

# A Decision Procedure for String to Code Point Conversion

**Andrew Reynolds**

Cesare Tinelli



Andres Noetzli

Clark Barrett

IJCAR 2020

July 3, 2020

**Stanford**  
University

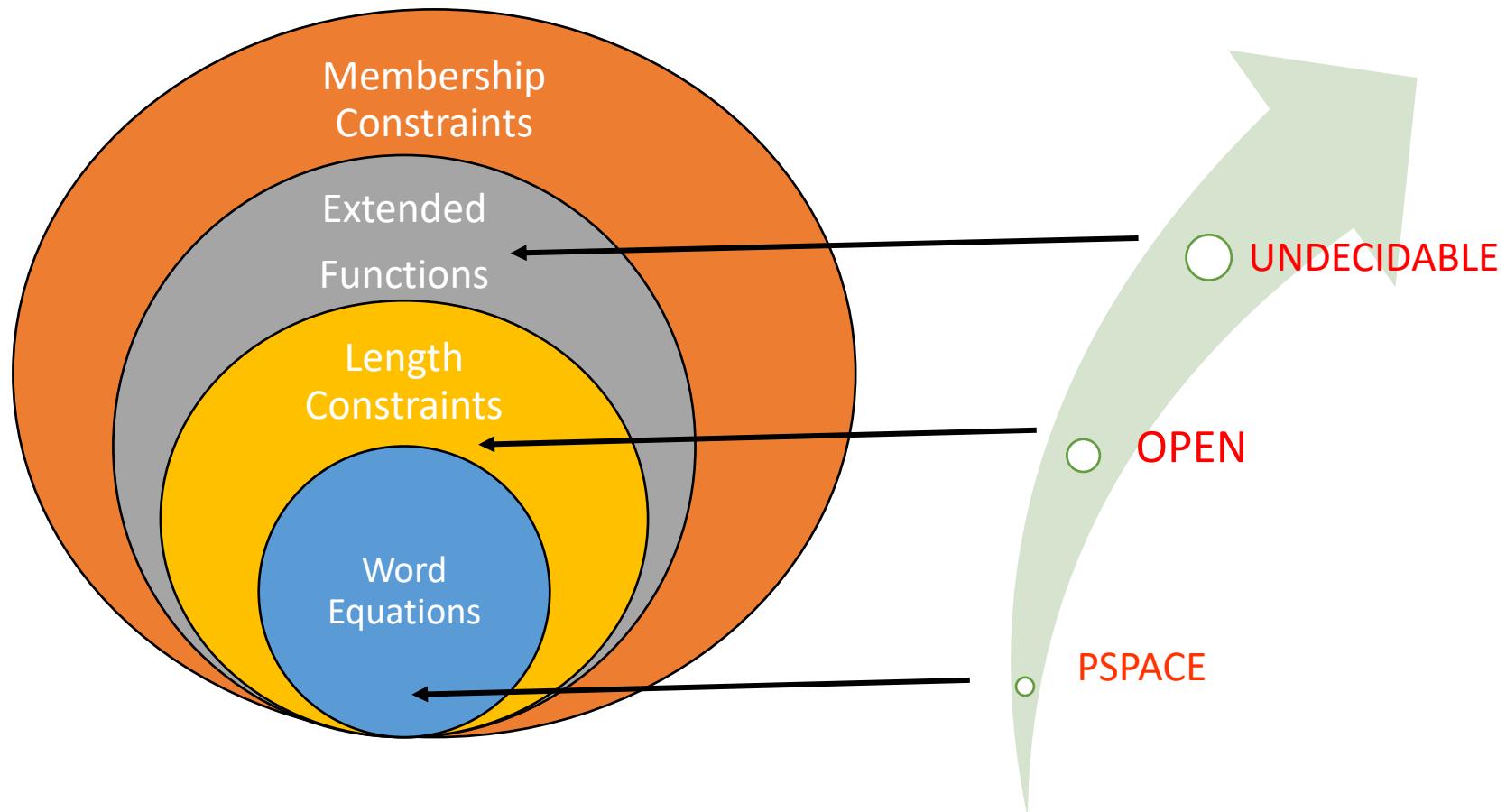
# SMT Solvers for Strings

- Strings are important in formal methods applications:
  - SQL injections, XSS vulnerabilities in web applications
  - Reasoning about access policies in cloud applications
  - Data processing
- In the past decade, efficient SMT solvers for strings have emerged
  - **cvc4**, Norn, S3 , Trau, Z3, Z3str, ...

