

# Z3-Noodler: Shepherding Decision Procedures



**Vojtěch Havlena, et al.**

Brno University of Technology, Czech Republic

MOSCA'25

# SMT String constraint solving

- Checking **satisfiability** of formulae with **string variables** and operations

$$\underbrace{x = yz \wedge y \neq u}_{(dis)equations} \wedge \overbrace{x \in (ab)^* a^+ (b|c)}^{\text{regular constraints}} \wedge \overbrace{|x| = 2|u| + 1}^{\text{length constraints}} \wedge \underbrace{\text{contains}(u, \text{replace}(z, b, c)) \wedge \dots}_{\text{more complex operations}}$$

# SMT String constraint solving

- Checking **satisfiability** of formulae with **string variables** and operations

$$\underbrace{x = yz \wedge y \neq u \wedge x \in (ab)^* a^+ (b|c)}_{(dis)equations} \wedge \overbrace{(ab)^* a^+ (b|c)}^{\text{regular constraints}} \wedge \overbrace{|x| = 2|u| + 1}^{\text{length constraints}} \wedge \underbrace{\text{contains}(u, \text{replace}(z, b, c)) \wedge \dots}_{\text{more complex operations}}$$

- **Motivation:** **large and complex real-world programs** need security guarantees

- **analysis** of string manipulating programs (**vulnerabilities** of web applications)

```

let x = y.substring(1, y.length - 1);    x0 = substr(y, 1, |y| - 1) ∧
let z = y.concat(x);                    z0 = y · x0 ∧
assert(x == z);                          x0 ≠ z0
  
```

- analysing **access control policies** (AWS/Rego)

```

action: deactivate,
resource: (a1, a2),
condition: {StringLike, s3:prefix, home*}
  
```

$$A = \text{"deactivate"} \wedge (R = \text{"a1"} \vee R = \text{"a2"}) \wedge \text{prefix} \in \text{home}^*$$

- **verification** of cockpit systems (Boeing), etc.

- $\leadsto$  **efficient and expressive** SMT string solvers are **needed**

# Z3-Noodler: Highlight

- Based on SMT solver **Z3**
  - Z3-Noodler replaces Z3's **string theory solver**
  - low-level interaction with Z3 (given by the theory interface)
    - gathers conjunction of string constraints (provided by SMT core)
    - core Noodler's procedure
    - provide theory lemma
  - uses **Z3's linear arithmetic** (LIA) theory solver

## Z3-Noodler: Highlight

- Based on SMT solver **Z3**
  - Z3-Noodler replaces Z3's **string theory solver**
  - low-level interaction with Z3 (given by the theory interface)
    - gathers conjunction of string constraints (provided by SMT core)
    - core Noodler's procedure
    - provide theory lemma
  - uses **Z3's linear arithmetic** (LIA) theory solver
- Uses **Nondeterministic finite automata** (**Mata** automata library)

[TACAS'24]

## Z3-Noodler: Highlight

- Based on SMT solver **Z3**
  - Z3-Noodler replaces Z3's **string theory solver**
  - low-level interaction with Z3 (given by the theory interface)
    - gathers conjunction of string constraints (provided by SMT core)
    - core Noodler's procedure
    - provide theory lemma
  - uses **Z3's linear arithmetic** (LIA) theory solver
- Uses **Nondeterministic finite automata** (**Mata** automata library)
- Combines various **decision procedures** for different fragments

[TACAS'24]

## Z3-Noodler: Highlight

- Based on SMT solver **Z3**
  - Z3-Noodler replaces Z3's **string theory solver**
  - low-level interaction with Z3 (given by the theory interface)
    - gathers conjunction of string constraints (provided by SMT core)
    - core Noodler's procedure
    - provide theory lemma
  - uses **Z3's linear arithmetic** (LIA) theory solver
- Uses **Nondeterministic finite automata** (**Mata** automata library) [TACAS'24]
- Combines various **decision procedures** for different fragments
- Support for various **string functions/predicates**
  - more **complex** string functions/predicates (`replace`, `substr`, `indexof` ...)
  - string-integer conversions, transducer constraints (`replace_all`)

## Z3-Noodler: Highlight

- Based on SMT solver **Z3**
  - Z3-Noodler replaces Z3's **string theory solver**
  - low-level interaction with Z3 (given by the theory interface)
    - gathers conjunction of string constraints (provided by SMT core)
    - core Noodler's procedure
    - provide theory lemma
  - uses **Z3's linear arithmetic** (LIA) theory solver
- Uses **Nondeterministic finite automata** (**Mata** automata library) [TACAS'24]
- Combines various **decision procedures** for different fragments
- Support for various **string functions/predicates**
  - more **complex** string functions/predicates (`replace`, `substr`, `indexof` ...)
  - string-integer conversions, transducer constraints (`replace_all`)
- the **winner** of SMT-COMP'24 string division



# Combining decision procedures

## ■ **interface** of decision procedure

- $\text{isSuitable}(\psi) \rightarrow \mathbb{B}$
- $\text{init}(\psi)$
- $\text{preprocess}()$
- $\text{nextSolution}() \rightarrow \mathbb{B}_3$
- $\text{getLIA}() \rightarrow \Phi_{\text{LIA}} \times \{\text{precise}, \text{underapprox}\}$
- $\text{getModel}(\theta, x) \rightarrow \Sigma^*$

# Combining decision procedures

## ■ **interface** of decision procedure

- $\text{isSuitable}(\psi) \rightarrow \mathbb{B}$
- $\text{init}(\psi)$
- $\text{preprocess}()$
- $\text{nextSolution}() \rightarrow \mathbb{B}_3$
- $\text{getLIA}() \rightarrow \Phi_{\text{LIA}} \times \{\text{precise}, \text{underapprox}\}$
- $\text{getModel}(\theta, x) \rightarrow \Sigma^*$

