

CS:4980 Topics in Computer Science II

Introduction to Automated Reasoning

Theory Solvers III

Cesare Tinelli

Spring 2024



Credits

These slides are based on slides originally developed by **Cesare Tinelli** at the University of Iowa, and by **Clark Barrett, Caroline Trippel**, and **Andrew (Haoze) Wu** at Stanford University. Adapted by permission.

Roadmap for Today

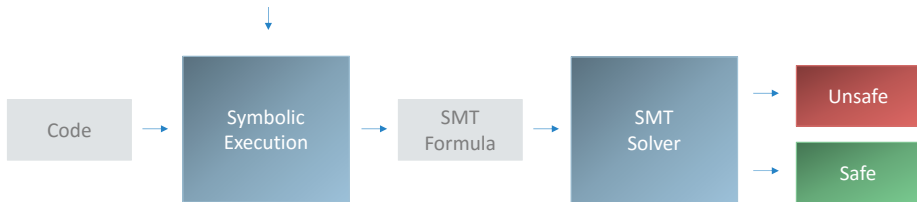
Theory Solvers

- Strings

Motivation: Symbolic Execution

Symbolic Execution

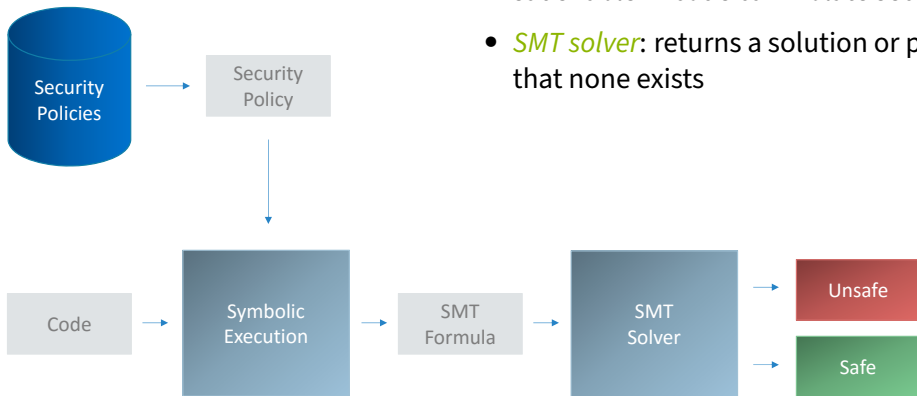
- *Enumerate program paths* that end in a bad state
 - (e.g., invalid memory access)
- Represent program *inputs* as SMT *variables*
- Translate *statements* in the path into *constraints* on the variables
- Constraints represent *all possible executions* along the path
- Solving the constraints determines whether the path is *feasible*



Example: Symbolic Execution for Security

Security Vulnerabilities

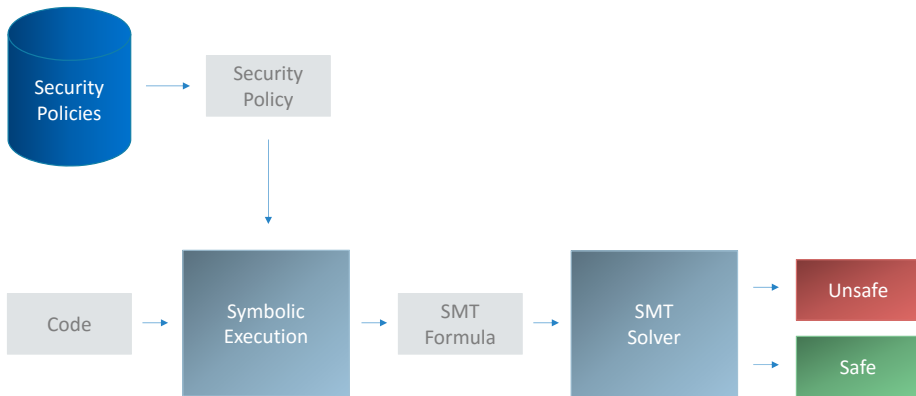
- Input: *code* and *security policy*
- *Symbolic execution*: generates formula satisfiable iff code can violate security policy
- *SMT solver*: returns a solution or proves that none exists



String Analysis

Strings in Symbolic Execution

- Input code may manipulate *strings*



Basic Theory of Strings

Alphabet

\mathcal{A} fixed *finite* set of characters

Constants

Empty string $\epsilon : \text{String}$ (i.e., $\text{rank}(\epsilon) = \langle \text{String} \rangle$)

Character string $c : \text{String}$ for all $c \in \mathcal{A}$

Integer numeral $n : \text{Int}$ for all $n \geq 0$

Operators

Concatenation $_ \cdot _ : \text{String} \times \text{String} \rightarrow \text{String}$ (i.e., $\text{rank}(\cdot) = \langle \text{String}, \text{String}, \text{String} \rangle$)

Length $|_| : \text{String} \rightarrow \text{Int}$

Membership $_ \in _ : \text{String} \times \text{RegEx} \rightarrow \text{Bool}$

Addition $_ + _ : \text{Int} \times \text{Int} \rightarrow \text{Int}$

Comparison $_ > _ : \text{Int} \times \text{Int} \rightarrow \text{Int}$

Basic Theory of Strings

Alphabet

\mathcal{A} fixed *finite* set of characters

Constants

Empty string $\epsilon : \text{String}$ (i.e., $\text{rank}(\epsilon) = \langle \text{String} \rangle$)

Character string $c : \text{String}$ for all $c \in \mathcal{A}$

Integer numeral $n : \text{Int}$ for all $n \geq 0$

Operators

Concatenation $_ \cdot _ : \text{String} \times \text{String} \rightarrow \text{String}$ (i.e., $\text{rank}(\cdot) = \langle \text{String}, \text{String}, \text{String} \rangle$)

Length $|_| : \text{String} \rightarrow \text{Int}$

Membership $_ \in _ : \text{String} \times \text{RegEx} \rightarrow \text{Bool}$

Addition $_ + _ : \text{Int} \times \text{Int} \rightarrow \text{Int}$

Comparison $_ > _ : \text{Int} \times \text{Int} \rightarrow \text{Int}$

Basic Theory of Strings

Alphabet

\mathcal{A} fixed *finite* set of characters

Constants

Empty string $\epsilon : \text{String}$ (i.e., $\text{rank}(\epsilon) = \langle \text{String} \rangle$)

Character string $c : \text{String}$ for all $c \in \mathcal{A}$

Integer numeral $n : \text{Int}$ for all $n \geq 0$

Operators

Concatenation $_ \cdot _ : \text{String} \times \text{String} \rightarrow \text{String}$ (i.e., $\text{rank}(\cdot) = \langle \text{String}, \text{String}, \text{String} \rangle$)

