

On Regular Constraint Propagation for Solving String Constraints

Philipp Rümmer
University of Regensburg
Uppsala University

EPFL, 2024-06-15

*Joint work with
Artur Jež, Matt Hague, Anthony W. Lin, Oliver Markgraf, and others*

What are String Constraints?

Strings =

- finite sequences of letters over a finite alphabet (e.g., Unicode)

String constraints =

- quantifier-free formulas over string variables, including operations like:
 - String concatenation (word equations)
 - String length
 - Substring, character access
 - Regular expressions

Strings in Verification

```
// Pre = (true)
String s= '';
// P1 = (s ∈ ε)
while(*){
    // P2 = (s = u · v ∧ u ∈ a* ∧ v ∈ b* ∧ |u| = |v|)
    s= 'a' + s + 'b';
}
// P3 = P2
assert(!s.contains('ba') && (s.length() % 2) == 0);
// Post = P3
```

Strings in Verification

ASCII, Unicode

```
// Pre = (true)
String s = '';
// P1 = (s ∈ ε)
while (*) {
    // P2 = (s = u · v ∧ u ∈ a* ∧ v ∈ b* ∧ |u| = |v|)
    s = 'a' + s + 'b';
}
// P3 = P2
assert(!s.contains('ba') && (s.length() % 2) == 0);
// Post = P3
```

Strings in Verification

Regular expression

assertion: $s = \epsilon$

```
// Pre = ( )
String s = '';
// P1 = (s ∈  $\epsilon$ )
while(*){
    // P2 = (s = u · v ∧ u ∈ a* ∧ v ∈ b* ∧ |u| = |v|)
    s = 'a' + s + 'b';
}
// P3 = P2
assert(!s.contains('ba') && (s.length() % 2) == 0);
// Post = P3
```

Strings in Verification

```
// Pre = (true)
String s = '';
// P1 = (s ∈ ε)
while (*) {
    // P2 = (s = u · v ∧ u ∈ a* ∧ v ∈ b* ∧ |u| = |v|)
    s = 'a' + s + 'b';
}
// P3 = P2
assert(!s.contains('ba') && (s.length() % 2) == 0);
// Post = P3
```

Word/string
concatenation

Strings in Verification

```
// Pre = (true)
String s = '';
// P1 = (s ∈ ε)
while(*){
    // P2 = (s = u · v ∧ u ∈ a* ∧ v ∈ b* ∧ |u| = |v|)
    s = 'a' + s + 'b';
}
// P3 = P2
assert(!s.contains('ba') && (s.length() % 2) == 0);
// Post = P3
```

Loop invariant combining
word equations,
regex constraints,
length constraints

Strings in Verification

```
// Pre = (true)
String s = '';
// P1 = (s ∈ ε)
while (*) {
    // P2 = (s = u · v ∧ u ∈ a* ∧ v ∈ b* ∧ |u| = |v|)
    s = 'a' + s + 'b';
}
// P3 = P2
assert(!s.contains('ba') && (s.length() % 2) == 0);
// Post = P3
```

**Substring
constraint**

Strings in Verification

```
// Pre = (true)
String s = '';
// P1 = (s ∈ ε)
while (*) {
    // P2 = (s = u · v ∧ u ∈ a* ∧ v ∈ b* ∧ |u| = |v|)
    s = 'a' + s + 'b';
}
// P3 = P2
assert (!s.contains('ba') && (s.length() % 2) == 0);
// Post = P3
```

Or regex:

$s \notin \Sigma^* \cdot ba \cdot \Sigma^*$

Strings in Verification

```
// Pre = (true)
String s = '';
// P1 = (s ∈ ε)
while (*) {
    // P2 = (s = u · v ∧ u ∈ a* ∧ v ∈ b* ∧ |u| = |v|)
    s = 'a' + s + 'b';
}
// P3 = P2
assert(!s.contains('ba') && (s.length() % 2) == 0);
// Post = P3
```

**Presburger
length constraint**

Strings in Verification

```
// Pre = (true)
String s= '';
// P1 = (s ∈ ε)
while(*){
    // P2 = (s = u · v ∧ u ∈ a* ∧ v ∈ b* ∧ |u| = |v|)
    s= 'a' + s + 'b';
}
// P3 = P2
assert(!s.contains('ba') && (s.length() % 2) == 0);
// Post = P3
```

Strings in Verification

```
// Pre = (true)
String s = '';
// P1 = (s ∈ ε)
while(*){
    // P2 = (s = u · v ∧ u ∈ a* ∧ v ∈ b* ∧ |u| = |v|)
    s = 'a' + s + 'b';
}
// P3 = P2
assert(!s.contains('ba') && (s.length() % 2) == 0);
// Post = P3
```

Main application area of string solving:
security analysis

A billion SMT queries a day

CAV keynote lecture by the director of applied science for AWS Identity explains how AWS is making the power of automated reasoning available to all customers.

By [Neha Rungta](#)

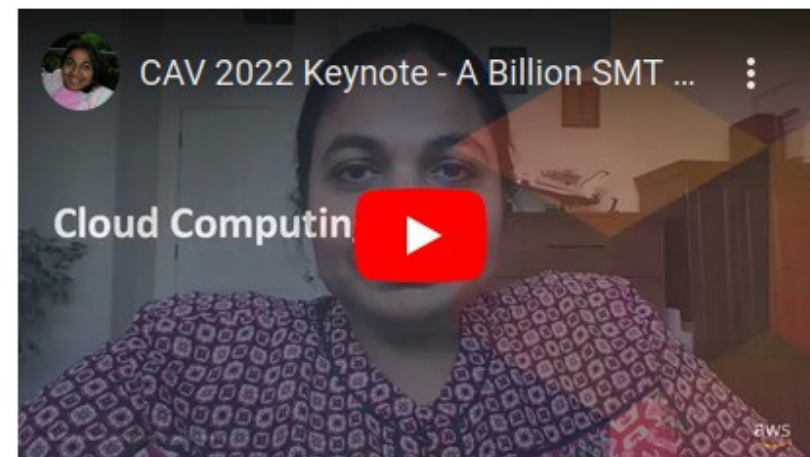
August 18, 2022

 Share

At this year's Computer-Aided Verification (CAV) conference — a leading automated-reasoning conference collocated with the Federated Logic Conferences (FLoC) — Amazon's Neha Rungta delivered a keynote talk in which she suggested that innovations at Amazon have "ushered in the golden age of automated reasoning".

Amazon scientists and engineers are using automated reasoning to prove the correctness of critical internal systems and to help customers prove the security of their cloud infrastructures. Many of these innovations are being driven by powerful reasoning engines called SMT solvers.

Satisfiability problems, or SAT, ask whether it's possible to assign variables true/false values that satisfy a set of



A billion SMT queries a day

Neha Rungta's 2022 CAV keynote

Satisfiability Modulo Theories

- Many new SMT solvers for strings have emerged in the last years
- SMT-LIB theory of strings
- SMT-COMP has a `QF_Strings` division

Some String Solvers

- Hampi
- Kaluza
- Stranger
- Gecode+S
- GStrings
- Z3
- Z3-str/2/3/4/alpha
- CVC4/cvc5
- Norn
- S3/p/#
- TRAU
- nfa2sat
- OSTRICH
- Z3-Noodler
- Woorpje
- BEK, REX
- SLOG, SLENT
- *(many more)*

What's Decidable about ...

$$z = w_1 \cdot y \cdot w_2 \cdot x \cdot w_3$$

Word Equations

What's Decidable about ...

$$z = w_1 \cdot y \cdot w_2 \cdot x \cdot w_3$$

Word Equations

This is already in
solved form

What's Decidable about ...

In general, two main decision procedures:
Makanin & Recompression

This is already in
solved form

$$z = w_1 \cdot y \cdot w_2 \cdot x \cdot w_3$$

Word Equations

What's Decidable about ...