# Overview

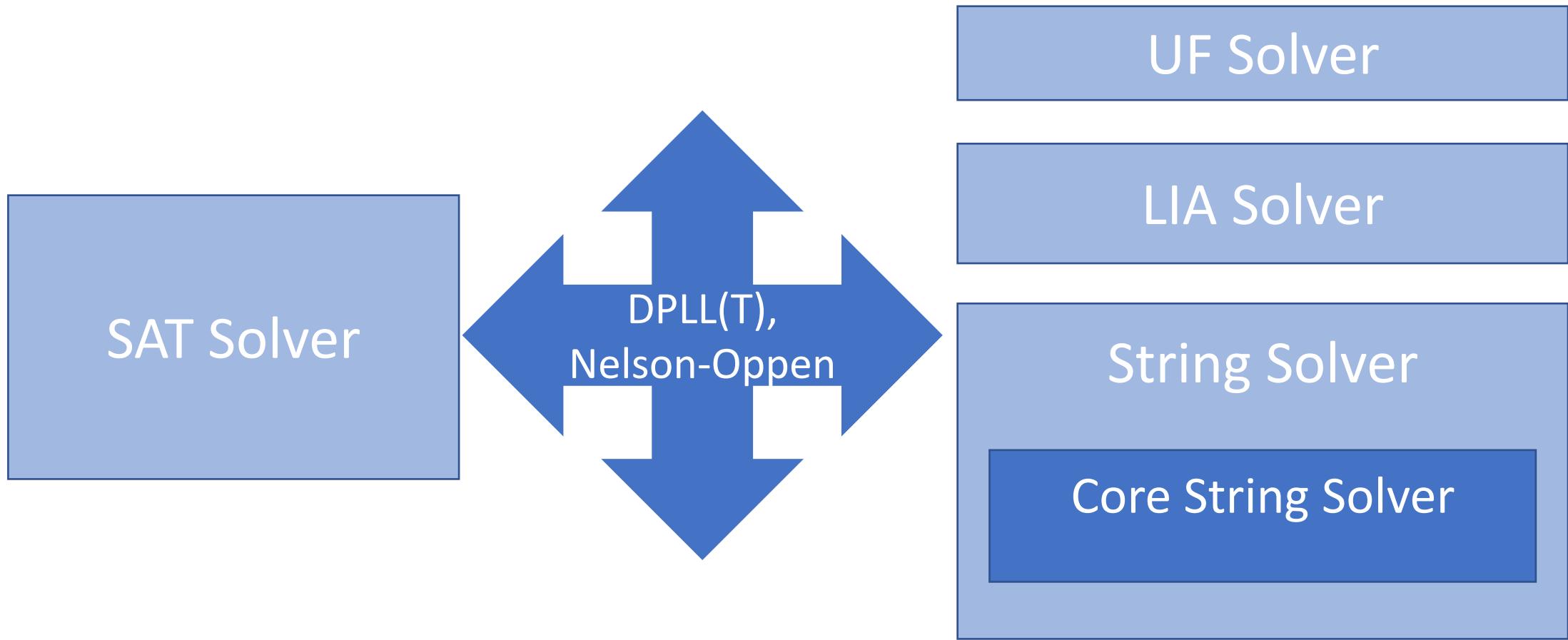
- A **decision procedure** for a simple theory of strings  $T_{AS}$ 
  - Strings + length constraints + *string-to-code conversion*
- **Improved reductions** for extended string functions
- Implementation in CVC4
- Experimental Results
  - Significant improvement in state-of-the-art for SMT string solvers

# Strings and RegExp: Theoretical Challenges



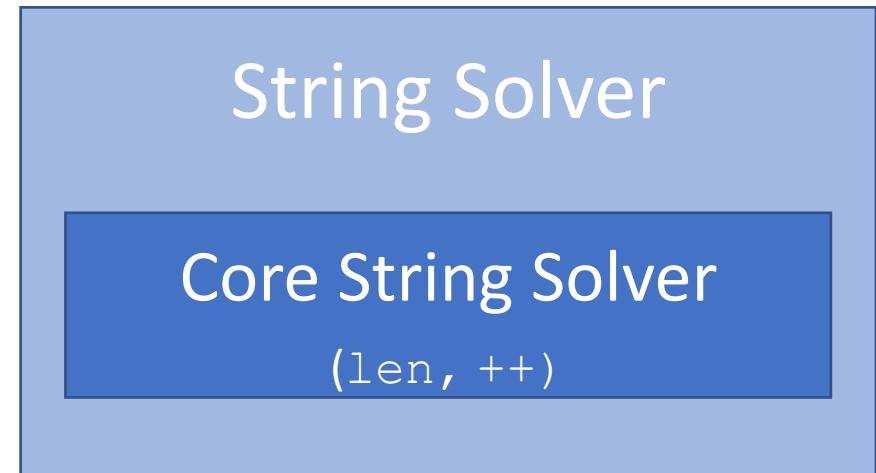
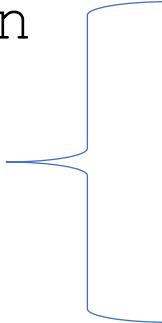
- Many applications require *extended string functions* and *RegEx memberships*
  - `to_int(x)≠44`, `to_lower(x) = "abc"`, `x ∈ range("A", "Z")`

# SMT Solvers for Strings + Extended Functions



# SMT Solvers for Strings + Extended Functions

- *Core string solver [Liang et al CAV 2014]*
  - (Dis)equalities over strings:
    - Variables, constants, concatenation `++`
  - (Dis)equalities with string length `len`

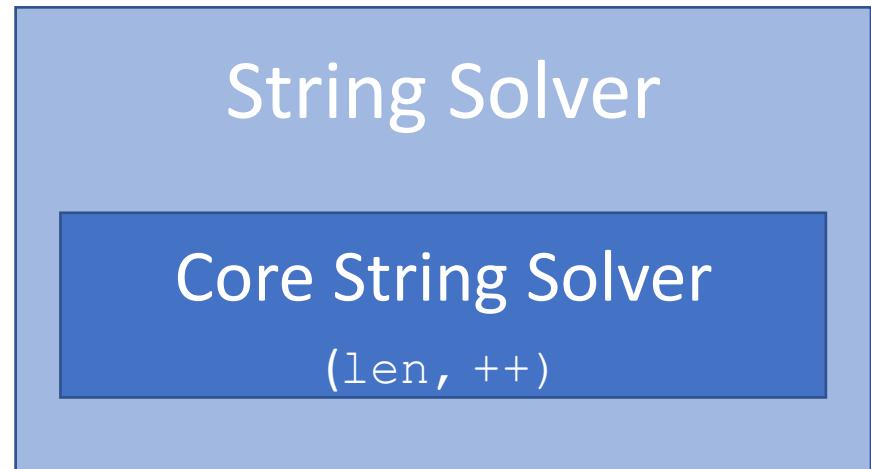


# SMT Solvers for Strings + Extended Functions

```
len(x) > 5 ∧  
x++y = z++"A" ∧  
to_int(x) ≠ z
```

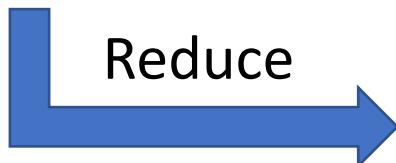
- Extended functions can be reduced to word equations

[Reynolds et al CAV 2017]



# SMT Solvers for Strings + Extended Functions

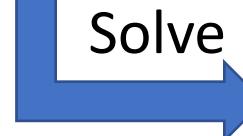
```
len(x)>5 ∧  
x++y=z++"A" ∧  
to_int(x)≠z
```



```
... ∧ i_x ≠ z ∧  
¬F_is_num ⇒ i_x = -1 ∧  
F_is_num ⇒ (i_x = sti(len(x)) ∧ sti(0) = 0 ∧  
∀0 ≤ i < len(x). sti(i+1) = 10 * sti(i) +  
(ite x[i] = "9" 9  
(ite x[i] = "8" 8 ...  
(ite x[i] = "1" 1 0) ...)
```