Length $|_| : \text{String} \rightarrow \text{Int}$

Membership $_ \in _ : \text{String} \times \text{Regex} \rightarrow \text{Bool}$

Addition $_ + _ : \text{Int} \times \text{Int} \rightarrow \text{Int}$

Comparison $_ > _ : \text{Int} \times \text{Int} \rightarrow \text{Int}$

Basic Theory of Strings

Alphabet

\mathcal{A} fixed *finite* set of

Constants

Empty string $\epsilon : \text{String}$

Character string $c : \text{String}$

Integer numeral $n : \text{Int}$

Operators

Concatenation $_ \cdot _ : \text{String} \times \text{String} \rightarrow \text{String}$ (i.e., $\text{rank}(\cdot) = \langle \text{String}, \text{String}, \text{String} \rangle$)

Length $|_ : \text{String} \rightarrow \text{Int}$

Membership $_ \in _ : \text{String} \times \text{Regex} \rightarrow \text{Bool}$

Addition $_ + _ : \text{Int} \times \text{Int} \rightarrow \text{Int}$

Comparison $_ > _ : \text{Int} \times \text{Int} \rightarrow \text{Bool}$

Challenge: complexity

concatenation + equality: *word equations problem*

- Decidable in PSPACE

+ length

- Decidability open

+ replace (all instances of some substring)

- Undecidable

Basic Theory of Strings

Alphabet

\mathcal{A} fixed *finite* set of characters

Pragmatic approach

- Rule-based proof system
- Use existing arithmetic theory solver
- Embrace incompleteness

Constants

Empty string $\epsilon : \text{String}$

Character string $c : \text{String}$

Integer numeral $n : \text{Int}$ for

Operators

Concatenation $_ \cdot _ : \text{String} \times \text{String} \rightarrow \text{String}$ (i.e., $\text{rank}(\cdot) = \langle \text{String}, \text{String}, \text{String} \rangle$)

Length $|_| : \text{String} \rightarrow \text{Int}$

Membership $_ \in _ : \text{String} \times \text{Regex} \rightarrow \text{Bool}$

Addition $_ + _ : \text{Int} \times \text{Int} \rightarrow \text{Int}$

Comparison $_ > _ : \text{Int} \times \text{Int} \rightarrow \text{Int}$

Satisfiability Proof System for Strings

Proof States

A *proof state* is either:

- One of the distinguished states **SAT**, **UNSAT**
- A pair $\langle S; A \rangle$, where S contains *string* constraints and A contains *arithmetic* constraints

Assumptions

- All literals are flat
- For every string variable x in S , there exists a variable ℓ_x , such that $\ell_x = |x| \in S$
- Ignore regular expression membership for now

Notation

Definitions

- $\mathcal{T}(S)$ denotes all terms in S
- $S \models \alpha$ means that α follows from S using the rules of QF_UF
- $A \models_{LIA} \alpha$ means that α follows from A in the theory of linear integer arithmetic

Normalization function for length

- $|\epsilon|_{\downarrow} = 0$
- $|c|_{\downarrow} = 1$ for all $c \in \mathcal{A}$
- $|s_1 \cdot \dots \cdot s_n|_{\downarrow} = |s_1|_{\downarrow} + \dots + |s_n|_{\downarrow}$

Core Rules

$$\text{A-CONF} \frac{A \models_{LIA} \perp}{\text{UNSAT}}$$

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

$$\text{S-CONF} \frac{S \models \perp}{\text{UNSAT}}$$

$$\text{S-PROP} \frac{S \models s \doteq t \quad s, t \in \mathcal{T}(S) \quad s, t \text{ are } \Sigma_{LIA}\text{-terms}}{A := A, s \doteq t}$$

$$\text{S-A} \frac{x, y \in \mathcal{T}(S) \cap \mathcal{T}(A) \quad x, y : \text{Int}}{A := A, x \doteq y \quad A := A, x \doteq y}$$

$$\text{L-INTRO} \frac{s \in \mathcal{T}(S) \quad s : \text{String}}{S := S, |s| \doteq |s| \downarrow}$$

$$\text{L-VALID} \frac{x \in \mathcal{T}(S) \quad x : \text{String}}{S := S, x \doteq \epsilon \quad A \doteq A, \ell_x > 0}$$

$$\text{CONST-CONF} \frac{S \models c \doteq d \quad c \in \mathcal{A} \quad d \in \mathcal{A} \setminus \{c\}}{\text{UNSAT}}$$

$$\text{SAT} \frac{\text{no other rule applies}}{\text{SAT}}$$

Example Derivation

Let $S_0 = \{x \doteq y \cdot x \cdot z, y \doteq "a", \ell_x \doteq |x|, \ell_y \doteq |y|, \ell_z \doteq |z|\}$

$A_0 = \emptyset$

$\langle S_0; A_0 \rangle$	
$\langle y \cdot x \cdot z \doteq y + x + z ; \emptyset \rangle$	
$\langle "a" \doteq 1; \emptyset \rangle$	
$\langle \emptyset; \ell_x \doteq \ell_y + \ell_x + \ell_z \rangle$	
$\langle \emptyset; \ell_y \doteq 1 \rangle$	
$\langle z \doteq \epsilon; \emptyset \rangle$	$\langle \emptyset; \ell_z > 0 \rangle$
$\langle c \doteq 0; \emptyset \rangle$	UNSAT
$\langle \emptyset; \ell_z \doteq 0 \rangle$	
UNSAT	

For each derivation step, we show only the **difference** between the derived state and the previous one

Example Derivation

Let $S_0 = \{x \doteq y \cdot x \cdot z, y \doteq "a", \ell_x \doteq |x|, \ell_y \doteq |y|, \ell_z \doteq |z|\}$
 $A_0 = \emptyset$

$$\begin{array}{c}
 \langle S_0; A_0 \rangle \\
 \hline
 \langle |y \cdot x \cdot z| \doteq |y| + |x| + |z|; \emptyset \rangle \\
 \hline
 \langle |"a"| \doteq 1; \emptyset \rangle \\
 \hline
 \langle \emptyset; \ell_x \doteq \ell_y + \ell_x + \ell_z \rangle \\
 \hline
 \langle \emptyset; \ell_y \doteq 1 \rangle \\
 \hline
 \begin{array}{cc}
 \langle z \doteq \epsilon; \emptyset \rangle & \langle \emptyset; \ell_z > 0 \rangle \\
 \hline
 \langle |\epsilon| \doteq 0; \emptyset \rangle & \text{UNSAT} \\
 \hline
 \langle \emptyset; \ell_z \doteq 0 \rangle & \\
 \hline
 \text{UNSAT} &
 \end{array}
 \end{array}$$