## ■ procedure **selection**

- suitability check
- **ordered** by the most specific to the most general

# Combining decision procedures

## ■ **interface** of decision procedure

- $\text{isSuitable}(\psi) \rightarrow \mathbb{B}$
- $\text{init}(\psi)$
- $\text{preprocess}()$
- $\text{nextSolution}() \rightarrow \mathbb{B}_3$
- $\text{getLIA}() \rightarrow \Phi_{\text{LIA}} \times \{\text{precise, underapprox}\}$
- $\text{getModel}(\theta, x) \rightarrow \Sigma^*$

## ■ procedure **selection**

- suitability check
- **ordered** by the most specific to the most general

## ■ **execution**

- iteratively call  $\text{nextSolution}()$
- check if  $\text{getLIA}()$  is satisfiable

# Stabilization-based procedure

[FM'23], [OOPSLA'23]

$$\underbrace{x = yz \wedge y \neq u}_{\text{(dis)equations}} \wedge \underbrace{x \in (ab)^* a^+ (b|c)}_{\text{regular constraints}} \wedge \underbrace{|x| = 2|u| + 1}_{\text{length constraints}}$$

# Stabilization-based procedure

[FM'23], [OOPSLA'23]

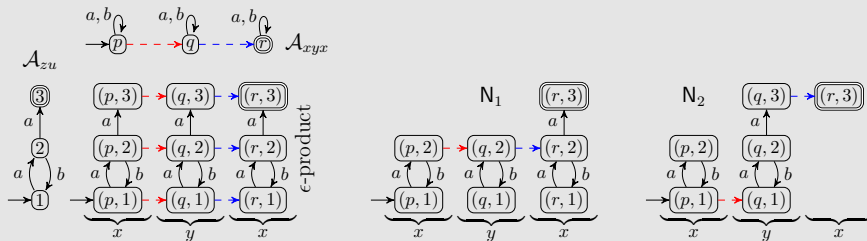
$$\underbrace{x = yz \wedge y \neq u}_{\text{(dis)equations}} \wedge \underbrace{x \in (ab)^* a^+ (b|c)}_{\text{regular constraints}} \wedge \underbrace{|x| = 2|u| + 1}_{\text{length constraints}}$$

- tight integration of **equations** with **regular constraints**
- complete on **chain-free constraints**
- iteratively **refining the languages** of variables  $Lang$ ; until **stable form**
- feasible equation splits refining languages computed by **noodlification**
- create LIA formula **encoding possible lengths of words** in each language in  $Lang$

# Noodlification (No-lengths)

$$xyx = zu \quad \begin{array}{l} u \mapsto (baba)^* a \quad z \mapsto a(ba)^* \quad x \mapsto \Sigma^* \quad y \mapsto \Sigma^* \end{array}$$

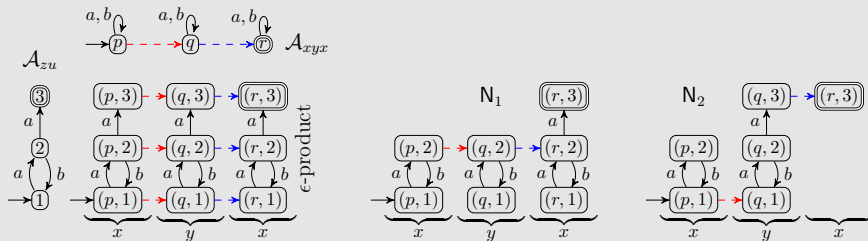
- Use right side to **refine languages** of variables  $x, y$  on the left side
- Leads to two **feasible noodles**:



# Noodlification (No-lengths)

$$xyx = zu \quad u \mapsto (baba)^*a \quad z \mapsto a(ba)^* \quad x \mapsto \Sigma^* \quad y \mapsto \Sigma^*$$

- Use right side to **refine languages** of variables  $x, y$  on the left side
- Leads to two **feasible noodles**:



- Two refinements

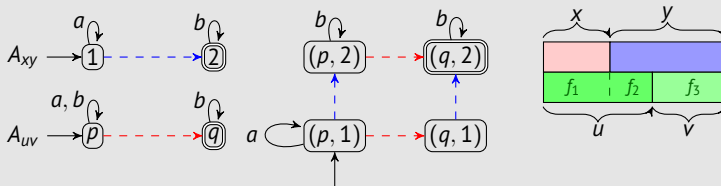
$$u \mapsto (baba)^*a \quad z \mapsto a(ba)^* \quad x \mapsto \mathbf{a} \quad y \mapsto \mathbf{(ba)^*}$$

$$u \mapsto (baba)^*a \quad z \mapsto a(ba)^* \quad x \mapsto \epsilon \quad y \mapsto \mathbf{a(ba)^*a}$$

# Noodlification (Lengths)

$$xy = uv \quad x \mapsto a^* \quad y \mapsto b^* \quad u \mapsto \Sigma^* \quad v \mapsto b^*$$

- **Align&Split** augmented with **noodlification** [Abdulla-CAV'14]
- **epsilon transitions** possibly on right side (only between relevant variables)
- epsilon defines **alignments** (derive substitutions with fresh variables)
- generate only **language-feasible splits**



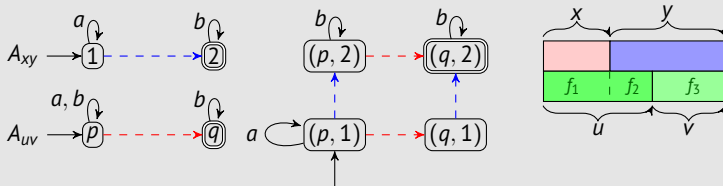


## Noodlification (Lengths)

$$xy = uv \quad x \mapsto a^* \quad y \mapsto b^* \quad u \mapsto \Sigma^* \quad v \mapsto b^*$$

