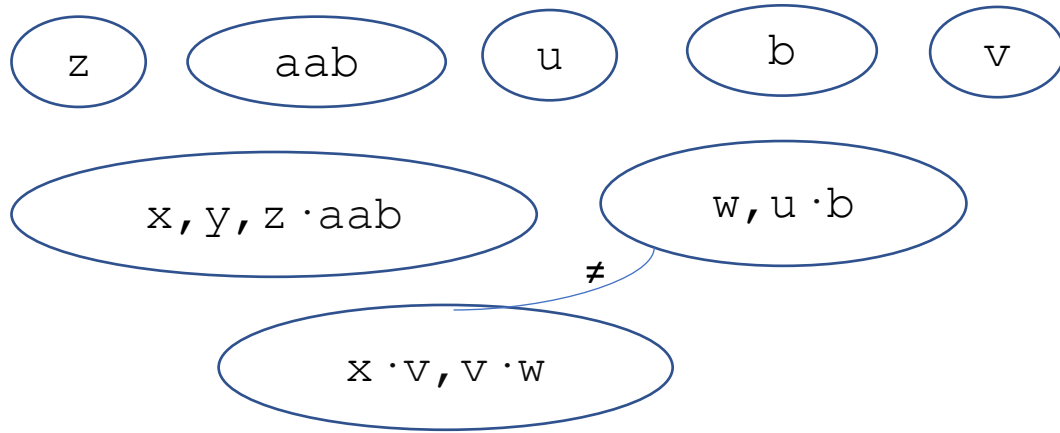


Group terms by equivalence classes:

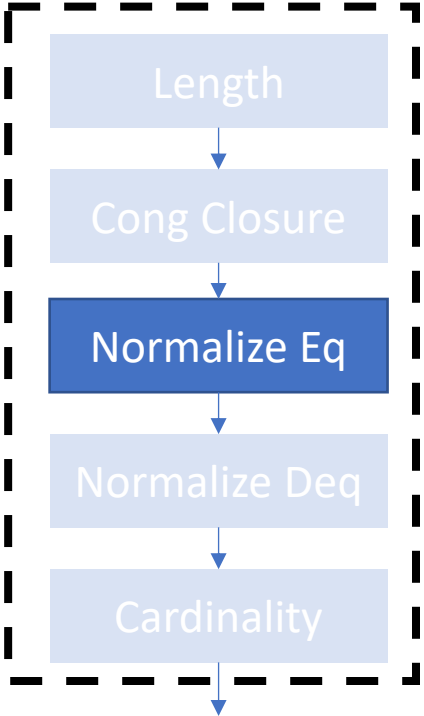


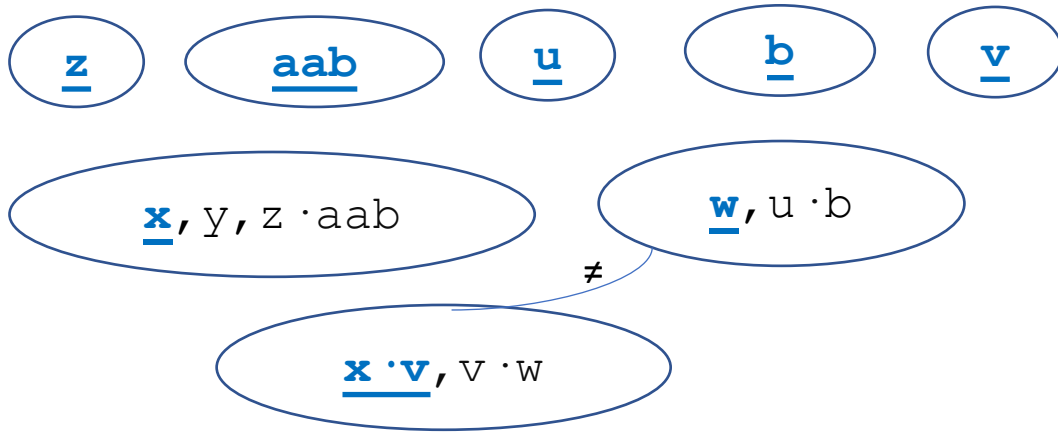
Return lemma corresponding to T_S -conflict if disequal terms end up in the same eq class

3. Normalize Equalities

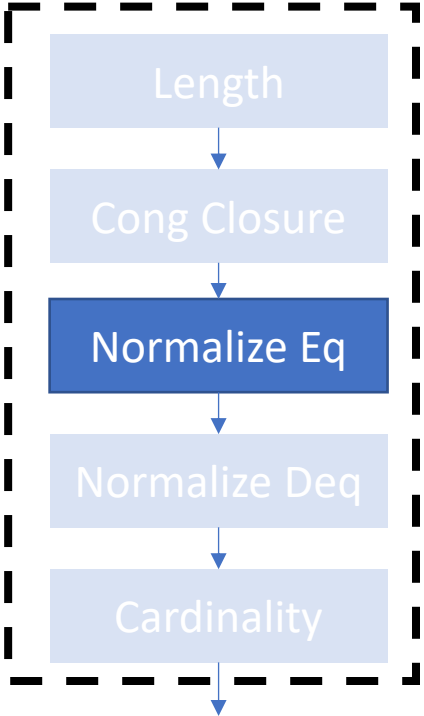


$$\begin{aligned}x &= z \cdot \mathbf{aab} \\y &= x \\w &= u \cdot \mathbf{b} \\x \cdot v &= v \cdot w \\x \cdot v &\neq w\end{aligned}$$



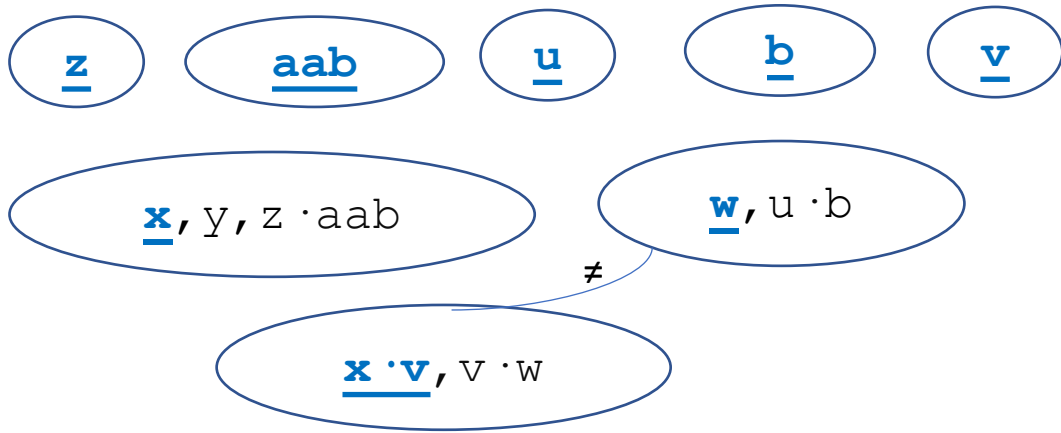


$$\begin{aligned}
 x &= z \cdot aab \\
 y &= x \\
 w &= u \cdot b \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$



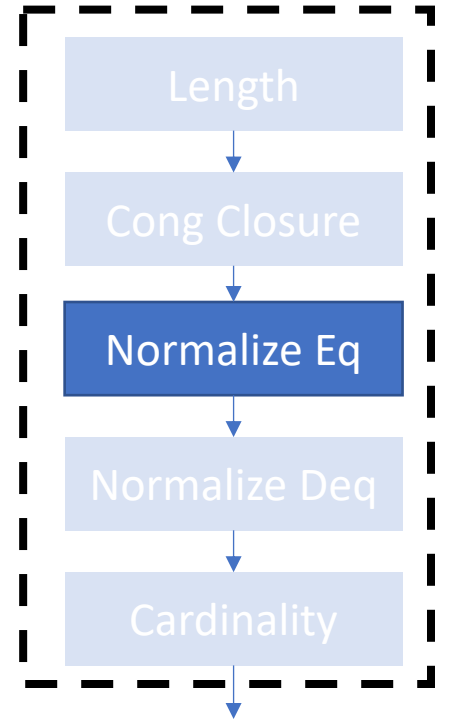
Compute *normal forms* for equivalence classes

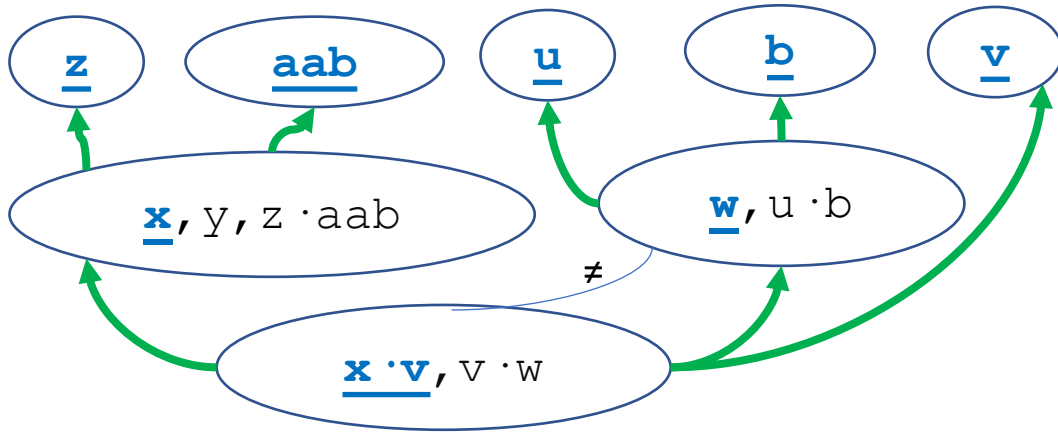
- A normal form is a concatenation of string terms $r_1 \cdot \dots \cdot r_n$ where each r_i is the **representative** of its equivalence class
Restriction: string constants must be chosen as representatives
- An equivalence class can be **assigned** a normal form $r_1 \cdot \dots \cdot r_n$ if:
 Each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



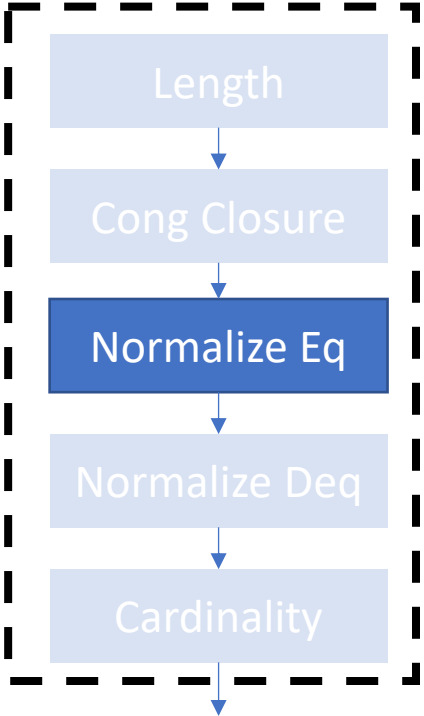
Normal forms computed **bottom-up**

$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$



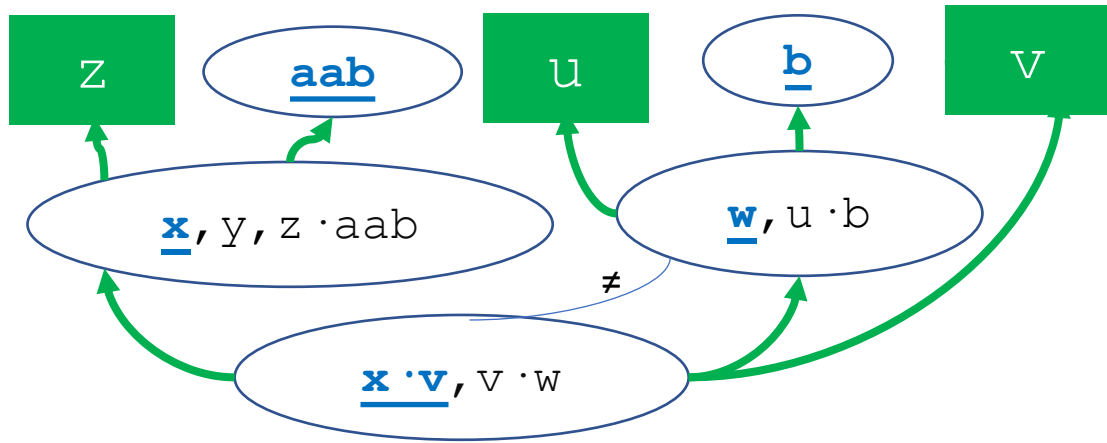


$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$

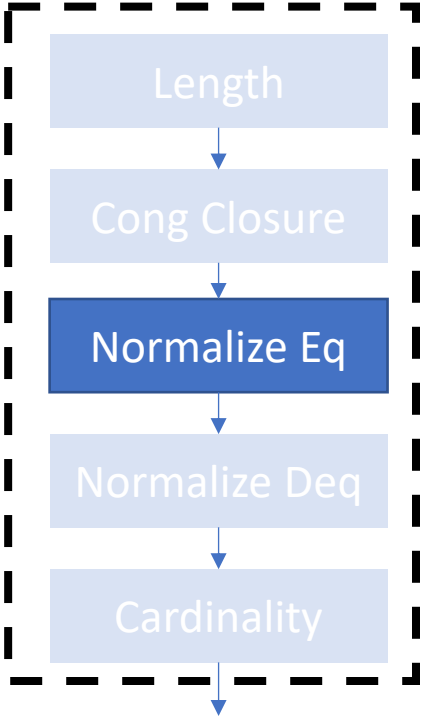


Normal forms computed bottom-up

- First, compute **containment relation** induced by concatenation terms
 This relation is guaranteed to be acyclic due to length elaboration step (cycle \Rightarrow LIA-conflict)

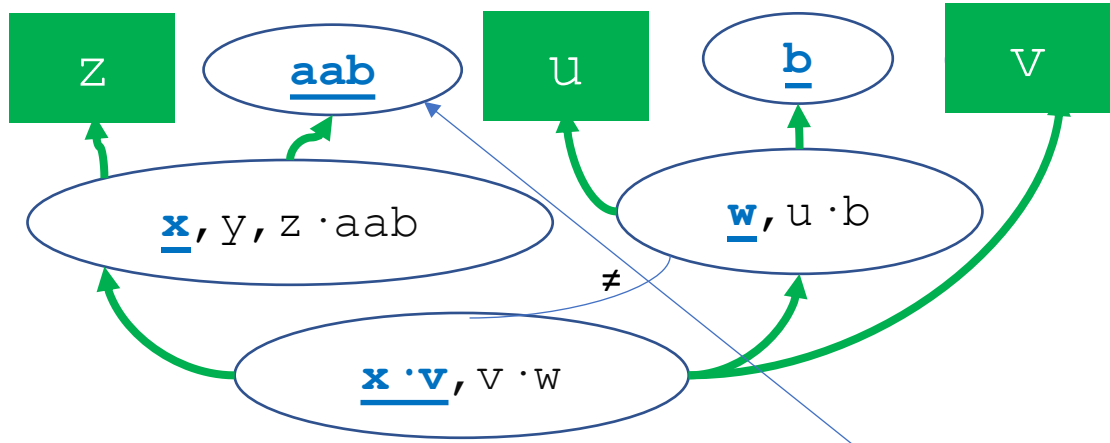


$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$



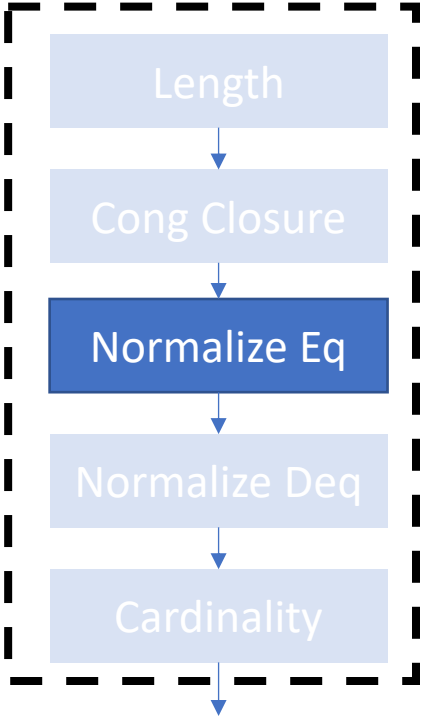
Normal forms computed bottom-up

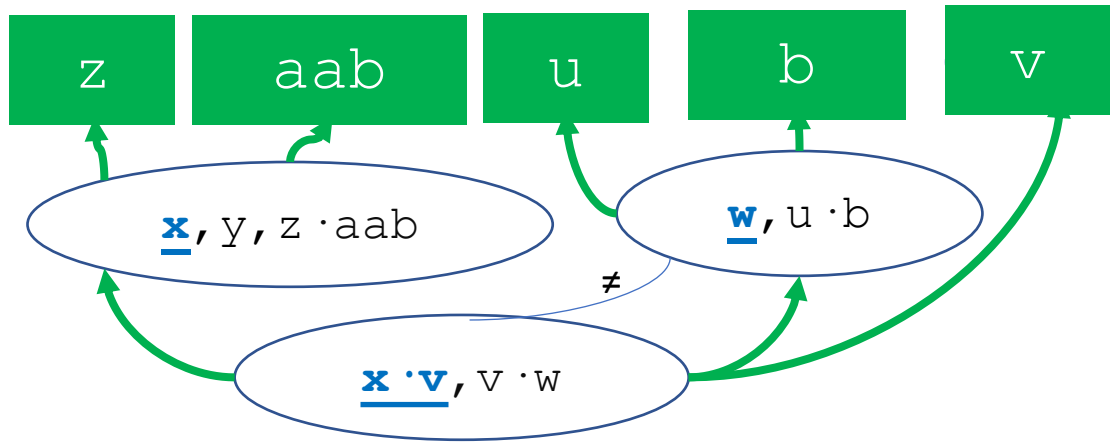
- First, compute containment relation induced by concatenation terms
 - This relation is guaranteed to be acyclic due to length elaboration step (cycle \Rightarrow LIA-conflict)
- **Base** case: eq classes with just variables can be assigned representative as a normal form
- **Inductive** case: compare the *expanded* forms t_1, \dots, t_n of each non-variable
 - If $t_1 \cong \dots \cong t_n$, assign one. If there exists distinct t_i, t_j , then try to equate them



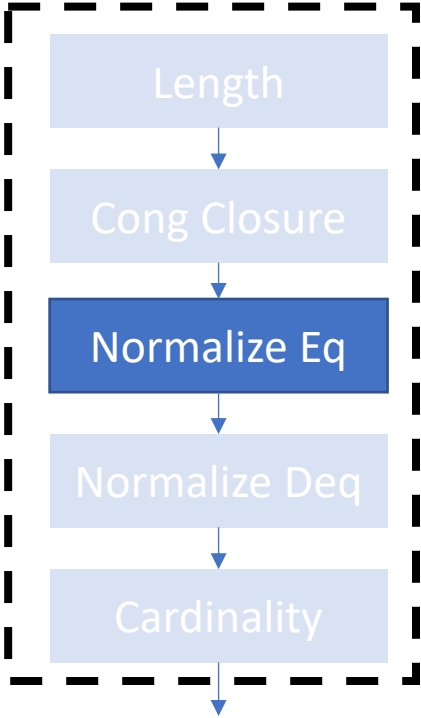
$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$

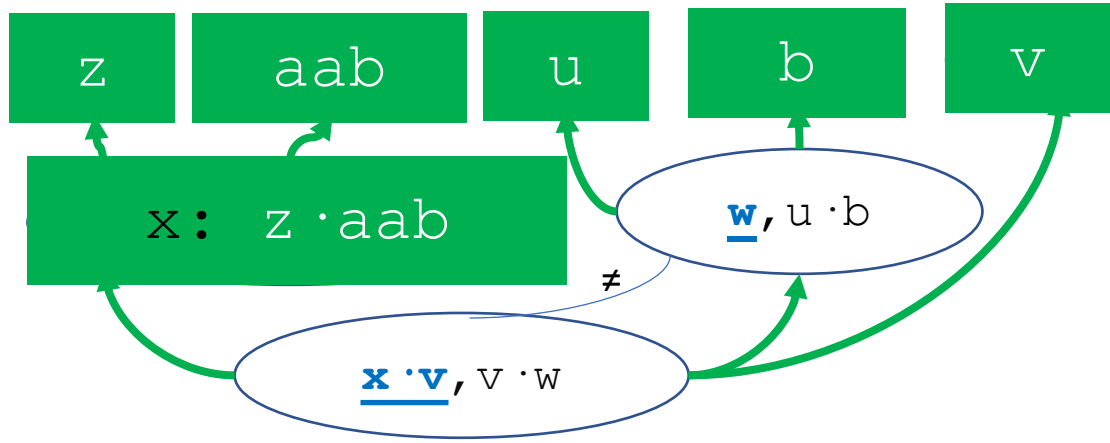
Single non-variable string term \Rightarrow assign