In general, two main
decision procedures:
Makanin & Recompression

This is already in
solved form

$$z = w_1 \cdot y \cdot w_2 \cdot x \cdot w_3$$



Word Equations

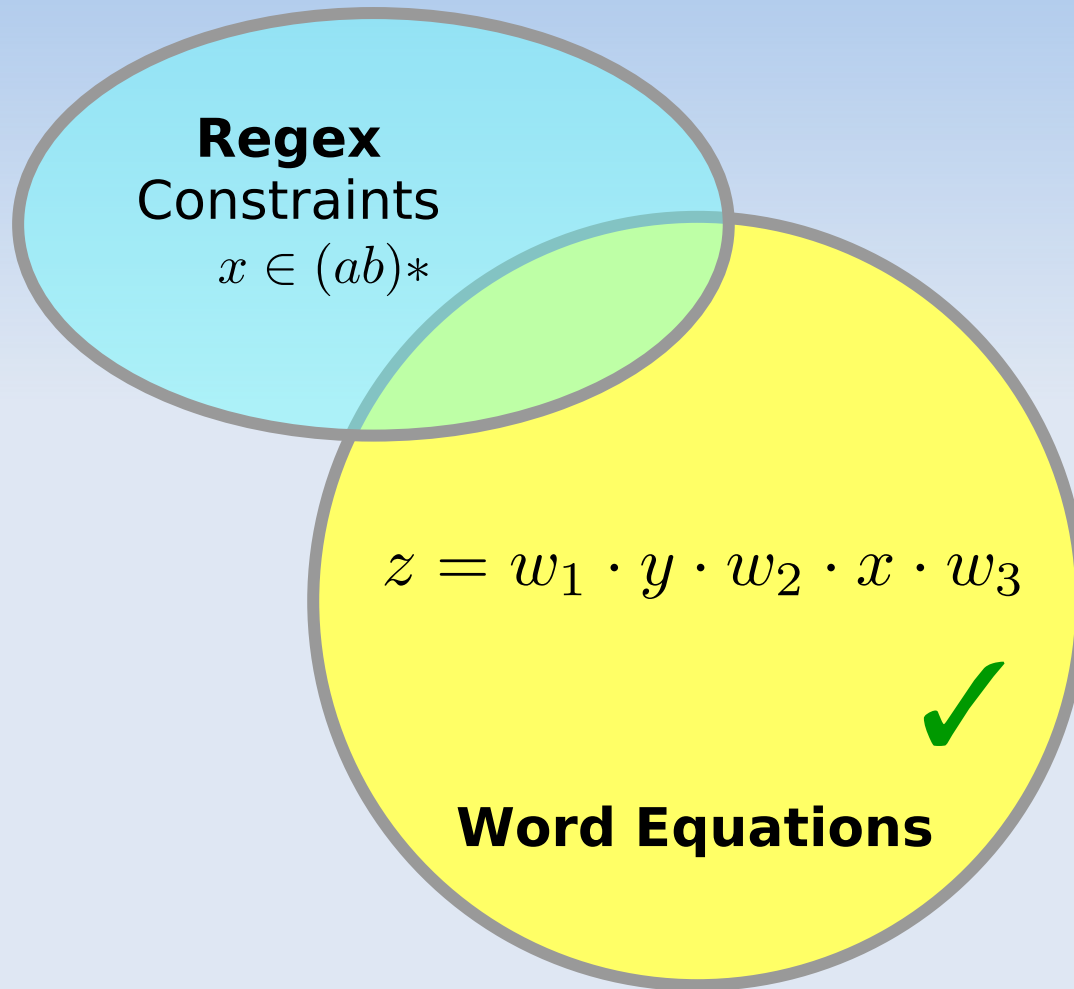
What's Decidable about ...

$$z = w_1 \cdot y \cdot w_2 \cdot x \cdot w_3$$

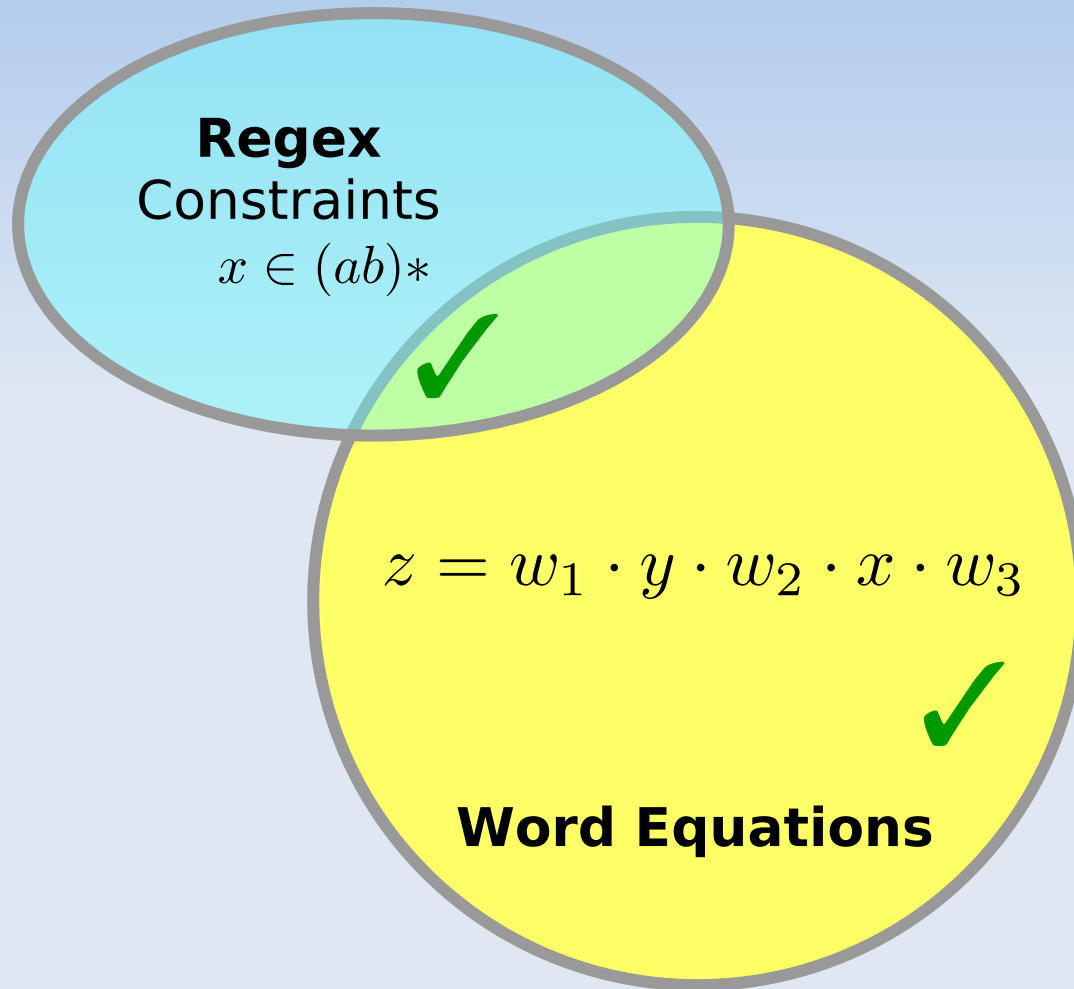


Word Equations

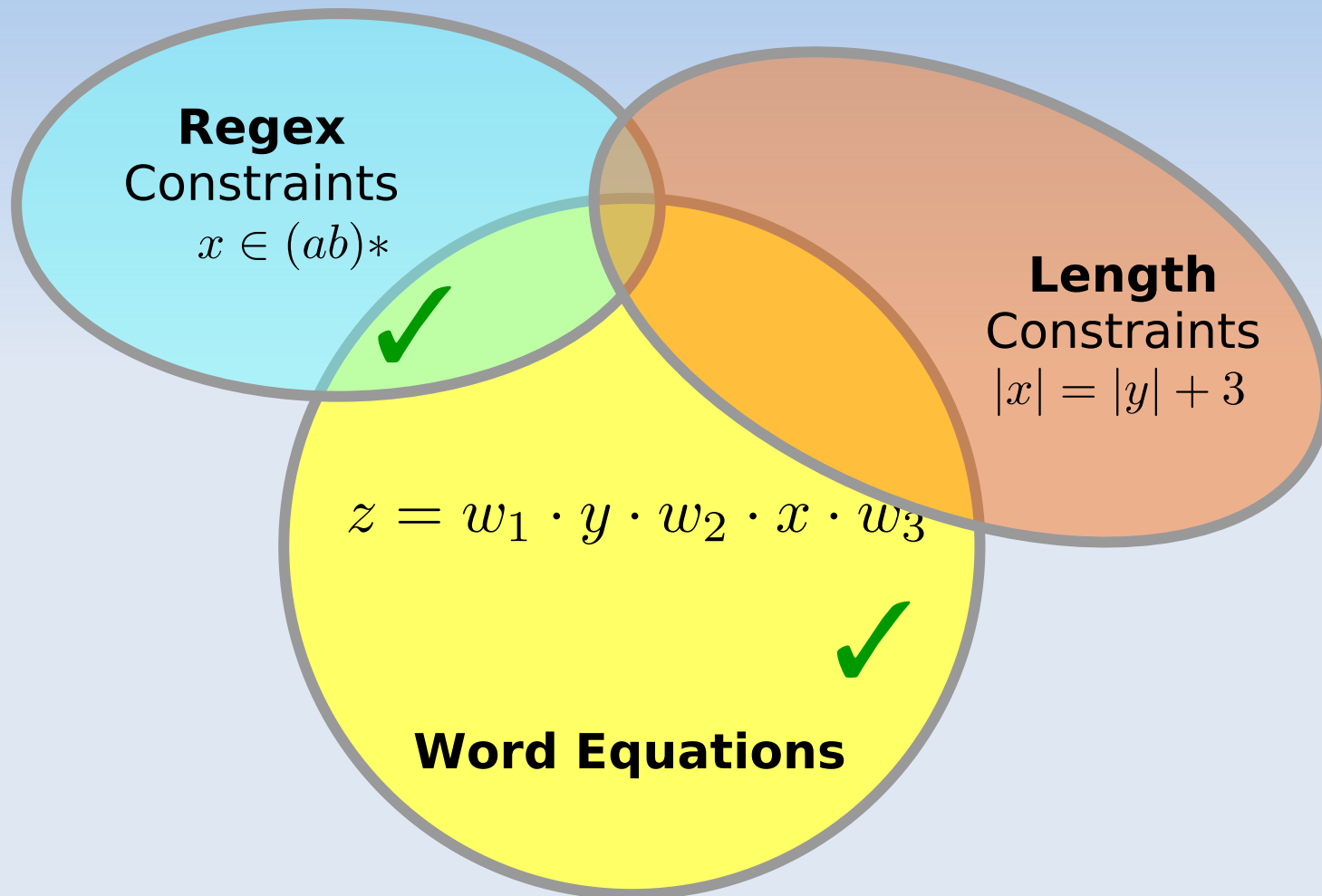
What's Decidable about ...