- **Align&Split** augmented with **noodlification** [Abdulla-CAV'14]
- **epsilon transitions** possibly on right side (only between relevant variables)
- epsilon defines **alignments** (derive substitutions with fresh variables)
- generate only **language-feasible splits**

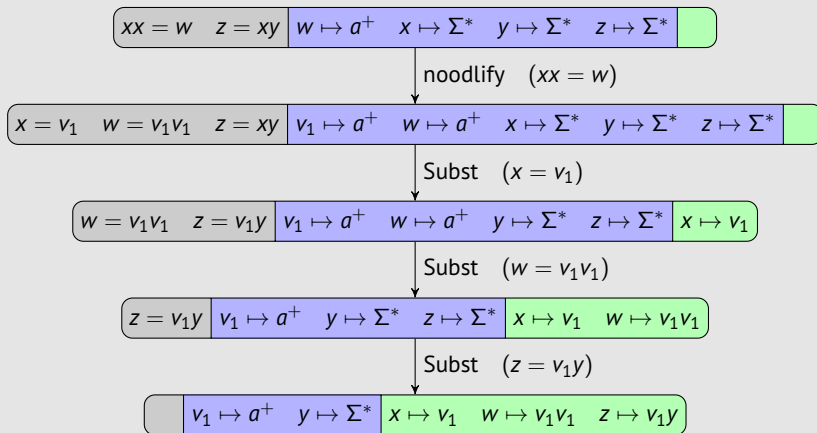
[Abdulla-CAV'14]



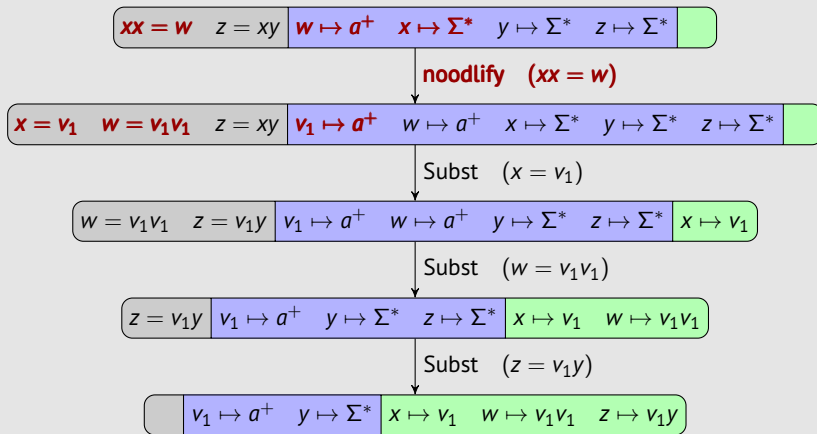
$$x = f_1 \quad y = f_2 f_3 \quad v = f_3 \quad u = f_1 f_2 \quad x \mapsto a^* \quad y \mapsto b^* \quad u \mapsto \Sigma^* \quad v \mapsto b^* \quad f_1 \mapsto a^* \quad f_2 \mapsto b^* \quad f_3 \mapsto b^*$$

$$x = f_1 \quad y = f_2 \quad v = f_2 \quad u = f_1 \quad x \mapsto a^* \quad y \mapsto b^* \quad u \mapsto \Sigma^* \quad v \mapsto b^* \quad f_1 \mapsto a^* \quad f_2 \mapsto b^*$$

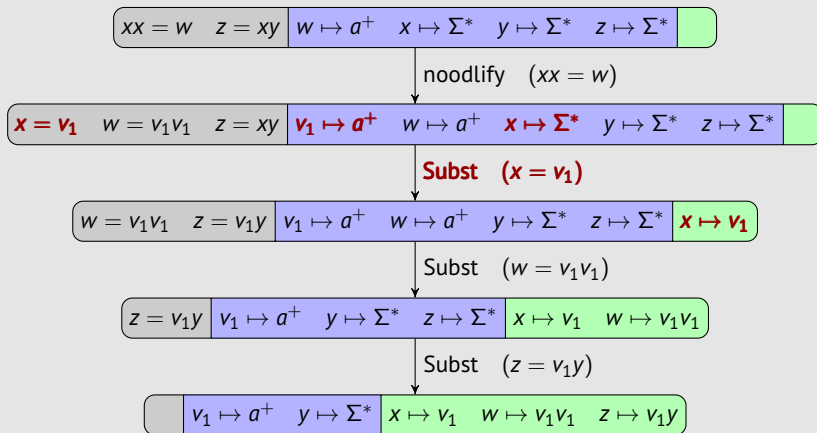
**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$



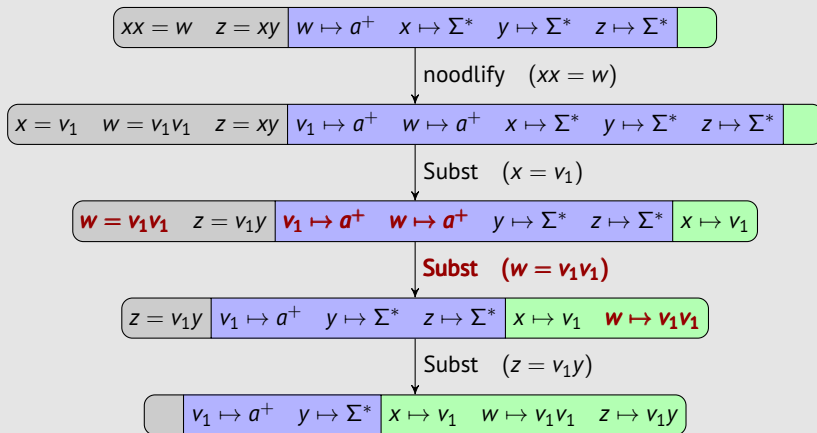
**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$