For each derivation step, we show only the **difference** between the derived state and the previous one

Concatenation Rules

If x is a variable of S , we can recursively expand x by substituting using equalities from S whose right-hand sides are concatenation terms until this is no longer possible

If t is the result, we write $S \models^* x = t$

We write z as a short-hand for a concatenation of zero or more variables ($z = x_1 \cdot x_2 \cdot \dots \cdot x_n$, with $z = \epsilon$ when $n = 0$)

$$\text{C-EQ} \quad \frac{S \models^* x \doteq z \quad S \models^* y \doteq z}{S := S, x \doteq y}$$

$$\text{C-SPLIT} \quad \frac{S \models^* x \doteq w \cdot u \cdot z \quad S \models^* x \doteq w \cdot v \cdot z'}{\begin{array}{l} A := A, l_u > l_v; S := S, u \doteq v \cdot k \\ A := A, l_u < l_v; S := S, v \doteq u \cdot k \\ A := A, l_u \doteq l_v; S := S, u \doteq v \end{array}}$$

Note: k is a fresh variable

Concatenation Rules

If x is a variable of S , we can recursively expand x by substituting using equalities from S whose right-hand sides are concatenation terms until this is no longer possible

If t is the result, we write $S \models^* x = t$

We write x as a short-hand for a concatenation of zero or more variables ($x = x_1 \cdot x_2 \cdot \dots \cdot x_n$, with $x = \epsilon$ when $n = 0$)

$$\text{C-EQ} \quad \frac{S \models^* x \doteq z \quad S \models^* y \doteq z}{S := S, x \doteq y}$$

$$\text{C-SPLIT} \quad \frac{S \models^* x \doteq w \cdot u \cdot z \quad S \models^* x \doteq w \cdot v \cdot z'}{A := A, \ell_u > \ell_v; S := S, u \doteq v \cdot k \\ A := A, \ell_u < \ell_v; S := S, v \doteq u \cdot k \\ A := A, \ell_u \doteq \ell_v; S := S, u \doteq v}$$

Note: k is a fresh variable

Concatenation Rules

If x is a variable of S , we can recursively expand x by substituting using equalities from S whose right-hand sides are concatenation terms until this is no longer possible

If t is the result, we write $S \models^* x = t$

We write \mathbf{z} as a short-hand for a concatenation of zero or more variables ($\mathbf{z} = z_1 \cdot z_2 \cdot \dots \cdot z_n$, with $\mathbf{z} = \epsilon$ when $n = 0$)

$$\text{C-EQ} \quad \frac{S \models^* x \doteq \mathbf{z} \quad S \models^* y \doteq \mathbf{z}}{S := S, x \doteq y}$$

$$\text{C-SPLIT} \quad \frac{S \models^* x \doteq w \cdot u \cdot \mathbf{z} \quad S \models^* x \doteq w \cdot v \cdot \mathbf{z}'}{A := A, l_u > l_v; S := S, u \doteq v \cdot k \\ A := A, l_u < l_v; S := S, v \doteq u \cdot k \\ A := A, l_u \doteq l_v; S := S, u \doteq v}$$

Note: k is a fresh variable

Concatenation Rules

If x is a variable of S , we can recursively expand x by substituting using equalities from S whose right-hand sides are concatenation terms until this is no longer possible

If t is the result, we write $S \models^* x = t$

We write \mathbf{z} as a short-hand for a concatenation of zero or more variables ($\mathbf{z} = z_1 \cdot z_2 \cdot \dots \cdot z_n$, with $\mathbf{z} = \epsilon$ when $n = 0$)

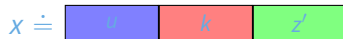
$$\text{C-EQ} \frac{S \models^* x \doteq \mathbf{z} \quad S \models^* y \doteq \mathbf{z}}{S := S, x \doteq y}$$

$$\text{C-SPLIT} \frac{S \models^* x \doteq \mathbf{w} \cdot u \cdot \mathbf{z} \quad S \models^* x \doteq \mathbf{w} \cdot v \cdot \mathbf{z}'}{\begin{array}{l} A := A, \ell_u > \ell_v; S := S, u \doteq v \cdot k \\ A := A, \ell_u < \ell_v; S := S, v \doteq u \cdot k \\ A := A, \ell_u \doteq \ell_v; S := S, u \doteq v \end{array}}$$

Note: k is a fresh variable

Example of C-Split

$$\begin{array}{c}
 \text{C-SPLIT} \quad \frac{S \models^* x \dot{=} \mathbf{w} \cdot u \cdot \mathbf{z} \quad S \models^* x \dot{=} \mathbf{w} \cdot v \cdot \mathbf{z}'}{
 \begin{array}{l}
 A := A, \ell_u > \ell_v; S := S, u \dot{=} v \cdot k \\
 A := A, \ell_u < \ell_v; S := S, v \dot{=} u \cdot k \\
 A := A, \ell_u \dot{=} \ell_v; S := S, u \dot{=} v
 \end{array}
 }
 \end{array}$$



Core and Concat Rules

$$\text{A-CONF} \frac{A \models_{LIA} \perp}{\text{UNSAT}}$$

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

$$\text{S-CONF} \frac{S \models \perp}{\text{UNSAT}}$$

$$\text{S-PROP} \frac{S \models s \doteq t \quad s, t \in \mathcal{T}(S) \quad s, t \text{ are } \Sigma_{LIA}\text{-terms}}{A := A, s \doteq t}$$

$$\text{CONST-CONF} \frac{S \models c \doteq d \quad c \in \mathcal{A} \quad d \in \mathcal{A} \setminus \{c\}}{\text{UNSAT}}$$

$$\text{S-A} \frac{x, y \in \mathcal{T}(S) \cap \mathcal{T}(A) \quad x, y : \text{Int}}{A := A, x \doteq y \quad A := A, x \doteq y}$$

$$\text{L-INTRO} \frac{s \in \mathcal{T}(S) \quad s : \text{String}}{S := S, |s| \doteq |s| \downarrow}$$

$$\text{L-VALID} \frac{x \in \mathcal{T}(S) \quad x : \text{String}}{S := S, x \doteq \epsilon \quad A \doteq A, \ell_x > 0}$$

$$\text{SAT} \frac{\text{no other rule applies}}{\text{SAT}}$$

$$\text{C-EQ} \frac{S \models^* x \doteq \mathbf{z} \quad S \models^* y \doteq \mathbf{z}}{S := S, x \doteq y}$$

$$\text{C-SPLIT} \frac{S \models^* x \doteq \mathbf{w} \cdot \mathbf{u} \cdot \mathbf{z} \quad S \models^* x \doteq \mathbf{w} \cdot \mathbf{v} \cdot \mathbf{z}'}{A := A, \ell_u > \ell_v; S := S, u \doteq v \cdot k \quad A := A, \ell_u < \ell_v; S := S, v \doteq u \cdot k \quad A := A, \ell_u \doteq \ell_v; S := S, u \doteq v}$$

Properties of the proof system

Is the proof system sound? terminating?