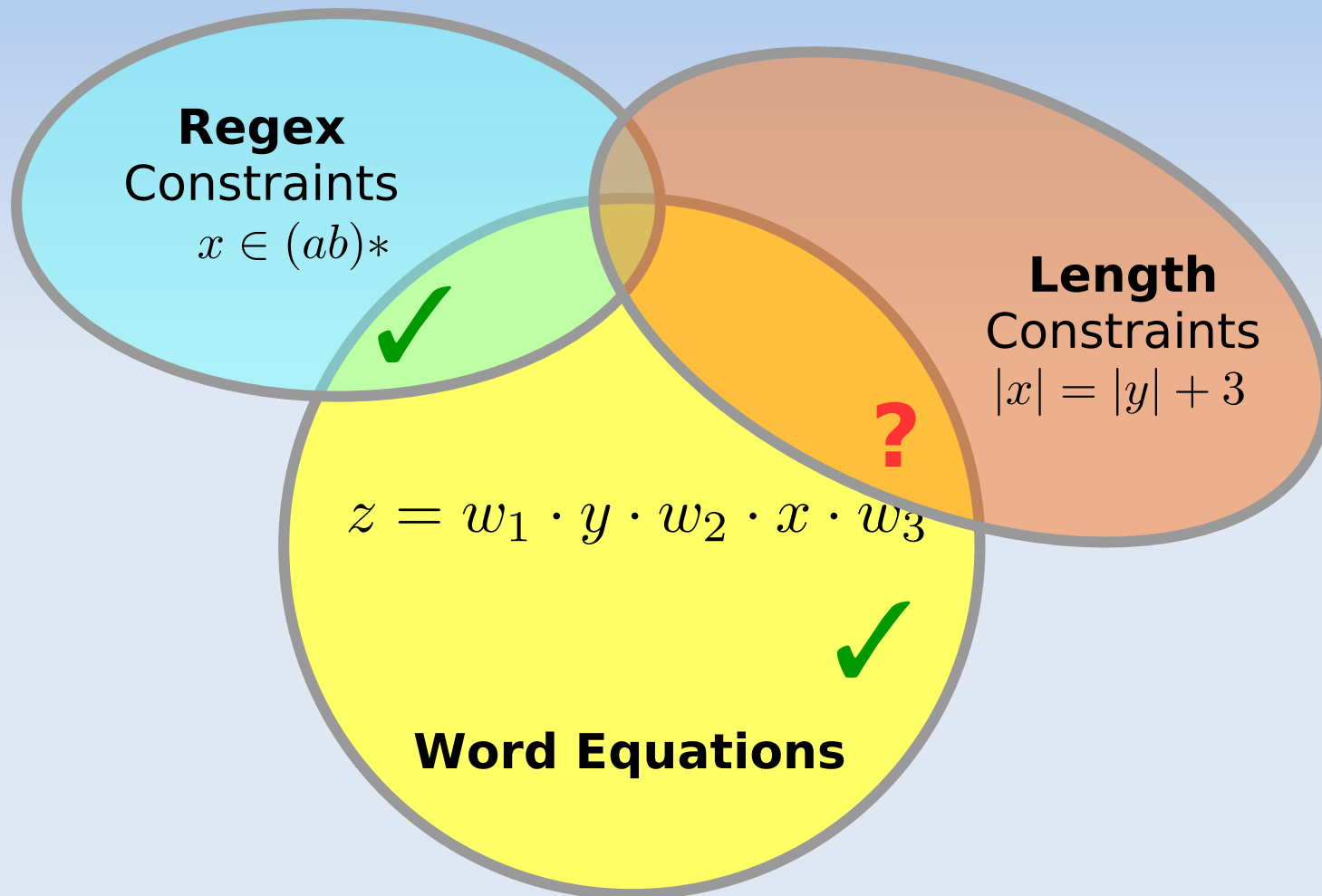
What's Decidable about ...



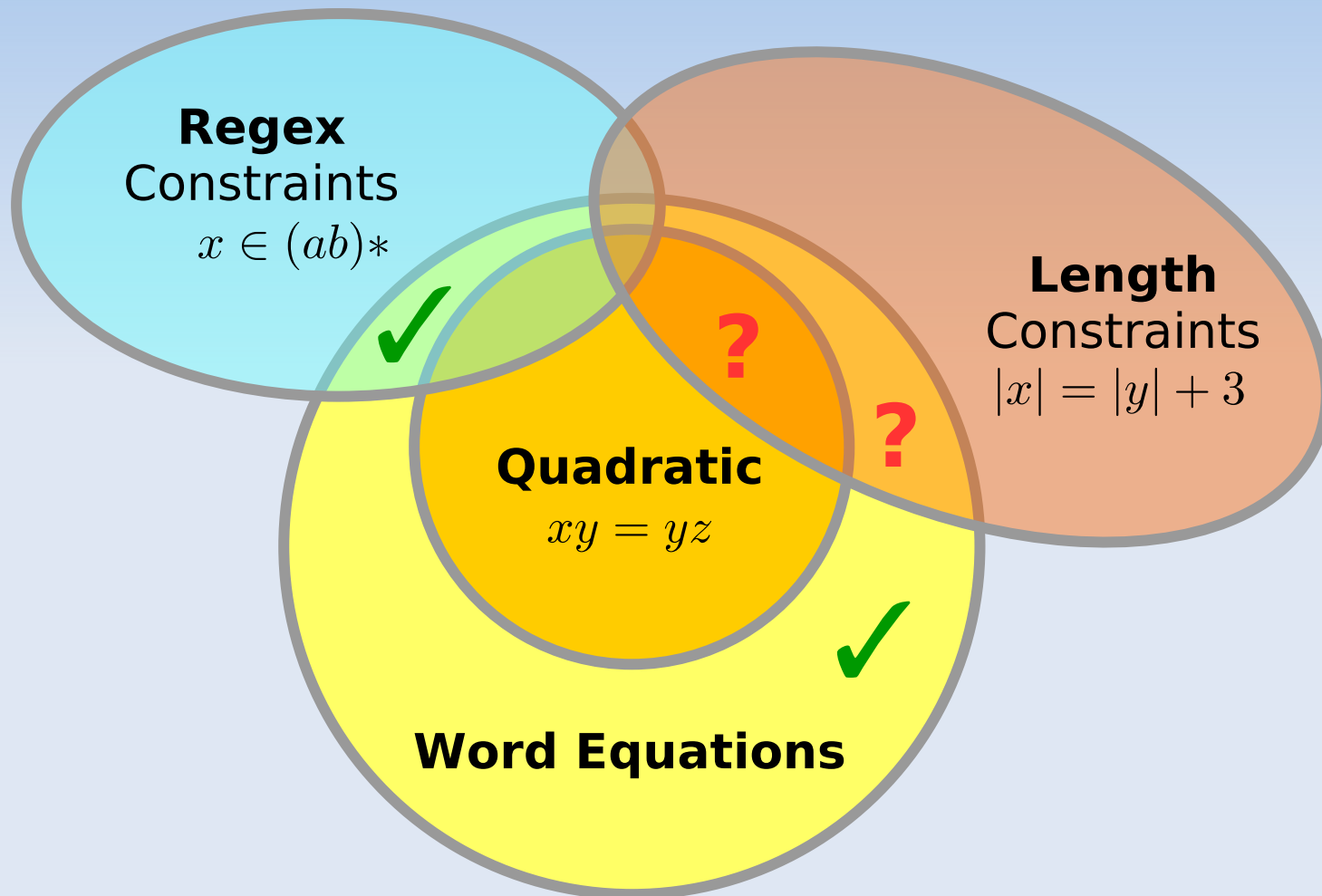
What's Decidable about ...