Core String Solver

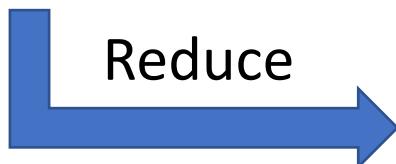
(len, ++)



SAT or UNSAT

# SMT Solvers for Strings + Extended Functions (via code)

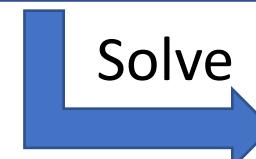
```
len(x) > 5 ∧  
x++y = z++"A" ∧  
to_int(x) ≠ z
```



```
... ∧ i_x ≠ z ∧  
¬F_is_num ⇒ i_x = -1 ∧  
F_is_num ⇒ (i_x = sti(len(x)) ∧ sti(0) = 0 ∧  
∀0 ≤ i < len(x). sti(i+1) = 10 * sti(i) +  
(code(x[i]) - 48)
```

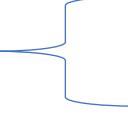
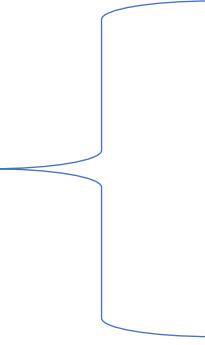
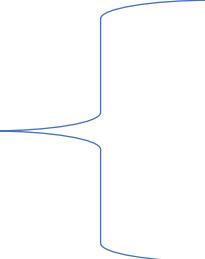
- In this paper:
  - New procedure for string-to-code point
  - Enables **efficient reductions** for constraints in common applications

Core String Solver  
(`len`, `++`, `code`)



SAT or UNSAT

# Signature of Theory of Strings $\Sigma_{AS}$

Sorts		Int                    Str                    Bool
$\Sigma_A$		$n : \text{Int}$ for all $n \in \mathbb{Z}$ $+ : \text{Int} \times \text{Int} \rightarrow \text{Int}$ $- : \text{Int} \rightarrow \text{Int}$ $\geq : \text{Int} \times \text{Int} \rightarrow \text{Bool}$
$\Sigma_S$		$l : \text{Str}$ for all $l \in \mathcal{A}^*$ $\text{len} : \text{Str} \rightarrow \text{Int}$ <b>code</b> : $\text{Str} \rightarrow \text{Int}$

... (Dis)equalities  $=, \neq$  over all sorts

$\Rightarrow$  Intentionally omit string concatenation

# String to code operator **code**

- Assume ordering on characters of alphabet  $\mathcal{A}$ :
  - $c_1 < \dots < c_{|\mathcal{A}|-1}$
  - For each  $c_i$ , we call  $i$  its code point
- $\text{code} : \text{Str} \rightarrow \text{Int}$  is interpreted as:
  1. For  $w$  in  $\mathcal{A}^1$ ,  $\text{code}(w)$  is the code point of the single character in  $w$
  2. For all other  $w$ ,  $\text{code}(w)$  is  $-1$
- Notes:
  - In practice, we are interested in reasoning about Unicode strings
    - E.g.  $\text{code}("A")=65$ ,  $\text{code}("5")=53$
  - Alphabet can have any cardinality ( $\geq 1$ )

# Rule-based calculus for strings

- Rules operate on configurations  $\langle A, S \rangle$  where:
  - **A**: a set of arithmetic constraints
  - **S**: a set of string constraints
- Each rule in guarded assignment form:

	required conditions	or	required conditions
<b>RuleName</b>			
	modifications to the configuration		“unsat”

- Some rules have multiple conclusions, delimited by “||”:

	required conditions			
modification #1		...		modification #n

# Rule-based calculus for strings

- A run of the procedure generates a ***derivation tree***:
  - Nodes are configurations
  - Child nodes are obtained by legally applying a rule to its parent
- Given derivation tree with root  $\langle A, S \rangle$ 
  - A closed derivation tree (where all leaves are “**unsat**”)  
 $\Rightarrow \langle A, S \rangle$  is  $T_{AS}$ -unsat
  - A derivation tree with a saturated configuration  
 $\Rightarrow \langle A, S \rangle$  is  $T_{AS}$ -sat

# Core rules

$$\text{A-Conf} \frac{\text{A} \models_{\text{LIA}} \perp}{\text{unsat}}$$

$$\text{A-Prop} \frac{\text{A} \models_{\text{LIA}} s \approx t \quad s, t \in \mathcal{T}(\text{A} \cup \text{S})}{\text{S} := \text{S}, s \approx t}$$

$$\text{S-Conf} \frac{\text{S} \models \perp}{\text{unsat}}$$

$$\text{S-Prop} \frac{\text{S} \models s \approx t \quad s, t : \text{Int} \quad s, t \in \mathcal{T}(\text{A} \cup \text{S})}{\text{A} := \text{A}, s \approx t}$$

$$\text{L-Intro} \frac{\text{S} \models x \approx l \quad x : \text{Str}}{\text{S} := \text{S}, \text{len}(x) \approx (\text{len}(l)) \downarrow}$$

$$\text{L-Valid} \frac{x \in \mathcal{T}(\text{A} \cup \text{S}) \quad x : \text{Str}}{\text{S} := \text{S}, x \approx \epsilon \quad \| \quad \text{A} := \text{A}, \text{len}(x) > 0}$$

$$\text{Card} \frac{\text{S} \models \text{len}(x_1) \approx \dots \approx \text{len}(x_n) \quad n > 1}{\|_{1 \leq i < j \leq n} \quad \text{S} := \text{S}, x_i \approx x_j \quad \| \quad \text{A} := \text{A}, \text{len}(x_1) > \lfloor \log_{|\mathcal{A}|} (n - 1) \rfloor}$$

- Rules adopted from previous work [\[Liang et al CAV 2014\]](#)

# Core rules

$$\boxed{\begin{array}{c} \text{A-Conf} \frac{A \models_{\text{LIA}} \perp}{\text{unsat}} \quad \text{A-Prop} \frac{A \models_{\text{LIA}} s \approx t \quad s, t \in \mathcal{T}(A \cup S)}{S := S, s \approx t} \end{array}}$$