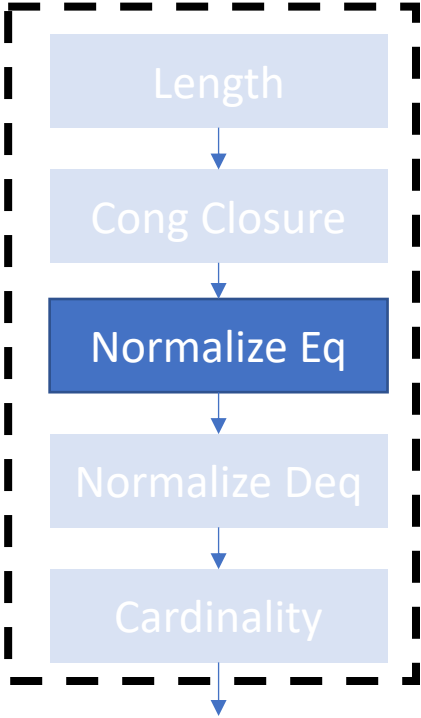


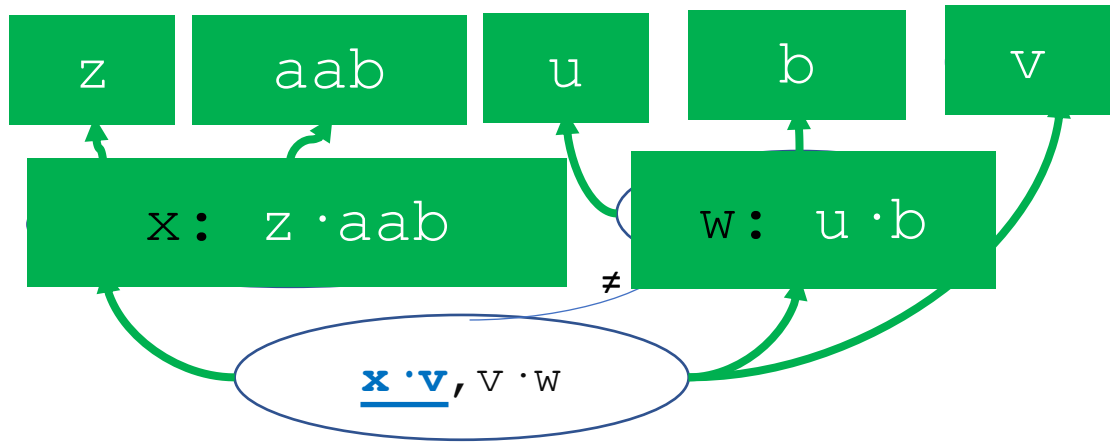
$$\begin{aligned}
 x &= z \cdot aab \\
 y &= x \\
 w &= u \cdot b \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$



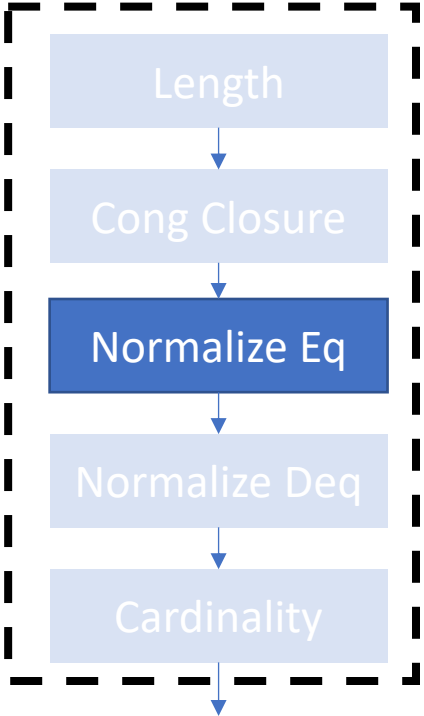


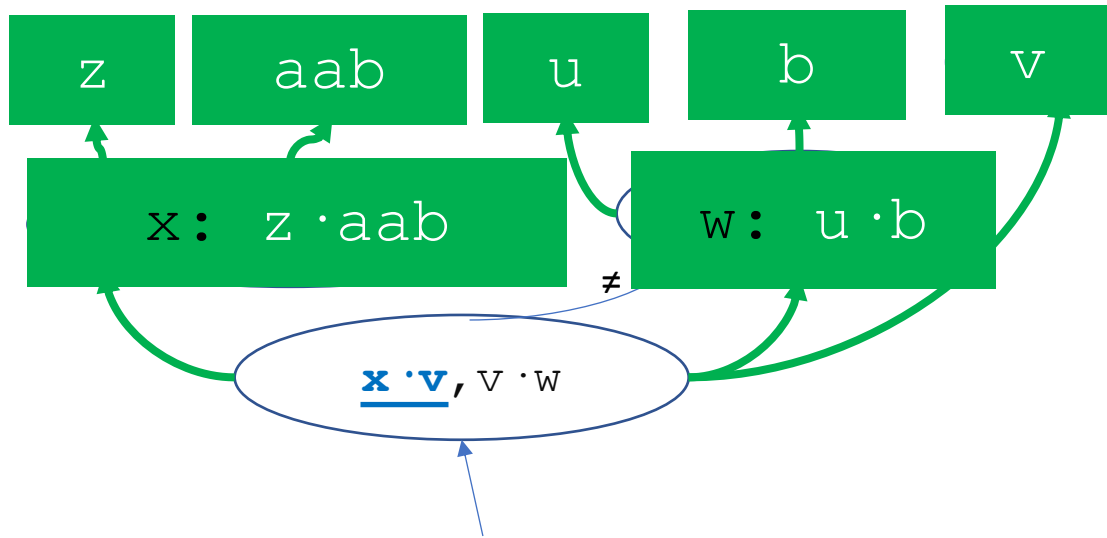
$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$



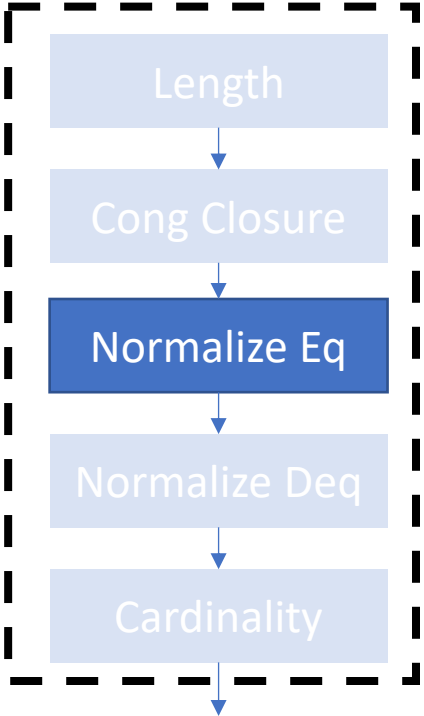


$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$



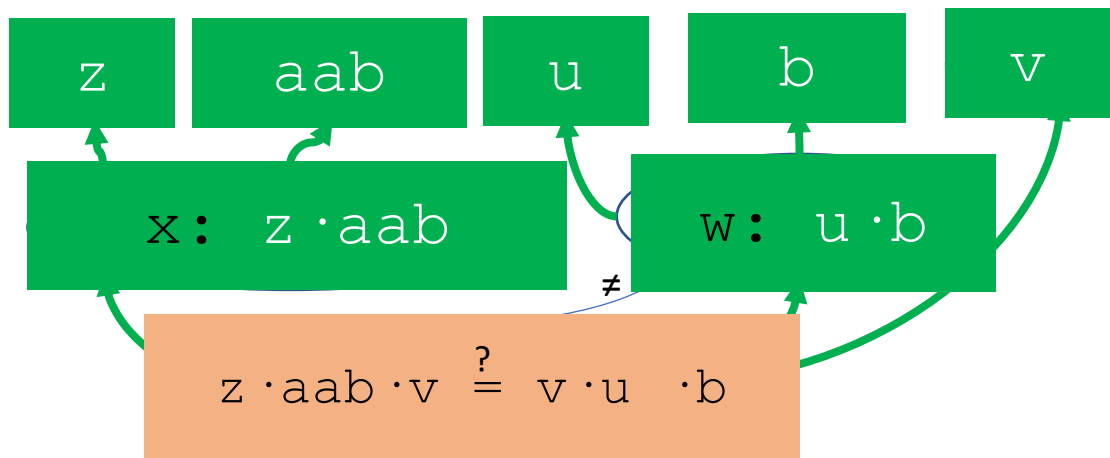


$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$

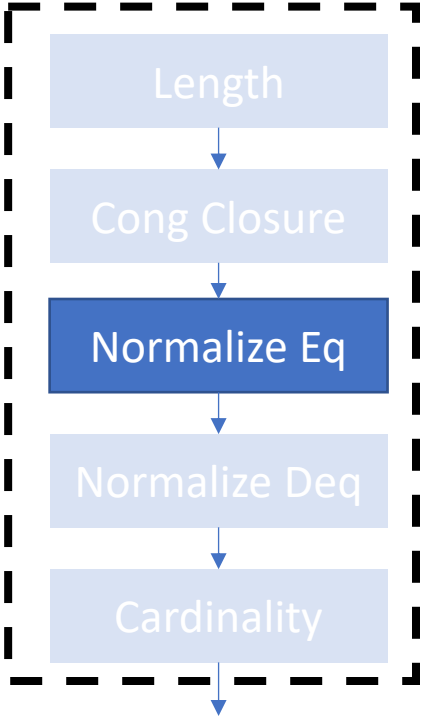


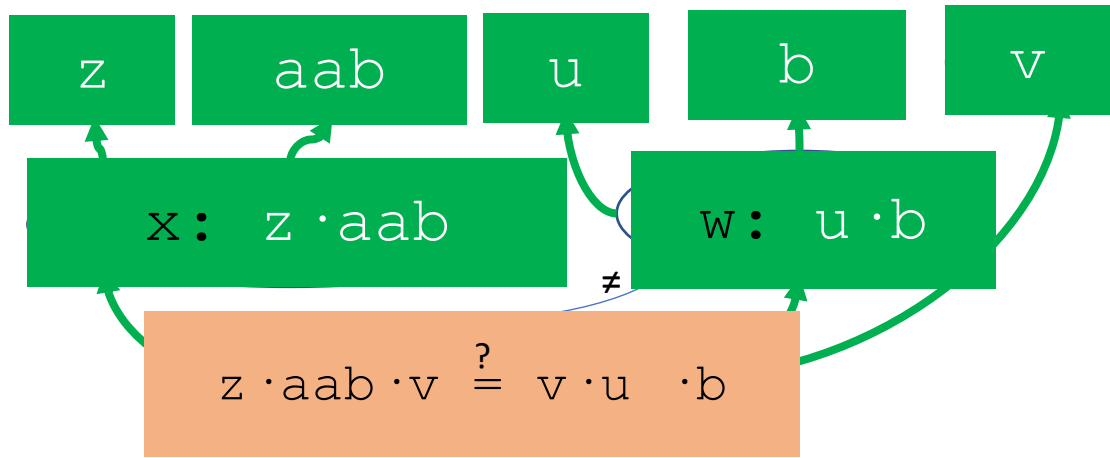
• Equivalence class with two non-variable terms with **distinct** expanded forms:

- $x \cdot v = (z \cdot \mathbf{aab}) \cdot v = \mathbf{z \cdot aab \cdot v}$
- $v \cdot w = v \cdot (u \cdot \mathbf{b}) = \mathbf{v \cdot u \cdot b}$

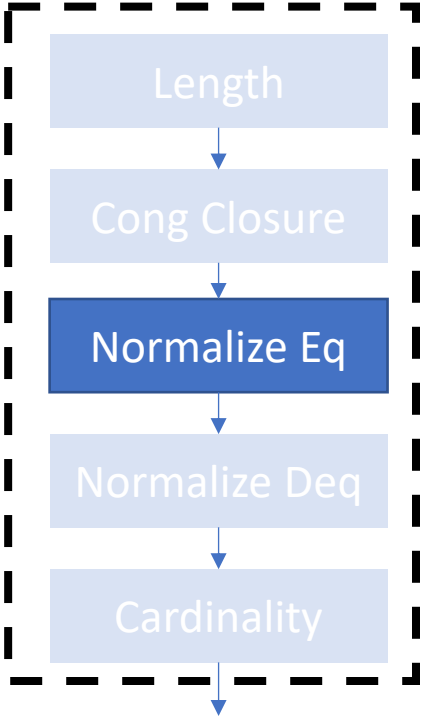


$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$



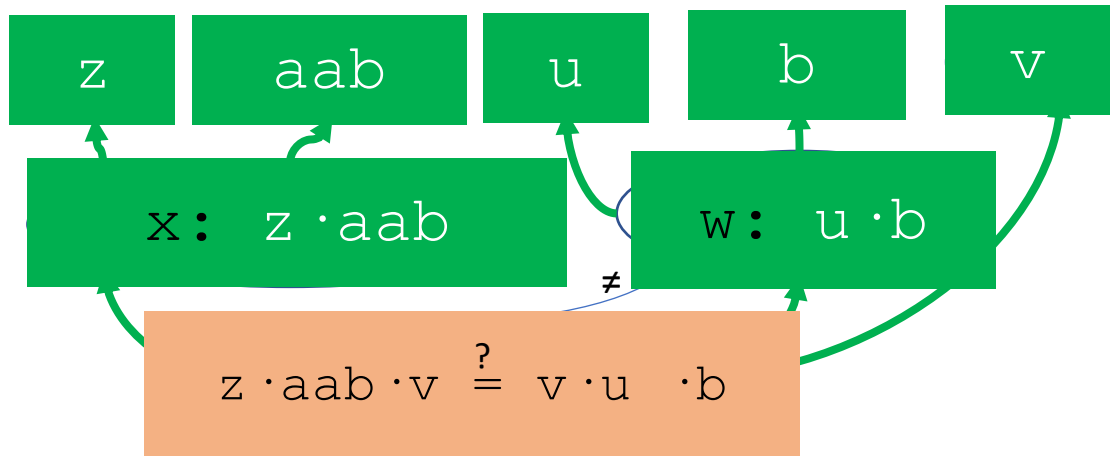


$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$

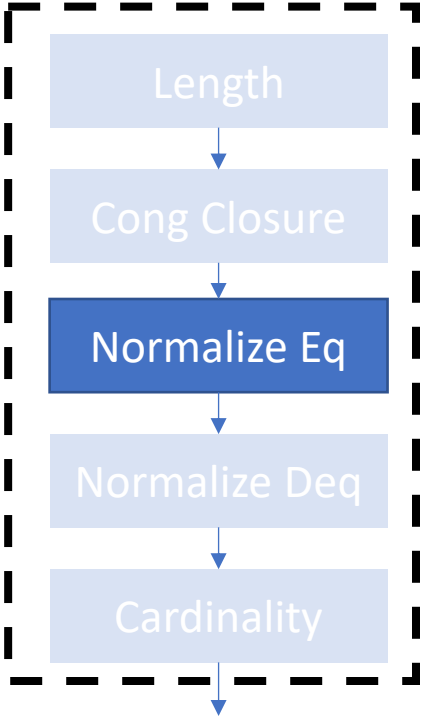


Goal: split strings so that **all** aligning components are equal





$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$

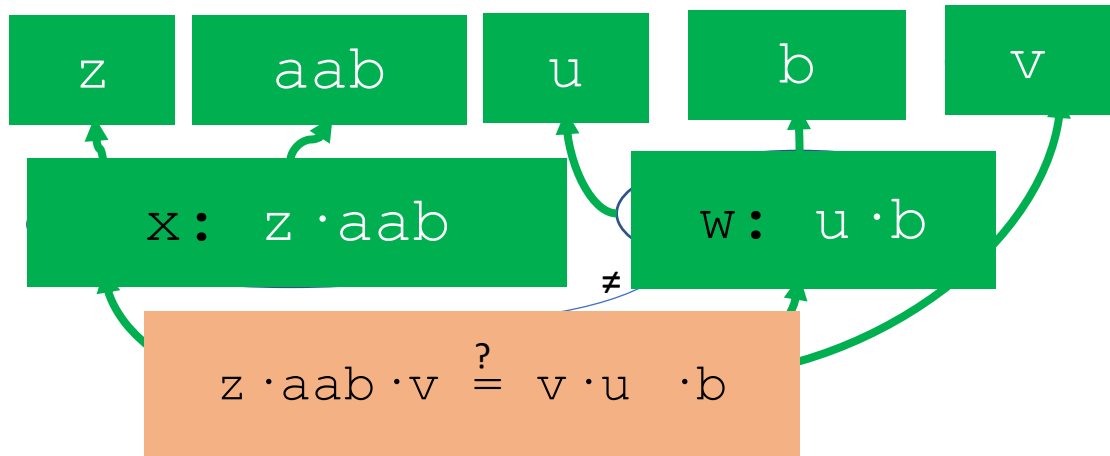


Consider **three cases** for making these two terms equal:

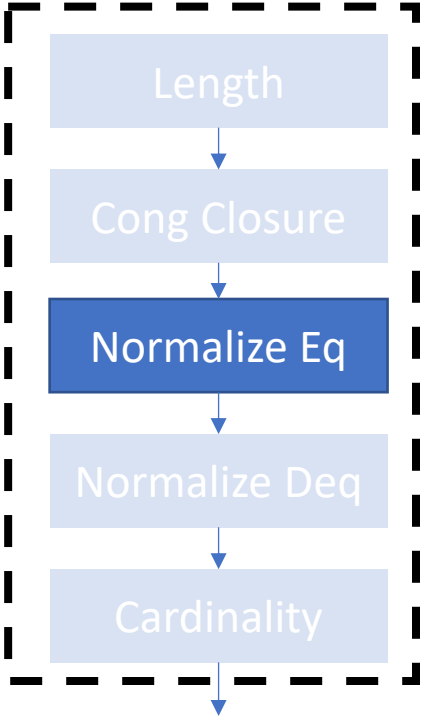


|| Case: $|z| = |v|$

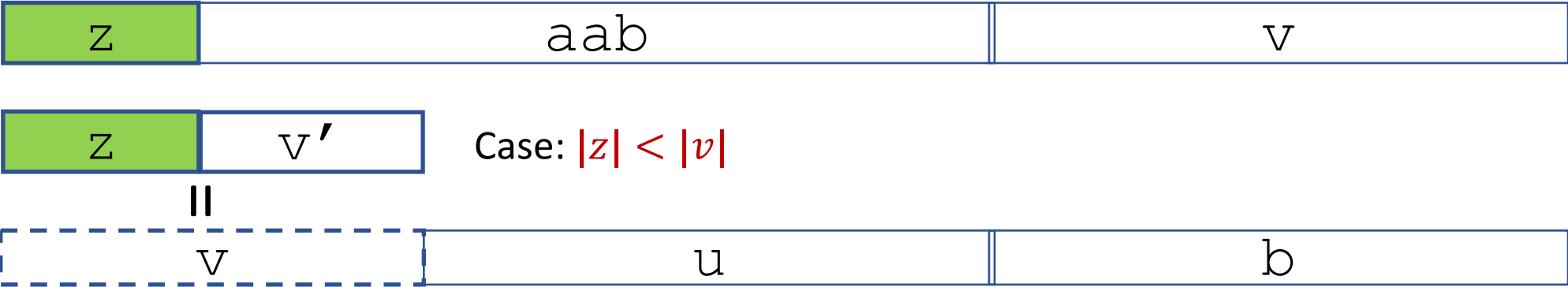


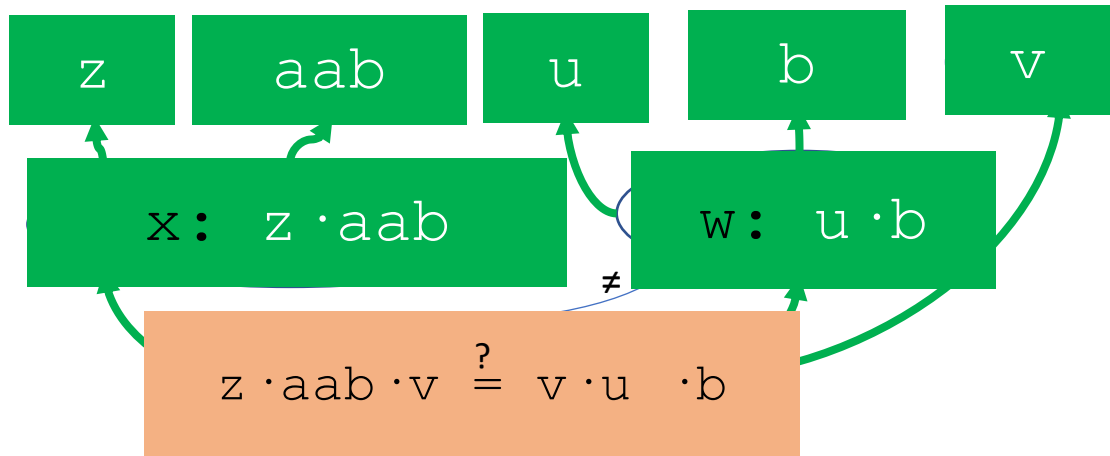


$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$

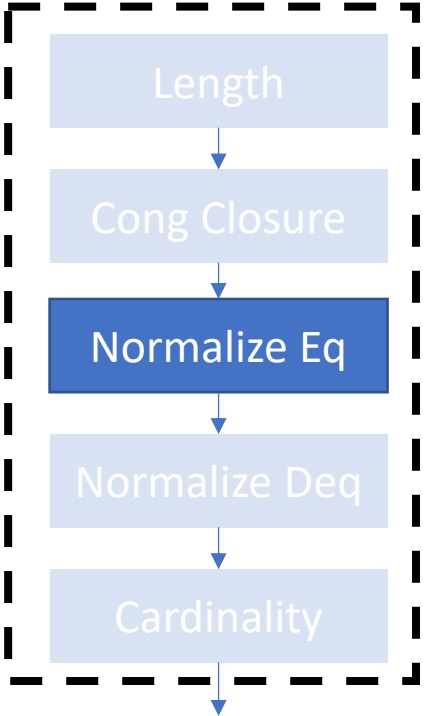


Consider **three cases** for making these two terms equal:

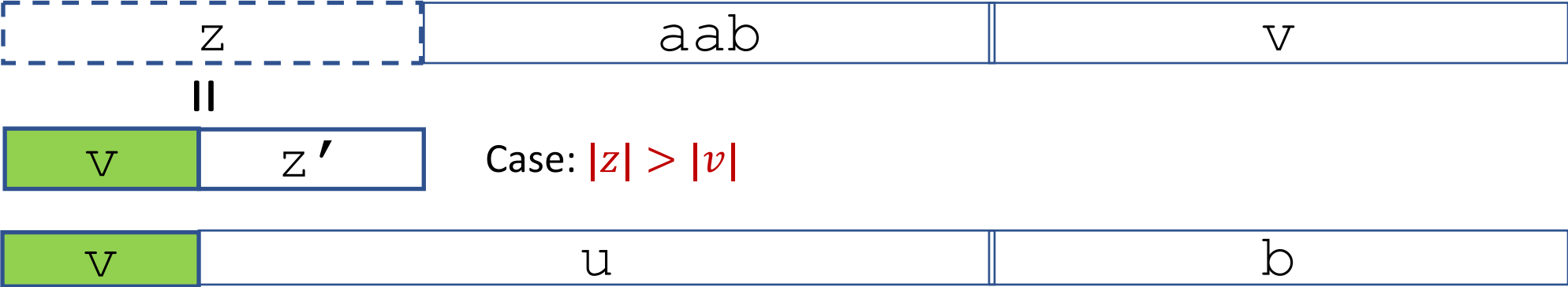




$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w
 \end{aligned}$$



Consider **three cases** for making these two terms equal:



$$x = z \cdot \mathbf{aab}$$

$$y = x$$

$$w = u \cdot \mathbf{b}$$

$$x \cdot v = v \cdot w$$

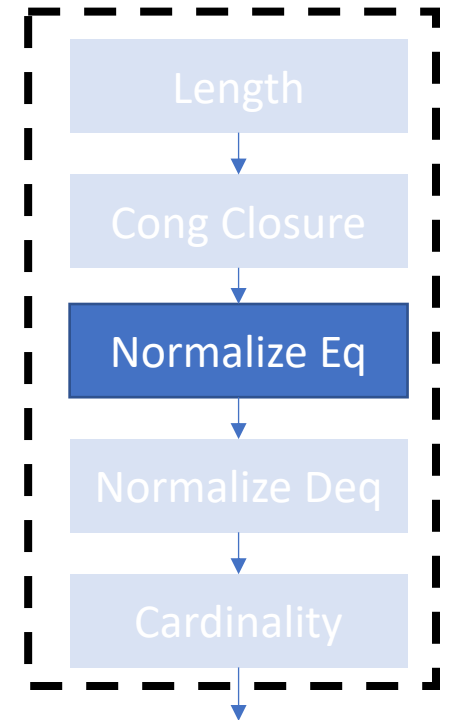
$$x \cdot v \neq w$$

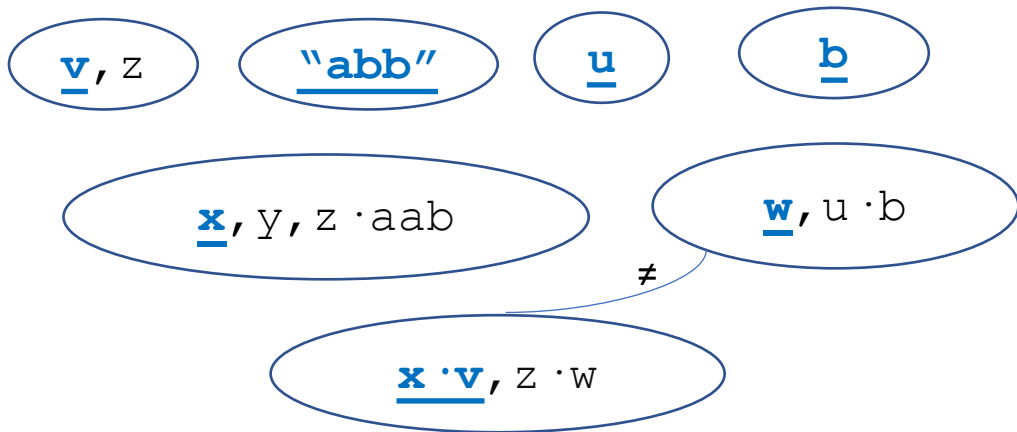
$$z = v$$

Equal case:



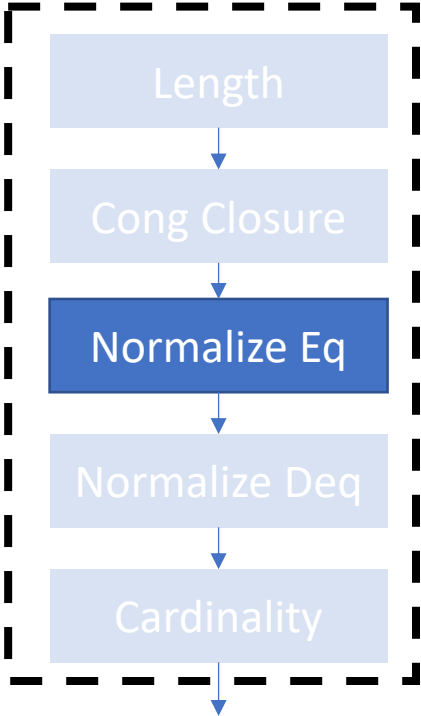
||



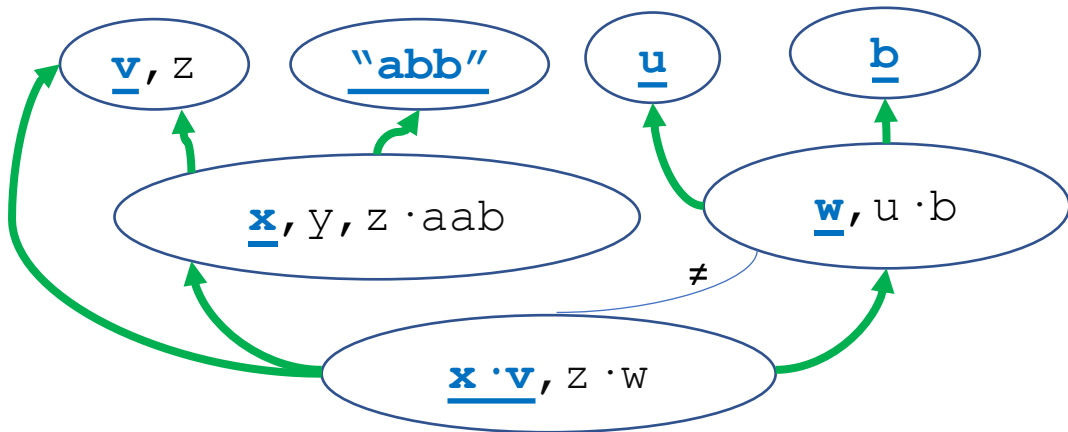


Equal case:

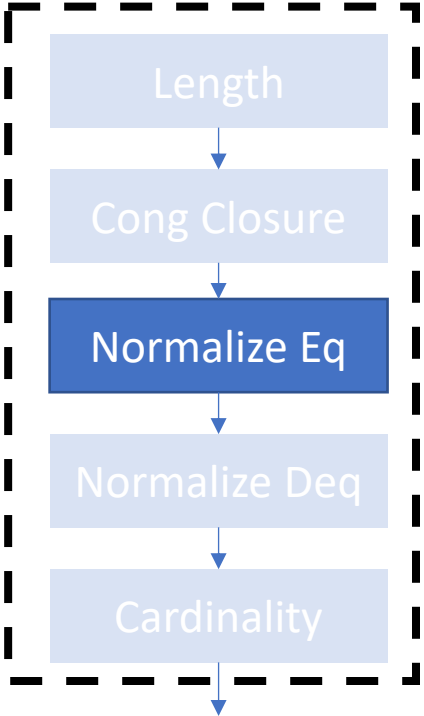
$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w \\
 z &= v
 \end{aligned}$$



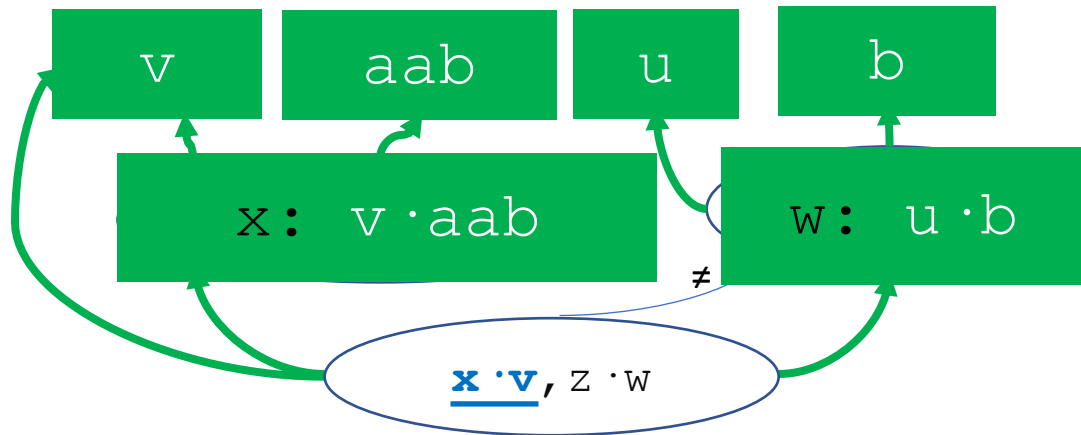
Recompute **congruence closure**



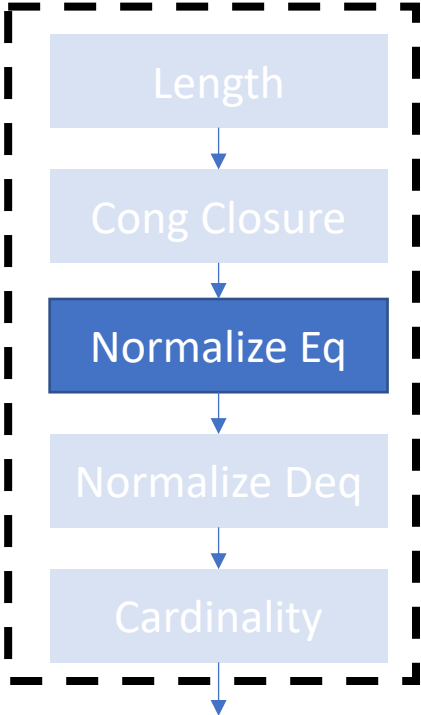
$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w \\
 z &= v
 \end{aligned}$$



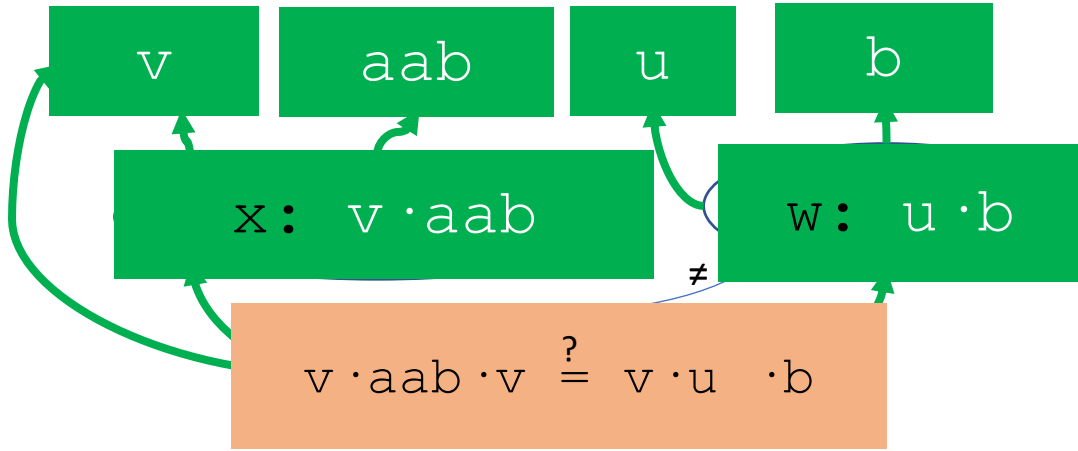
Recompute congruence closure and **normal forms**



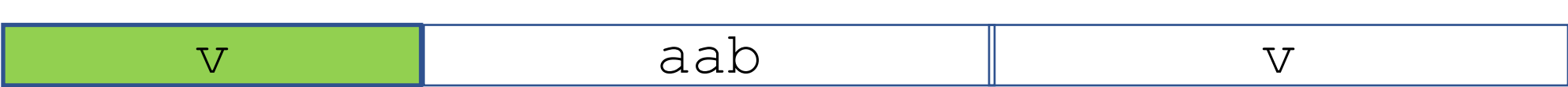
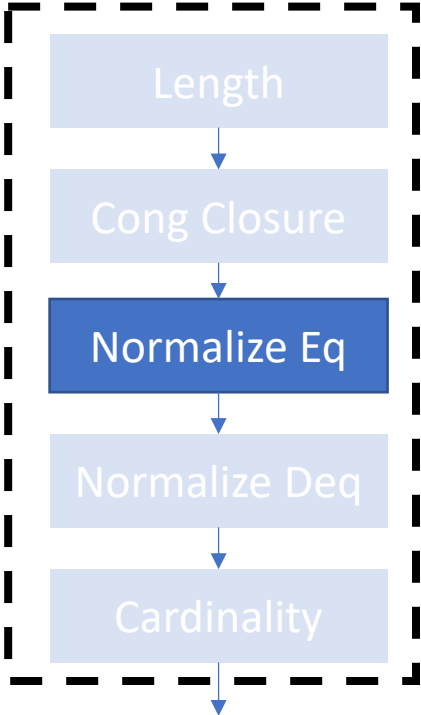
$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w \\
 z &= v
 \end{aligned}$$

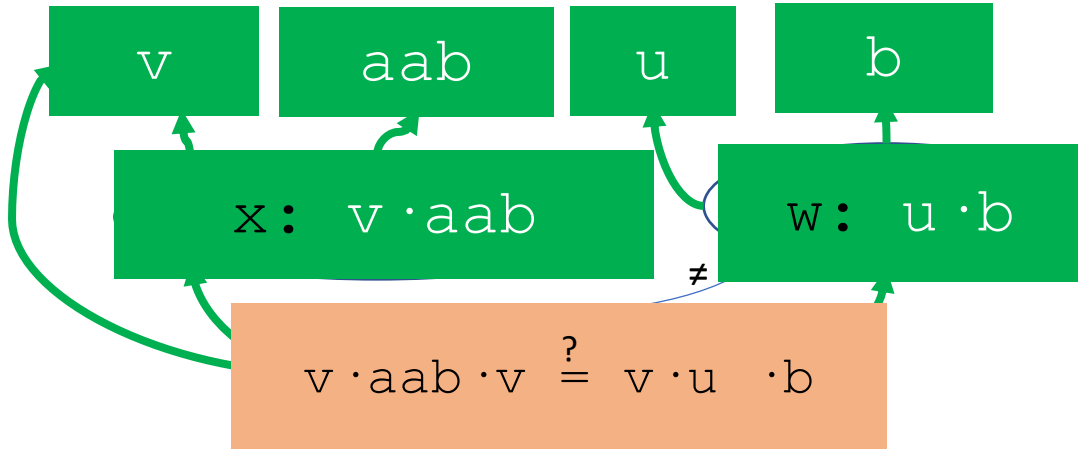


Recompute congruence closure and **normal forms**

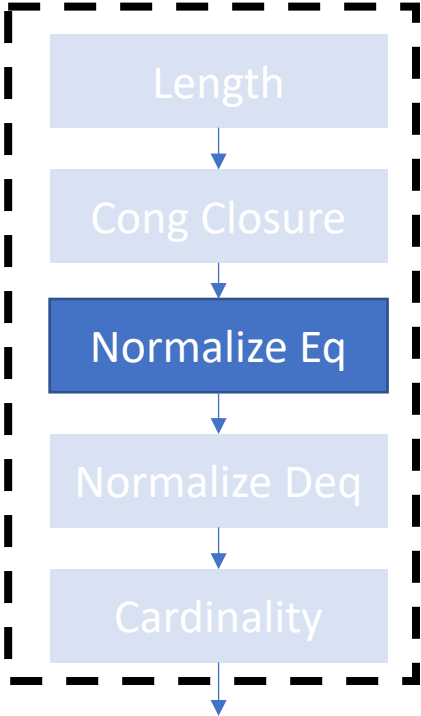


$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w \\
 z &= v
 \end{aligned}$$





$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &= v \cdot w \\
 x \cdot v &\neq w \\
 z &= v
 \end{aligned}$$



Repeat the process on these components



Splitting on String Equalities

Choice of equalities is quite **sophisticated** and **critical** to performance:

1. Prefers propagations over splits

E.g., $x \cdot w = y \cdot w \Rightarrow x = y$ over $x \cdot w = z \cdot v \Rightarrow (x = z \cdot x' \vee z = x \cdot z')$

2. Considers both the prefix and suffix of strings

E.g., $w \cdot x = w \cdot y \Rightarrow x = y$

3. Exploits length entailment [\[Zheng et al., 2015\]](#)

If $|x| > |y|$ according to the arithmetic solver,
then $x \cdot w = y \cdot v \wedge |x| > |y| \Rightarrow x = y \cdot x'$

Splitting on String Equalities

Choice of equalities is quite **sophisticated** and **critical** to performance:

4. Propagates constraints based on adjacent constants

E.g., $x \cdot \mathbf{b} = \mathbf{aab} \cdot y \Rightarrow x = \mathbf{aa} \cdot x'$, since **b** cannot overlap with prefix **aa**

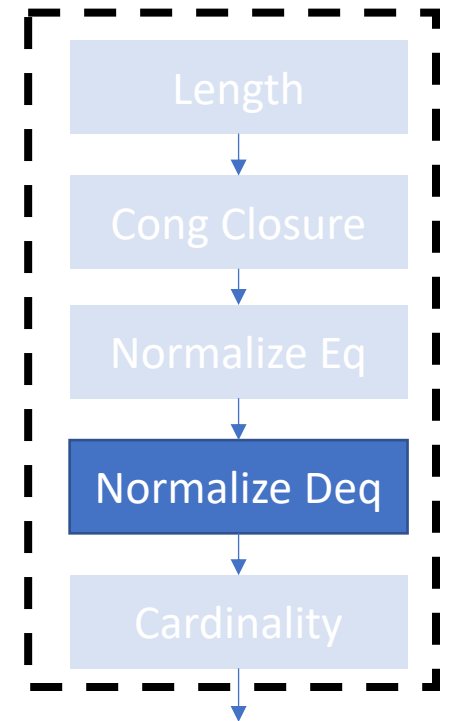
5. Treats looping word equations specially [[Liang et al., 2014](#)]

Splitting leads to non-termination; instead, reduce to RE membership

E.g., $x \cdot \mathbf{ba} = \mathbf{ab} \cdot x \Rightarrow x \in (\mathbf{ab})^* \mathbf{a}$

String Solver: Normalize Disequalities

$$\begin{aligned}x &= z \cdot \mathbf{aab} \\y &= x \\w &= u \cdot \mathbf{b} \\x \cdot v &\neq v \cdot w\end{aligned}$$