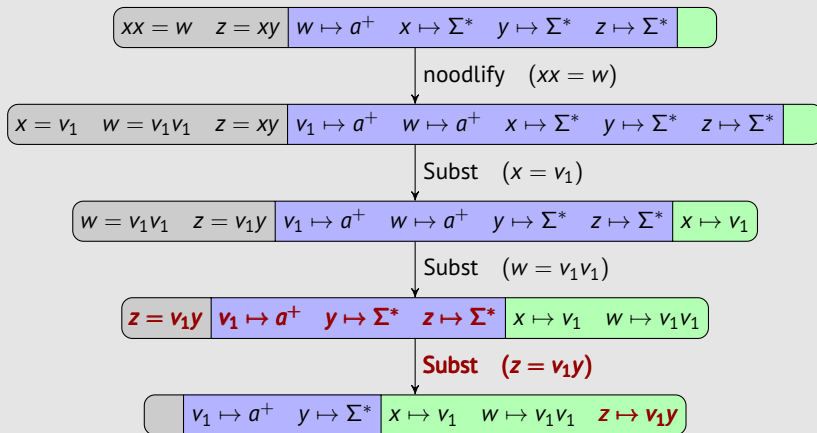
**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$



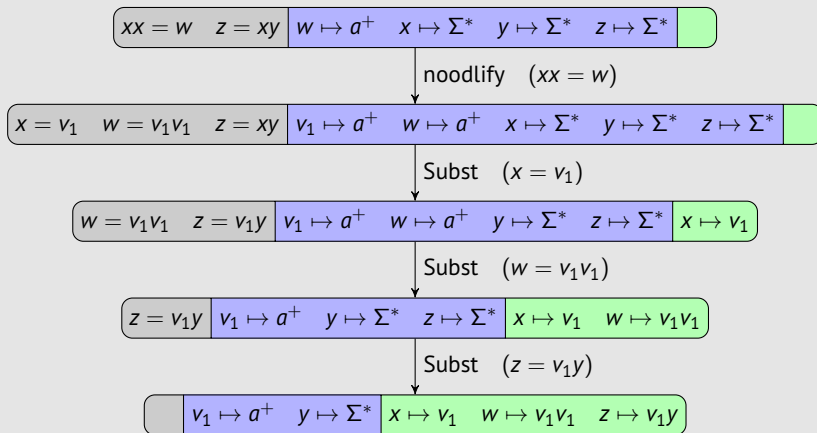
**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$



**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$



**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$



**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$

	$v_1 \mapsto a^+ \quad y \mapsto \Sigma^*$	$x \mapsto v_1 \quad w \mapsto v_1 v_1 \quad z \mapsto v_1 y$
--	--	---

■ stable solution **(*Lang*,  $\sigma$ )**:

- **language assignment** *Lang*:  $v_1 \mapsto a^+, y \mapsto \Sigma^*$
- **substitution map**  $\sigma$ :  $x \mapsto v_1, w \mapsto v_1 v_1, z \mapsto v_1 y$



**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$

	$v_1 \mapsto a^+ \quad y \mapsto \Sigma^*$	$x \mapsto v_1 \quad w \mapsto v_1 v_1 \quad z \mapsto v_1 y$
--	--	---

■ stable solution **(Lang,  $\sigma$ )**:

■ **language assignment**  $Lang: v_1 \mapsto a^+, y \mapsto \Sigma^*$

■ **substitution map**  $\sigma: x \mapsto v_1, w \mapsto v_1 v_1, z \mapsto v_1 y$

■ **LIA formula** encoding **possible lengths** of variables:

$$\varphi_{\text{len}} \stackrel{\text{def.}}{\iff} \quad \wedge \quad \wedge \quad \wedge \quad \wedge$$

**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$

	$v_1 \mapsto a^+ \quad y \mapsto \Sigma^*$	$x \mapsto v_1 \quad w \mapsto v_1 v_1 \quad z \mapsto v_1 y$
--	--	---

■ stable solution **(Lang,  $\sigma$ )**:

■ language assignment *Lang*:  $v_1 \mapsto a^+, y \mapsto \Sigma^*$

■ substitution map  $\sigma$ :  $x \mapsto v_1, w \mapsto v_1 v_1, z \mapsto v_1 y$

■ **LIA formula** encoding **possible lengths** of variables:

$$\varphi_{\text{len}} \stackrel{\text{def.}}{\iff} |v_1| \geq 1 \wedge \quad \wedge \quad \wedge \quad \wedge$$

**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$

	$v_1 \mapsto a^+ \quad y \mapsto \Sigma^*$	$x \mapsto v_1 \quad w \mapsto v_1 v_1 \quad z \mapsto v_1 y$
--	--	---

■ stable solution **(Lang,  $\sigma$ )**:

■ **language assignment**  $Lang: v_1 \mapsto a^+, y \mapsto \Sigma^*$

■ **substitution map**  $\sigma: x \mapsto v_1, w \mapsto v_1 v_1, z \mapsto v_1 y$

■ **LIA formula** encoding **possible lengths** of variables:

$$\varphi_{\text{len}} \stackrel{\text{def.}}{\iff} |v_1| \geq 1 \wedge |y| \geq 0 \wedge \quad \wedge \quad \wedge$$

**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$

	$v_1 \mapsto a^+ \quad y \mapsto \Sigma^*$	$x \mapsto v_1 \quad w \mapsto v_1 v_1 \quad z \mapsto v_1 y$
--	--	---

■ stable solution **(Lang,  $\sigma$ )**:

■ **language assignment**  $Lang: v_1 \mapsto a^+, y \mapsto \Sigma^*$

■ **substitution map**  $\sigma: x \mapsto v_1, w \mapsto v_1 v_1, z \mapsto v_1 y$