Models LIA solver :

- Conflicts
- Propagate Int equalities from A to S

$$\boxed{\begin{array}{c} \text{S-Conf} \frac{S \models \perp}{\text{unsat}} \quad \text{S-Prop} \frac{S \models s \approx t \quad s, t : \text{Int} \quad s, t \in \mathcal{T}(A \cup S)}{A := A, s \approx t} \end{array}}$$

$$\text{L-Intro} \frac{S \models x \approx l \quad x : \text{Str}}{S := S, \text{len}(x) \approx (\text{len}(l)) \downarrow} \quad \text{L-Valid} \frac{x \in \mathcal{T}(A \cup S) \quad x : \text{Str}}{S := S, x \approx \epsilon \quad \| \quad A := A, \text{len}(x) > 0}$$

$$\text{Card} \frac{S \models \text{len}(x_1) \approx \dots \approx \text{len}(x_n) \quad n > 1}{\|_{1 \leq i < j \leq n} \quad S := S, x_i \approx x_j \quad \| \quad A := A, \text{len}(x_1) > [\log_{|\mathcal{A}|} (n - 1)]}$$

- Rules adopted from previous work [Liang et al CAV 2014]

# Core rules

$$\begin{array}{c}
 \text{A-Conf} \frac{A \models_{\text{LIA}} \perp}{\text{unsat}} \quad \text{A-Prop} \frac{A \models_{\text{LIA}} s \approx t \quad s, t \in \mathcal{T}(A \cup S)}{S := S, s \approx t} \\
 \boxed{\text{S-Conf} \frac{S \models \perp}{\text{unsat}} \quad \text{S-Prop} \frac{S \models s \approx t \quad s, t : \text{Int} \quad s, t \in \mathcal{T}(A \cup S)}{A := A, s \approx t}} \quad \boxed{\text{Models core string solver :} \\
 \bullet \text{ Conflicts by pure equality reasoning} \\
 \bullet \text{ Propagate Int equalities from } S \text{ to } A} \\
 \text{L-Intro} \frac{S \models x \approx l \quad x : \text{Str}}{S := S, \text{len}(x) \approx (\text{len}(l)) \downarrow} \quad \text{L-Valid} \frac{x \in \mathcal{T}(A \cup S) \quad x : \text{Str}}{S := S, x \approx \epsilon \quad || \quad A := A, \text{len}(x) > 0} \\
 \text{Card} \frac{S \models \text{len}(x_1) \approx \dots \approx \text{len}(x_n) \quad n > 1}{\parallel_{1 \leq i < j \leq n} S := S, x_i \approx x_j \quad || \quad A := A, \text{len}(x_1) > [\log_{|\mathcal{A}|} (n - 1)]}
 \end{array}$$

- Rules adopted from previous work [Liang et al CAV 2014]

# Core rules

$$\text{A-Conf} \frac{A \models_{\text{LIA}} \perp}{\text{unsat}}$$

$$\text{A-Prop} \frac{A \models_{\text{LIA}} s \approx t \quad s, t \in \mathcal{T}(A \cup S)}{S := S, s \approx t}$$

$$\text{S-Conf} \frac{S \models \perp}{\text{unsat}}$$

$$\text{S-Prop} \frac{S \models s \approx t \quad s, t : \text{Int} \quad s, t \in \mathcal{T}(A \cup S)}{A := A, s \approx t}$$

$$\text{L-Intro} \frac{S \models x \approx l \quad x : \text{Str}}{S := S, \text{len}(x) \approx (\text{len}(l)) \downarrow}$$

$$\text{L-Valid} \frac{x \in \mathcal{T}(A \cup S) \quad x : \text{Str}}{S := S, x \approx \epsilon \quad || \quad A := A, \text{len}(x) > 0}$$

$$\text{Card} \frac{\quad \quad \quad S \models \text{len}(x_1) \approx \dots \approx \text{len}(x_n) \quad n > 1}{\parallel_{1 \leq i < j \leq n} \quad S := S, x_i \approx x_j \quad || \quad A := A, \text{len}(x_1) > \lfloor \log_{|\mathcal{A}|} (n - 1) \rfloor}$$

- Reasoning about length:
- Fixed for constant strings
  - Length is non-negative

- Rules adopted from previous work [Liang et al CAV 2014]

# Core rules

$$\text{A-Conf} \frac{\text{A} \models_{\text{LIA}} \perp}{\text{unsat}}$$

$$\text{A-Prop} \frac{\text{A} \models_{\text{LIA}} s \approx t \quad s, t \in \mathcal{T}(\text{A} \cup \text{S})}{\text{S} := \text{S}, s \approx t}$$

$$\text{S-Conf} \frac{\text{S} \models \perp}{\text{unsat}}$$

$$\text{S-Prop} \frac{\text{S} \models s \approx t \quad s, t : \text{Int} \quad s, t \in \mathcal{T}(\text{A} \cup \text{S})}{\text{A} := \text{A}, s \approx t}$$

$$\text{L-Intro} \frac{\text{S} \models x \approx l \quad x : \text{Str}}{\text{S} := \text{S}, \text{len}(x) \approx (\text{len}(l)) \downarrow}$$

$$\text{L-Valid} \frac{x \in \mathcal{T}(\text{A} \cup \text{S}) \quad x : \text{Str}}{\text{S} := \text{S}, x \approx \epsilon \quad \| \quad \text{A} := \text{A}, \text{len}(x) > 0}$$

$$\text{Card} \frac{\text{S} \models \text{len}(x_1) \approx \dots \approx \text{len}(x_n) \quad n > 1}{\|_{1 \leq i < j \leq n} \quad \text{S} := \text{S}, x_i \approx x_j \quad \| \quad \text{A} := \text{A}, \text{len}(x_1) > [\log_{|\mathcal{A}|} (n - 1)]}$$

Reasoning about **cardinality** :

- Only finitely many strings exist for each length, since alphabet  $\mathcal{A}$  is finite.