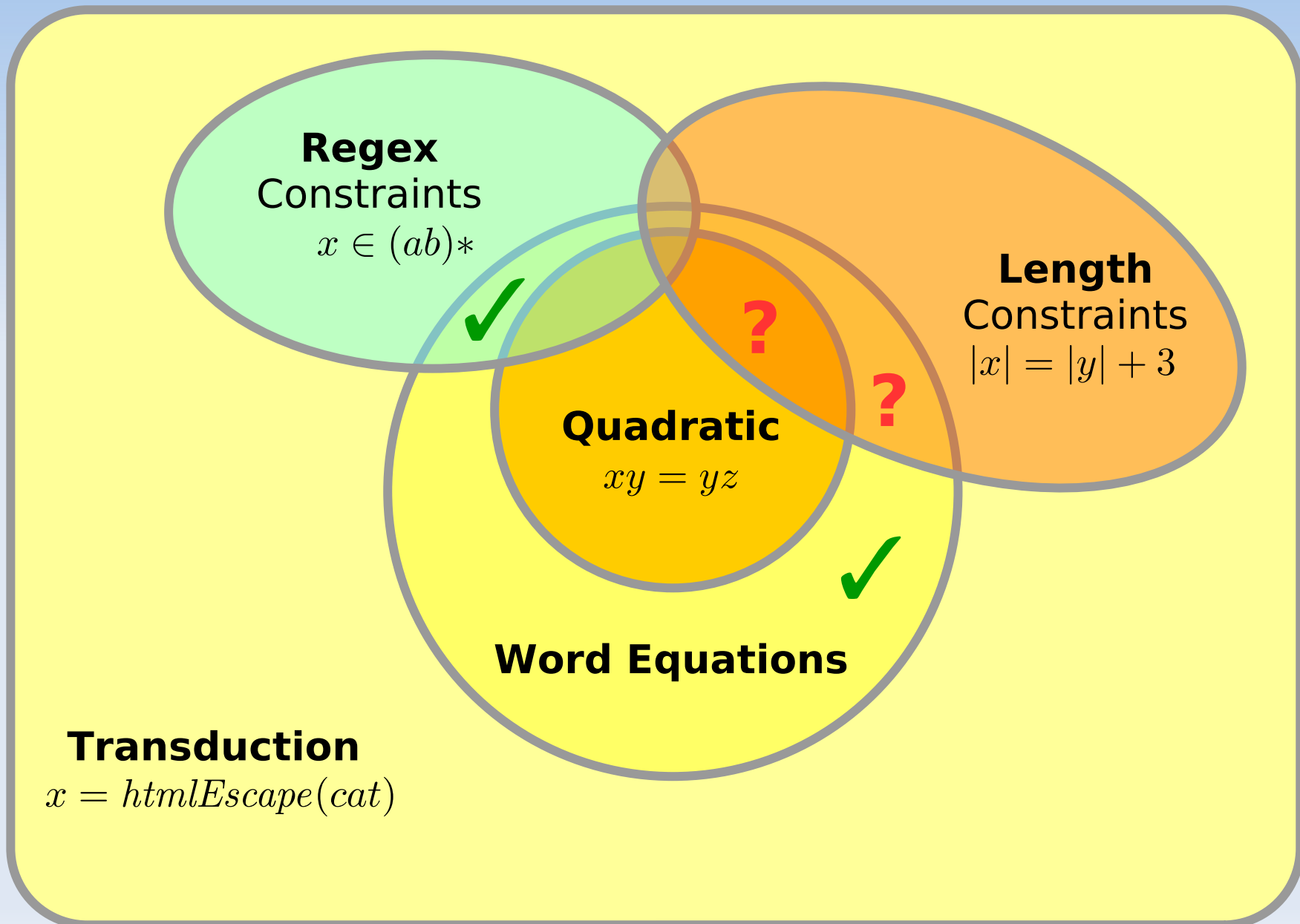
What's Decidable about ...



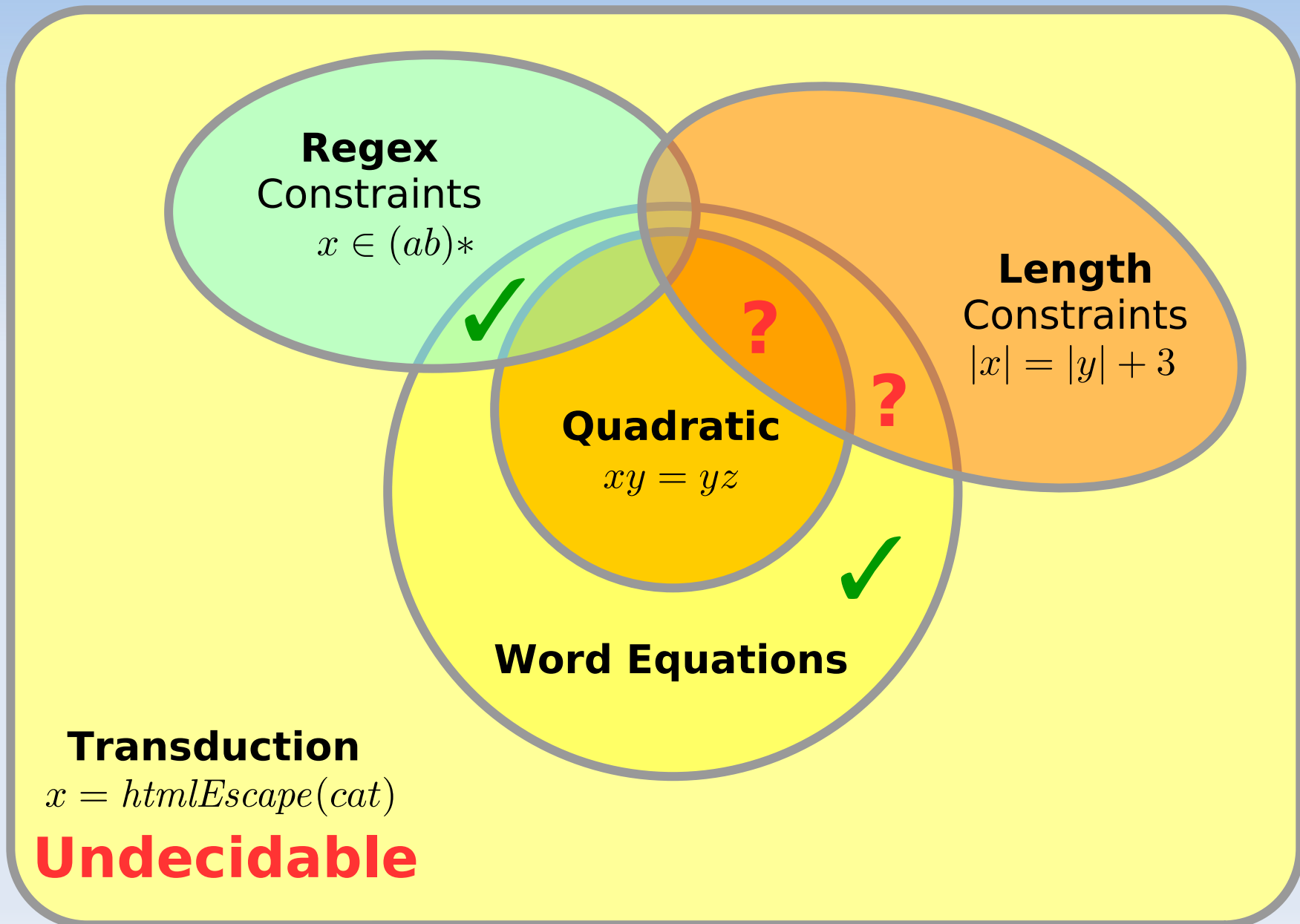
What's Decidable about ...



What's Decidable about ...



What's Decidable about ...



In Practice

Solvers use wide variety of techniques:

- Encoding as bit-vectors
- Encoding as SAT problem
- Rewriting/simplification rules
- Automata methods, derivatives
- Splitting rules for word equations
- (Re)Compression
- Propagation and CP methods
- *etc.*

Completeness

- Already for just word equations, decision procedures are **complicated** and **impractical** (not implementable?)
- String solvers are generally **incomplete** (for proving word equations unsat)

Completeness


- Already for just word equations, decision procedures are **complicated** and **impractical** (not implementable?)
- String solvers are generally **incomplete** (for proving word equations unsat)
- *Simpler decision procedures exist for various fragments (and are implemented by some solvers)*

Some Identified Fragments

- Quadratic word equations
- Tree-shaped
- Acyclic
- Chainfree
- Straightline
- Cost-enriched straightline
- Weakly chaining

Some Identified Fragments

- Quadratic word equations
- Tree-shaped
- Acyclic
- Chainfree
- Straightline
- Cost-enriched straightline
- Weakly chaining



Sometimes
complicated
definitions

Ad-hoc decision
procedures

Relationship
often unclear

Some Identified Fragments

- Quadratic word equations
- Tree shaped
- Chain
- Straightline
- Cost-enriched straightline
- Weakly chaining

Is it possible to unify definitions?