Properties of the proof system

The proof system is

- **refutation sound**
 - easily checkable by examining each proof rule
- **solution sound**
 - proving this is highly non-trivial
- **not terminating**
 - for pathological unsat cases, C-SPLIT can be applied infinitely often
- **incomplete**
 - a consequence of non-termination

Properties of the proof system

The proof system is

- **refutation sound**
 - easily checkable by examining each proof rule
- **solution sound**
 - proving this is highly non-trivial
- **not terminating**
 - for pathological unsat cases, C-SPLIT can be applied infinitely often
- **incomplete**
 - a consequence of non-termination

Properties of the proof system

The proof system is

- **refutation sound**
 - easily checkable by examining each proof rule
- **solution sound**
 - proving this is highly non-trivial
- **not terminating**
 - for pathological unsat cases, C-SPLIT can be applied infinitely often
- **incomplete**
 - a consequence of non-termination

Properties of the proof system

The proof system is

- **refutation sound**
 - easily checkable by examining each proof rule
- **solution sound**
 - proving this is highly non-trivial
- **not terminating**
 - for pathological unsat cases, **C-SPLIT** can be applied infinitely often
- **incomplete**
 - a consequence of non-termination

Properties of the proof system

The proof system is

- **refutation sound**
 - easily checkable by examining each proof rule
- **solution sound**
 - proving this is highly non-trivial
- **not terminating**
 - for pathological unsat cases, **C-SPLIT** can be applied infinitely often
- **incomplete**
 - a consequence of non-termination

Iterating to Improve the Solver: More String Operators

1. Extend the theory by adding *new operators*

- $\text{substr}(x, n, m) : \text{String}$, the maximal substring of x , starting at position n , with length $\leq m$
- $\text{contains}(x, y) : \text{Bool}$, true iff x contains y as a substring
- $\text{index_of}(x, y, n) : \text{Int}$, position of the first occurrence of y in x , starting from position n
- $\text{replace}(x, y, z) : \text{String}$, the result of replacing the first occurrence of x in y by z

2. Implement them by *reduction* to the core theory

Iterating to Improve the Solver: More String Operators

1. Extend the theory by adding *new operators*

- $\text{substr}(x, n, m) : \text{String}$, the maximal substring of x , starting at position n , with length $\leq m$
- $\text{contains}(x, y) : \text{Bool}$, true iff x contains y as a substring
- $\text{index_of}(x, y, n) : \text{Int}$, position of the first occurrence of y in x , starting from position n
- $\text{replace}(x, y, z) : \text{String}$, the result of replacing the first occurrence of x in y by z

2. Implement them by reduction to the core theory

Iterating to Improve the Solver: More String Operators

1. Extend the theory by adding *new operators*

- $\text{substr}(x, n, m) : \text{String}$, the maximal substring of x , starting at position n , with length $\leq m$
- $\text{contains}(x, y) : \text{Bool}$, true iff x contains y as a substring
- $\text{index_of}(x, y, n) : \text{Int}$, position of the first occurrence of y in x , starting from position n
- $\text{replace}(x, y, z) : \text{String}$, the result of replacing the first occurrence of x in y by z

2. Implement them by **reduction** to the core theory

New Operators as Macros

$$x \dot{=} y \equiv \max(x - y, 0)$$

$$x \dot{=} \text{substr}(y, n, m) \equiv \text{ite}(0 \leq n < |y| \wedge 0 \leq m, \\ y \dot{=} z_1 \cdot x \cdot z_2 \wedge |z_1| \dot{=} n \wedge |z_2| \dot{=} |y| - (m + n), \\ x \dot{=} \epsilon)$$

$$\text{contains}(y, z) \equiv \exists k. 0 \leq k \leq |y| - |z| \wedge \text{substr}(y, k, |z|) \dot{=} z$$

$$x \dot{=} \text{index_of}(y, z, n) \equiv \text{ite}(0 \leq n \leq |y| \wedge \text{contains}(y', z), \\ \text{substr}(y', x', |z|) \dot{=} z \wedge \neg \text{contains}(\text{substr}(y', 0, x' + |z| - 1), z), \\ x \dot{=} -1) \\ \text{with } y' \dot{=} \text{substr}(y, n, |y| - n) \text{ and } x' \dot{=} x - n$$

$$x \dot{=} \text{replace}(y, z, w) \equiv \text{ite}(\text{contains}(y, z) \wedge z \dot{=} \epsilon, \\ x \dot{=} z_1 \cdot w \cdot z_2 \wedge y \dot{=} z_1 \cdot z \cdot z_2 \wedge \text{index_of}(y, z, 0) \dot{=} |z_1|, \\ x \dot{=} y)$$

New Operators as Macros

$$x \dot{=} y \equiv \max(x - y, 0)$$

$$x \dot{=} \text{substr}(y, n, m) \equiv \text{ite}(0 \leq n < |y| \wedge 0 < m, \\ y \dot{=} z_1 \cdot x \cdot z_2 \wedge |z_1| \dot{=} n \wedge |z_2| \dot{=} |y| \dot{-} (m + n), \\ x \dot{=} \epsilon)$$

$$\text{contains}(y, z) \equiv \exists k. 0 \leq k \leq |y| - |z| \wedge \text{substr}(y, k, |z|) \dot{=} z$$

$$x \dot{=} \text{index_of}(y, z, n) \equiv \text{ite}(0 \leq n \leq |y| \wedge \text{contains}(y', z), \\ \text{substr}(y', x', |z|) \dot{=} z \wedge \neg \text{contains}(\text{substr}(y', 0, x' + |z| - 1), z), \\ x \dot{=} -1) \\ \text{with } y' \dot{=} \text{substr}(y, n, |y| - n) \text{ and } x' \dot{=} x - n$$

$$x \dot{=} \text{replace}(y, z, w) \equiv \text{ite}(\text{contains}(y, z) \wedge z \dot{=} \epsilon, \\ x \dot{=} z_1 \cdot w \cdot z_2 \wedge y \dot{=} z_1 \cdot z \cdot z_2 \wedge \text{index_of}(y, z, 0) \dot{=} |z_1|, \\ x \dot{=} y)$$