- Rules adopted from previous work [[Liang et al CAV 2014](#)]

# Rules for code points (1/3)

$$\text{C-Intro} \frac{S \models x \approx l \quad l \in A^1}{A := A, \text{code}(x) \approx (\text{code}(l))\downarrow}$$

$$\text{C-Collapse} \frac{S \models x \approx l \quad \text{code}(x) \in T(S)}{A := A, \text{code}(x) \approx (\text{code}(l))\downarrow}$$

- Code of string  $x$  must be equal to its value, if it is forced by an equality
  - If  $x$  is equal to a constant string of length one:
    - $x = "A" \in S$  ...  $\text{code}(x) = 65 \in A$
  - If  $x$  is equal to any constant string, and  $\text{code}(x)$  is a term in  $S$ 
    - $x = "abc"$ ,  $\text{code}(x) = \text{code}(y) \in S$  ...  $\text{code}(x) = -1 \in A$

# Rules for code points (2/3)

$$\text{C-Valid} \frac{\text{code}(x) \in \mathcal{T}(\mathcal{S})}{\mathbf{A} := \mathbf{A}, \text{len}(x) \approx 1, \text{code}(x) \approx -1 \quad \parallel \quad \mathbf{A} := \mathbf{A}, \text{len}(x) \approx 1, 0 \leq \text{code}(x) < |\mathcal{A}|}$$

- For each string  $x$  :
  - The length of  $x$  is not 1 and its code point is  $-1$
  - The length of  $x$  is 1 and its code point is in range  $[0...|\mathcal{A}|-1]$

# Rules for code points (3/3)

$$\text{C-Inj} \frac{\text{code}(x), \text{code}(y) \in \mathcal{T}(S) \quad x, y \text{ distinct}}{A := A, \text{code}(x) \approx -1 \quad || \quad A := A, \text{code}(x) \not\approx \text{code}(y) \quad || \quad S := S, x \approx y}$$

- Code is injective over the domain of strings of length 1
  - For each  $\text{code}(x), \text{code}(y)$ , either:
    1. One of these strings, w.l.o.g. assume  $x$ , is not length 1 ...  $\text{code}(x) = -1$
    2.  $x$  and  $y$  have different code points ...  $\text{code}(x) \neq \text{code}(y)$
    3.  $x$  and  $y$  are the same string ...  $x = y$

# Example #1

$\langle A_0, S_0 \rangle$

- $A = \{len(x) > len(y), code(x) = code(y), code(x) \geq 0\}$
- $S = \{\}$

# Example #1

$$\frac{\langle A_0, S_0 \rangle}{\text{code}(x) = \text{code}(y) \in S} \text{-Prop}$$

- $A = \{len(x) > len(y), \text{code}(x) = \text{code}(y), \text{code}(x) \geq 0\}$
- $S = \{\text{code}(x) = \text{code}(y)\}$

# Example #1

$$\frac{\begin{array}{c} \langle A_0, S_0 \rangle \\ \hline \text{code}(x) = \text{code}(y) \in S \end{array}}{\text{code}(x) = -1 \in A \quad \text{code}(x) \neq \text{code}(y) \in A \quad x = y \in S} \quad \begin{array}{c} \text{A-Prop} \\ \hline \text{C-Inj} \end{array}$$

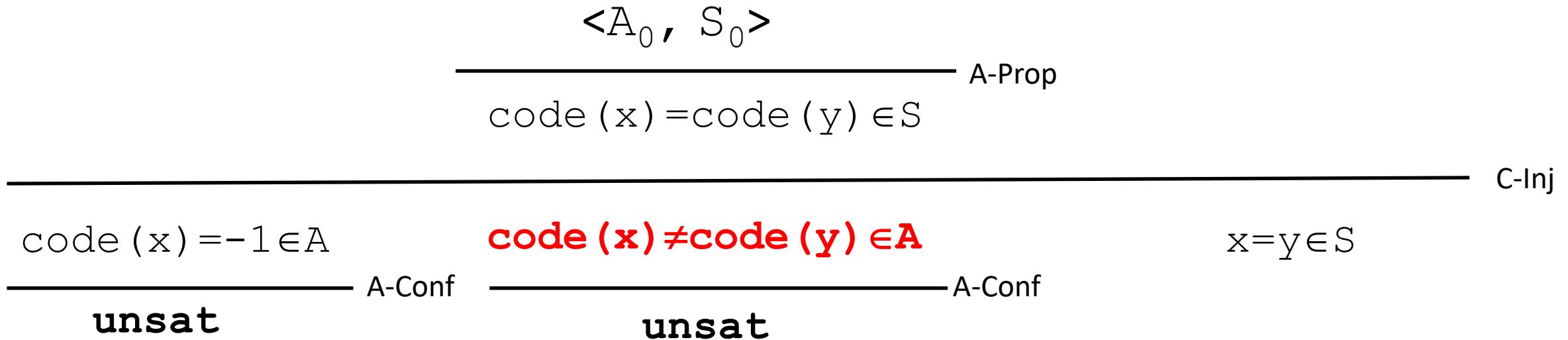
- $A = \{len(x) > len(y), \text{code}(x) = \text{code}(y), \text{code}(x) \geq 0\}$
- $S = \{\text{code}(x) = \text{code}(y)\}$

# Example #1

$$\frac{\begin{array}{c} \frac{\begin{array}{c} \langle A_0, S_0 \rangle \\ \hline \text{A-Prop} \end{array}}{\text{code}(x) = \text{code}(y) \in S} \\ \hline \text{C-Inj} \end{array}}{\begin{array}{c} \text{code}(x) = -1 \in A \\ \text{code}(x) \neq \text{code}(y) \in A \\ x = y \in S \end{array}} \\ \hline \text{A-Conf} \\ \text{unsat} \end{array}$$