■ **LIA formula** encoding **possible lengths** of variables:

$$\varphi_{\text{len}} \stackrel{\text{def.}}{\iff} |v_1| \geq 1 \wedge |y| \geq 0 \wedge |x| = |v_1| \wedge \quad \wedge$$

**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$

	$v_1 \mapsto a^+ \quad y \mapsto \Sigma^*$	$x \mapsto v_1 \quad w \mapsto v_1 v_1 \quad z \mapsto v_1 y$
--	--	---

■ stable solution **(Lang,  $\sigma$ )**:

■ **language assignment**  $Lang: v_1 \mapsto a^+, y \mapsto \Sigma^*$

■ **substitution map**  $\sigma: x \mapsto v_1, w \mapsto v_1 v_1, z \mapsto v_1 y$

■ **LIA formula** encoding **possible lengths** of variables:

$$\varphi_{\text{len}} \stackrel{\text{def.}}{\iff} |v_1| \geq 1 \wedge |y| \geq 0 \wedge |x| = |v_1| \wedge |w| = |v_1| + |v_1| \wedge$$

**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$

	$v_1 \mapsto a^+ \quad y \mapsto \Sigma^*$	$x \mapsto v_1 \quad w \mapsto v_1 v_1 \quad z \mapsto v_1 y$
--	--	---

■ stable solution **(Lang,  $\sigma$ )**:

■ **language assignment**  $Lang: v_1 \mapsto a^+, y \mapsto \Sigma^*$

■ **substitution map**  $\sigma: x \mapsto v_1, w \mapsto v_1 v_1, z \mapsto v_1 y$

■ **LIA formula** encoding **possible lengths** of variables:

$$\varphi_{\text{len}} \stackrel{\text{def.}}{\iff} |v_1| \geq 1 \wedge |y| \geq 0 \wedge |x| = |v_1| \wedge |w| = |v_1| + |v_1| \wedge |z| = |v_1| + |y|$$

**Example:**  $xx = w \wedge z = xy \wedge w \in a^+ \wedge |z| = 2|w| - |x|$

	$v_1 \mapsto a^+ \quad y \mapsto \Sigma^*$	$x \mapsto v_1 \quad w \mapsto v_1 v_1 \quad z \mapsto v_1 y$
--	--	---

■ stable solution (**Lang**,  $\sigma$ ):

■ **language assignment**  $Lang: v_1 \mapsto a^+, y \mapsto \Sigma^*$

■ **substitution map**  $\sigma: x \mapsto v_1, w \mapsto v_1 v_1, z \mapsto v_1 y$

■ **LIA formula** encoding **possible lengths** of variables:

$$\varphi_{\text{len}} \stackrel{\text{def.}}{\iff} |v_1| \geq 1 \wedge |y| \geq 0 \wedge |x| = |v_1| \wedge |w| = |v_1| + |v_1| \wedge |z| = |v_1| + |y|$$

■ ask LIA solver if  $|z| = 2|w| - |x| \wedge \varphi_{\text{len}}$  is satisfiable

■ it is, we have model  $|v_1| = |x| = 1, |w| = |y| = 2, |z| = 3$

■ we can choose any word from  $Lang(v_1)$  and  $Lang(y)$  with correct lengths:

$$v_1 = a \text{ and } y = bc$$

■ models for  $x, w$ , and  $z$  are computed using the substitution map  $\sigma$ :

$$x = v_1 = a, w = v_1 v_1 = aa, \text{ and } z = v_1 y = abc$$

# Pure regular constraints

[TACAS'25]

$$\bigwedge_{1 \leq i \leq n} x \in \mathcal{S}_i \quad \wedge \quad \bigwedge_{1 \leq i \leq m} x \notin \mathcal{R}_i$$

- $\text{aut}(\mathcal{S})$ : **Regex**  $\rightsquigarrow$  **NFA**
- $L(P) \cap L(U) \neq \emptyset$  where  $P = \bigcap_{1 \leq i \leq n} \text{aut}(\mathcal{S}_i)$ ,  $U = \bigcap_{1 \leq i \leq m} \text{aut}(\mathcal{R}_i)^c$
- **Problem:** **Expensive complement computation** (determinization) for negations



# Pure regular constraints

[TACAS'25]

$$\bigwedge_{1 \leq i \leq n} x \in \mathcal{S}_i \quad \wedge \quad \bigwedge_{1 \leq i \leq m} x \notin \mathcal{R}_i$$

- $\text{aut}(\mathcal{S})$ : **Regex**  $\rightsquigarrow$  **NFA**
- $L(P) \cap L(U) \neq \emptyset$  where  $P = \bigcap_{1 \leq i \leq n} \text{aut}(\mathcal{S}_i)$ ,  $U = \bigcap_{1 \leq i \leq m} \text{aut}(\mathcal{R}_i)^c$
- **Problem: Expensive complement computation** (determinization) for negations
- **eager simulation reduction/determinisation** (if necessary) for aut
- sort  $\mathcal{S}_i$  by **estimated NFA size**
- **inclusion/universality checking**:  $L(P) \not\subseteq L(V)$  where  $V = \bigcup_{1 \leq i \leq m} \text{aut}(\mathcal{R}_i)$ 
  - **antichain-based algorithms**: perform well on real-world problems
- avoid construction for simpler cases

## Pure regular constraints: Syntactic check

- Analyze **single regexes** ( $x \in \mathcal{R}, x \notin \mathcal{R}$ ) to **extract syntactic properties**
- used to decide emptiness/universality
- **Propagate flags**  $(e, u, \ell)$  through operations:

- $e \in \mathbb{B}_3$ : the regex includes the empty word

$$\mathbb{B}_3 = \{\top, \perp, \text{undef}\}$$

- $u \in \mathbb{B}_3$ : the regex is universal

- $\ell \in \mathbb{N} \cup \{\text{undef}\}$ : the minimum length of words recognized by the regex

$$R_1 : (e_1, u_1, \ell_1) \quad R_2 : (e_2, u_2, \ell_2) \quad \text{re} . ++(R_1, R_2)$$

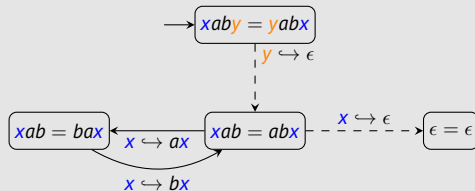
$$(e_1 \wedge e_2, u, \ell_1 + \ell_2), \ell_1 + \ell_2 > 0 \rightsquigarrow u = \perp, \text{otherwise } u = u_1 \wedge u_2$$

- Completely **avoid the NFA construction by reasoning about the flags**
- undef: when flags are insufficient  $\rightsquigarrow$  construct NFAs

# Nielsen transformation

- **Quadratic** (each variable has at most two occurrences) and **beyond**
- Create a **Nielsen graph** (finite for a quadratic system)
  - node: set of equations, Nielsen transformation metarules
  - priority queue; **pruning** of the state space
- lengths: **counter abstraction system**
  - heuristics for LIA formulae generation
  - iteratively generating LIA for **extended runs** (under-approximation)
  - **self-loop saturation**

[LIN-LMCS'21]



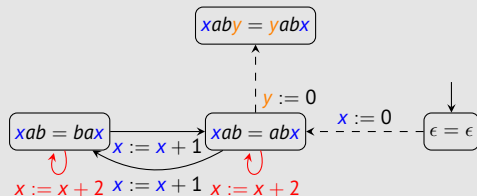
$$(x \hookrightarrow \alpha x) : \frac{(xu = \alpha v) \in \mathcal{E}}{\text{trim}(\mathcal{E}[x/\alpha x])}$$

$$(x \hookrightarrow \epsilon) : \frac{(xu = v) \in \mathcal{E}}{\text{trim}(\mathcal{E}[x/\epsilon])}$$

# Nielsen transformation

- **Quadratic** (each variable has at most two occurrences) and **beyond**
- Create a **Nielsen graph** (finite for a quadratic system)
  - node: set of equations, Nielsen transformation metarules
  - priority queue; **pruning** of the state space
- lengths: **counter abstraction system**
  - heuristics for LIA formulae generation
  - iteratively generating LIA for **extended runs** (under-approximation)
  - **self-loop saturation**

[LIN-LMCS'21]



$$(x \hookrightarrow \alpha x) : \frac{(xu = \alpha v) \in \mathcal{E}}{\text{trim}(\mathcal{E}[x/\alpha x])}$$

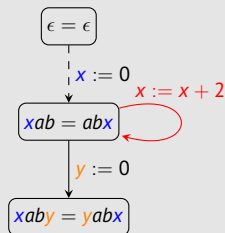
$$(x \hookrightarrow \epsilon) : \frac{(xu = v) \in \mathcal{E}}{\text{trim}(\mathcal{E}[x/\epsilon])}$$

# Nielsen transformation: Derivation of a LIA formula

$$xaby = yabx \wedge |x| \geq 50$$

# Nielsen transformation: Derivation of a LIA formula

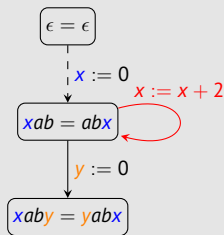
$$xaby = yabx \wedge |x| \geq 50$$



# Nielsen transformation: Derivation of a LIA formula

$$xaby = yabx \wedge |x| \geq 50$$

- Fresh counter variables for each step

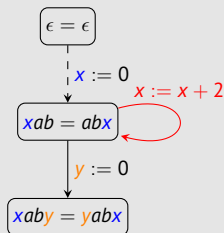


$$\begin{aligned} \varphi(x, y) \Leftrightarrow & x_0 = 0 \wedge y_0 = 0 \wedge \\ & x_1 = 0 \wedge y_1 = y_0 \wedge \\ & x_2 = x_1 + 2k \wedge y_2 = y_1 \wedge \\ & y_3 = 0 \wedge x_3 = x_2 \wedge \\ & x = x_3 \wedge y = y_3 \end{aligned}$$

# Nielsen transformation: Derivation of a LIA formula

$$xaby = yabx \wedge |x| \geq 50$$

- Fresh counter variables for each step



$$\begin{aligned} \varphi(x, y) \Leftrightarrow & x_0 = 0 \wedge y_0 = 0 \wedge \\ & x_1 = 0 \wedge y_1 = y_0 \wedge \\ & x_2 = x_1 + 2k \wedge y_2 = y_1 \wedge \\ & y_3 = 0 \wedge x_3 = x_2 \wedge \\ & x = x_3 \wedge y = y_3 \end{aligned}$$

Is  $\varphi(|x|, |y|) \wedge |x| \geq 50$  satisfiable?