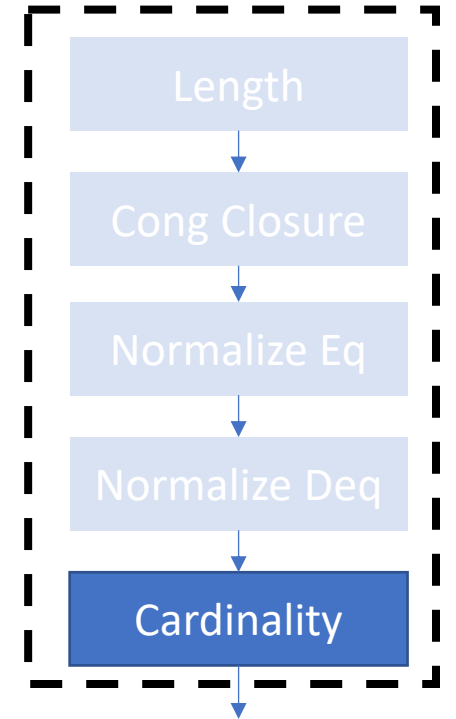


Disequalities are handled analogously to equalities

- 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
- **Goal:** find **any** two aligning components that are disequal

5. Check Cardinality Constraints

$$\begin{aligned}x &= z \cdot \mathbf{aab} \\y &= x \\w &= u \cdot \mathbf{b} \\x \cdot v &\neq v \cdot w \\v &\neq z\end{aligned}$$



5. Check Cardinality Constraints

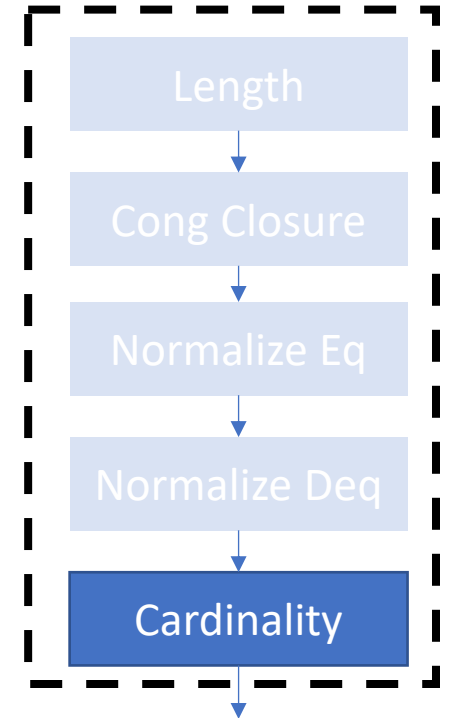
\mathcal{M}_Σ may be unsatisfiable
because Σ is **finite**

$$\begin{aligned}x &= z \cdot \mathbf{aab} \\y &= x \\w &= u \cdot \mathbf{b} \\x \cdot v &\neq v \cdot w \\v &\neq z\end{aligned}$$

Example:

- Σ consists of 256 characters, and
- \mathcal{M}_Σ entails that 257 distinct strings of length 1 exist

$$\text{distinct}(s_1, \dots, s_{257}), |s_1| = 1, \dots, |s_{257}| = 1 \models \perp$$

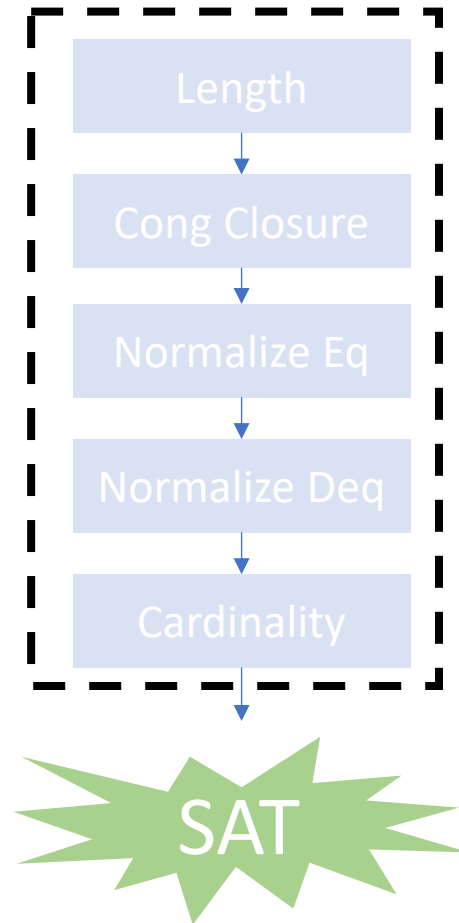


Finally: Compute Model

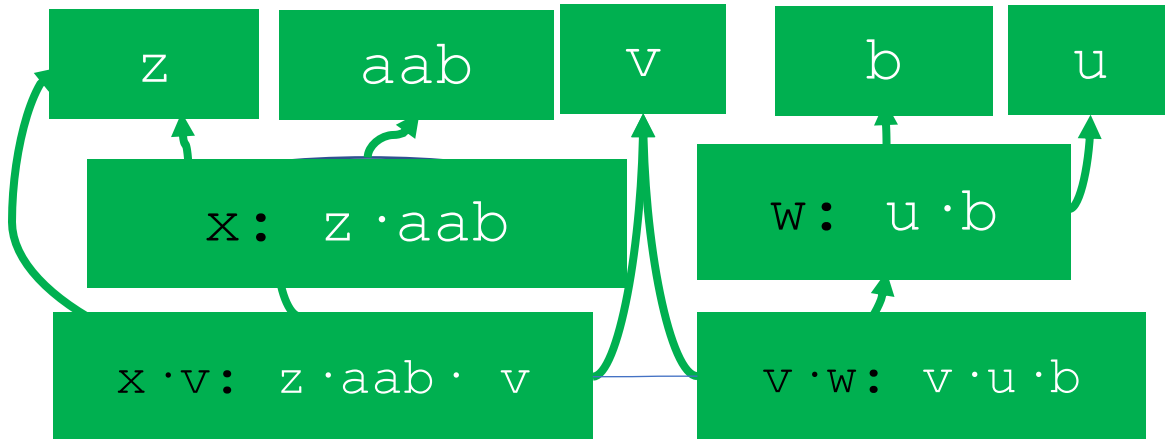
If all steps finish with
no new lemmas:

- \mathcal{M}_s is T_s -satisfiable
- Compute **model based on normal forms**
 - assign string constants to eq classes whose normal form is a variable
 - Length fixed by model from arithmetic solver
 - Interpret each var as the value of its eq class' normal form

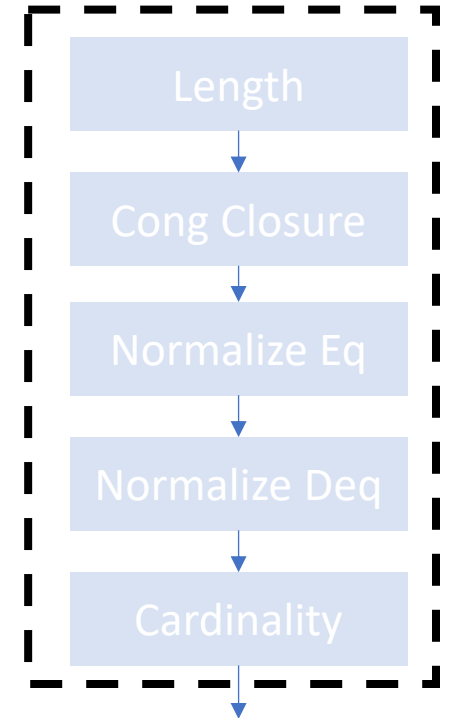
$$\begin{aligned}x &= z \cdot \mathbf{aab} \\y &= x \\w &= u \cdot \mathbf{b} \\x \cdot v &\neq v \cdot w \\v &\neq z\end{aligned}$$



Compute Model

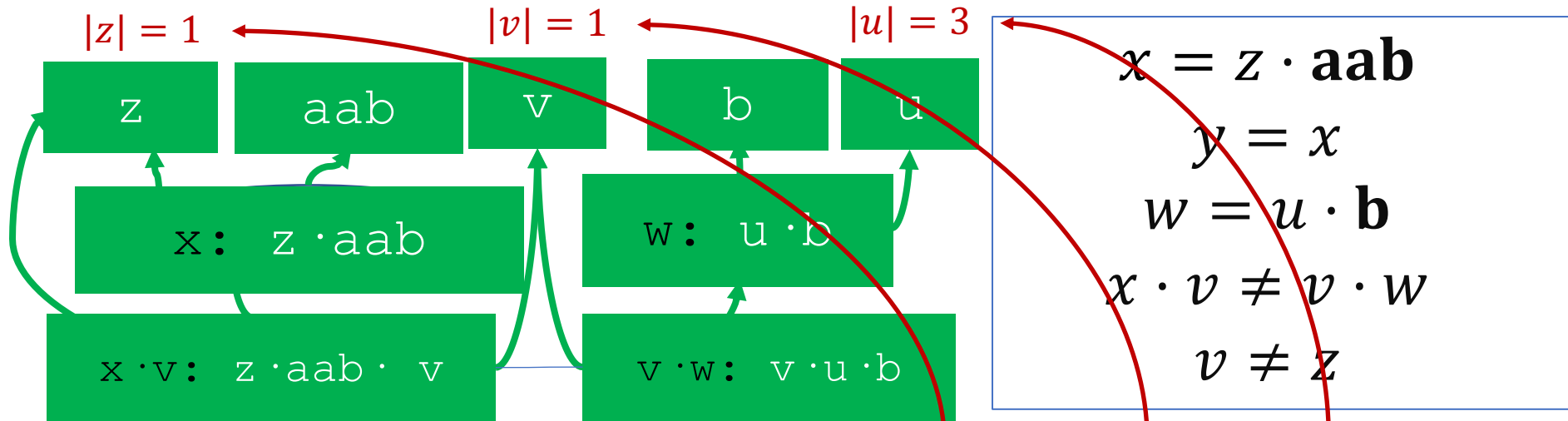


$$\begin{aligned}x &= z \cdot \mathbf{aab} \\y &= x \\w &= u \cdot \mathbf{b} \\x \cdot v &\neq v \cdot w \\v &\neq z\end{aligned}$$



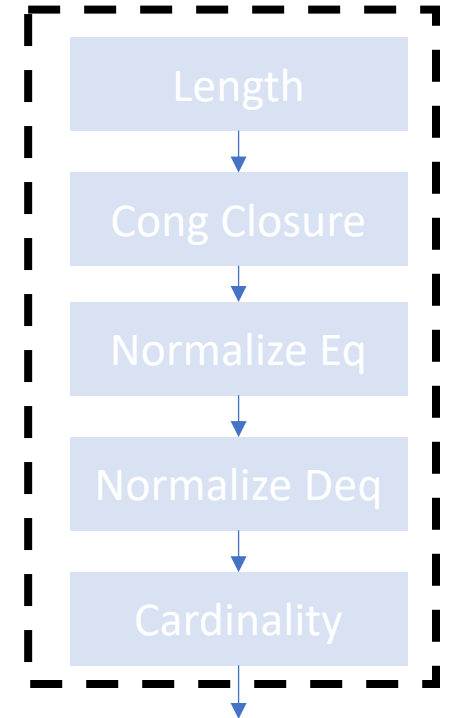
SAT

Compute Model

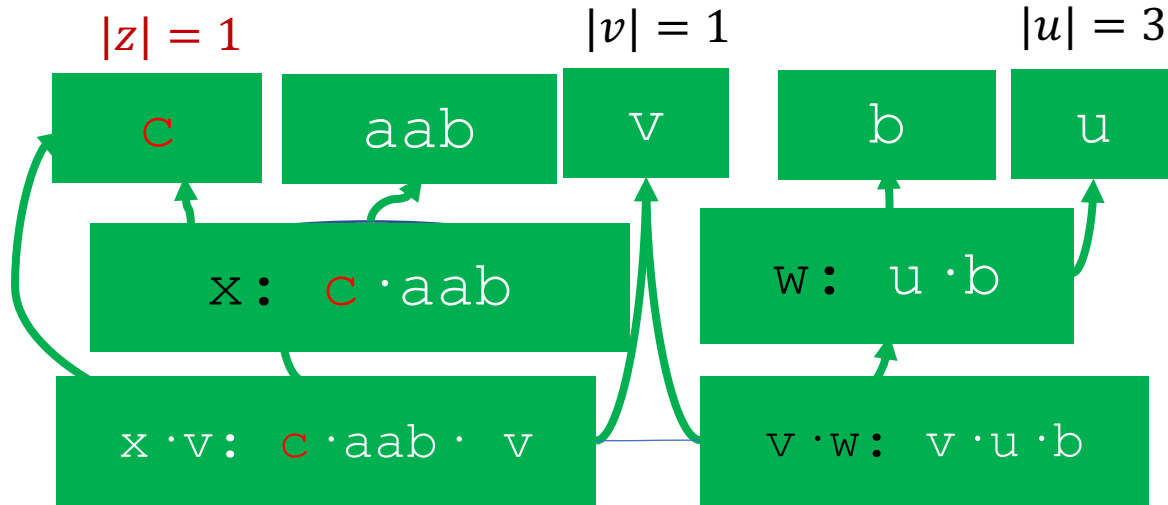


Example:

arith model



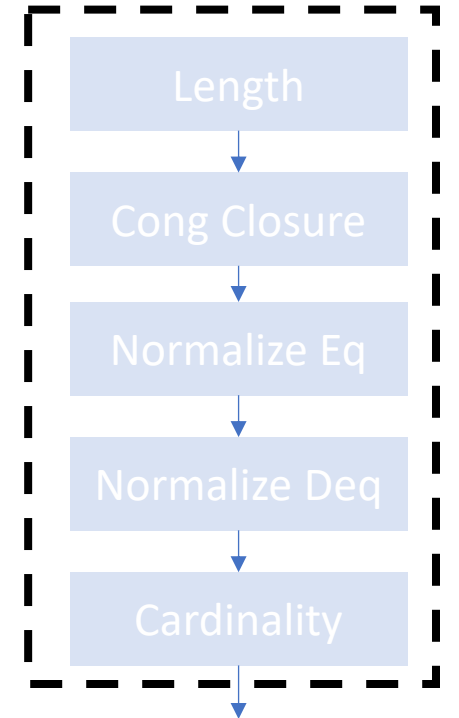
Compute Model



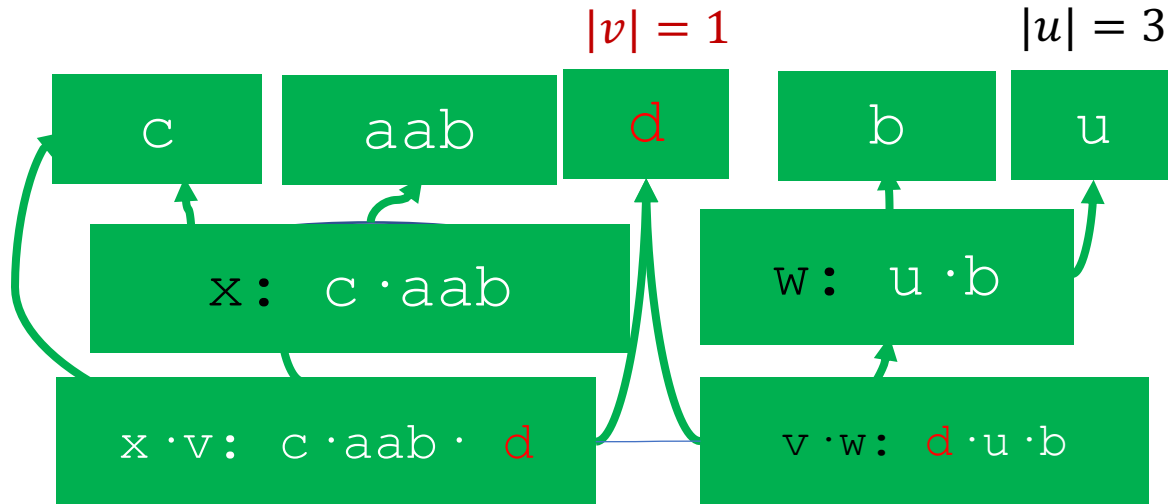
$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &\neq v \cdot w \\
 v &\neq z
 \end{aligned}$$

Example:

- $z \mapsto c$



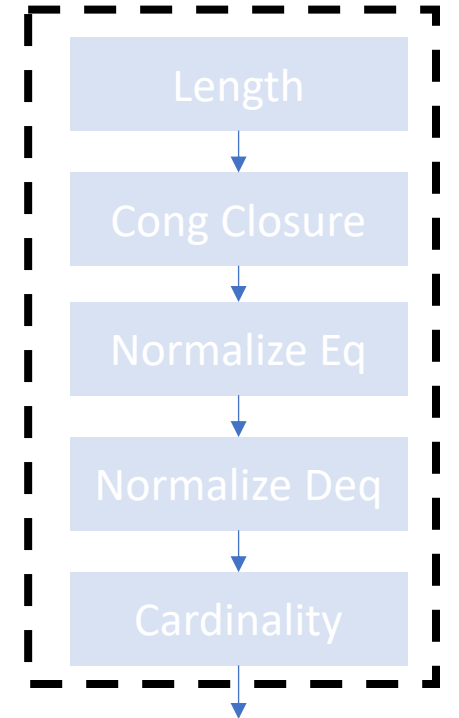
Compute Model



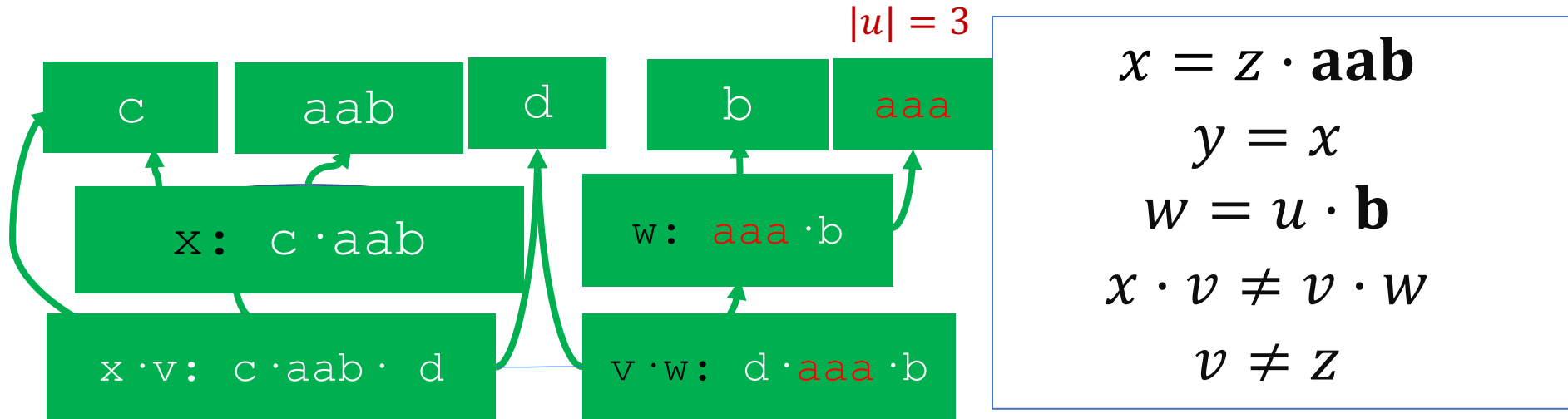
$$\begin{aligned}x &= z \cdot \mathbf{aab} \\y &= x \\w &= u \cdot \mathbf{b} \\x \cdot v &\neq v \cdot w \\v &\neq z\end{aligned}$$

Example:

- $z \mapsto \mathbf{c}$
- $v \mapsto \mathbf{d}$



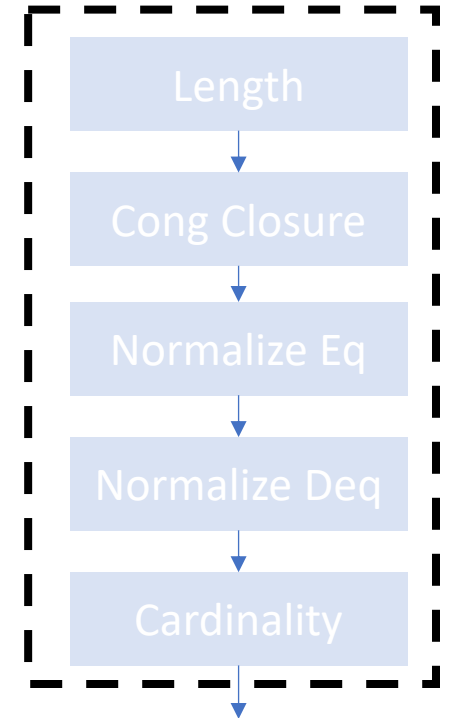
Compute Model



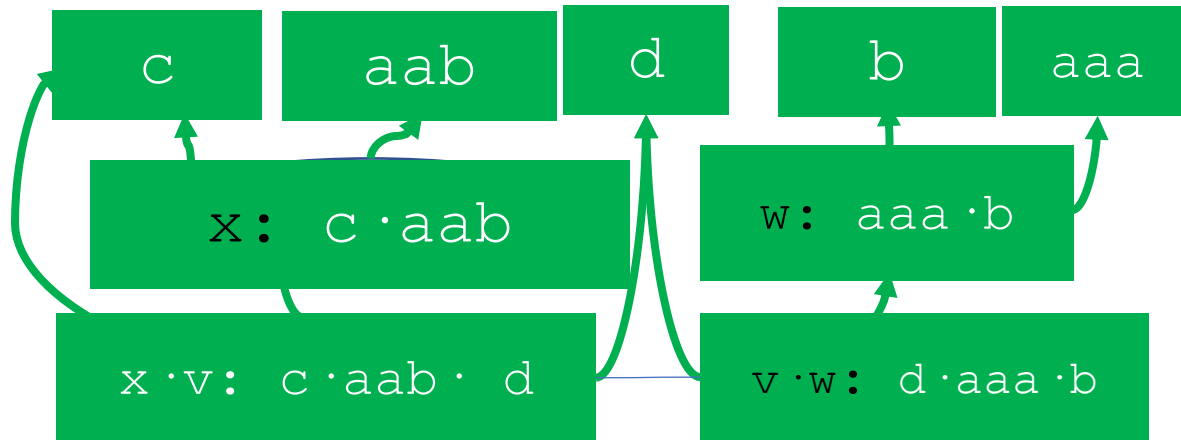
Example:

- $z \mapsto \mathbf{c}$
- $v \mapsto \mathbf{d}$
- $u \mapsto \mathbf{aaa}$

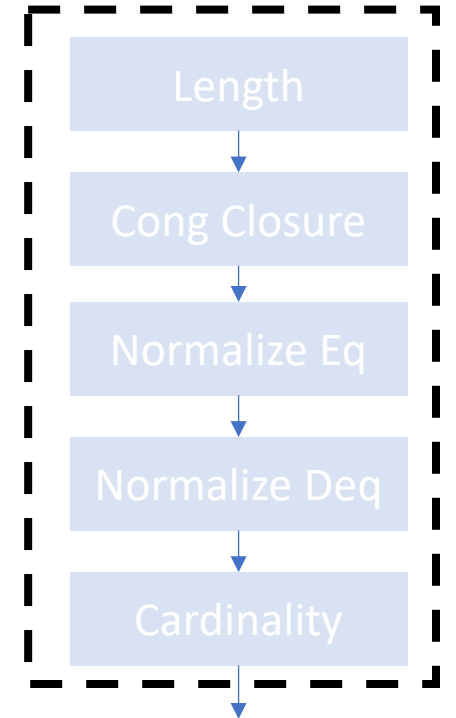
Check-cardinality step ensures there are enough constants



Compute Model



$$\begin{aligned}
 x &= z \cdot \mathbf{aab} \\
 y &= x \\
 w &= u \cdot \mathbf{b} \\
 x \cdot v &\neq v \cdot w \\
 v &\neq z
 \end{aligned}$$



SAT

Example:

- $z \mapsto \mathbf{c}$
- $v \mapsto \mathbf{d}$
- $u \mapsto \mathbf{aaa}$
- Other vars assigned to value of the normal form of their eq classes
 $x \mapsto \mathbf{caab}$ $y \mapsto \mathbf{caab}$ $w \mapsto \mathbf{aaab}$

Saturation criterion for procedure ensures this model satisfies \mathcal{M}_s

Techniques for **Fast** String Solving in cvc5

- Finite model finding
- Context-dependent simplification for extended constraints
- Witness sharing
- Regular expression elimination
- String to code point conversion

Finite Model Finding for Strings

Finite Model Finding for Strings

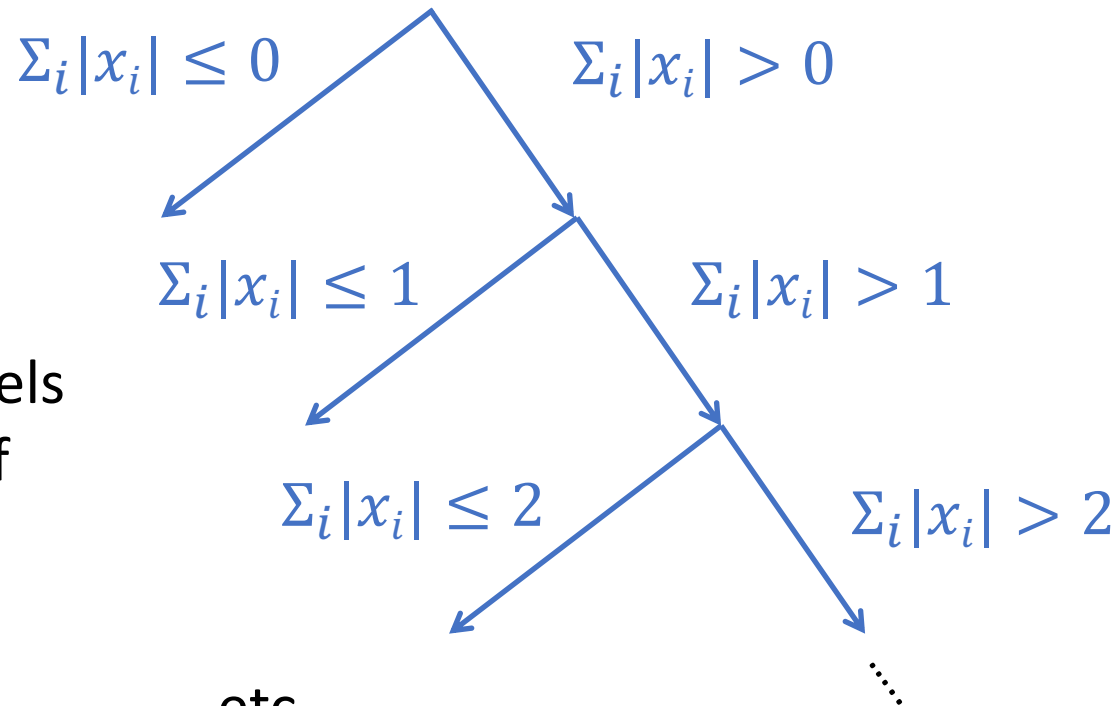
Idea: Incrementally bound the lengths of input string variables x_1, \dots, x_n

⇒ Improves solver's solving time for problems with small models

Search for models
where sum of
lengths is 0

Search for models
where sum of
lengths is 1

etc.



Context-Dependent Simplification for Extended String Constraints

[Reynolds, Woo, Barrett, Brumley, Liang and Tinelli, CAV'17]

Extended String Constraint Language

Substring: $\text{substr}(x, n, l)$

- the substring of string x starting at position n of length at most l

String contains: $\text{contains}(x, y)$

- true iff string x contains y as a substring

Find index: $\text{indexof}(x, n, p)$

- the position of the first occurrence of string y in x , starting from position n if any; -1 otherwise

String replace: $\text{replace}(x, y, y')$

- the result of replacing the first occurrence of string y in x (if any) with y'

Example:

$$\neg \text{contains}(\text{substr}(x, 0, 3), \mathbf{a}) \wedge 0 \leq \text{indexof}(x, \mathbf{ab}, 0) < 4$$

How do we handle **Extended String Constraints**?

\neg contains(x , \mathbf{a})

How do we handle Extended String Constraints?

Naively, by **reduction** to basic constraints + bounded \forall

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

How do we handle Extended String Constraints?

Naively, by **reduction** to basic constraints + bounded \forall

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

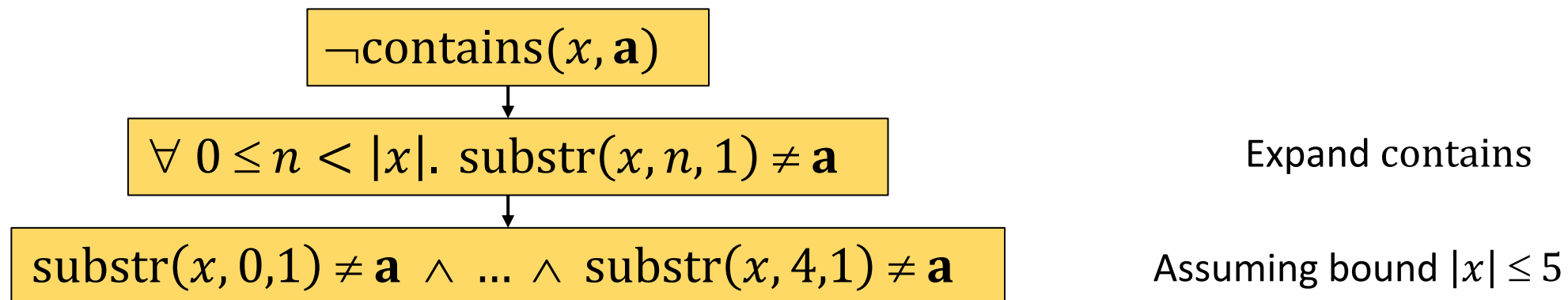


$\forall 0 \leq n < |x|. \text{ substr}(x, n, 1) \neq \mathbf{a}$

Expand contains

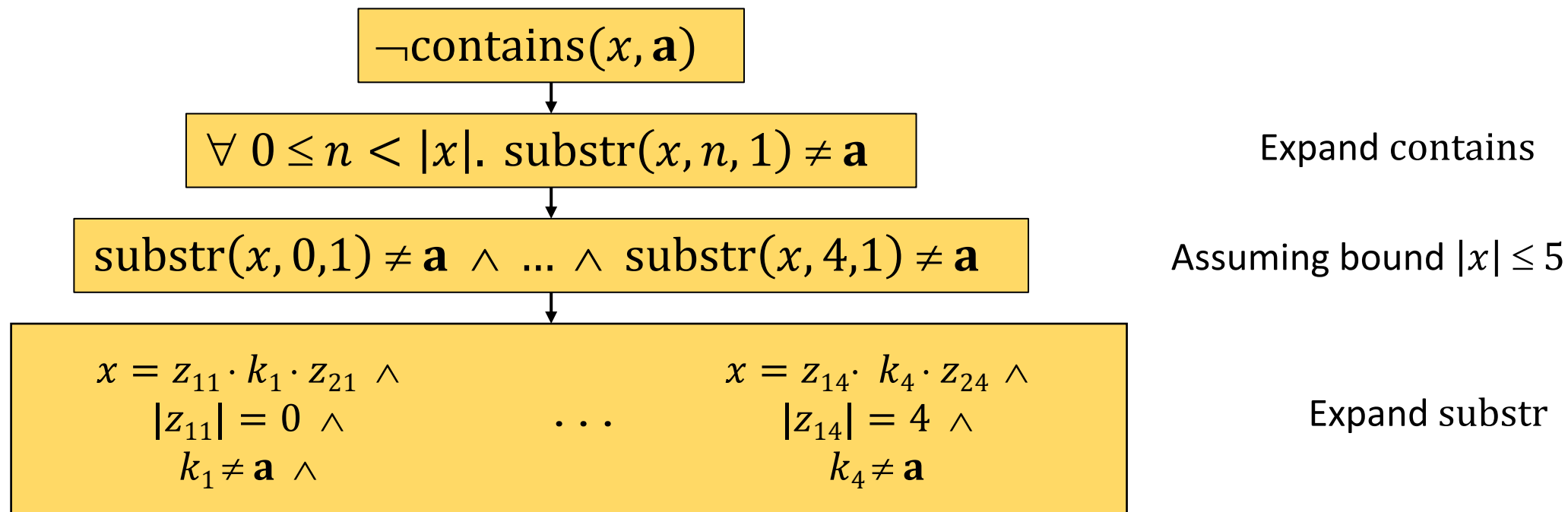
How do we handle Extended String Constraints?

Naively, by **reduction** to basic constraints + bounded \forall



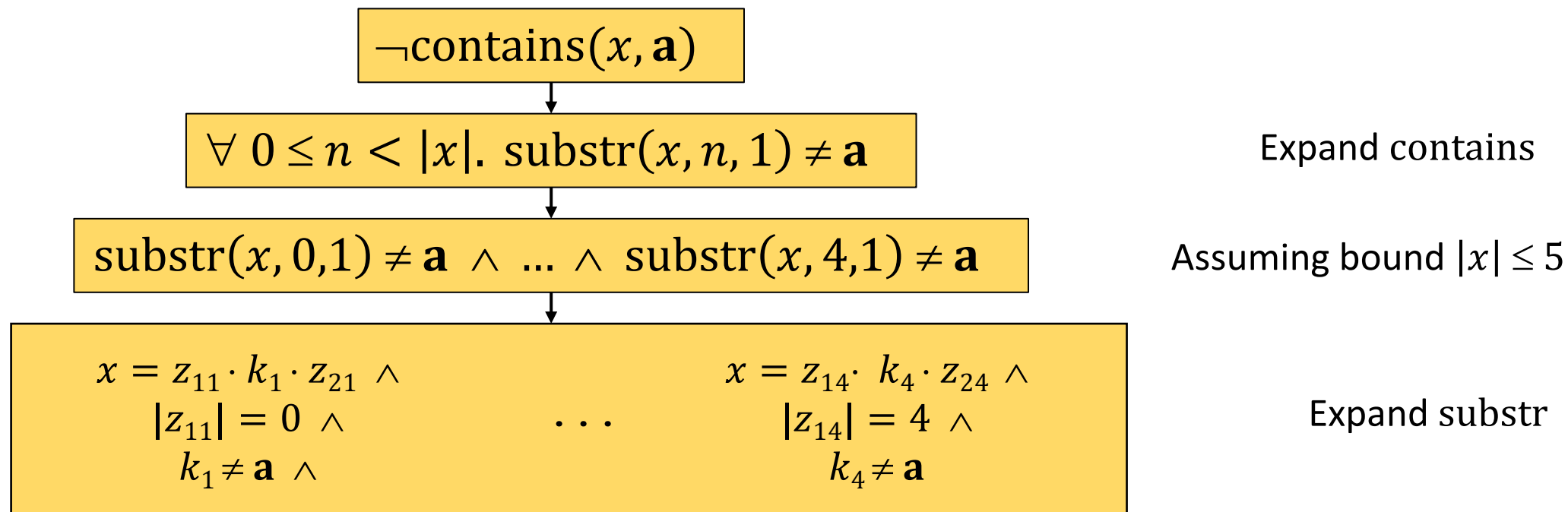
How do we handle Extended String Constraints?

Naively, by **reduction** to basic constraints + bounded \forall



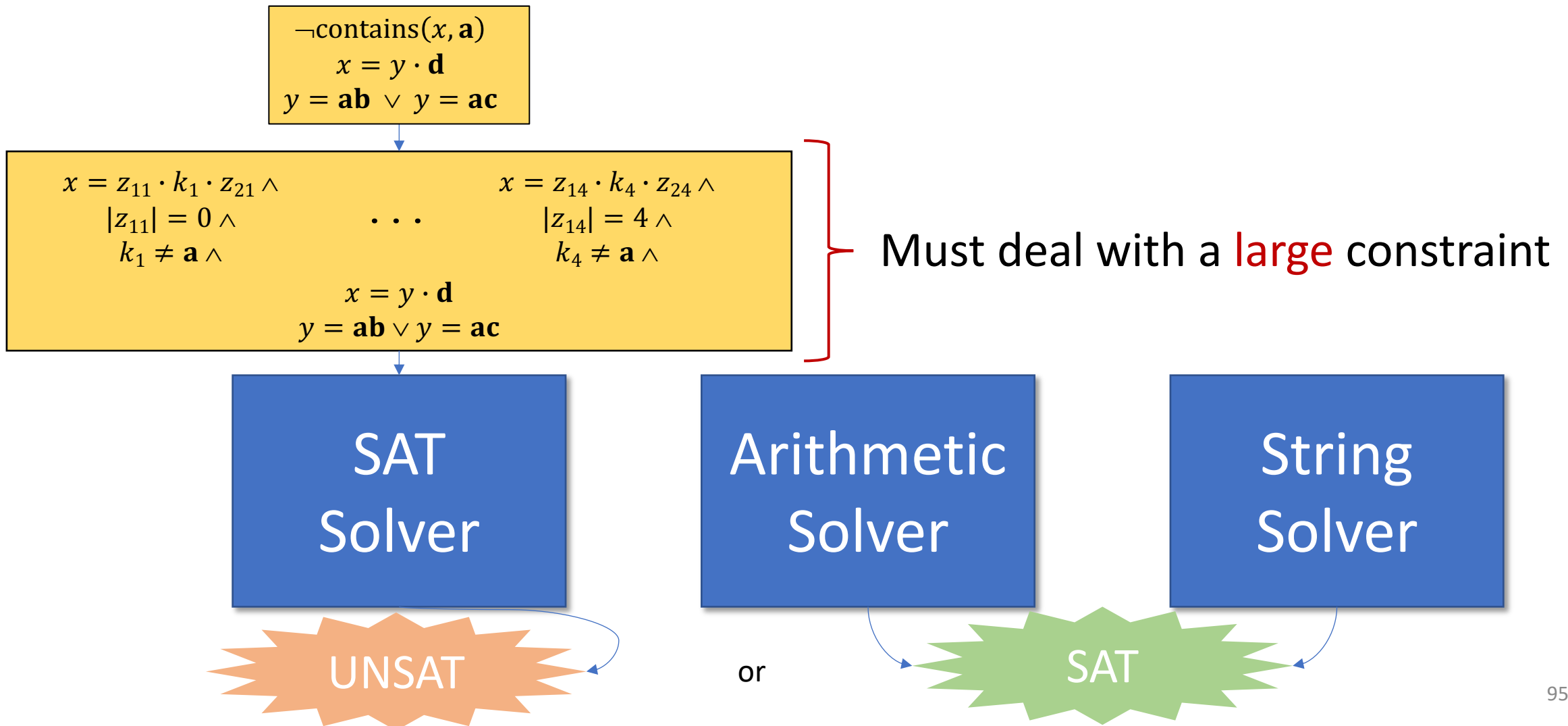
How do we handle Extended String Constraints?

Naively, by **reduction** to basic constraints + bounded \forall

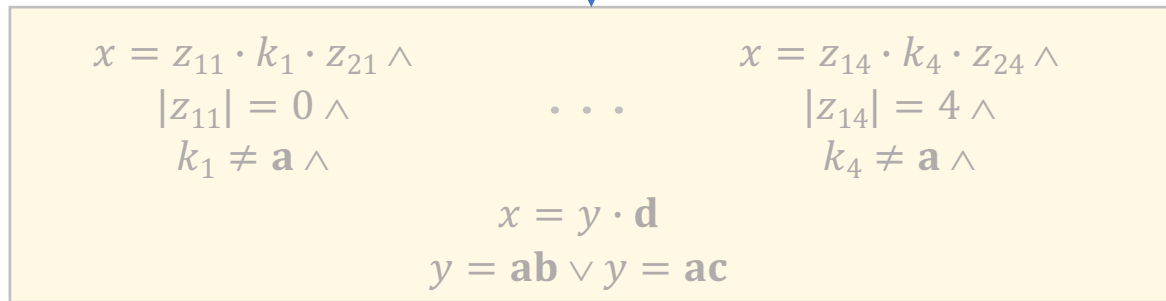
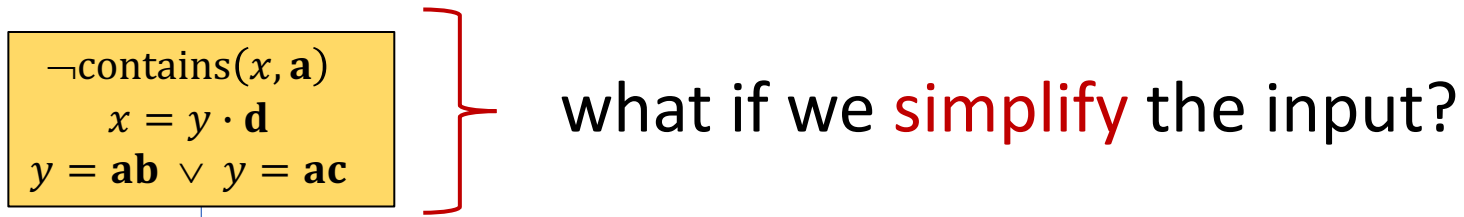


Approach followed by many solvers [Bjorner et al. 2009, Zheng et al. 2013, Li et al. 2013, Trinh et al. 2014]

(Eager) Expansion of Extended Constraints



(Eager) Expansion of Extended Constraints



or



SMT Solvers + Simplification

All SMT solvers implement **simplification** techniques

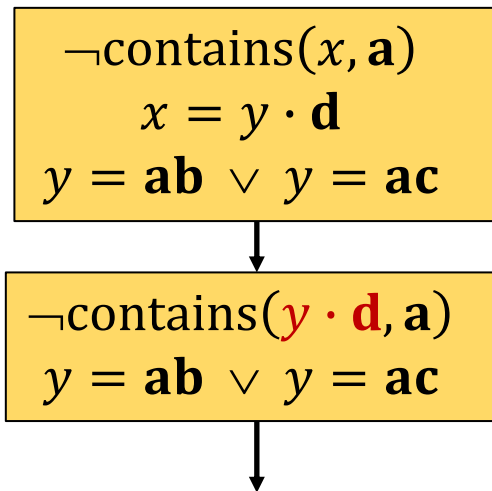
(also called *normalization* or *rewrite rules*)