Possible to have a **common proof format?**

Ad-hoc decision
procedures

Relationship
often unclear

Our Working Hypothesis

(Most/many/some/...) decision procedures can be understood as a combination of:

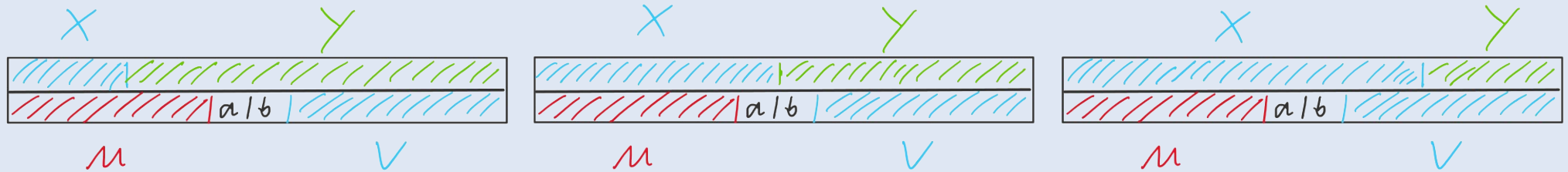
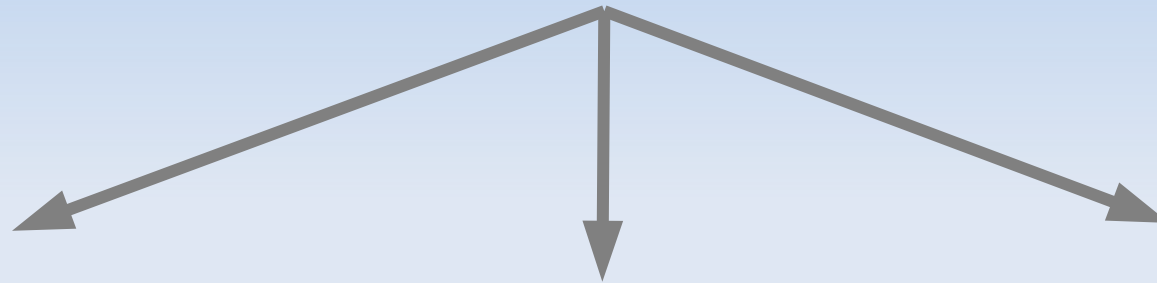
- Nielsen's transformation
- Propagation of regular constraints
- (Simplification rules)

Nielsen Transformation

$$x \cdot y = u \cdot 'ab' \cdot v$$

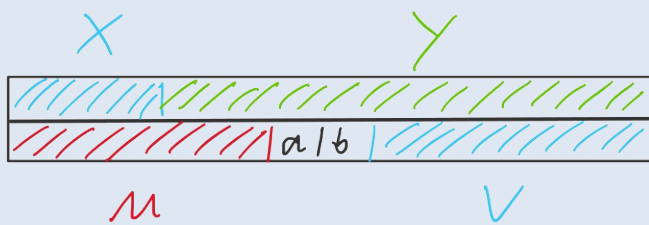
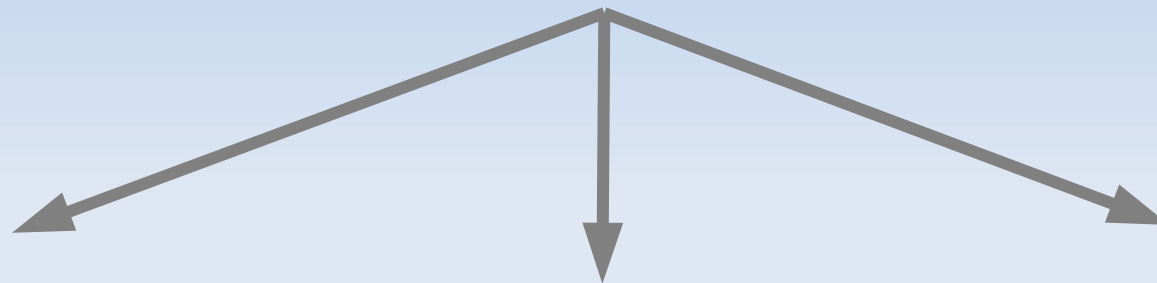
Nielsen Transformation

$$x \cdot y = u \cdot 'ab' \cdot v$$

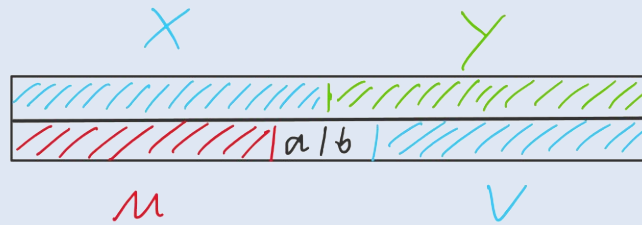


Nielsen Transformation

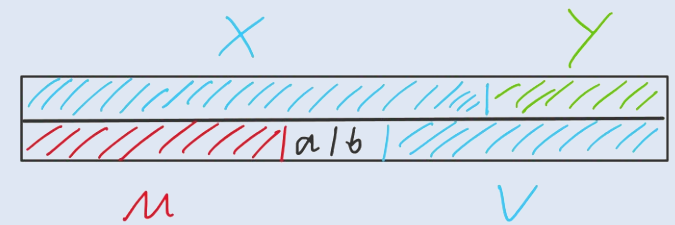
$$x \cdot y = u \cdot 'ab' \cdot v$$



$$\begin{aligned} x &= u_1 \\ \wedge y &= u_2 \cdot 'ab' \cdot v \\ \wedge u &= u_1 \cdot u_2 \end{aligned}$$



$$\begin{aligned} x &= u \cdot 'a' \\ \wedge y &= 'b' \cdot v \end{aligned}$$

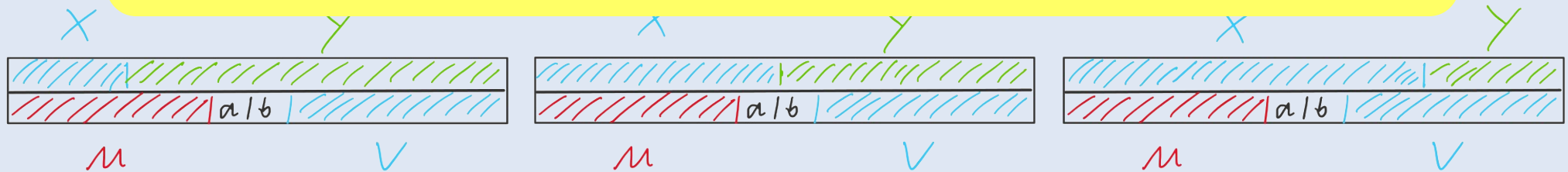


$$\begin{aligned} x &= u \cdot 'ab' \cdot v_1 \\ \wedge y &= v_2 \\ \wedge v &= v_1 \cdot v_2 \end{aligned}$$

Nielsen Transformation

$$x \cdot y = u \cdot 'ab' \cdot v$$

In general, this style of reasoning does not terminate. Most solvers split equations anyway.



$$\begin{aligned} x &= u_1 \\ \wedge y &= u_2 \cdot 'ab' \cdot v \\ \wedge u &= u_1 \cdot u_2 \end{aligned}$$

$$\begin{aligned} x &= u \cdot 'a' \\ \wedge y &= 'b' \cdot v \end{aligned}$$

$$\begin{aligned} x &= u \cdot 'ab' \cdot v_1 \\ \wedge y &= v_2 \\ \wedge v &= v_1 \cdot v_2 \end{aligned}$$

Our Working Hypothesis

(Most/many/some/...) decision procedures can be understood as a combination of:

- Nielsen's transformation
- Propagation of regular constraints
- (Simplification rules)

Constraint Normalization

String constraints can be rewritten to Boolean combinations of:

- Equations $x = y$
- Function applications $x = f(\bar{y})$
- Membership predicates $x \in \mathcal{L}$
- (ignoring length)

$x \cdot y = 'ab'$ becomes $z = \text{concat}(x, y) \wedge z \in ab$

Regular Constraint Prop.