# Length-based decision procedure

$$x = \text{abyc} \wedge x = zw \wedge x = \text{uddc} \wedge y = \text{vad} \wedge y = \text{as} \wedge |x| + |y| \geq 20$$

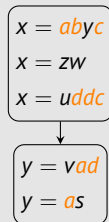
- Large systems (many equations, unrestricted variables and literals)  
 $\rightsquigarrow$  **noodles explosion**
- positions of literals matters
- Symbolically **encode all possible alignments of literals** (their positions) into LIA formulae
- String formula reduced into a **LIA formula**

# Length-based decision procedure

$$x = \text{abyc} \wedge x = zw \wedge x = \text{uddc} \wedge y = \text{vad} \wedge y = \text{as} \wedge |x| + |y| \geq 20$$

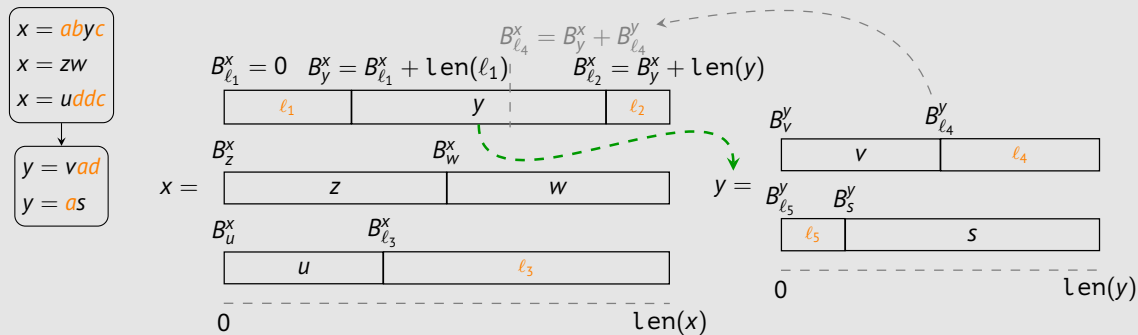
- Large systems (many equations, unrestricted variables and literals)  
 $\rightsquigarrow$  **noodles explosion**
- positions of literals matters
- Symbolically **encode all possible alignments of literals** (their positions) into LIA formulae
- String formula reduced into a **LIA formula**
- **Block graph**: dependencies between blocks
- **Block-acyclic string constraint**: acyclic block graph

$$\bigwedge_{1 \leq i \leq n} x = R_i$$



# Length-based decision procedure: Alignments to LIA formula

$$x = \text{abyc} \wedge x = zw \wedge x = \text{uddc} \wedge y = \text{vad} \wedge y = \text{as}$$



# Extended Constraints

- **string-integer conversions**: `to_int`, `from_code`, etc.
  - stable **solution**  $(Lang, \sigma)$ ; set of **conversion constraints**  
 $\mathcal{C} = \{k = \text{to\_int}(x), y = \text{from\_code}(l), \dots\}$ , assume finite  $Lang(x)$
  - create **concise**  $\varphi_c$  for  $c \in \mathcal{C}$

[SAT'24]

# Extended Constraints

## ■ string-integer conversions: `to_int`, `from_code`, etc.

[SAT'24]

- stable **solution**  $(Lang, \sigma)$ ; set of **conversion constraints**

$\mathcal{C} = \{k = \text{to\_int}(x), y = \text{from\_code}(l), \dots\}$ , assume finite  $Lang(x)$

- create **concise**  $\varphi_c$  for  $c \in \mathcal{C}$
- combining **lengths and conversions**  $|x| \geq 3 \wedge \text{to\_int}(x) < 100$
- handle **substitutions**  $x \mapsto x_1 \cdots x_n \in \sigma, \ell_i$ : possible length of  $x_i$  (constant)

$$\text{to\_int}(x) = \sum_{1 \leq i \leq n} \left( \text{to\_int}(x_i) \cdot 10^{\ell_{i+1} + \dots + \ell_n} \right) \wedge \bigwedge_{1 \leq i \leq n} (|x_i| = \ell_i)$$

# Extended Constraints

- **string-integer conversions:** `to_int`, `from_code`, etc.

[SAT'24]

- stable **solution**  $(Lang, \sigma)$ ; set of **conversion constraints**

$\mathcal{C} = \{k = \text{to\_int}(x), y = \text{from\_code}(l), \dots\}$ , assume finite  $Lang(x)$

- create **concise**  $\varphi_c$  for  $c \in \mathcal{C}$

- combining **lengths and conversions**  $|x| \geq 3 \wedge \text{to\_int}(x) < 100$

- handle **substitutions**  $x \mapsto x_1 \cdots x_n \in \sigma, \ell_i$ : possible length of  $x_i$  (constant)

$$\text{to\_int}(x) = \sum_{1 \leq i \leq n} \left( \text{to\_int}(x_i) \cdot 10^{\ell_{i+1} + \dots + \ell_n} \right) \wedge \bigwedge_{1 \leq i \leq n} (|x_i| = \ell_i)$$

- **intervals**  $(0 \leq \text{to\_int}(x) \leq 7 \wedge |x| = 1) \vee (20 \leq \text{to\_int}(x) \leq 59 \wedge |x| = 2)$
- encoding invalid strings
- **underapproximation** of infinite languages

# Extended Constraints

- **transducer constraints** ( $T(xy, uv)$ ): `replace_all`
  - **transducer noodlification**: sequence of transducers
  - stable solution  $\rightsquigarrow$  construct **multitape transducer** composing the same variables substitutions
  - Parikh image for lengths

# Extended Constraints

- **transducer constraints** ( $T(xy, uv)$ ): `replace_all`
  - **transducer noodlification**: sequence of transducers
  - stable solution  $\rightsquigarrow$  construct **multitape transducer** composing the same variables substitutions
  - Parikh image for lengths
- **hard-positional constraints**:  $\neg$ contains [PLDI'25]
  - conversion to **quantified LIA** after stable solution found



# Experimental evaluation

- **SMT-LIB benchmarks**, split into 3 categories:
  - **Regex** (mainly regular and length constraints): **AutomatArk**, **Denghang**, **Redos**, **StringFuzz**, **Sygus-qgen**
  - **Equations** (mostly word equations and length constraints with some small number of more complex constraints): **Kaluza**, **Kepler**, **Norn**, **Omark**, **Slent**, **Slog**, **Webapp**, **Woorpje**
  - **Predicates** (complex predicates): **FullStrInt**, **LeetCode**, **PyEx**, **StrSmallRw**, **Transducer+**
- Timeout: 120 s, memory limit: 8 GiB