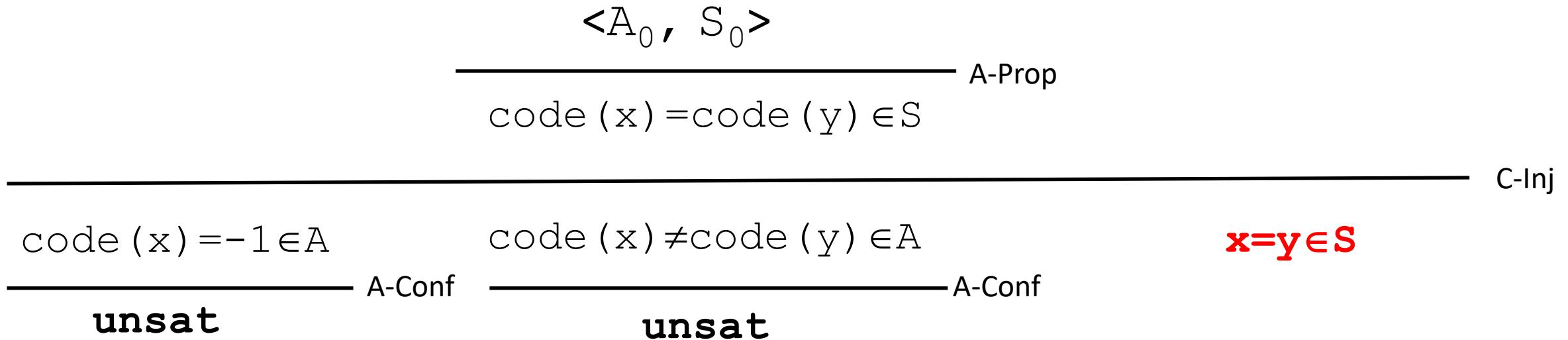
- $A = \{\text{len}(x) > \text{len}(y), \text{code}(x) = \text{code}(y), \underline{\text{code}(x) \geq 0}, \underline{\text{code}(x) = -1}\}$
- $S = \{\text{code}(x) = \text{code}(y)\}$

# Example #1



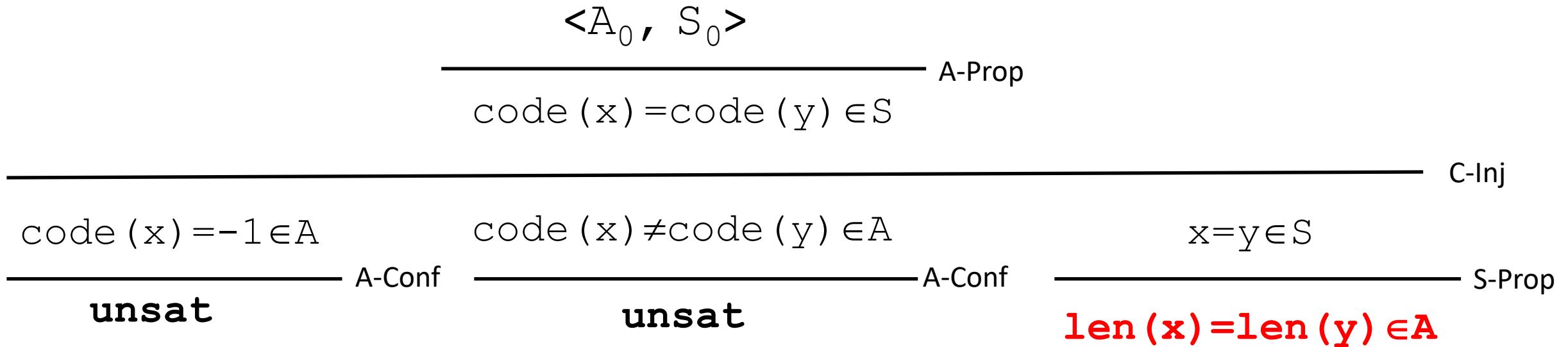
- $A = \{\text{len}(x) > \text{len}(y), \underline{\text{code}(x) = \text{code}(y)}, \text{code}(x) \geq 0, \underline{\text{code}(x) \neq \text{code}(y)}\}$
- $S = \{\text{code}(x) = \text{code}(y)\}$

# Example #1



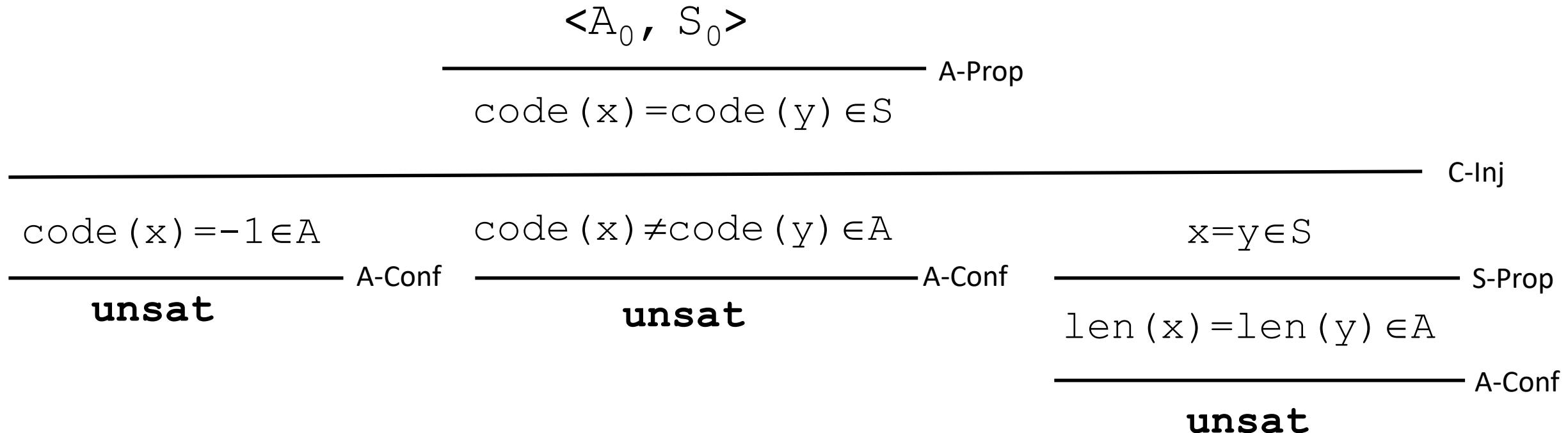
- $A = \{\text{len}(x) > \text{len}(y), \text{ code}(x) = \text{code}(y), \text{ code}(x) \geq 0\}$
- $S = \{\text{code}(x) = \text{code}(y), \text{ x=y}\}$