New Operators as Macros

$$x \dot{=} y \equiv \max(x - y, 0)$$

$$x \dot{=} \text{substr}(y, n, m) \equiv \text{ite}(0 \leq n < |y| \wedge 0 < m, \\ y \dot{=} z_1 \cdot x \cdot z_2 \wedge |z_1| \dot{=} n \wedge |z_2| \dot{=} |y| \dot{-} (m + n), \\ x \dot{=} \epsilon)$$

$$\text{contains}(y, z) \equiv \exists k. 0 \leq k \leq |y| - |z| \wedge \text{substr}(y, k, |z|) \dot{=} z$$

$$x \dot{=} \text{index_of}(y, z, n) \equiv \text{ite}(0 \leq n \leq |y| \wedge \text{contains}(y', z), \\ \text{substr}(y', x', |z|) \dot{=} z \wedge \neg \text{contains}(\text{substr}(y', 0, x' + |z| - 1), z), \\ x \dot{=} -1) \\ \text{with } y' \dot{=} \text{substr}(y, n, |y| - n) \text{ and } x' \dot{=} x - n$$

$$x \dot{=} \text{replace}(y, z, w) \equiv \text{ite}(\text{contains}(y, z) \wedge z \dot{=} \epsilon, \\ x \dot{=} z_1 \cdot w \cdot z_2 \wedge y \dot{=} z_1 \cdot z \cdot z_2 \wedge \text{index_of}(y, z, 0) \dot{=} |z_1|, \\ x \dot{=} y)$$

New Operators as Macros

$$x \dot{=} y \equiv \max(x - y, 0)$$

$$x \dot{=} \text{substr}(y, n, m) \equiv \text{ite}(0 \leq n < |y| \wedge 0 < m, \\ y \dot{=} z_1 \cdot x \cdot z_2 \wedge |z_1| \dot{=} n \wedge |z_2| \dot{=} |y| \dot{-} (m + n), \\ x \dot{=} \epsilon)$$

$$\text{contains}(y, z) \equiv \exists k. 0 \leq k \leq |y| - |z| \wedge \text{substr}(y, k, |z|) \dot{=} z$$

$$x \dot{=} \text{index_of}(y, z, n) \equiv \text{ite}(0 \leq n \leq |y| \wedge \text{contains}(y', z), \\ \text{substr}(y', x', |z|) \dot{=} z \wedge \neg \text{contains}(\text{substr}(y', 0, x' + |z| - 1), z), \\ x \dot{=} -1) \\ \text{with } y' \dot{=} \text{substr}(y, n, |y| - n) \text{ and } x' \dot{=} x - n$$

$$x \dot{=} \text{replace}(y, z, w) \equiv \text{ite}(\text{contains}(y, z) \wedge z \dot{=} \epsilon, \\ x \dot{=} z_1 \cdot w \cdot z_2 \wedge y \dot{=} z_1 \cdot z \cdot z_2 \wedge \text{index_of}(y, z, 0) \dot{=} |z_1|, \\ x \dot{=} y)$$

New Operators as Macros

$$x \dot{=} y \equiv \max(x - y, 0)$$

$$x \dot{=} \text{substr}(y, n, m) \equiv \text{ite}(0 \leq n < |y| \wedge 0 < m, \\ y \dot{=} z_1 \cdot x \cdot z_2 \wedge |z_1| \dot{=} n \wedge |z_2| \dot{=} |y| \dot{-} (m + n), \\ x \dot{=} \epsilon)$$

$$\text{contains}(y, z) \equiv \exists k. 0 \leq k \leq |y| - |z| \wedge \text{substr}(y, k, |z|) \dot{=} z$$

$$x \dot{=} \text{index_of}(y, z, n) \equiv \text{ite}(0 \leq n \leq |y| \wedge \text{contains}(y', z), \\ \text{substr}(y', x', |z|) \dot{=} z \wedge \neg \text{contains}(\text{substr}(y', 0, x' + |z| - 1), z), \\ x \dot{=} -1) \\ \text{with } y' \dot{=} \text{substr}(y, n, |y| - n) \text{ and } x' \dot{=} x - n$$

$$x \dot{=} \text{replace}(y, z, w) \equiv \text{ite}(\text{contains}(y, z) \wedge z \dot{=} \epsilon, \\ x \dot{=} z_1 \cdot w \cdot z_2 \wedge y \dot{=} z_1 \cdot z \cdot z_2 \wedge \text{index_of}(y, z, 0) \dot{=} |z_1|, \\ x \dot{=} y)$$

Reasoning about New Operators: Performance

Iterate and Improve

- Extend the implementation to **reason directly** on the new operators
- How?
 - Keep formulas with **original new operators**
 - Periodically try to **simplify** them based on new knowledge

Reasoning about New Operators: Performance

Iterate and Improve

- Extend the implementation to **reason directly** on the new operators
- How?
 - Keep formulas with **original new operators**
 - Periodically try to **simplify** them based on new knowledge

Reasoning about New Operators: Performance

Iterate and Improve

- Extend the implementation to **reason directly** on the new operators
- How?
 - Keep formulas with **original new operators**
 - Periodically try to **simplify** them based on new knowledge

Reasoning about New Operators: Performance

Iterate and Improve

- Extend the implementation to **reason directly** on the new operators
- How?
 - Keep formulas with **original new operators**
 - Periodically try to **simplify** them based on new knowledge

Simplification rules for New Operators

Example: `contains`

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot t) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot t) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, t) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot t) \longrightarrow \perp$	if $\text{contains}(l_1, t) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot t, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot t, s) \longrightarrow \top$	if $\text{contains}(t, s) \longrightarrow^* \top$
$\text{contains}(t \cdot s, t \cdot u) \longrightarrow \top$	if $\text{contains}(s, u) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot t, l_2) \longrightarrow \text{contains}(t, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(t \cdot l_1, l_2) \longrightarrow \text{contains}(t, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, t) = \top \longrightarrow \epsilon = t$	
...	

Simplification rules for New Operators

Example: contains

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot t) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot t) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, t) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot t) \longrightarrow \perp$	if $\text{contains}(l_1, t) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot t, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot t, s) \longrightarrow \top$	if $\text{contains}(t, s) \longrightarrow^* \top$
$\text{contains}(t \cdot s, t \cdot u) \longrightarrow \top$	if $\text{contains}(s, u) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot t, l_2) \longrightarrow \text{contains}(t, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(t \cdot l_1, l_2) \longrightarrow \text{contains}(t, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, t) = \top \longrightarrow \epsilon = t$	
...	