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

SMT Solvers + Simplification

All SMT solvers implement **simplification** techniques

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

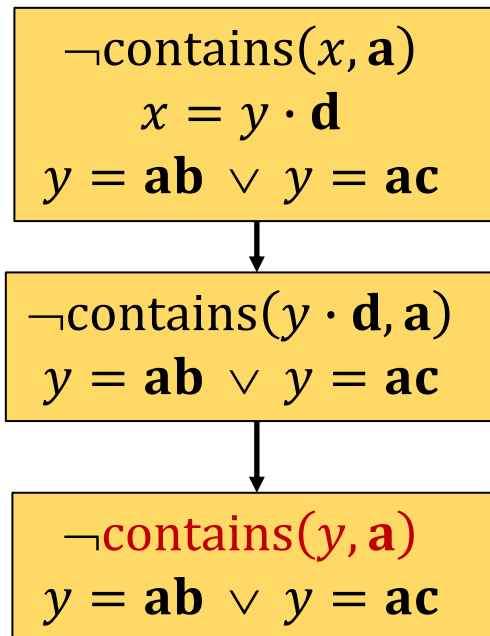


since $x = y \cdot \mathbf{d}$

SMT Solvers + Simplification

All SMT solvers implement **simplification** techniques

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



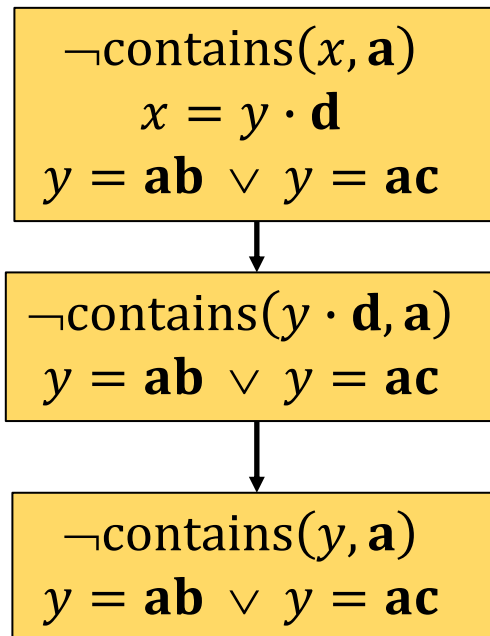
since $x = y \cdot \mathbf{d}$

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

SMT Solvers + Simplification

All SMT solvers implement **simplification** techniques

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



since $x = y \cdot \mathbf{d}$

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

- Leads to smaller inputs

Some problems can be solved by simplification alone

(Lazy) Expansion + Simplification

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

SAT
Solver

Arithmetic
Solver

String
Solver

(Lazy) Expansion + Simplification

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

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

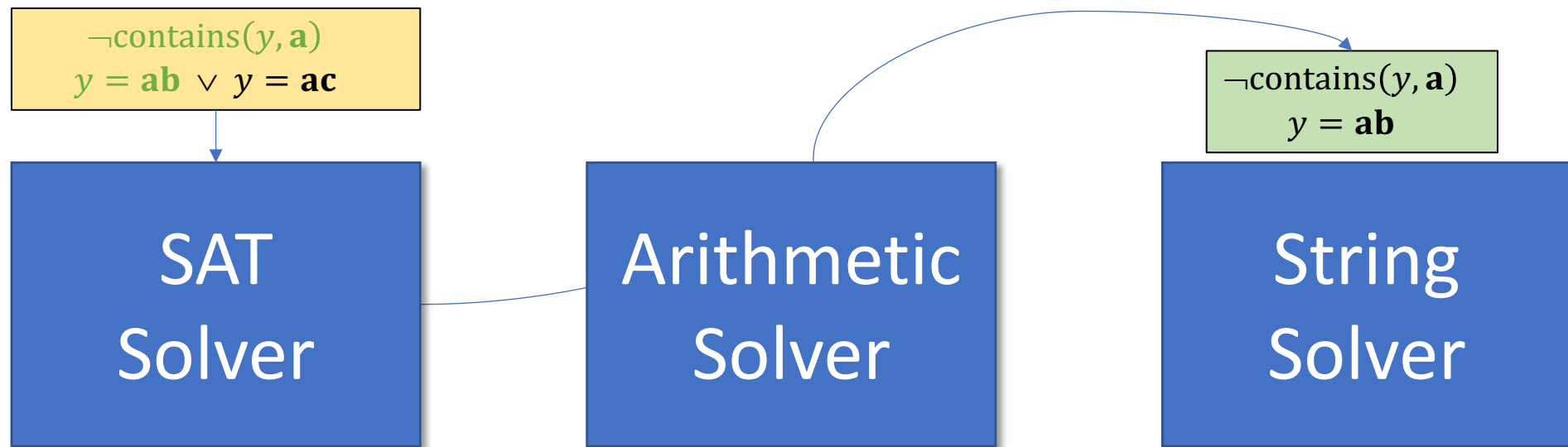
Simplify the input

SAT
Solver

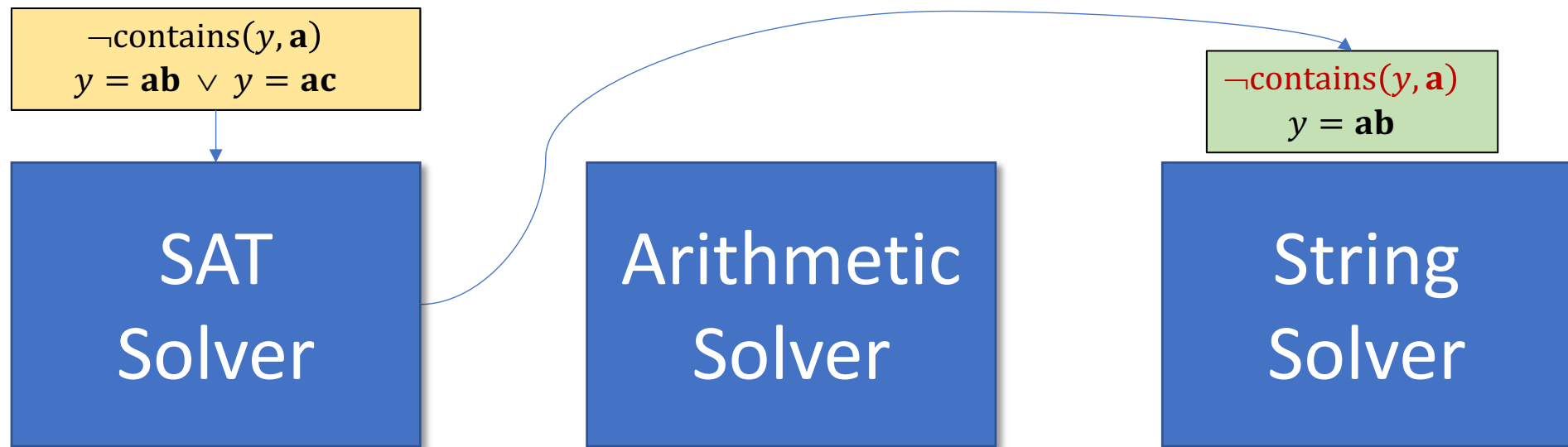
Arithmetic
Solver

String
Solver

(Lazy) Expansion + Simplification



(Lazy) Expansion + Simplification

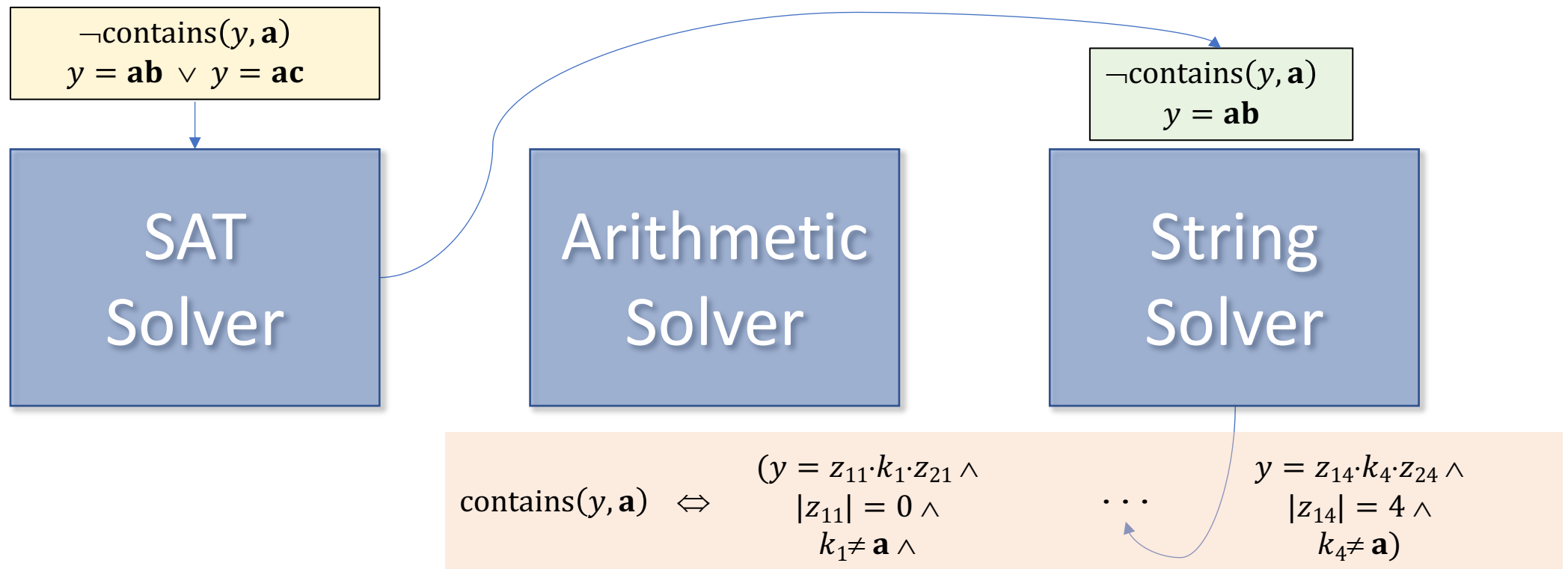


Still have a large constraint!

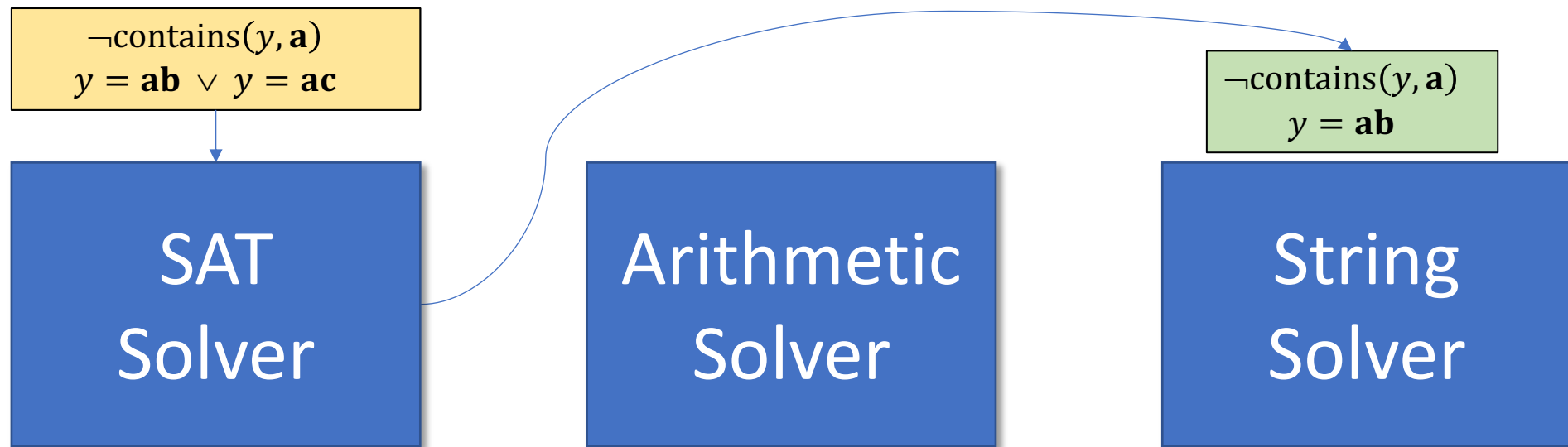
$$\text{contains}(y, \mathbf{a}) \Leftrightarrow (y = z_{11} \cdot k_1 \cdot z_{21} \wedge |z_{11}| = 0 \wedge k_1 \neq \mathbf{a} \wedge \dots \wedge y = z_{14} \cdot k_4 \cdot z_{24} \wedge |z_{14}| = 4 \wedge k_4 \neq \mathbf{a})$$

(Lazy) Expansion + Simplification

What if we simplify based on the **context**?

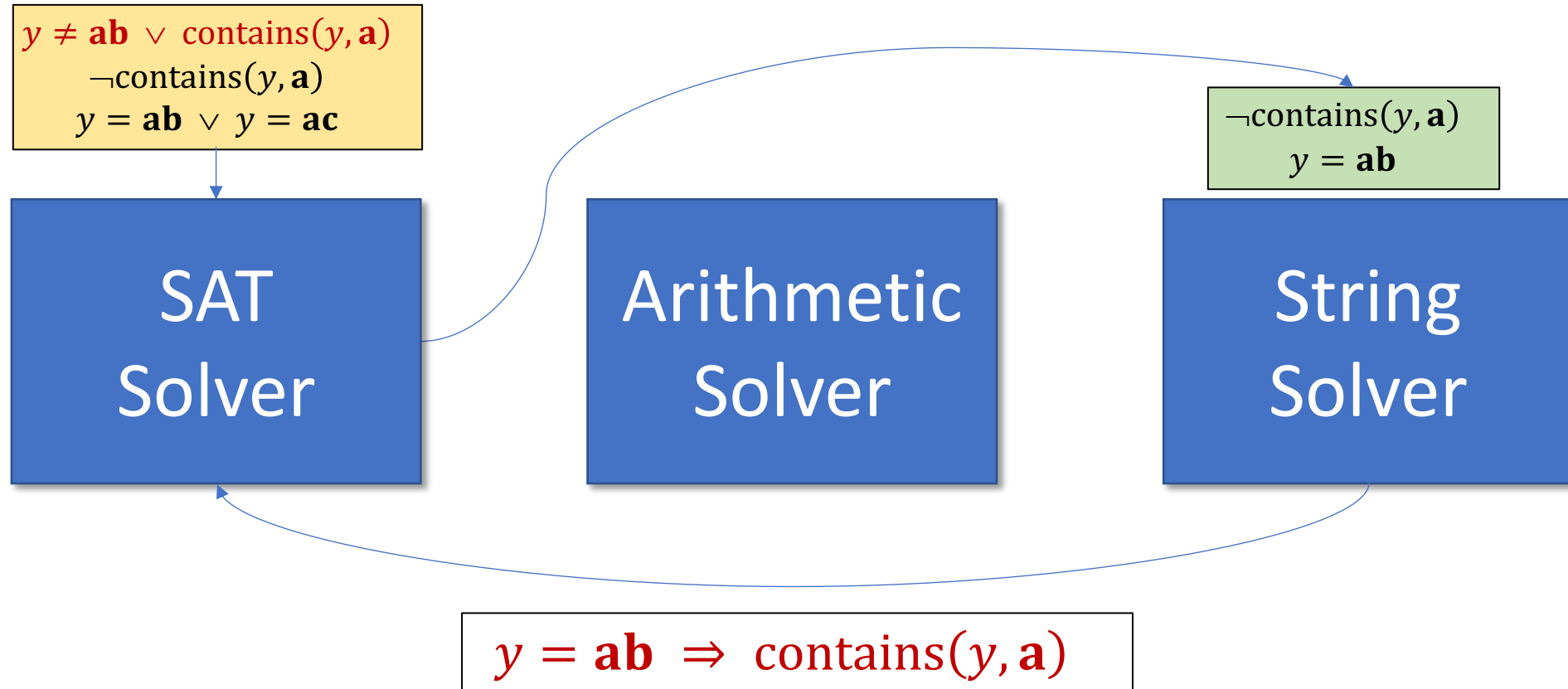


(Lazy) Expansion + **Context-Dependent** Simplification [Reynolds et al., CAV'17]

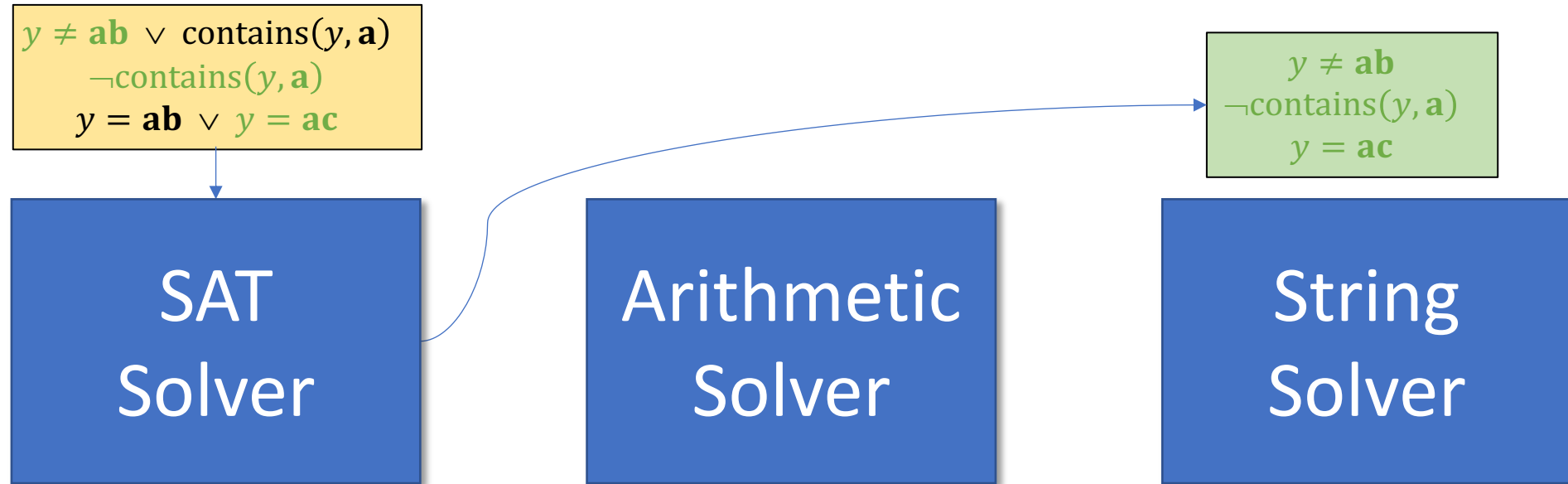


Since $\text{contains}(y, a)$ is true when $y = \mathbf{ab}$...