$$\begin{array}{l}
 \notin \frac{\Gamma, x \in e^c}{\Gamma, x \notin e} \quad \neq \frac{\Gamma, x \neq y, y = f(x_1, \dots, x_n)}{\Gamma, x \neq f(x_1, \dots, x_n)} \text{ where } y \text{ is fresh} \quad \text{CUT} \frac{\Gamma, x \in e}{\Gamma} \quad \frac{\Gamma, x \in e^c}{\Gamma} \\
 \\
 =\text{-PROP} \frac{\Gamma, x \in e, x = y, y \in e}{\Gamma, x \in e, x = y} \quad \neq\text{-SUBSUME} \frac{\Gamma, x \in e_1, y \in e_2}{\Gamma, x \in e_1, x \neq y, y \in e_2} \text{ if } \mathcal{L}(e_1) \cap \mathcal{L}(e_2) = \emptyset \\
 \\
 =\text{-PROP-ELIM} \frac{\Gamma, x \in e, y \in e}{\Gamma, x \in e, x = y} \text{ if } |\mathcal{L}(e)| = 1 \quad \neq\text{-PROP-ELIM} \frac{\Gamma, x \in e, y \in e^c}{\Gamma, x \in e, x \neq y} \text{ if } |\mathcal{L}(e)| = 1 \\
 \\
 \text{CLOSE} \frac{}{\Gamma, x \in e_1, \dots, x \in e_n} \text{ if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \emptyset \\
 \\
 \text{SUBSUME} \frac{\Gamma, x \in e_1, \dots, x \in e_n}{\Gamma, x \in e, x \in e_1, \dots, x \in e_n} \text{ if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) \subseteq \mathcal{L}(e) \\
 \\
 \text{INTERSECT} \frac{\Gamma, x \in e}{\Gamma, x \in e_1, \dots, x \in e_n} \text{ if } n > 1 \text{ and } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \mathcal{L}(e) \\
 \\
 \text{FWD-PROP} \frac{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \text{ if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n)) \\
 \\
 \text{FWD-PROP-ELIM} \frac{\Gamma, x \in e, x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \text{ if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n)) \text{ and } |\mathcal{L}(e)| = 1 \\
 \\
 \text{BWD-PROP} \frac{\{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1^i, \dots, x_n \in e_n^i\}_{i=1}^k}{\Gamma, x \in e, x = f(x_1, \dots, x_n)} \text{ if } f^{-1}(\mathcal{L}(e)) = \bigcup_{i=1}^k (\mathcal{L}(e_1^i) \times \dots \times \mathcal{L}(e_n^i))
 \end{array}$$

String Constraints with Real-World Regular Expressions. Taolue Chen, Matthew Hague, Zhilei Han, Denghang Hu, Alejandro Flores-Lamas, Anthony W. Lin, Shuanglong Kan, PR, Zhilin Wu, POPL'22

$$\text{CLOSE} \frac{}{\Gamma, x \in e_1, \dots, x \in e_n}$$

$$\text{if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \emptyset$$

$$\notin \frac{\Gamma, x \in e^c}{\Gamma, x \notin e} \quad \neq \frac{\Gamma, x \neq y, y = f(x_1, \dots, x_n)}{\Gamma, x \neq f(x_1, \dots, x_n)} \quad \text{where } y \text{ is fresh} \quad \text{CUT} \frac{\Gamma, x \in e \quad \Gamma, x \in e^c}{\Gamma}$$

$$\text{=-PROP} \frac{\Gamma, x \in e, x = y, y \in e}{\Gamma, x \in e, x = y} \quad \neq\text{-SUBSUME} \frac{\Gamma, x \in e_1, y \in e_2}{\Gamma, x \in e_1, x \neq y, y \in e_2} \quad \text{if } \mathcal{L}(e_1) \cap \mathcal{L}(e_2) = \emptyset$$

$$\text{=-PROP-ELIM} \frac{\Gamma, x \in e, y \in e}{\Gamma, x \in e, x = y} \quad \text{if } |\mathcal{L}(e)| = 1 \quad \neq\text{-PROP-ELIM} \frac{\Gamma, x \in e, y \in e^c}{\Gamma, x \in e, x \neq y} \quad \text{if } |\mathcal{L}(e)| = 1$$

$$\text{CLOSE} \frac{}{\Gamma, x \in e_1, \dots, x \in e_n} \quad \text{if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \emptyset$$

$$\text{SUBSUME} \frac{\Gamma, x \in e_1, \dots, x \in e_n}{\Gamma, x \in e, x \in e_1, \dots, x \in e_n} \quad \text{if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) \subseteq \mathcal{L}(e)$$

$$\text{INTERSECT} \frac{\Gamma, x \in e}{\Gamma, x \in e_1, \dots, x \in e_n} \quad \text{if } n > 1 \text{ and } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \mathcal{L}(e)$$

$$\text{FWD-PROP} \frac{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \quad \text{if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n))$$

$$\text{FWD-PROP-ELIM} \frac{\Gamma, x \in e, x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \quad \text{if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n)) \text{ and } |\mathcal{L}(e)| = 1$$

$$\text{BWD-PROP} \frac{\{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1^i, \dots, x_n \in e_n^i\}_{i=1}^k}{\Gamma, x \in e, x = f(x_1, \dots, x_n)} \quad \text{if } f^{-1}(\mathcal{L}(e)) = \bigcup_{i=1}^k (\mathcal{L}(e_1^i) \times \dots \times \mathcal{L}(e_n^i))$$

String Constraints with Real-World Regular Expressions. Taolue Chen, Matthew Hague, Zhilei Han, Denghang Hu, Alejandro Flores-Lamas, Anthony W. Lin, Shuanglong Kan, PR, Zhilin Wu, POPL'22

$$\text{FWD-PROP} \frac{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \quad \text{if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n))$$

$$\notin \frac{\Gamma, x \in e^c}{\Gamma, x \notin e} \quad \neq \frac{\Gamma, x \neq y, y = f(x_1, \dots, x_n)}{\Gamma, x \neq f(x_1, \dots, x_n)} \quad \text{where } y \text{ is fresh} \quad \text{CUT} \frac{\Gamma, x \in e \quad \Gamma, x \in e^c}{\Gamma}$$

$$\text{=-PROP} \frac{\Gamma, x \in e, x = y, y \in e}{\Gamma, x \in e, x = y} \quad \neq\text{-SUBSUME} \frac{\Gamma, x \in e_1, y \in e_2}{\Gamma, x \in e_1, x \neq y, y \in e_2} \quad \text{if } \mathcal{L}(e_1) \cap \mathcal{L}(e_2) = \emptyset$$

$$\text{=-PROP-ELIM} \frac{\Gamma, x \in e, y \in e}{\Gamma, x \in e, x = y} \quad \text{if } |\mathcal{L}(e)| = 1 \quad \neq\text{-PROP-ELIM} \frac{\Gamma, x \in e, y \in e^c}{\Gamma, x \in e, x \neq y} \quad \text{if } |\mathcal{L}(e)| = 1$$

$$\text{CLOSE} \frac{}{\Gamma, x \in e_1, \dots, x \in e_n} \quad \text{if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \emptyset$$

$$\text{SUBSUME} \frac{\Gamma, x \in e_1, \dots, x \in e_n}{\Gamma, x \in e, x \in e_1, \dots, x \in e_n} \quad \text{if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) \subseteq \mathcal{L}(e)$$

$$\text{INTERSECT} \frac{\Gamma, x \in e}{\Gamma, x \in e_1, \dots, x \in e_n} \quad \text{if } n > 1 \text{ and } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \mathcal{L}(e)$$

$$\text{FWD-PROP} \frac{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \quad \text{if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n))$$

$$\text{FWD-PROP-ELIM} \frac{\Gamma, x \in e, x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \quad \text{if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n)) \text{ and } |\mathcal{L}(e)| = 1$$

$$\text{BWD-PROP} \frac{\{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1^i, \dots, x_n \in e_n^i\}_{i=1}^k}{\Gamma, x \in e, x = f(x_1, \dots, x_n)} \quad \text{if } f^{-1}(\mathcal{L}(e)) = \bigcup_{i=1}^k (\mathcal{L}(e_1^i) \times \dots \times \mathcal{L}(e_n^i))$$

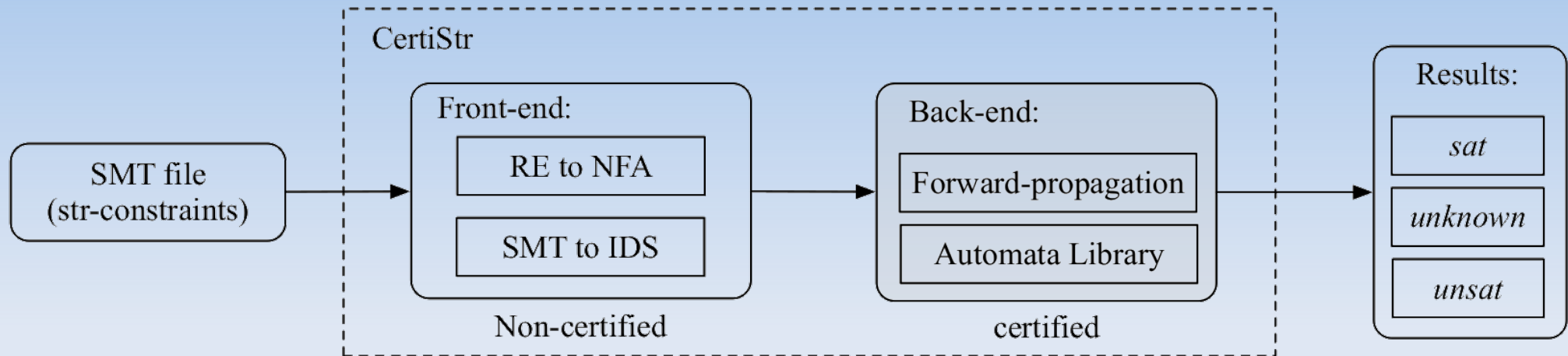
String Constraints with Real-World Regular Expressions. Taolue Chen, Matthew Hague, Zhilei Han, Denghang Hu, Alejandro Flores-Lamas, Anthony W. Lin, Shuanglong Kan, PR, Zhilin Wu, POPL'22