Simplification rules for New Operators

Example: contains

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot t) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot t) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, t) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot t) \longrightarrow \perp$	if $\text{contains}(l_1, t) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot t, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot t, s) \longrightarrow \top$	if $\text{contains}(t, s) \longrightarrow^* \top$
$\text{contains}(t \cdot s, t \cdot u) \longrightarrow \top$	if $\text{contains}(s, u) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot t, l_2) \longrightarrow \text{contains}(t, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(t \cdot l_1, l_2) \longrightarrow \text{contains}(t, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(c, t) = \top \longrightarrow \epsilon = t$	
...	

Simplification rules for New Operators

Example: contains

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot t) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot t) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, t) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot t) \longrightarrow \perp$	if $\text{contains}(l_1, t) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot t, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot t, s) \longrightarrow \top$	if $\text{contains}(t, s) \longrightarrow^* \top$
$\text{contains}(t \cdot s, t \cdot u) \longrightarrow \top$	if $\text{contains}(s, u) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot t, l_2) \longrightarrow \text{contains}(t, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(t \cdot l_1, l_2) \longrightarrow \text{contains}(t, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, t) = \top \longrightarrow \epsilon = t$	
...	

Simplification rules for New Operators

Example: **contains**

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot \mathbf{t}, s) \longrightarrow \top$	if $\text{contains}(\mathbf{t}, s) \longrightarrow^* \top$
$\text{contains}(\mathbf{t} \cdot s, \mathbf{t} \cdot u) \longrightarrow \top$	if $\text{contains}(s, u) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(\mathbf{t} \cdot l_1, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, \mathbf{t}) = \top \longrightarrow \epsilon = \mathbf{t}$	
...	

Simplification rules for New Operators

Example: contains

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot \mathbf{t}, s) \longrightarrow \top$	if $\text{contains}(\mathbf{t}, s) \longrightarrow^* \top$
$\text{contains}(\mathbf{t} \cdot s, \mathbf{t} \cdot u) \longrightarrow \top$	if $\text{contains}(s, u) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(\mathbf{t} \cdot l_1, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, \mathbf{t}) = \top \longrightarrow \epsilon = \mathbf{t}$	
...	

Simplification rules for New Operators

Example: **contains**

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot \mathbf{t}, s) \longrightarrow \top$	if $\text{contains}(\mathbf{t}, s) \longrightarrow^* \top$
$\text{contains}(\mathbf{t} \cdot s, \mathbf{t} \cdot u) \longrightarrow \top$	if $\text{contains}(s, u) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(\mathbf{t} \cdot l_1, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, \mathbf{t}) = \top \longrightarrow \epsilon = \mathbf{t}$	
...	

Simplification rules for New Operators

Example: **contains**

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot \mathbf{t}, s) \longrightarrow \top$	if $\text{contains}(\mathbf{t}, s) \longrightarrow^* \top$
$\text{contains}(\mathbf{t} \cdot s, \mathbf{t} \cdot u) \longrightarrow \top$	if $\text{contains}(s, u) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(\mathbf{t} \cdot l_1, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, \mathbf{t}) \longrightarrow \top$	$\epsilon = \mathbf{t}$
...	...

Simplification rules for New Operators

Example: **contains**

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot \mathbf{t}, s) \longrightarrow \top$	if $\text{contains}(\mathbf{t}, s) \longrightarrow^* \top$
$\text{contains}(t \cdot \mathbf{s}, t \cdot \mathbf{u}) \longrightarrow \top$	if $\text{contains}(\mathbf{s}, \mathbf{u}) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(\mathbf{t} \cdot l_1, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, t) = \top \longrightarrow \epsilon = t$	
...	

Simplification rules for New Operators

Example: **contains**

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot \mathbf{t}, s) \longrightarrow \top$	if $\text{contains}(\mathbf{t}, s) \longrightarrow^* \top$
$\text{contains}(t \cdot \mathbf{s}, t \cdot \mathbf{u}) \longrightarrow \top$	if $\text{contains}(\mathbf{s}, \mathbf{u}) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(t \cdot l_1, l_2) \longrightarrow \text{contains}(t, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, t) = \top \longrightarrow \epsilon = t$	
...	

Simplification rules for New Operators

Example: **contains**

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot \mathbf{t}, s) \longrightarrow \top$	if $\text{contains}(\mathbf{t}, s) \longrightarrow^* \top$
$\text{contains}(t \cdot \mathbf{s}, t \cdot \mathbf{u}) \longrightarrow \top$	if $\text{contains}(\mathbf{s}, \mathbf{u}) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(\mathbf{t} \cdot l_1, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$

$\text{contains}(\epsilon, t) = \top \longrightarrow \epsilon = t$

...

Simplification rules for New Operators

Example: **contains**

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot \mathbf{t}, s) \longrightarrow \top$	if $\text{contains}(\mathbf{t}, s) \longrightarrow^* \top$
$\text{contains}(t \cdot \mathbf{s}, t \cdot \mathbf{u}) \longrightarrow \top$	if $\text{contains}(\mathbf{s}, \mathbf{u}) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(\mathbf{t} \cdot l_1, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, t) = \top \longrightarrow \epsilon = t$	

Simplification rules for New Operators

Example: **contains**

(l_1, l_2 denote string constants)

$\text{contains}(l_1, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(l_1, l_2) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if l_1 does not contain l_2
$\text{contains}(l_1, l_2 \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1 \setminus l_2, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1, x \cdot \mathbf{t}) \longrightarrow \perp$	if $\text{contains}(l_1, \mathbf{t}) \longrightarrow^* \perp$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \top$	if l_1 contains l_2
$\text{contains}(x \cdot \mathbf{t}, s) \longrightarrow \top$	if $\text{contains}(\mathbf{t}, s) \longrightarrow^* \top$
$\text{contains}(t \cdot \mathbf{s}, t \cdot \mathbf{u}) \longrightarrow \top$	if $\text{contains}(\mathbf{s}, \mathbf{u}) \longrightarrow^* \top$
$\text{contains}(l_1 \cdot \mathbf{t}, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_l l_2 = \epsilon$
$\text{contains}(\mathbf{t} \cdot l_1, l_2) \longrightarrow \text{contains}(\mathbf{t}, l_2)$	if $l_1 \sqcup_r l_2 = \epsilon$
$\text{contains}(\epsilon, t) = \top \longrightarrow \epsilon = t$	

...