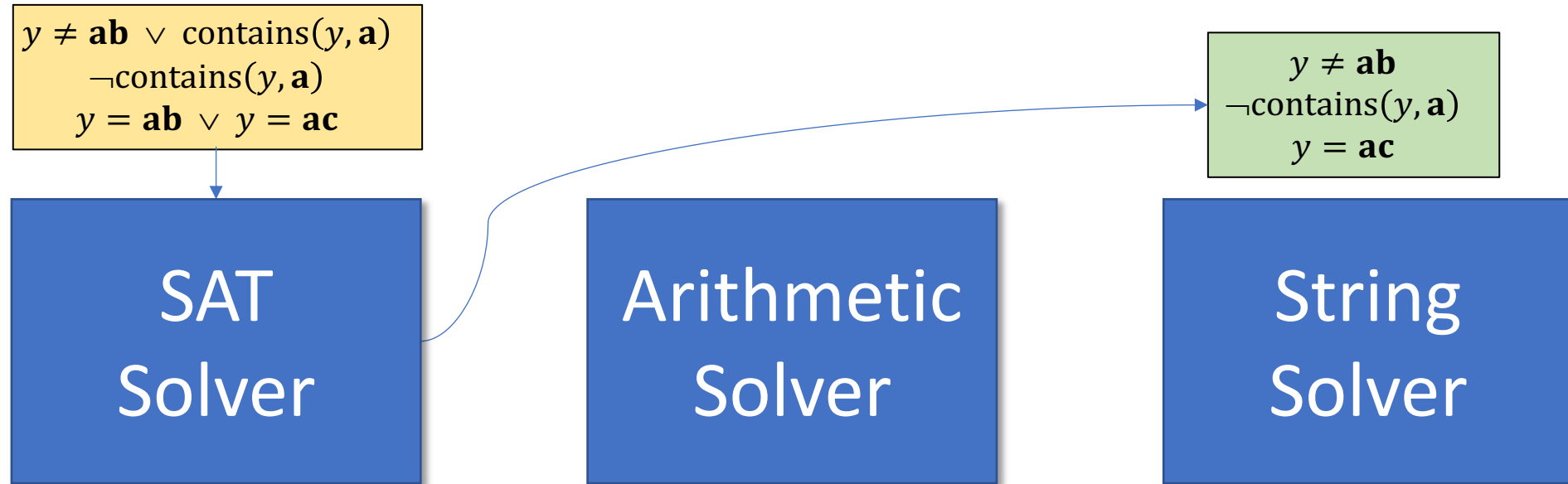
(Lazy) Expansion + Context-Dependent Simplification



(Lazy) Expansion + Context-Dependent Simplification

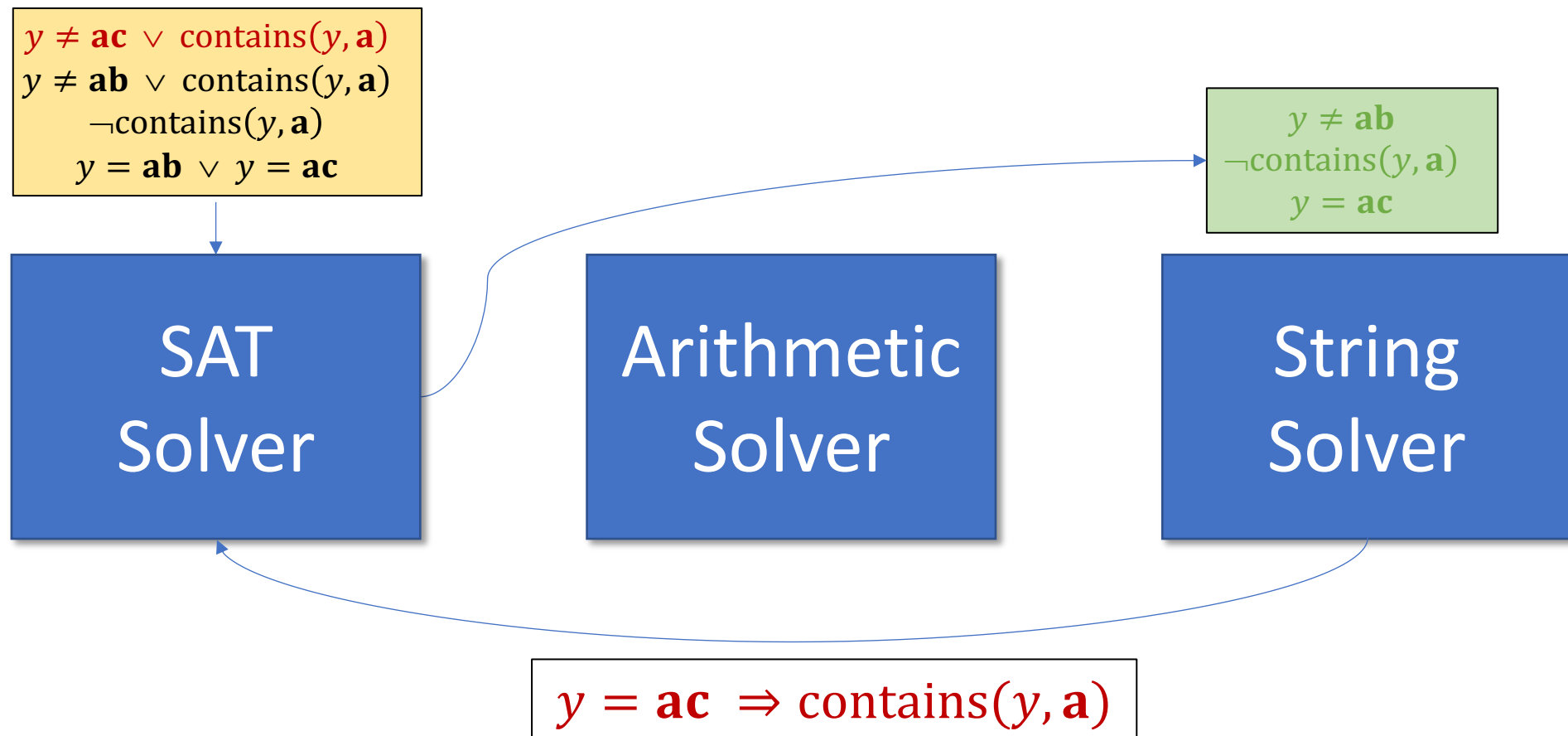


(Lazy) Expansion + Context-Dependent Simplification



$\text{contains}(y, \mathbf{a})$ is true also when $y = \mathbf{ac}$...

(Lazy) Expansion + Context-Dependent Simplification



(Lazy) Expansion + Context-Dependent Simplification

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

Did not need to expand
contains at all!

$y \neq \mathbf{ab}$
 $\neg \text{contains}(y, \mathbf{a})$
 $y = \mathbf{ac}$

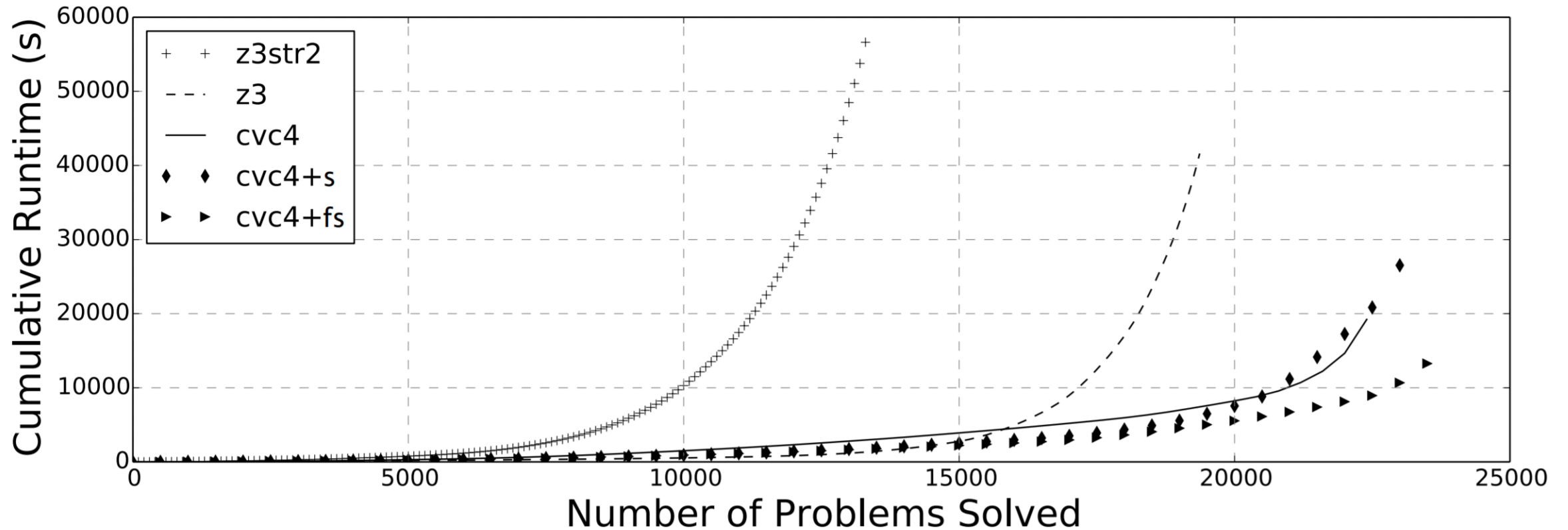
SAT
Solver

Arithmetic
Solver

String
Solver

UNSAT

Results on Symbolic Execution [\[Reynolds et al., CAV'17\]](#)



cvc4+fs (context-dependent simplification + finite model finding) solves **23,802** benchmarks in **5h8m**

- Without finite model finding, solves **23,266** in **8h46m**

- Without either finite model finding or cd-simplification, solves **22,607** in **6h38m**

Aggressive Simplifications for Strings

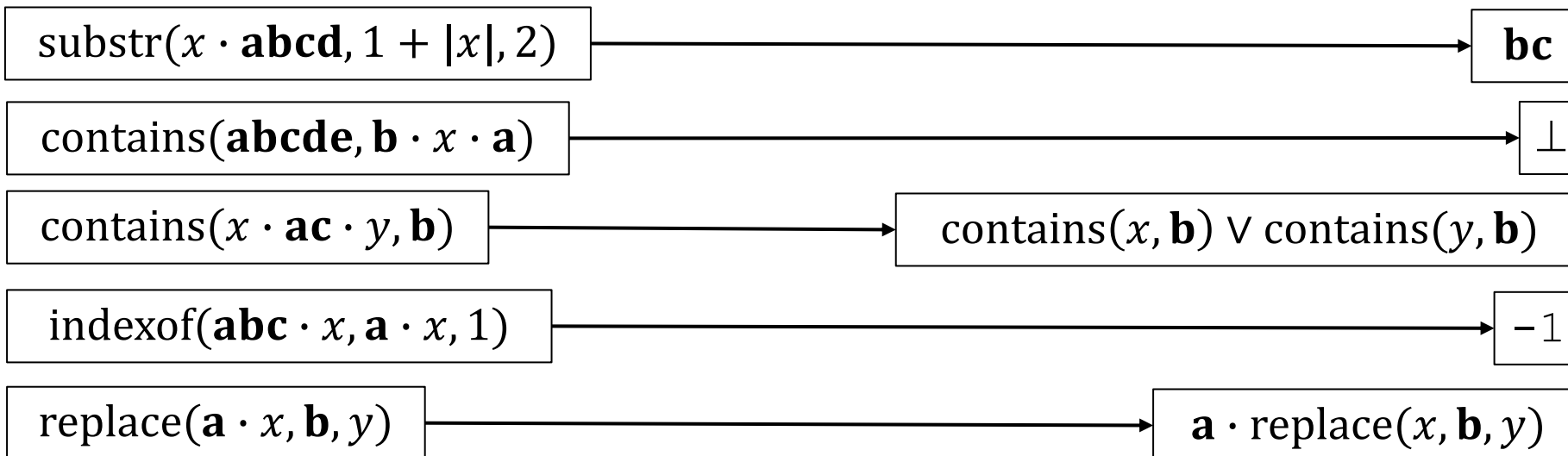
[Reynolds, Noetzli, Tinelli and Barrett, CAV'19]

Many Simplification Rules for Strings

Unlike arithmetic:

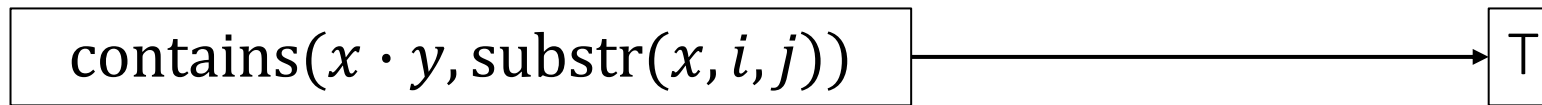
$$x + x + 7y = y - 4 \longrightarrow 2x + 6y + 4 = 0$$

... **simplification** rules for **strings** can be **quite complex**:



Abstraction-based **Rewriting**

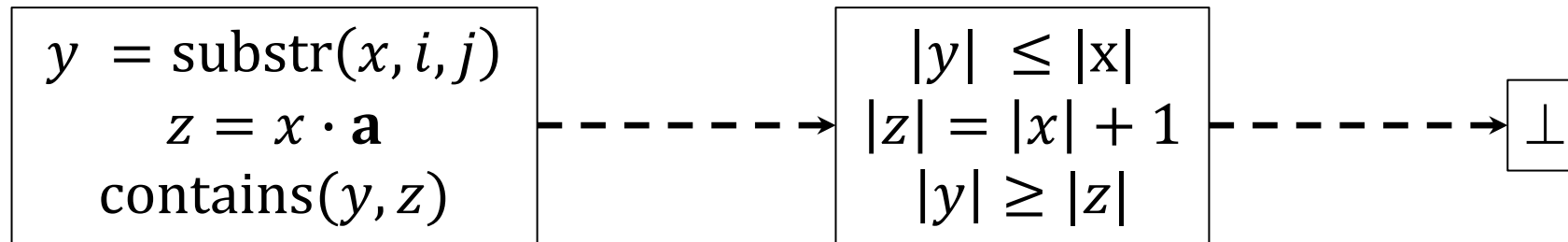
Considering the **string containment lattice**



(since $x \cdot y$ contains x , which contains $\text{substr}(x, \dots)$)

Abstraction-based Propagators

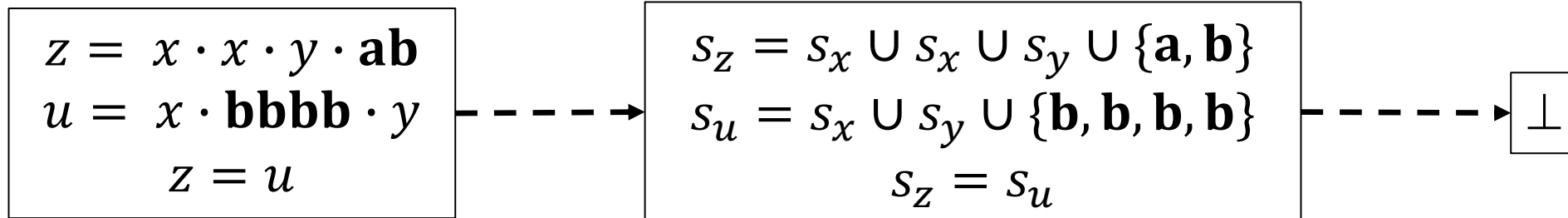
1. Abstracting strings by their length



$$|y| \geq |z| = |x| + 1 \geq |y| + 1 > |y|$$

Abstraction-based Propagators

2. Abstracting strings by their multiset of characters



(s_z contains one extra occurrence of \mathbf{a} than s_u)

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

- > 3,000 lines of C++ (and growing) for simplification rules in cvc5
- important aspect of modern string solving

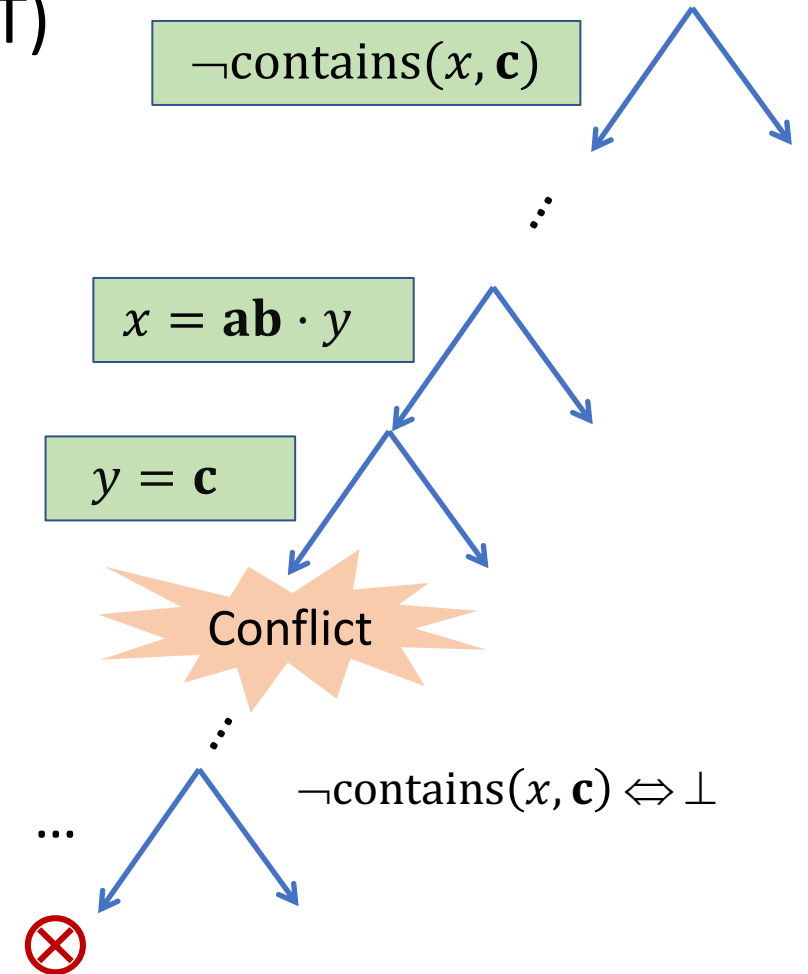
Even Faster Conflicts and Lazier Reductions

[Noetzli, Reynolds, Barbosa, Barrett and Tinelli, CAV'22]

Even **Faster Conflicts** and Lazier Reductions

Idea: apply simplifications **eagerly** during CDCL(T) search

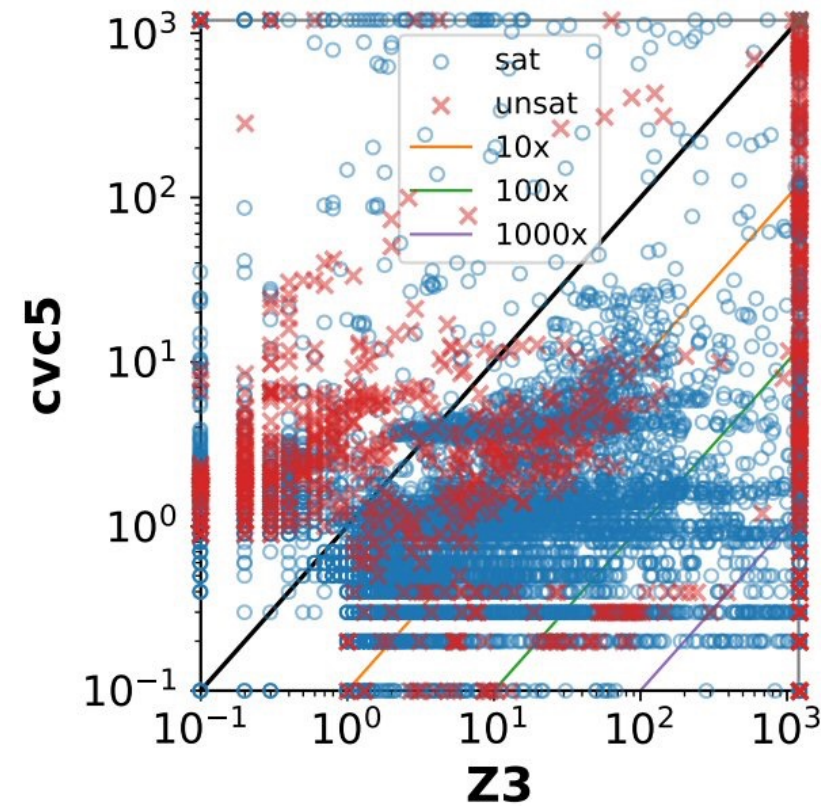
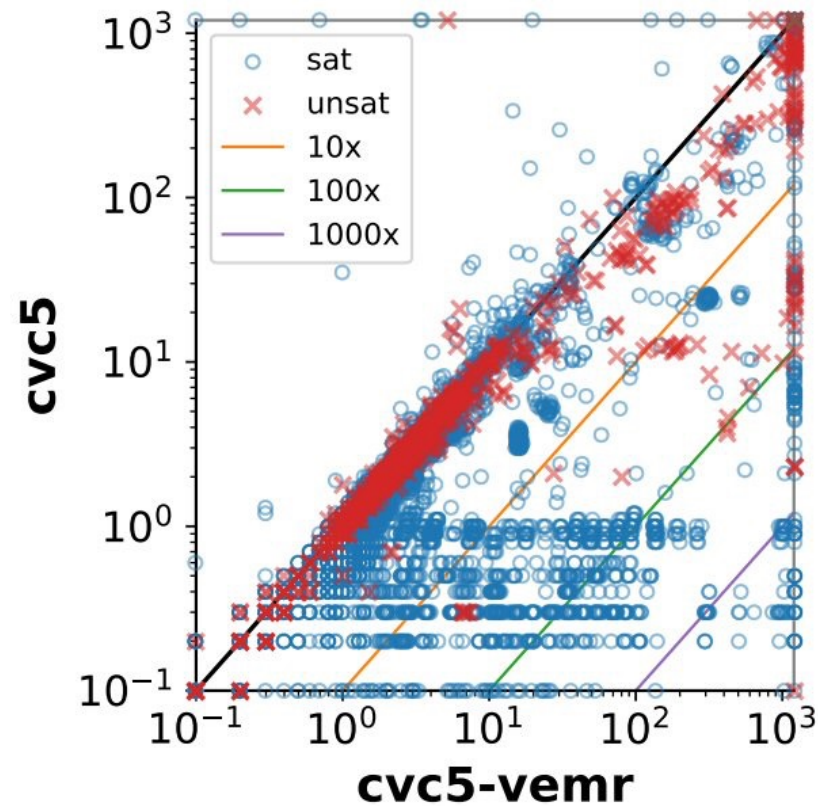
- Instrument congruence closure to detect conflicts via:
 - evaluation of concrete terms
 - inferred **properties of equivalence classes**
 - Upper/lower bounds for integer equivalence classes
 - Prefix and suffix approximations for string equivalence classes
- Report conflicts as soon as they arise
 - Avoids unnecessary expansion of extended functions