# Example #1



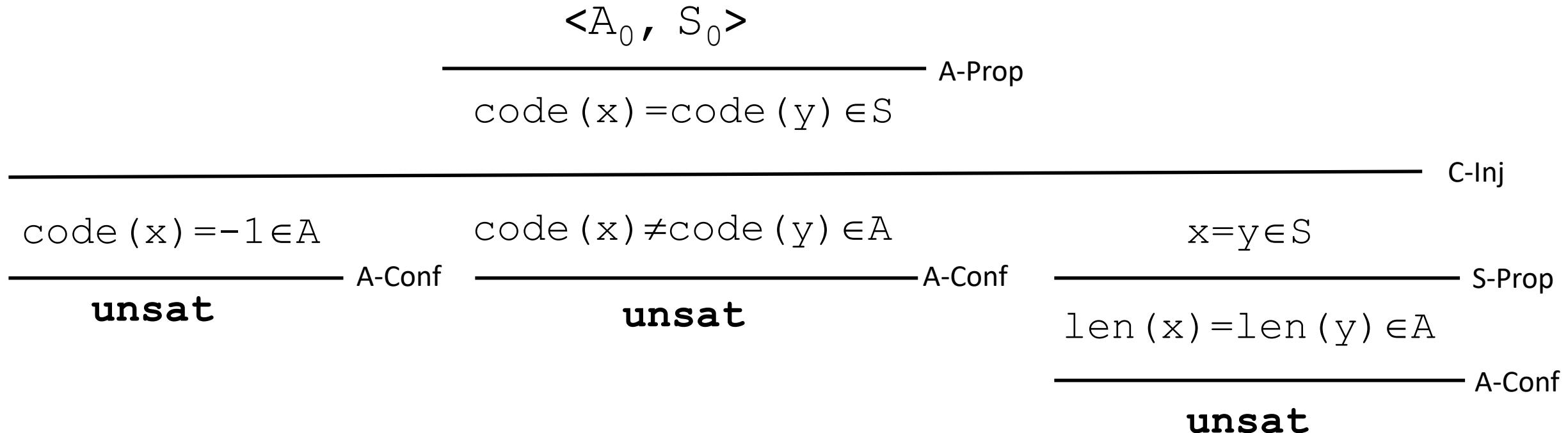
- $A = \{\text{len}(x) > \text{len}(y), \text{ code}(x) = \text{code}(y), \text{ code}(x) \geq 0, \text{ len}(x) = \text{len}(y)\}$
- $S = \{\text{code}(x) = \text{code}(y), x = y\}$

# Example #1



- $A = \{\underline{\text{len}(x) > \text{len}(y)}, \text{ code}(x) = \text{code}(y), \text{ code}(x) \geq 0, \underline{\text{len}(x) = \text{len}(y)}\}$
- $S = \{\text{code}(x) = \text{code}(y), x = y\}$

# Example #1



- Closed derivation tree ...  $\therefore A_0 \wedge S_0$  is  $T_{AS}$ -unsatisfiable

## Example #2

$\langle A_0, S_0 \rangle$

- $A = \{106 \geq \text{code}(x) \geq 97\}$
- $S = \{x \neq y, x \neq z, y = c_{97}, z = c_{106}\}$

## Example #2

$\langle A_0, S_0 \rangle$

---

C-Intro (2)

**code(y)=97, code(z)=106 ∈ A**

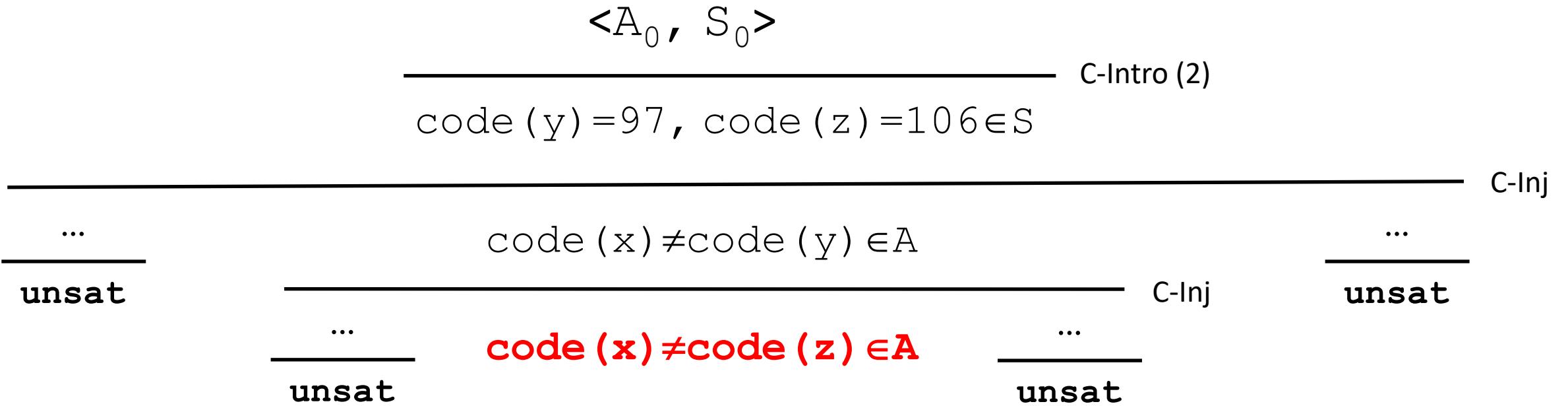
- $A = \{106 \geq \text{code}(x) \geq 97, \text{code}(y) = 97, \text{code}(z) = 106\}$
- $S = \{x \neq y, x \neq z, y = c_{97}, z = c_{106}\}$

## Example #2

	$\langle A_0, S_0 \rangle$	
		C-Intro (2)
	$\text{code}(y) = 97, \text{code}(z) = 106 \in S$	
		C-Inj
...	$\text{code}(x) \neq \text{code}(y) \in A$	...
unsat		unsat
		C-Inj
...	$\text{code}(x) \neq \text{code}(z) \in A$	...
unsat		unsat

- $A = \{106 \geq \text{code}(x) \geq 97, \text{code}(y) = 97, \text{code}(z) = 106, \text{code}(x) \neq \text{code}(y), \text{code}(x) \neq \text{code}(z)\}$
- $S = \{x \neq y, x \neq z, y = c_{97}, z = c_{106}\}$