Reasoning about New Operators: More Performance

Iterate and Improve

- Supercharge the simplifier
- Many simplifications are conditional
- Build a mini-inference engine inside the simplifier to verify *simplification conditions*

Conditional Simplifications based on String Length

Notation: $\vdash C$ states that simplifier can prove simplification condition C

$$t \doteq s \longrightarrow \perp$$

$$\text{if } \vdash |t| > |s|$$

$$t \doteq s \cdot r \cdot q \longrightarrow t \doteq s \cdot q \wedge r \doteq \epsilon$$

$$\text{if } \vdash |s| + |q| \geq |t|$$

$$\text{contains}(t, s) \longrightarrow t \doteq s$$

$$\text{if } \vdash |s| \geq |t|$$

$$\text{substr}(t, v, w) \longrightarrow \epsilon$$

$$\text{if } \vdash 0 > v \vee v \geq |t| \vee 0 \geq w$$

$$\text{substr}(t \cdot s, v, w) \longrightarrow \text{substr}(s, v - |t|, w)$$

$$\text{if } \vdash v \geq |t|$$

$$\text{substr}(s \cdot t, v, w) \longrightarrow \text{substr}(s, v, w)$$

$$\text{if } \vdash |s| \geq v + w$$

$$\text{substr}(t \cdot s, 0, w) \longrightarrow t \cdot \text{substr}(s, 0, w - |t|)$$

$$\text{if } \vdash w \geq |t|$$

$$\text{index_of}(t, s, v) \longrightarrow \text{ite}(\text{substr}(t, v) \doteq s, v, -1) \quad \text{if } \vdash v + |s| \geq |t|$$

...

Conditional Simplifications based on String Length

Notation: $\vdash C$ states that simplifier can prove simplification condition C

$t \doteq s \longrightarrow \perp$	if $\vdash t > s $
$t \doteq s \cdot r \cdot q \longrightarrow t \doteq s \cdot q \wedge r \doteq \epsilon$	if $\vdash s + q \geq t $
$\text{contains}(t, s) \longrightarrow t \doteq s$	if $\vdash s \geq t $
$\text{substr}(t, v, w) \longrightarrow \epsilon$	if $\vdash 0 > v \vee v \geq t \vee 0 \geq w$
$\text{substr}(t \cdot s, v, w) \longrightarrow \text{substr}(s, v - t , w)$	if $\vdash v \geq t $
$\text{substr}(s \cdot t, v, w) \longrightarrow \text{substr}(s, v, w)$	if $\vdash s \geq v + w$
$\text{substr}(t \cdot s, 0, w) \longrightarrow t \cdot \text{substr}(s, 0, w - t)$	if $\vdash w \geq t $
$\text{index_of}(t, s, v) \longrightarrow \text{ite}(\text{substr}(t, v) \doteq s, v, -1)$	if $\vdash v + s \geq t $
...	

Conditional Simplifications based on String Length

Notation: $\vdash C$ states that simplifier can prove simplification condition C

$t \doteq s \longrightarrow \perp$	if $\vdash t > s $
$t \doteq s \cdot r \cdot q \longrightarrow t \doteq s \cdot q \wedge r \doteq \epsilon$	if $\vdash s + q \geq t $
$\text{contains}(t, s) \longrightarrow t \doteq s$	if $\vdash s \geq t $
$\text{substr}(t, v, w) \longrightarrow \epsilon$	if $\vdash 0 > v \vee v \geq t \vee 0 \geq w$
$\text{substr}(t \cdot s, v, w) \longrightarrow \text{substr}(s, v - t , w)$	if $\vdash v \geq t $
$\text{substr}(s \cdot t, v, w) \longrightarrow \text{substr}(s, v, w)$	if $\vdash s \geq v + w$
$\text{substr}(t \cdot s, 0, w) \longrightarrow t \cdot \text{substr}(s, 0, w - t)$	if $\vdash w \geq t $
$\text{index_of}(t, s, v) \longrightarrow \text{ite}(\text{substr}(t, v) \doteq s, v, -1)$	if $\vdash v + s \geq t $
...	

Conditional Simplifications based on String Length

Notation: $\vdash C$ states that simplifier can prove simplification condition C

$$t \doteq s \longrightarrow \perp$$

$$\text{if } \vdash |t| > |s|$$

$$t \doteq s \cdot r \cdot q \longrightarrow t \doteq s \cdot q \wedge r \doteq \epsilon$$

$$\text{if } \vdash |s| + |q| \geq |t|$$

$$\text{contains}(t, s) \longrightarrow t \doteq s$$

$$\text{if } \vdash |s| \geq |t|$$

$$\text{substr}(t, v, w) \longrightarrow \epsilon$$

$$\text{if } \vdash 0 > v \vee v \geq |t| \vee 0 \geq w$$

$$\text{substr}(t \cdot s, v, w) \longrightarrow \text{substr}(s, v - |t|, w)$$

$$\text{if } \vdash v \geq |t|$$

$$\text{substr}(s \cdot t, v, w) \longrightarrow \text{substr}(s, v, w)$$

$$\text{if } \vdash |s| \geq v + w$$

$$\text{substr}(t \cdot s, 0, w) \longrightarrow t \cdot \text{substr}(s, 0, w - |t|)$$

$$\text{if } \vdash w \geq |t|$$

$$\text{index_of}(t, s, v) \longrightarrow \text{ite}(\text{substr}(t, v) \doteq s, v, -1)$$

$$\text{if } \vdash v + |s| \geq |t|$$

...

Conditional Simplifications based on String Length

Notation: $\vdash C$ states that simplifier can prove simplification condition C

$t \doteq s \longrightarrow \perp$	if $\vdash t > s $
$t \doteq s \cdot r \cdot q \longrightarrow t \doteq s \cdot q \wedge r \doteq \epsilon$	if $\vdash s + q \geq t $
$\text{contains}(t, s) \longrightarrow t \doteq s$	if $\vdash s \geq t $
$\text{substr}(t, v, w) \longrightarrow \epsilon$	if $\vdash 0 > v \vee v \geq t \vee 0 \geq w$
$\text{substr}(t \cdot s, v, w) \longrightarrow \text{substr}(s, v - t , w)$	if $\vdash v \geq t $
$\text{substr}(s \cdot t, v, w) \longrightarrow \text{substr}(s, v, w)$	if $\vdash s \geq v + w$
$\text{substr}(t \cdot s, 0, w) \longrightarrow t \cdot \text{substr}(s, 0, w - t)$	if $\vdash w \geq t $
$\text{index_of}(t, s, v) \longrightarrow \text{ite}(\text{substr}(t, v) \doteq s, v, -1)$	if $\vdash v + s \geq t $

Conditional Simplifications based on String Length

Notation: $\vdash C$ states that simplifier can prove simplification condition C

$t \doteq s \longrightarrow \perp$	if $\vdash t > s $
$t \doteq s \cdot r \cdot q \longrightarrow t \doteq s \cdot q \wedge r \doteq \epsilon$	if $\vdash s + q \geq t $
$\text{contains}(t, s) \longrightarrow t \doteq s$	if $\vdash s \geq t $
$\text{substr}(t, v, w) \longrightarrow \epsilon$	if $\vdash 0 > v \vee v \geq t \vee 0 \geq w$
$\text{substr}(t \cdot s, v, w) \longrightarrow \text{substr}(s, v - t , w)$	if $\vdash v \geq t $
$\text{substr}(s \cdot t, v, w) \longrightarrow \text{substr}(s, v, w)$	if $\vdash s \geq v + w$
$\text{substr}(t \cdot s, 0, w) \longrightarrow t \cdot \text{substr}(s, 0, w - t)$	if $\vdash w \geq t $
$\text{index_of}(t, s, v) \longrightarrow \text{ite}(\text{substr}(t, v) \doteq s, v, -1)$	if $\vdash v + s \geq t $

...