Forward Propagation

$$\text{FWD-PROP} \frac{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \quad \text{if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n))$$

- “Image of regular languages is regular”
- Covers functions concat, transducers, etc.
- Over-approximate if argument variables not pairwise distinct
- Yields decision procedure for **tree-shaped constraints**

String Solver *CertiStr*



- Verified solver for tree-shaped constraints, written in Isabelle/HOL
- The hardest part: implementation of symbolic automata

Regular Constraint Prop.

$$\begin{array}{l}
 \notin \frac{\Gamma, x \in e^c}{\Gamma, x \notin e} \quad \neq \frac{\Gamma, x \neq y, y = f(x_1, \dots, x_n)}{\Gamma, x \neq f(x_1, \dots, x_n)} \text{ where } y \text{ is fresh} \quad \text{CUT} \frac{\Gamma, x \in e \quad \Gamma, x \in e^c}{\Gamma} \\
 \\
 =\text{-PROP} \frac{\Gamma, x \in e, x = y, y \in e}{\Gamma, x \in e, x = y} \quad \neq\text{-SUBSUME} \frac{\Gamma, x \in e_1, y \in e_2}{\Gamma, x \in e_1, x \neq y, y \in e_2} \text{ if } \mathcal{L}(e_1) \cap \mathcal{L}(e_2) = \emptyset \\
 \\
 =\text{-PROP-ELIM} \frac{\Gamma, x \in e, y \in e}{\Gamma, x \in e, x = y} \text{ if } |\mathcal{L}(e)| = 1 \quad \neq\text{-PROP-ELIM} \frac{\Gamma, x \in e, y \in e^c}{\Gamma, x \in e, x \neq y} \text{ if } |\mathcal{L}(e)| = 1 \\
 \\
 \text{CLOSE} \frac{}{\Gamma, x \in e_1, \dots, x \in e_n} \text{ if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \emptyset \\
 \\
 \text{SUBSUME} \frac{\Gamma, x \in e_1, \dots, x \in e_n}{\Gamma, x \in e, x \in e_1, \dots, x \in e_n} \text{ if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) \subseteq \mathcal{L}(e) \\
 \\
 \text{INTERSECT} \frac{\Gamma, x \in e}{\Gamma, x \in e_1, \dots, x \in e_n} \text{ if } n > 1 \text{ and } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \mathcal{L}(e) \\
 \\
 \text{FWD-PROP} \frac{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \text{ if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n)) \\
 \\
 \text{FWD-PROP-ELIM} \frac{\Gamma, x \in e, x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \text{ if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n)) \text{ and } |\mathcal{L}(e)| = 1 \\
 \\
 \text{BWD-PROP} \frac{\{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1^i, \dots, x_n \in e_n^i\}_{i=1}^k}{\Gamma, x \in e, x = f(x_1, \dots, x_n)} \text{ if } f^{-1}(\mathcal{L}(e)) = \bigcup_{i=1}^k (\mathcal{L}(e_1^i) \times \dots \times \mathcal{L}(e_n^i))
 \end{array}$$

String Constraints with Real-World Regular Expressions. Taolue Chen, Matthew Hague, Zhilei Han, Denghang Hu, Alejandro Flores-Lamas, Anthony W. Lin, Shuanglong Kan, PR, Zhilin Wu, POPL'22

$$\text{BWD-PROP} \frac{\{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1^i, \dots, x_n \in e_n^i\}_{i=1}^k}{\Gamma, x \in e, x = f(x_1, \dots, x_n)} \text{ if } f^{-1}(\mathcal{L}(e)) = \bigcup_{i=1}^k (\mathcal{L}(e_1^i) \times \dots \times \mathcal{L}(e_n^i))$$

$$\notin \frac{\Gamma, x \in e^c}{\Gamma, x \notin e} \quad \neq \frac{\Gamma, x \neq y, y = f(x_1, \dots, x_n)}{\Gamma, x \neq f(x_1, \dots, x_n)} \text{ where } y \text{ is fresh} \quad \text{CUT} \frac{\Gamma, x \in e \quad \Gamma, x \in e^c}{\Gamma}$$

$$\text{=-PROP} \frac{\Gamma, x \in e, x = y, y \in e}{\Gamma, x \in e, x = y} \quad \neq\text{-SUBSUME} \frac{\Gamma, x \in e_1, y \in e_2}{\Gamma, x \in e_1, x \neq y, y \in e_2} \text{ if } \mathcal{L}(e_1) \cap \mathcal{L}(e_2) = \emptyset$$

$$\text{=-PROP-ELIM} \frac{\Gamma, x \in e, y \in e}{\Gamma, x \in e, x = y} \text{ if } |\mathcal{L}(e)| = 1 \quad \neq\text{-PROP-ELIM} \frac{\Gamma, x \in e, y \in e^c}{\Gamma, x \in e, x \neq y} \text{ if } |\mathcal{L}(e)| = 1$$

$$\text{CLOSE} \frac{}{\Gamma, x \in e_1, \dots, x \in e_n} \text{ if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \emptyset$$

$$\text{SUBSUME} \frac{\Gamma, x \in e_1, \dots, x \in e_n}{\Gamma, x \in e, x \in e_1, \dots, x \in e_n} \text{ if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) \subseteq \mathcal{L}(e)$$

$$\text{INTERSECT} \frac{\Gamma, x \in e}{\Gamma, x \in e_1, \dots, x \in e_n} \text{ if } n > 1 \text{ and } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \mathcal{L}(e)$$

$$\text{FWD-PROP} \frac{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \text{ if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n))$$

$$\text{FWD-PROP-ELIM} \frac{\Gamma, x \in e, x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \text{ if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n)) \text{ and } |\mathcal{L}(e)| = 1$$

$$\text{BWD-PROP} \frac{\{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1^i, \dots, x_n \in e_n^i\}_{i=1}^k}{\Gamma, x \in e, x = f(x_1, \dots, x_n)} \text{ if } f^{-1}(\mathcal{L}(e)) = \bigcup_{i=1}^k (\mathcal{L}(e_1^i) \times \dots \times \mathcal{L}(e_n^i))$$

String Constraints with Real-World Regular Expressions. Taolue Chen, Matthew Hague, Zhilei Han, Denghang Hu, Alejandro Flores-Lamas, Anthony W. Lin, Shuanglong Kan, PR, Zhilin Wu, POPL'22

Backward Propagation

$$\text{BWD-PROP} \frac{\{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1^i, \dots, x_n \in e_n^i\}_{i=1}^k}{\Gamma, x \in e, x = f(x_1, \dots, x_n)} \quad \text{if} \quad f^{-1}(\mathcal{L}(e)) = \bigcup_{i=1}^k (\mathcal{L}(e_1^i) \times \dots \times \mathcal{L}(e_n^i))$$

- “Pre-image of regular language is a recognizable relation”
- More general than forward prop.:
 - Also covers replace-all, etc.
 - Also precise when arguments coincide
- Yields decision procedure for **straightline constraints**

Some Fragments

- Quadratic word equations
- Tree-shaped
- Acyclic
- Chainfree
- Straightline
- Cost-enriched straightline
- Weakly chaining

Some Fragments

- Quadratic word equations
- Tree-shaped
- Acyclic
- Chainfree
- Straightline
- Cost-enriched straightline
- Weakly chaining



Forward + backward

Some Fragments

- Quadratic word equations
- Tree-shaped
- Acyclic
- Chainfree
- Straightline
- Cost-enriched straightline
- Weakly chaining

Nielsen + backward,
propagate
cost-enriched regexes
(aka Parikh automata)

Forward + backward

Some Fragments

- Quadratic word equations

- Tree-shaped

- Acyclic

Nielsen + backward,
propagate
cost-enriched regexes
(aka Parikh automata)

- Chainfree

- Straightline

Forward + backward

- Cost-enriched straightline

- Weakly chaining

???

Some Fragments

???

- Quadratic word equations
- Tree-shaped
- Acyclic
- Chainfree
- Straightline
- Cost-enriched straightline
- Weakly chaining

Nielsen + backward,
propagate
cost-enriched regexes
(aka Parikh automata)

Forward + backward

???

Where is the limit?