## Example #2



- Derivation tree with saturated configuration ...  $\therefore A_0 \wedge S_0$  is  $T_{AS}$ -satisfiable
  - Model for  $x$  determined by model for  $\text{code}(x)$

# Correctness Theorem

- Calculus is:
  - **Terminating** (all derivation trees are finite)
  - **Refutation-sound**
    - A closed derivation tree with root  $M \Rightarrow M$  is  $T_{AS}$ -unsat
  - **Solution-sound**
    - A derivation tree with root  $M$  with a saturated configuration  $\Rightarrow M$  is  $T_{AS}$ -sat
      - Proof is constructive, a model for  $M$  can be produced in all such cases

∴ Calculus is a decision procedure for  $T_{AS}$

# Implementation Details

- Implemented in CVC4 SMT solver
  - As part of its DPLL(T) solver for strings
- Combined with other procedures for strings
  - E.g. [[Liang et al CAV 2014](#), [Liang et al FroCoS 2015](#), [Reynolds et al CAV 2017](#)]
  - All rules of calculus with premises  $S \models \perp$  become  $S \models_T \perp$ 
    - Where  $\models_T$  denotes entailment modulo theory T
- Strategy crafted in these combinations:
  - C-Valid, C-Collapse applied eagerly
  - C-Intro, C-Inj applied lazily



# Reductions: Conversion Functions

- More efficient reductions that leverage `code`, including:
- Conversion between strings and integers `to_int(x)`:
  - ⊗ ... `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 `to_lower(x)`:
  - ⊗ ... `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: Ordering, RegExp Range

- Lexicographic ordering:

- $\otimes \quad x \leq y \Leftrightarrow \exists i \dots (x[i] = y[i]) \vee (x[i] = "A" \wedge y[i] = "B") \vee (x[i] = "A" \wedge y[i] = "C") \dots)$
- $\Rightarrow x \leq y \Leftrightarrow \exists i \dots \mathbf{code}(x[i]) \leq \mathbf{code}(y[i])$

- Regular expression ranges:

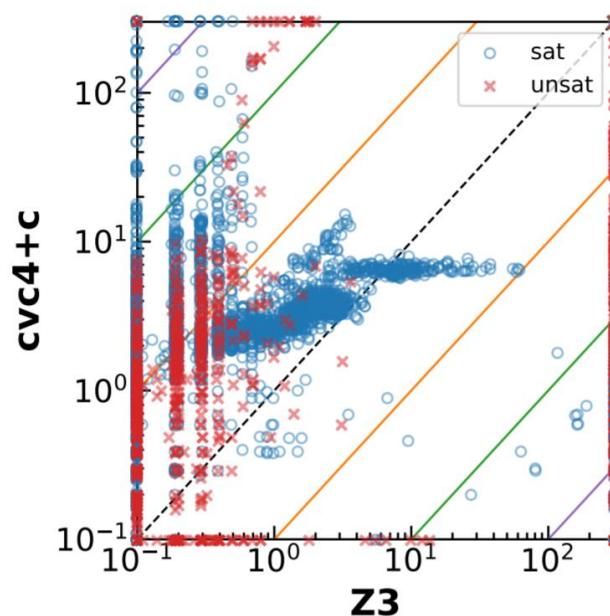
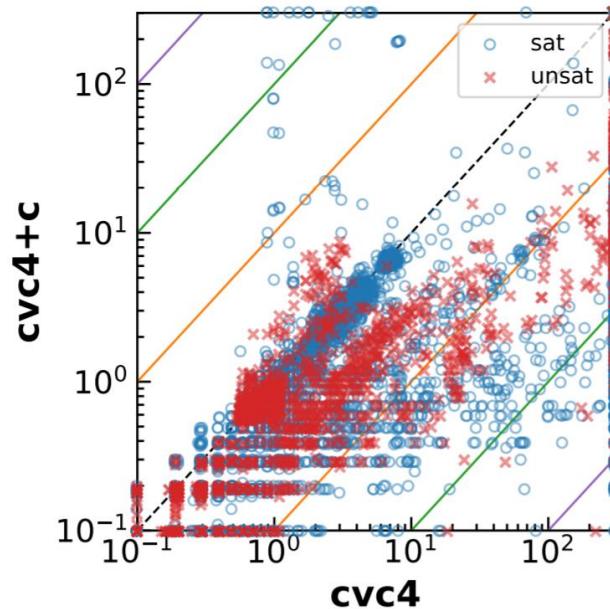
- $\otimes \quad x \in \text{range}(c_1, c_2) \Leftrightarrow \text{len}(x) = 1 \wedge (x = c_1 \vee \dots \vee x = c_2)$
- $\Rightarrow x \in \text{range}(c_1, c_2) \Leftrightarrow \mathbf{code}(c_1) \leq \mathbf{code}(x) \leq \mathbf{code}(c_2)$

# Experiments

- Compared solvers
  - **cvc4+c** : new version, includes techniques from this paper
  - **cvc4**: previous version
  - **z3**
- 21573 benchmarks from PyConByte
  - Correspond to concolic execution of Python code
  - Contain applications of `to_int`, `from_int`, regular expression ranges
- 300 second timeout

# Experimental Results

Benchmark Set		cvc4+c	cvc4	z3
py-conbyte_cvc4	sat	<b>1344</b>	1104	1187
	unsat	<b>8576</b>	8547	8482
	×	13	282	264
py-conbyte_trauc	sat	<b>1009</b>	929	697
	unsat	1424	1407	<b>1428</b>
	×	13	110	321
py-conbyte_z3seq	sat	<b>1354</b>	1126	1343
	unsat	<b>5864</b>	5797	5719
	×	35	330	191
py-conbyte_z3str	sat	<b>711</b>	652	692
	unsat	<b>1227</b>	1223	1223
	×	3	66	26
Total	sat	<b>4418</b>	3811	3919
	unsat	<b>17091</b>	16974	16852
	×	64	788	802



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

# Summary

- New **decision procedure** for string-to-code point conversion
- **Efficient reduction techniques**
  - Conversion functions, lexicographic ordering, regular expression ranges
- Implemented in CVC4
  - Combined with other procedures
- Experiments: significant improvement wrt SMT solvers for strings

# Future Work

- Support for more extended functions of interest to user applications
- Proof production for strings

