

# Solving String Constraints with Lengths by Stabilization

Yu-Fang Chen<sup>2</sup>, **David Chocholatý**<sup>1</sup>, Vojtěch Havlena<sup>1</sup>,  
Lukáš Holík<sup>1</sup>, Ondřej Lengál<sup>1</sup>, and Juraj Síč<sup>1</sup>

<sup>1</sup>Faculty of Information Technology, Brno University of Technology, Czech Republic

<sup>2</sup>Institute of Information Science, Academia Sinica, Taiwan



OOPSLA'23



# Motivation: SMT and string constraint solving

- SMT solvers are in demand in the industry
- Reasoning about strings is important:
  - Analysis of string manipulations in programs
  - Prevention of attacks such as XSS, SQL-injection, ... (scripting languages rely heavily on strings)
  - Regular-expression-based filters
  - AWS cloud access control policy analysis  
(Neha Rungta. *A Billion SMT Queries a Day*. CAV'22)
- Number of powerful SMT string constraint solvers.  
*Z3, cvc5, OSTRICH, Z3-Trau, Z3str3RE, Z3str4, ...*

# Motivation: SMT and string constraint solving

- SMT solvers are in demand in the industry
- Reasoning about strings is important:
  - Analysis of string manipulations in programs
  - Prevention of attacks such as XSS, SQL-injection, ... (scripting languages rely heavily on strings)
  - Regular-expression-based filters
  - AWS cloud access control policy analysis  
(Neha Rungta. *A Billion SMT Queries a Day*. CAV'22)
- Number of powerful SMT string constraint solvers.  
*Z3, cvc5, OSTRICH, Z3-Trau, Z3str3RE, Z3str4, ...*

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

# Our Approach

- Combination of two automata-based algorithms:

# Our Approach

- Combination of two automata-based algorithms:
  - 1 **STABILIZATION decision procedure** (FM'23).
    - Tightly integrating equations with regular constraints
    - Fast, but does not handle length constraints

- Supports:

- Our earlier versions

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

# Our Approach

- Combination of two automata-based algorithms:
  - 1 **STABILIZATION decision procedure** (FM'23).
    - Tightly integrating equations with regular constraints
    - Fast, but does not handle length constraints
  - 2 **ALIGN&SPLIT** (CAV'14)
    - Slower, but handles length constraints

- Supports:

- Our earlier versions

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

- Now

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

# ALIGN&SPLIT

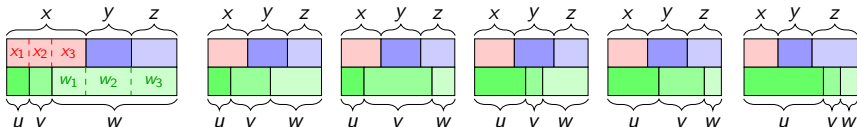
$$xyz = uvw \wedge x \in \mathcal{R}_x \wedge y \in \mathcal{R}_y \wedge z \in \mathcal{R}_z \wedge u \in \mathcal{R}_u \wedge v \in \mathcal{R}_v \wedge w \in \mathcal{R}_w \wedge \dots \wedge |x| = 100 + 4k$$

# ALIGN&SPLIT

$$xyz = uvw \wedge x \in \mathcal{R}_x \wedge y \in \mathcal{R}_y \wedge z \in \mathcal{R}_z \wedge u \in \mathcal{R}_u \wedge v \in \mathcal{R}_v \wedge w \in \mathcal{R}_w \wedge \dots \wedge |x| = 100 + 4k$$



Equation alignments (6)



$u/x_1$

$v/x_2$

$w/x_3yz$

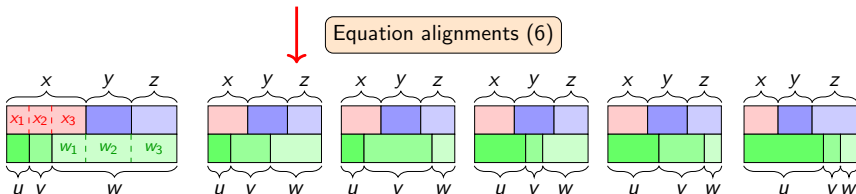
$x/x_1x_2x_3$

$x_1x_2x_3 \in \mathcal{R}_x$



# ALIGN&SPLIT

$$xyz = uvw \wedge x \in \mathcal{R}_x \wedge y \in \mathcal{R}_y \wedge z \in \mathcal{R}_z \wedge u \in \mathcal{R}_u \wedge v \in \mathcal{R}_v \wedge w \in \mathcal{R}_w \wedge \dots \wedge |x| = 100 + 4k$$