Conjecture:

- Nielsen + forward + backward + cut is **incomplete** for both quadratic and general word equations

$$\begin{array}{l}
\neq \frac{\Gamma, x \in e^c}{\Gamma, x \notin e} \quad \neq \frac{\Gamma, x \neq y, y = f(x_1, \dots, x_n)}{\Gamma, x \neq f(x_1, \dots, x_n)} \text{ where } y \text{ is fresh} \quad \text{CUT} \frac{\Gamma, x \in e \quad \Gamma, x \in e^c}{\Gamma} \\
\\
= \text{-PROP} \frac{\Gamma, x \in e, x = y, y \in e}{\Gamma, x \in e, x = y} \quad \neq \text{-SUBSUME} \frac{\Gamma, x \in e_1, y \in e_2}{\Gamma, x \in e_1, x \neq y, y \in e_2} \text{ if } \mathcal{L}(e_1) \cap \mathcal{L}(e_2) = \emptyset \\
\\
= \text{-PROP-ELIM} \frac{\Gamma, x \in e, y \in e}{\Gamma, x \in e, x = y} \text{ if } |\mathcal{L}(e)| = 1 \quad \neq \text{-PROP-ELIM} \frac{\Gamma, x \in e, y \in e^c}{\Gamma, x \in e, x \neq y} \text{ if } |\mathcal{L}(e)| = 1 \\
\\
\text{CLOSE} \frac{}{\Gamma, x \in e_1, \dots, x \in e_n} \text{ if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \emptyset \\
\\
\text{SUBSUME} \frac{\Gamma, x \in e_1, \dots, x \in e_n}{\Gamma, x \in e, x \in e_1, \dots, x \in e_n} \text{ if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) \subseteq \mathcal{L}(e) \\
\\
\text{INTERSECT} \frac{\Gamma, x \in e}{\Gamma, x \in e_1, \dots, x \in e_n} \text{ if } n > 1 \text{ and } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \mathcal{L}(e) \\
\\
\text{FWD-PROP} \frac{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \text{ if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n)) \\
\\
\text{FWD-PROP-ELIM} \frac{\Gamma, x \in e, x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \text{ if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n)) \text{ and } |\mathcal{L}(e)| = 1 \\
\\
\text{BWD-PROP} \frac{\{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1^i, \dots, x_n \in e_n^i\}_{i=1}^k}{\Gamma, x \in e, x = f(x_1, \dots, x_n)} \text{ if } f^{-1}(\mathcal{L}(e)) = \bigcup_{i=1}^k (\mathcal{L}(e_1^i) \times \dots \times \mathcal{L}(e_n^i))
\end{array}$$

String Constraints with Real-World Regular Expressions. Taolue Chen, Matthew Hague, Zhilei Han, Denghang Hu, Alejandro Flores-Lamas, Anthony W. Lin, Shuanglong Kan, PR, Zhilin Wu, POPL'22

$$\text{CUT} \frac{\Gamma, x \in e \quad \Gamma, x \in e^c}{\Gamma}$$

$$\notin \frac{\Gamma, x \in e^c}{\Gamma, x \notin e} \neq \frac{\Gamma, x \neq y, y = f(x_1, \dots, x_n)}{\Gamma, x \neq f(x_1, \dots, x_n)} \text{ where } y \text{ is fresh}$$

$$\text{CUT} \frac{\Gamma, x \in e \quad \Gamma, x \in e^c}{\Gamma}$$

$$\text{=-PROP} \frac{\Gamma, x \in e, x = y, y \in e}{\Gamma, x \in e, x = y} \quad \neq\text{-SUBSUME} \frac{\Gamma, x \in e_1, y \in e_2}{\Gamma, x \in e_1, x \neq y, y \in e_2} \text{ if } \mathcal{L}(e_1) \cap \mathcal{L}(e_2) = \emptyset$$

$$\text{=-PROP-ELIM} \frac{\Gamma, x \in e, y \in e}{\Gamma, x \in e, x = y} \text{ if } |\mathcal{L}(e)| = 1 \quad \neq\text{-PROP-ELIM} \frac{\Gamma, x \in e, y \in e^c}{\Gamma, x \in e, x \neq y} \text{ if } |\mathcal{L}(e)| = 1$$

$$\text{CLOSE} \frac{}{\Gamma, x \in e_1, \dots, x \in e_n} \text{ if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \emptyset$$

$$\text{SUBSUME} \frac{\Gamma, x \in e_1, \dots, x \in e_n}{\Gamma, x \in e, x \in e_1, \dots, x \in e_n} \text{ if } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) \subseteq \mathcal{L}(e)$$

$$\text{INTERSECT} \frac{\Gamma, x \in e}{\Gamma, x \in e_1, \dots, x \in e_n} \text{ if } n > 1 \text{ and } \mathcal{L}(e_1) \cap \dots \cap \mathcal{L}(e_n) = \mathcal{L}(e)$$

$$\text{FWD-PROP} \frac{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \text{ if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n))$$

$$\text{FWD-PROP-ELIM} \frac{\Gamma, x \in e, x_1 \in e_1, \dots, x_n \in e_n}{\Gamma, x = f(x_1, \dots, x_n), x_1 \in e_1, \dots, x_n \in e_n} \text{ if } \mathcal{L}(e) = f(\mathcal{L}(e_1), \dots, \mathcal{L}(e_n)) \text{ and } |\mathcal{L}(e)| = 1$$

$$\text{BWD-PROP} \frac{\{\Gamma, x \in e, x = f(x_1, \dots, x_n), x_1 \in e_1^i, \dots, x_n \in e_n^i\}_{i=1}^k}{\Gamma, x \in e, x = f(x_1, \dots, x_n)} \text{ if } f^{-1}(\mathcal{L}(e)) = \bigcup_{i=1}^k (\mathcal{L}(e_1^i) \times \dots \times \mathcal{L}(e_n^i))$$

String Constraints with Real-World Regular Expressions. Taolue Chen, Matthew Hague, Zhilei Han, Denghang Hu, Alejandro Flores-Lamas, Anthony W. Lin, Shuanglong Kan, PR, Zhilin Wu, POPL'22

The Next Steps

- **Properties** and **generalization** of regular constraint propagation
- What's **decidable** about strings constraints?
- From certified solver to proof checker: **proof format** for string solvers

For more details: see our POPL'24 tutorial
<https://eldarica.org/ostrich-popl24/>

Towards Alethe-Style Proofs

```
(regular-languages
  (! (re.from_automaton "automaton {init s0; ...;}") :id 1)
  (! (re.from_automaton "automaton {init s0; ...;}") :id 2)
  (! (re.from_automaton "automaton {init s0; ...;}") :id 3)
)

(assume h0 (str.in_re_id w 1))
(assume h1 (or (str.in_re_id w 2) (str.in_re_id w 3)))
(step t2 (cl (str.in_re_id w 2) (str.in_re_id w 3)) :rule or :premises (h1))

; start proof branch that spans until t3
(anchor :step t3)
(assume t3.h0 (str.in_re_id w 2))

(step t3.t1 (cl (not (str.in_re_id w 1)) (not (str.in_re_id w 2)))
           :rule re_empty_intersection)
(step t3.t2 (cl) :rule resolution :premises (h0 t3.h0 t3.t1))

(subproof t3 (cl (not (str.in_re_id w 2))))

(step t4 (cl (str.in_re_id w 3)) :rule resolution :premises (t2 t3))

(step t5 (cl (not (str.in_re_id w 1)) (not (str.in_re_id w 3)))
           :rule re_empty_intersection)
(step t6 (cl) :rule resolution :premises (h0 h2 t5))
```

Straightline Example

x in $a^*c^*b^*$

$y := \text{reverse}(x)$

y in b^*a^*

$z := \text{replaceAll}(y, a, b)$

z in b^*

Is there an input x satisfying all assertions?

Straightline Example

$x \text{ in } a^*c^*b^*$

$y := \text{reverse}(x)$

$y \text{ in } b^*a^*$

$z := \text{replaceAll}(y, a, b)$

$z \text{ in } b^*$

} $y \text{ in } (a \mid b)^*$

Is there an input x satisfying all assertions?

Straightline Example

$x \text{ in } a^*c^*b^*$

$y := \text{reverse}(x)$

$y \text{ in } b^*a^*$

$y \text{ in } (a \mid b)^*$

Is there an input x satisfying all assertions?

Straightline Example

$x \text{ in } a^*c^*b^*$

$y := \text{reverse}(x)$

$y \text{ in } b^*a^*$

$y \text{ in } (a \mid b)^*$

} $y \text{ in } (a \mid b)^* \ \& \ b^*a^*$

Is there an input x satisfying all assertions?

Straightline Example

$x \text{ in } a^*c^*b^*$

$y := \text{reverse}(x)$

$y \text{ in } b^*a^*$

Is there an input x satisfying all assertions?

Straightline Example

$x \text{ in } a^*c^*b^*$

$y := \text{reverse}(x)$

$y \text{ in } b^*a^*$

} $x \text{ in } a^*b^*$

Is there an input x satisfying all assertions?

Straightline Example

$x \text{ in } a^*c^*b^*$

$x \text{ in } a^*b^*$

Is there an input x satisfying all assertions?

Straightline Example

$x \text{ in } a^*c^*b^*$

$x \text{ in } a^*b^*$

} $x \text{ in } a^*b^*$

Is there an input x satisfying all assertions?

Straightline Example

$x \text{ in } a^*b^*$

Is there an input x satisfying all assertions?

Straightline Example

Easy to solve!

x in a^*b^*

Is there an input x satisfying all assertions?

Straightline Example

Easy to solve!

x in a^*b^*

Solution: $x = abb$