Even Faster Conflicts and Lazier Reductions

- Avoid reasoning about **unnecessary** reduction lemmas
- Regular expression inclusion tests
 - ⊗ E.g., do not reduce $x \in \Sigma^* \mathbf{a} \Sigma^*$ if already reduced $x \in \Sigma^* \mathbf{a} \Sigma^* \mathbf{b} \Sigma^*$ to \top
 - Since $\mathcal{L}(\Sigma^* \mathbf{a} \Sigma^* \mathbf{b} \Sigma^*) \subseteq \mathcal{L}(\Sigma^* \mathbf{a} \Sigma^*)$
 - Fast incomplete procedure for language inclusion
 - Can also be used for finding conflicts
- Model-based reductions
 - Construct candidate model \mathcal{M}
 - ⊗ Do not reduce, e.g., string predicates already satisfied by \mathcal{M}
 - Often, **negative** RE membership predicates are satisfied by current model

Even Faster Conflicts and Lazier Reductions



Results on 10,857 SMT-LIB string benchmarks; 1,200s timeout

- cvc5 solves 10,347; z3 solves 8,863

Witness Sharing + RE Elim

[Reynolds, Noetzli, Tinelli and Barrett, FMCAD'20]

Witness Sharing

Observation:

- There are often equivalent ways of expressing the same thing
 - E.g., string y is the result of removing the first character from string x :

$$\exists z. x = z \cdot y \wedge |z| = 1 \quad \text{substr}(x, 1, |x| - 1) = y \quad x \in \Sigma \cdot y$$

- Solving word equations, extended functions, and REs introduces **many fresh variables**

Idea:

- Formalize the definition for each introduced variable's **witness form**
- **Reuse variables** whose witness forms are semantically equivalent

Witness Sharing (Example)

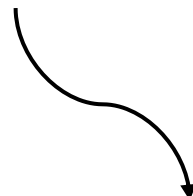
$$\frac{x \cdot w = \mathbf{a} \cdot u \quad |x| \neq 0}{x = \mathbf{a} \cdot k_1}$$

$$\frac{x \in \Sigma \cdot R}{x = k_2 \cdot k_3 \wedge k_2 \in \Sigma \wedge k_3 \in R}$$

Witness Sharing (Example)

$$\frac{x \cdot w = \mathbf{a} \cdot u \quad |x| \neq 0}{x = \mathbf{a} \cdot k_1}$$

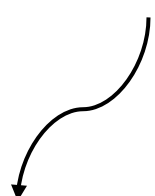
$$x = \mathbf{a} \cdot k_1$$



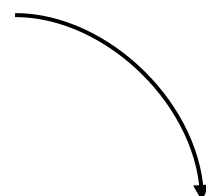
substr($x, 1, |x| - 1$)

$$\frac{x \in \Sigma \cdot R}{x = k_2 \cdot k_3 \wedge k_2 \in \Sigma \wedge k_3 \in R}$$

$$x = k_2 \cdot k_3 \wedge k_2 \in \Sigma \wedge k_3 \in R$$



substr($x, 0, 1$)



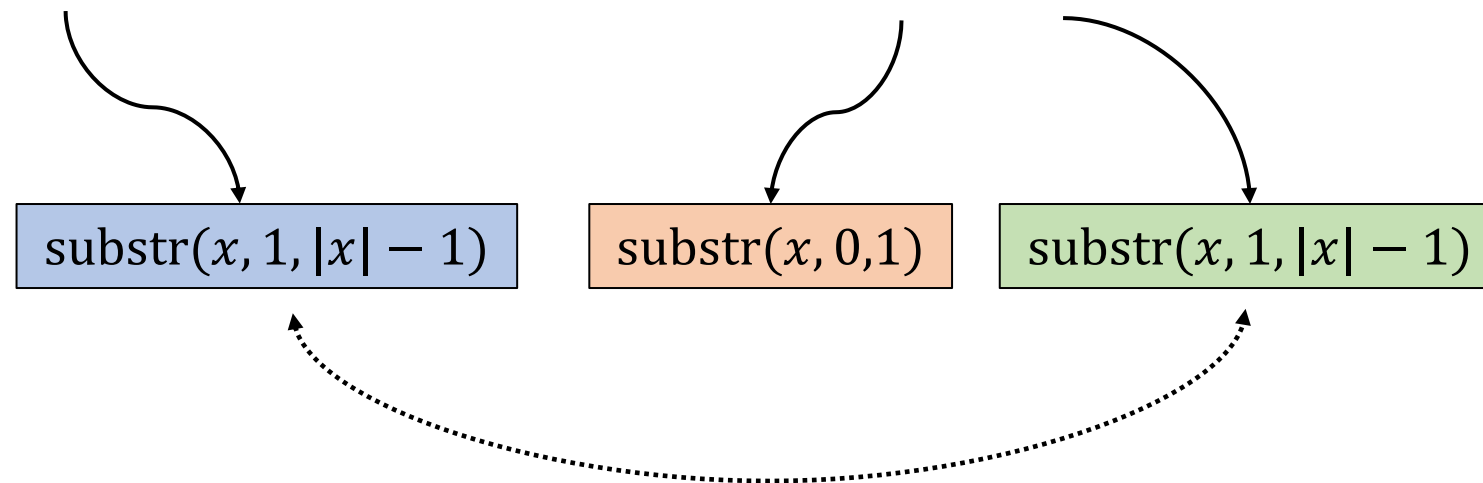
substr($x, 1, |x| - 1$)

witness forms

Witness Sharing (Example)

$$\frac{x \cdot w = \mathbf{a} \cdot u \quad |x| \neq 0}{x = \mathbf{a} \cdot k_1}$$

$$\frac{x \in \Sigma \cdot R}{x = k_2 \cdot k_3 \wedge k_2 \in \Sigma \wedge k_3 \in R}$$



Reuse variables whose witness form are (semantically) equivalent
 \Rightarrow Can use aggressive simplification to detect equivalent witness forms

Regular Expression Elimination

Idea: reduce REs to extended string constraints

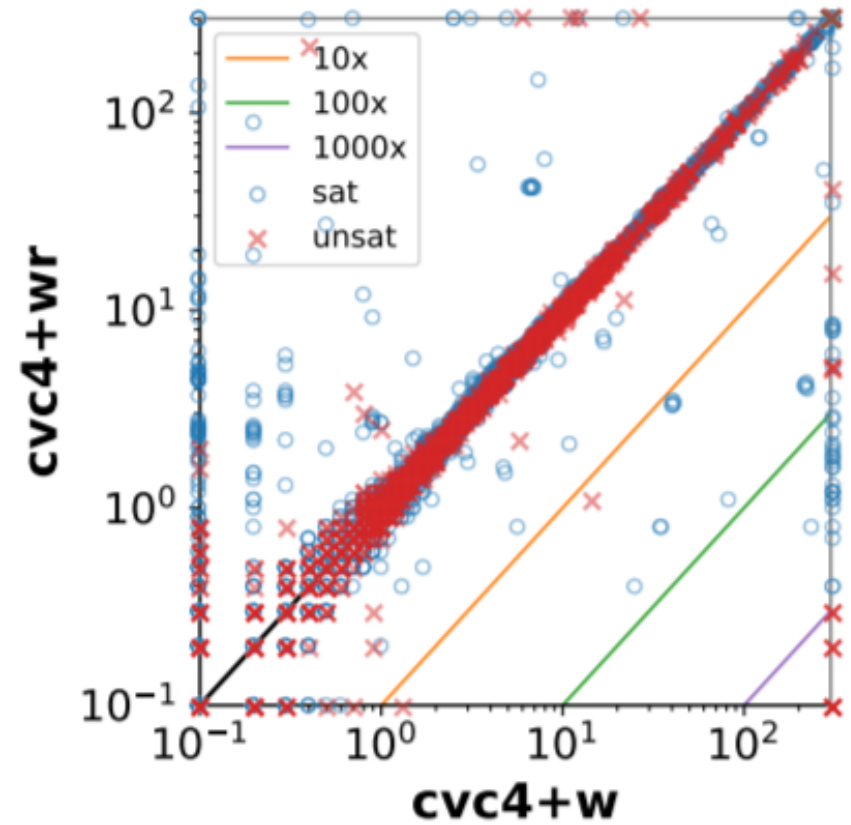
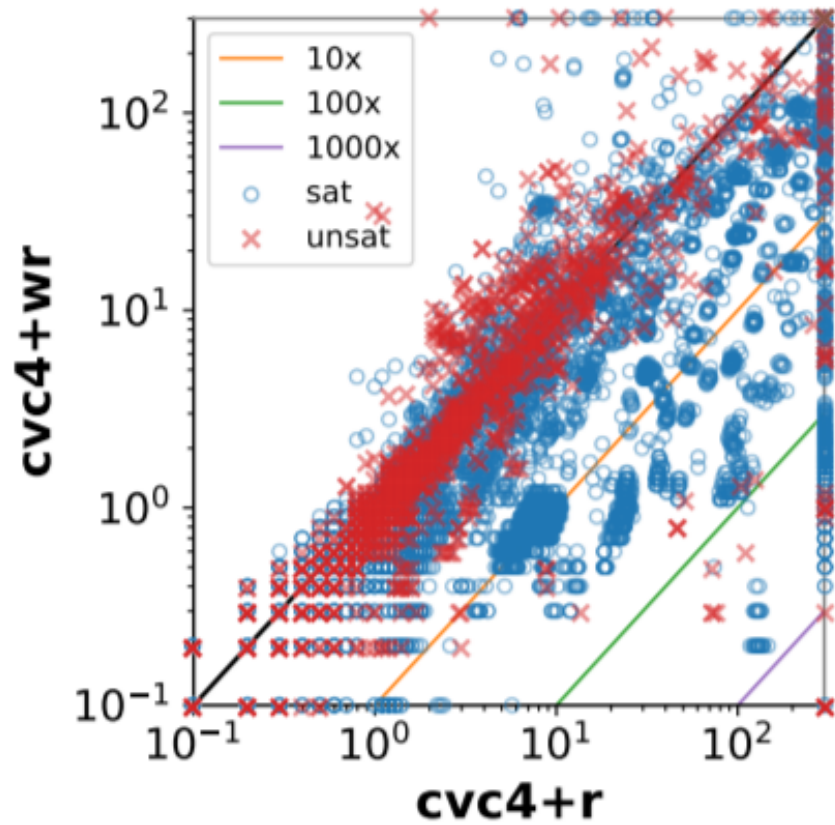
- Possible for many RE memberships occurring in practice

$$\boxed{x \in \Sigma \Sigma^* \Sigma} \iff \boxed{|x| \geq 2}$$

$$\boxed{x \in \Sigma^* \mathbf{abc} \Sigma^*} \iff \boxed{\text{contains}(x, \mathbf{abc})}$$

$$\boxed{x \in \Sigma^* \mathbf{a} \Sigma^* \mathbf{bcd} \Sigma^*} \iff \boxed{\begin{array}{l} \text{contains}(x, \mathbf{a}) \wedge \\ \text{contains}(\text{substr}(x, \text{indexof}(x, \mathbf{a}, 1) + 1, |x|), \mathbf{bcd}) \end{array}}$$

Impact of Witness Sharing + RE elim



String to Code Point Conversion

[Reynolds, Noetzli, Tinelli and Barrett, IJCAR'20]

Adding string-to-code operator `code`

Assume ordering on characters of alphabet Σ of size n :

- $c_1 < \dots < c_n$
- For each character c_i , we call i its *code point*

`code` : `String` \rightarrow `Int` is defined as follows:

1. `code(ci) = i` for all $c_i \in \Sigma$
2. `code(w) = -1` for all $w \in \Sigma^+$

Fragment with string length + code points (w/o concatenation):

- Devised a solving procedure that is **sound**, **complete**, and **terminating**

Reductions: Conversion Functions

Using `code` leads to efficient reductions, including:

- Conversion between strings and integers `toInt`:
 - ⊗ ... `ite(x[i] = 9, 9, ite(x[i] = 8, 8, ... ite(x[i] = 0, 0, -1) ...)`
 - ⇒ ... `ite(48 ≤ code(x[i]) ≤ 57, code(x[i]) - 48, -1)`
- Conversion between lowercase and uppercase strings `toLowerCase`:
 - ⊗ ... `ite(x[i] = A, a, ite(x[i] = B, b, ... ite(x[i] = Z, z, x[i]) ...)`
 - ⇒ ... `code(x[i]) + ite(65 ≤ code(x[i]) ≤ 90, 32, 0)`

Reductions: Conversion Functions

Using `code` leads to efficient reductions, including:

- Lexicographic ordering:

$$\otimes \quad x \leq y \Leftrightarrow \exists i \dots (x[i] = y[i] \vee (x[i] = \mathbf{a} \wedge y[i] = \mathbf{b}) \vee (x[i] = \mathbf{a} \wedge y[i] = \mathbf{c}) \dots)$$

$$\Rightarrow \quad x \leq y \Leftrightarrow \exists i \dots \mathbf{code}(x[i]) \leq \mathbf{code}(y[i])$$

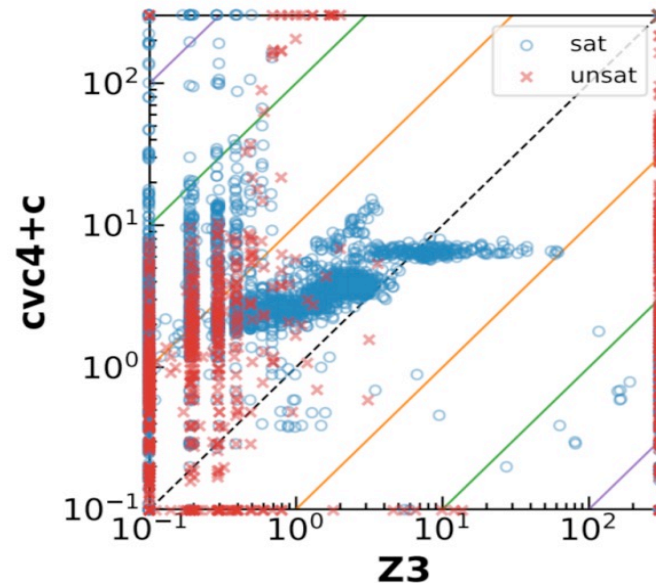
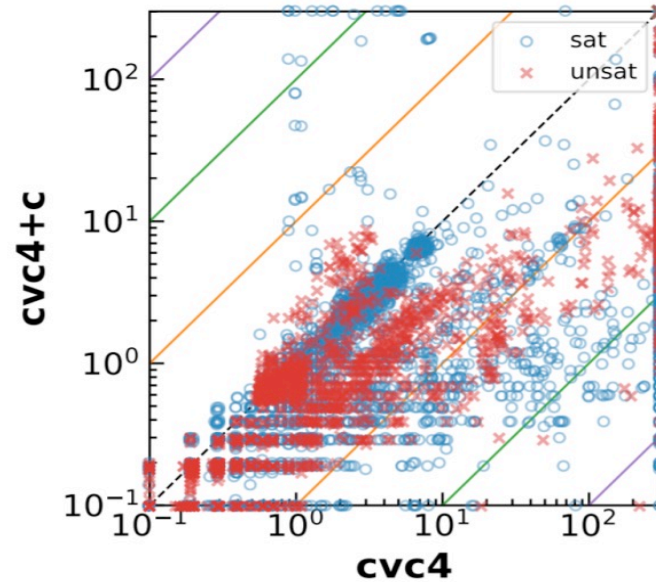
- Regular expression ranges:

$$\otimes \quad x \in \text{range}(c_1, c_2) \Leftrightarrow |x| = 1 \wedge (x = c_1 \vee \dots \vee x = c_2)$$

$$\Rightarrow \quad x \in \text{range}(c_1, c_2) \Leftrightarrow \mathbf{code}(c_1) \leq \mathbf{code}(x) \leq \mathbf{code}(c_2)$$

Experimental Results

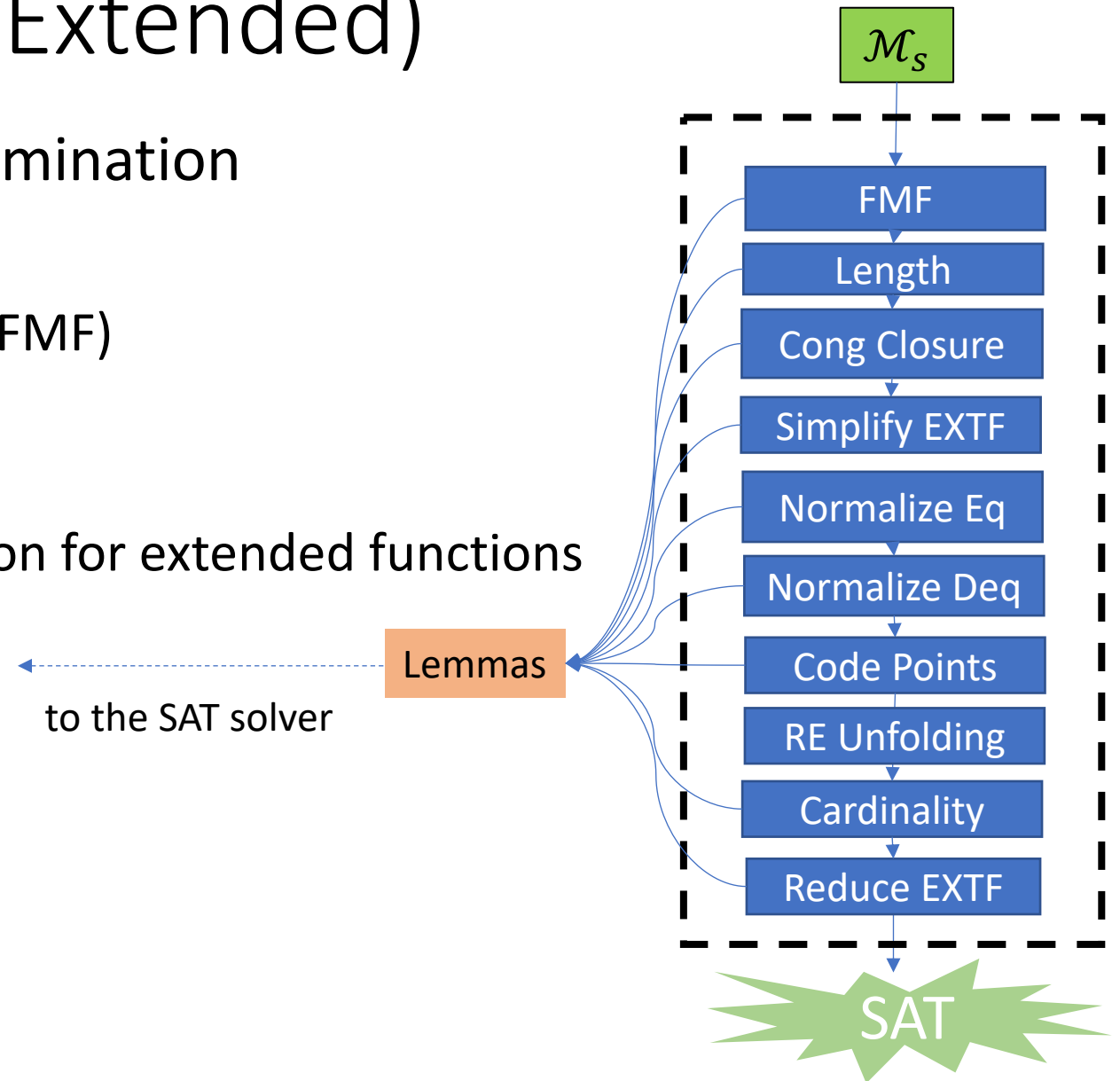
Benchmark Set		cvc4+c	cvc4	Z3
py-conbyte_cvc4	sat	1344	1104	1187
	unsat	8576	8547	8482
	×	13	282	264
py-conbyte_trauc	sat	1009	929	697
	unsat	1424	1407	1428
	×	13	110	321
py-conbyte_z3seq	sat	1354	1126	1343
	unsat	5864	5797	5719
	×	35	330	191
py-conbyte_z3str	sat	711	652	692
	unsat	1227	1223	1223
	×	3	66	26
Total	sat	4418	3811	3919
	unsat	17091	16974	16852
	×	64	788	802



- 10x t/o reduction
- Faster runtimes
- Improvement wrt state of the art

String Theory Solver (Extended)

- Preprocess based on reg-exp elimination
- Then, run inference strategy:
 1. Split on sum of lengths bound (FMF)
 2. Elaborate length constraints
 3. Congruence closure
 4. Context-dependent simplification for extended functions
 5. Normalize string equalities
 6. Normalize string disequalities
 7. Subprocedure for code points
 8. Regular expression unfolding
 9. Check cardinality constraints
 10. Reduce extended functions



Conclusions

SMT solvers can provide:

- Efficient (incomplete) procedure for word equations with length
- FMF, context-dependent simplification, RE elimination, witness sharing, ...

Ongoing work in cvc5:

- Proofs and proof certificates
 - Array-like reasoning (update + slices)
- cvc5 is open-source, available at <https://cvc5.github.io/>
 - Also supports theory of sequences, further extensions



Thanks for listening!