$u/x_1$   
 $v/x_2$   
 $w/x_3yz$   
 $x/x_1x_2x_3$

Automata  
splitting

Pure regular constraints formula:

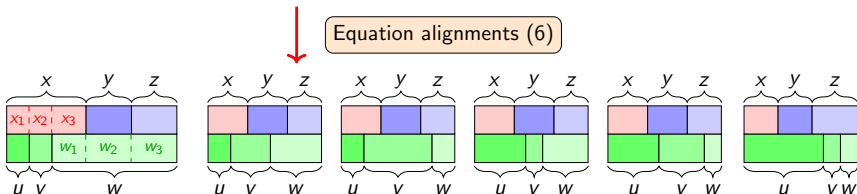
$$x_1 \in \mathcal{R}_{x_1} \wedge x_2 \in \mathcal{R}_{x_2} \wedge x_3 \in \mathcal{R}_{x_3}$$

(Monadic decomposition)

$$x_1x_2x_3 \in \mathcal{R}_x$$

# ALIGN&SPLIT

$$xyz = uvw \wedge x \in \mathcal{R}_x \wedge y \in \mathcal{R}_y \wedge z \in \mathcal{R}_z \wedge u \in \mathcal{R}_u \wedge v \in \mathcal{R}_v \wedge w \in \mathcal{R}_w \wedge \dots \wedge |x| = 100 + 4k$$



$u/x_1$   
 $v/x_2$   
 $w/x_3yz$   
 $x/x_1x_2x_3$

Automata  
splitting

Pure regular constraints formula:

$$x_1 \in \mathcal{R}_{x_1} \wedge x_2 \in \mathcal{R}_{x_2} \wedge x_3 \in \mathcal{R}_{x_3}$$

(Monadic decomposition)

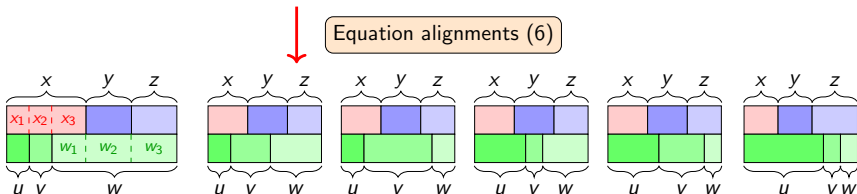
$$x_1x_2x_3 \in \mathcal{R}_x$$

Length  
extraction

$$\varphi[|x_1|, |x_2|, |x_3|, |y|, |z|] \wedge |x| = 100 + 4k$$

# ALIGN&SPLIT

$$xyz = uvw \wedge x \in \mathcal{R}_x \wedge y \in \mathcal{R}_y \wedge z \in \mathcal{R}_z \wedge u \in \mathcal{R}_u \wedge v \in \mathcal{R}_v \wedge w \in \mathcal{R}_w \wedge \dots \wedge |x| = 100 + 4k$$



$u/x_1$   
 $v/x_2$   
 $w/x_3yz$   
 $x/x_1x_2x_3$

Automata  
splitting

Pure regular constraints formula:

$$x_1 \in \mathcal{R}_{x_1} \wedge x_2 \in \mathcal{R}_{x_2} \wedge x_3 \in \mathcal{R}_{x_3}$$

(Monadic decomposition)

$$x_1x_2x_3 \in \mathcal{R}_x$$

Length  
extraction

$$\varphi[|x_1|, |x_2|, |x_3|, |y|, |z|] \wedge |x| = 100 + 4k$$

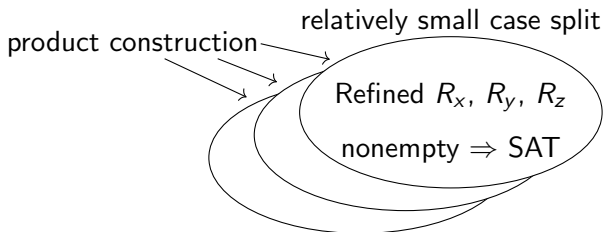
SMT LIA

# STABILIZATION

$$xyz = uvw \wedge x \in \mathcal{R}_x \wedge y \in \mathcal{R}_y \wedge z \in \mathcal{R}_z \wedge u \in \mathcal{R}_u \wedge v \in \mathcal{R}_v \wedge w \in \mathcal{R}_w \wedge \dots \wedge |x| = 100 + 4k$$