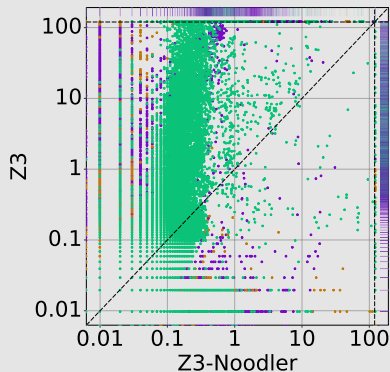
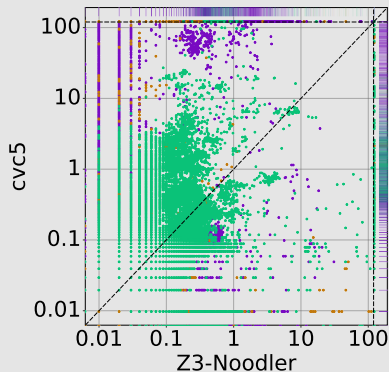
# Experimental evaluation: Procedures comparison

	number of calls	Regex proc.		Nielsen transf.		Length-based		Stabillization-based	
		called	solved	called	solved	called	solved	called	solved
<b>Sygu-s-qgen</b>	747	100%	100%	0%	0%	0%	0%	0%	0%
<b>Denghang</b>	999	0.10%	0.10%	0%	0%	96.10%	96.10%	3.80%	3.80%
<b>AutomatArk</b>	20,062	99.97%	99.97%	0%	0%	0.02%	0.02%	0.01%	0.01%
<b>StringFuzz</b>	9,941	46.45%	46.45%	0%	0%	27.98%	27.96%	25.58%	25.58%
<b>Redos</b>	2,952	70.02%	70.02%	0%	0%	11.21%	11.21%	18.77%	18.77%
<hr/>									
<b>Full Regex</b>	34,701	79.21%	79.21%	0%	0%	11.75%	11.74%	9.04%	9.04%
<hr/>									
<b>LeetCode</b>	874	1.37%	1.37%	0%	0%	59.27%	16.70%	81.92%	81.92%
<b>StrSmallRw</b>	6,327	0%	0%	0%	0%	4.85%	3.75%	96.25%	96.25%
<b>PyEx</b>	26,045	0.10%	0.10%	0%	0%	0.08%	0.08%	99.82%	99.82%
<b>FullStrInt</b>	9,003	0.04%	0.04%	0%	0%	0.26%	0.26%	99.70%	99.70%
<b>Transducer+</b>	0	-	-	-	-	-	-	-	-
<hr/>									
<b>Full Predicates</b>	42,249	0.10%	0.10%	0%	0%	2.06%	1.01%	98.89%	98.89%
<hr/>									
<b>Norn</b>	918	11.76%	11.76%	0%	0%	6.86%	6.86%	81.37%	81.37%
<b>Slog</b>	1,565	25.37%	25.37%	0%	0%	0.13%	0.13%	74.50%	74.50%
<b>Slent</b>	1,489	0.40%	0.40%	0%	0%	35.19%	30.09%	69.51%	69.51%
<b>Omark</b>	9	0%	0%	11.11%	11.11%	11.11%	0%	88.89%	88.89%
<b>Kepler</b>	579	0%	0%	99.83%	99.83%	0%	0%	0%	0%
<b>Woorpje</b>	478	0.84%	0.84%	43.10%	42.47%	30.96%	27.20%	20.50%	20.50%
<b>Webapp</b>	381	0.52%	0.52%	0%	0%	2.36%	0.26%	99.21%	99.21%
<b>Kaluza</b>	11,222	35.31%	35.31%	0%	0%	63.45%	61.78%	2.91%	2.91%
<hr/>									
<b>Full Equations</b>	16,641	26.92%	26.92%	4.72%	4.70%	47.27%	45.53%	22.59%	22.59%
<hr/>									
<b>All</b>	93,591	34.20%	34.20%	0.84%	0.84%	13.69%	12.91%	52.01%	52.01%

# Experimental evaluation: Running time

	Regex (32,242)		Equations (25,727)		Predicates (45,436)		All (103,405)	
	solved	time	solved	time	solved	time	solved	time
Z3-Noodler	32,232	3,688	25,301	1,147	45,035	6,353	<b>102,568</b>	<b>11,118</b>
Z3-Noodler <sup>M</sup>	32,228	4,010	25,299	1,456	45,035	7,321	<b>102,562</b>	<b>12,787</b>
cvc5	29,290	59,705	25,214	2,529	45,337	11,627	99,841	73,861
cvc5 <sup>M</sup>	29,287	59,892	25,214	2,756	45,337	12,220	99,838	74,868
Z3	29,075	51,379	24,569	3,240	44,101	74,094	97,745	128,712
Z3 <sup>M</sup>	29,064	51,830	24,571	4,013	44,096	74,708	97,731	130,551

# Experimental evaluation: Comparison with other solvers



Times are in seconds, axes are logarithmic, timeouts on side dashed lines (120 s)

● **Regex**, ● **Equations**, and ● **Predicates**.

# Conclusion and Acknowledgments

## ■ List of collaborators:

- František Blahoudek
- Yu-Fang Chen
- David Chocholatý
- Vojtěch Havlena
- Michal Hečko
- Lukáš Holík
- Jan Hranička
- Ondřej Lengál
- Juraj Síc



# Conclusion and Acknowledgments

## ■ List of collaborators:

- František Blahoudek
- Yu-Fang Chen
- David Chocholatý
- Vojtěch Havlena
- Michal Hečko
- Lukáš Holík
- Jan Hranička
- Ondřej Lengál
- Juraj Síč



## ■ Future work:

- optimization of **transducer constraints**
- application of Z3-Noodler on the **analysis of the security of web applications**