$$R_x \# R_y \# R_z \times R_u.R_v.R_w$$

- no new variables, no alignments 😊
- no monadic decomposition ☹️

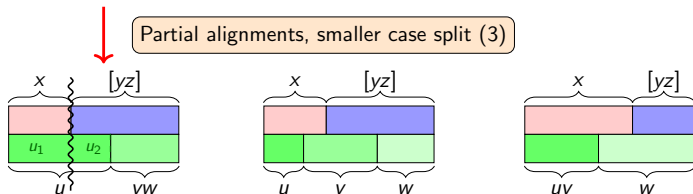


# Combined algorithm

$$xyz = uvw \wedge x \in \mathcal{R}_x \wedge y \in \mathcal{R}_y \wedge z \in \mathcal{R}_z \wedge u \in \mathcal{R}_u \wedge v \in \mathcal{R}_v \wedge w \in \mathcal{R}_w \wedge \dots \wedge |x| = 100 + 4k$$

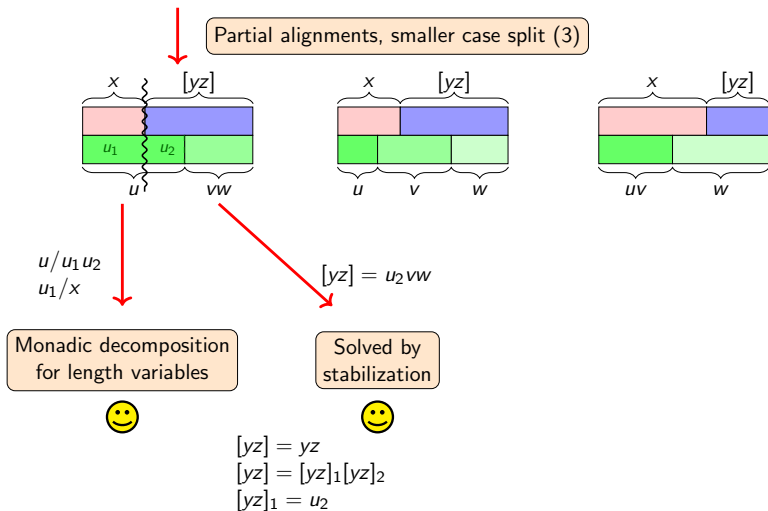
# Combined algorithm

$$xyz = uvw \wedge x \in \mathcal{R}_x \wedge y \in \mathcal{R}_y \wedge z \in \mathcal{R}_z \wedge u \in \mathcal{R}_u \wedge v \in \mathcal{R}_v \wedge w \in \mathcal{R}_w \wedge \dots \wedge |x| = 100 + 4k$$



# Combined algorithm

$$xyz = uvw \wedge x \in \mathcal{R}_x \wedge y \in \mathcal{R}_y \wedge z \in \mathcal{R}_z \wedge u \in \mathcal{R}_u \wedge v \in \mathcal{R}_v \wedge w \in \mathcal{R}_w \wedge \dots \wedge |x| = 100 + 4k$$



# Extensions and Completeness

- **Extended string constr** (`substr()`, `indexof()`, `replace()`, `contains()`, ...) rewritten to **basic constr** (equations, lengths, regular)
- **Saturation with basic constraints** further improves performance



# Extensions and Completeness

- **Extended string constr** (`substr()`, `indexof()`, `replace()`, `contains()`, ...) rewritten to **basic constr** (equations, lengths, regular)
- **Saturation with basic constraints** further improves performance
- A new reduction of **unrestricted disequations** to length constraints
- Complete on the **chain-free s.c. + unrestricted disequations**
  - (Abdulla, Atig, Diep, Holík, Janků. *Chain-Free String Constraints*. ATVA'19)
  - yields a new **largest known decidable fragment** of
    - (dis)equations,
    - regular and
    - length constraints

Z3-NOODLER<sup>1</sup>:

- Based on DPLL(T) SMT solver [Z3](#), replaces its theory of strings
- Relies on our fast automata library [MATA](#)<sup>2</sup>
- Very fast and competitive

---

<sup>1</sup><https://github.com/VeriFIT/z3-noodler>

<sup>2</sup><https://github.com/VeriFIT/mata>

# Experimental Evaluation<sup>3</sup>

## Unsolved instances

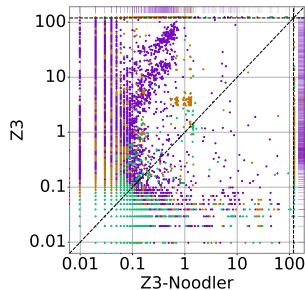
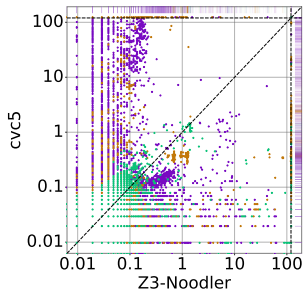
	Regex				Equations							Predicates			PyEx
	Aut	Den	StrFuzz	Syg	Kal	Kep	Norn	Slent	Slog	Web	Woo	StrInt	Leet	StrSm	
<i>Included</i>	15,995	999	10,050	343	19,432	587	1,027	1,128	1,976	267	809	11,669	2,652	1,670	23,845
<i>Unsupported</i>	0	0	1,568	0	0	0	0	0	0	414	0	5,299	0	210	0
Z3-NOODLER	62	0	0	0	259	4	0	5	0	0	243	4	4	55	4,424
CVC5	94	18	1037	0	0	240	85	22	0	40	54	0	0	4	34
Z3	113	118	340	0	164	313	124	74	71	61	25	4	0	32	1,071
Z3STR4	60	4	27	0	174	254	73	73	16	62	78	5	4	37	570
OSTRICH	55	15	229	0	288	387	1	130	7	65	53	37	26	*106	12,290
Z3STR3RE	66	27	*143	1	*144	311	133	87	55	*104	*118	64	192	*179	17,764

## Average runtimes

	Reg		Eq		Pred	
	avg	std	avg	std	avg	std
Z3-NOODLER	0.11	1.35	0.11	2.13	0.11	2.16
CVC5	1.17	8.51	0.11	2.15	0.03	0.15
Z3	1.92	9.71	0.18	2.83	0.04	0.42
Z3STR4	0.35	2.00	0.25	3.40	0.02	0.31
OSTRICH	4.29	8.67	4.28	9.28	12.71	15.08
Z3STR3RE	0.31	3.28	0.13	2.72	0.01	0.08

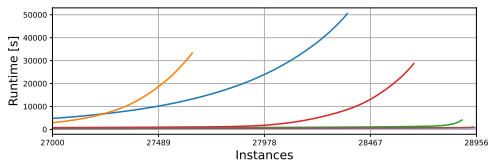
<sup>3</sup>newer results than in the paper

# Comparison with cvc5 and Z3

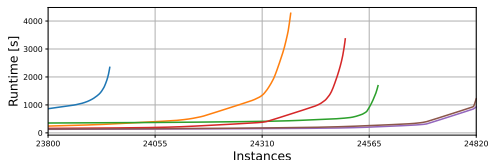


● **Regex**, ● **Equations**, ● **Predicates**.

# Virtual Best Solver

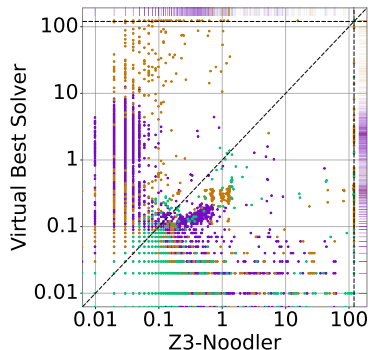


(a) VBS on regex-heavy



(b) VBS on equation-heavy

- z3
- cvc5
- Noodler
- cvc5+z3
- Noodler+cvc5+z3
- Noodler+cvc5



	Regex		Equations	
	Unsolved	Time	Unsolved	Time
VBS	1	427	19	1,304
<b>VBS - Z3-NOODLER</b>	<b>1</b>	<b>2,914</b>	<b>131</b>	<b>6,830</b>
VBS - cvc5	1	549	145	1,401
VBS - Z3	1	430	29	1,579
VBS - Z3STR4	1	473	19	1,416
VBS - OSTRICH	1	427	21	1,270
VBS - Z3STR3RE	1	510	20	1,307
cvc5 + Z3 + Z3-NOODLER	1	608	22	1,471
<b>cvc5 + Z3</b>	<b>278</b>	<b>27,916</b>	<b>303</b>	<b>2,805</b>

# Conclusion

- Combination of
  - **STABILIZATION** (efficient synergy of equations and regular constr.)
  - **ALIGN&SPLIT** (to handle lengths, disequations, and extended constr.)
- New SMT string solver **Z3-NOODLER**<sup>4</sup>
  - Fast and complementary to s.o.t.a. string solvers
  - Relying on our new fast automata library **MATA**<sup>5</sup>
- Ongoing work
  - other constraints – transducers, string/int conversion, ...
  - efficiency of automata
  - benchmark-specific heuristics
  - witness/certificate generation

---

<sup>4</sup><https://github.com/VeriFIT/z3-noodler>

<sup>5</sup><https://github.com/VeriFIT/mata>