Conditional Simplifications based on String Length

Notation: $\vdash C$ states that simplifier can prove simplification condition C

$t \doteq s \longrightarrow \perp$	if $\vdash t > s $
$t \doteq s \cdot r \cdot q \longrightarrow t \doteq s \cdot q \wedge r \doteq \epsilon$	if $\vdash s + q \geq t $
$\text{contains}(t, s) \longrightarrow t \doteq s$	if $\vdash s \geq t $
$\text{substr}(t, v, w) \longrightarrow \epsilon$	if $\vdash 0 > v \vee v \geq t \vee 0 \geq w$
$\text{substr}(t \cdot s, v, w) \longrightarrow \text{substr}(s, v - t , w)$	if $\vdash v \geq t $
$\text{substr}(s \cdot t, v, w) \longrightarrow \text{substr}(s, v, w)$	if $\vdash s \geq v + w$
$\text{substr}(t \cdot s, 0, w) \longrightarrow t \cdot \text{substr}(s, 0, w - t)$	if $\vdash w \geq t $
$\text{index_of}(t, s, v) \longrightarrow \text{ite}(\text{substr}(t, v) \doteq s, v, -1)$	if $\vdash v + s \geq t $

Conditional Simplifications based on String Length

Notation: $\vdash C$ states that simplifier can prove simplification condition C

$t \doteq s \longrightarrow \perp$	if $\vdash t > s $
$t \doteq s \cdot r \cdot q \longrightarrow t \doteq s \cdot q \wedge r \doteq \epsilon$	if $\vdash s + q \geq t $
$\text{contains}(t, s) \longrightarrow t \doteq s$	if $\vdash s \geq t $
$\text{substr}(t, v, w) \longrightarrow \epsilon$	if $\vdash 0 > v \vee v \geq t \vee 0 \geq w$
$\text{substr}(t \cdot s, v, w) \longrightarrow \text{substr}(s, v - t , w)$	if $\vdash v \geq t $
$\text{substr}(s \cdot t, v, w) \longrightarrow \text{substr}(s, v, w)$	if $\vdash s \geq v + w$
$\text{substr}(t \cdot s, 0, w) \longrightarrow t \cdot \text{substr}(s, 0, w - t)$	if $\vdash w \geq t $
$\text{index_of}(t, s, v) \longrightarrow \text{ite}(\text{substr}(t, v) \doteq s, v, -1)$	if $\vdash v + s \geq t $

Conditional Simplifications based on String Length

Notation: $\vdash C$ states that simplifier can prove simplification condition C

$t \doteq s \longrightarrow \perp$	if $\vdash t > s $
$t \doteq s \cdot r \cdot q \longrightarrow t \doteq s \cdot q \wedge r \doteq \epsilon$	if $\vdash s + q \geq t $
$\text{contains}(t, s) \longrightarrow t \doteq s$	if $\vdash s \geq t $
$\text{substr}(t, v, w) \longrightarrow \epsilon$	if $\vdash 0 > v \vee v \geq t \vee 0 \geq w$
$\text{substr}(t \cdot s, v, w) \longrightarrow \text{substr}(s, v - t , w)$	if $\vdash v \geq t $
$\text{substr}(s \cdot t, v, w) \longrightarrow \text{substr}(s, v, w)$	if $\vdash s \geq v + w$
$\text{substr}(t \cdot s, 0, w) \longrightarrow t \cdot \text{substr}(s, 0, w - t)$	if $\vdash w \geq t $
$\text{index_of}(t, s, v) \longrightarrow \text{ite}(\text{substr}(t, v) \doteq s, v, -1)$	if $\vdash v + s \geq t $
...	

Reasoning about Regular Expressions

Regular Expression Membership Example

$$\begin{aligned} x &\in [0..9]^* \cdot \text{"a"} \cdot \mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^* \wedge \\ x &\notin [0..9]^* \cdot \text{"a"} \cdot \mathcal{A}^* \end{aligned}$$

Automata-based approach

$$\frac{x \in R_1 \quad x \notin R_2}{x \in R_1 \cap \text{comp}(R_2)}$$

Problem:

Reasoning about Regular Expressions

Regular Expression Membership Example

$$\begin{aligned} x &\in [0..9]^* \cdot \text{"a"} \cdot \mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^* \wedge \\ x &\notin [0..9]^* \cdot \text{"a"} \cdot \mathcal{A}^* \end{aligned}$$

Automata-based approach

$$\frac{x \in R_1 \quad x \notin R_2}{x \in R_1 \cap \text{comp}(R_2)}$$

Problem: Complement and intersection are **expensive**

Reasoning about Regular Expressions

Regular Expression Membership Example

$$\begin{aligned} x &\in [0..9]^* \cdot \text{"a"} \cdot \mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^* \wedge \\ x &\notin [0..9]^* \cdot \text{"a"} \cdot \mathcal{A}^* \end{aligned}$$

Automata-based approach

$$\frac{x \in R_1 \quad x \notin R_2}{x \in R_1 \cap \text{comp}(R_2)}$$

Problem: Membership constraints may lead to non-terminating unfolding:

Reasoning about Regular Expressions

Regular Expression Membership Example

$$\begin{aligned} x &\in [0..9]^* \cdot \text{"a"} \cdot \mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^* \wedge \\ x &\notin [0..9]^* \cdot \text{"a"} \cdot \mathcal{A}^* \end{aligned}$$

Automata-based approach

$$\frac{x \in R_1 \quad x \notin R_2}{x \in R_1 \cap \text{comp}(R_2)}$$

Problem: Membership constraints may lead to non-terminating unfolding:

Example: $x_1 \in [0..9]^*$ is equivalent to

$$x \doteq \epsilon \vee x \in [0..9] \vee (\exists u, v, w. x \doteq u \cdot v \cdot w \wedge u \in [0..9] \wedge v \in [0..9]^* \wedge w \in [0..9])$$

Reasoning about Regular Expressions

Regular Expression Membership Example

$$\begin{aligned} x &\in [0..9]^* \cdot \text{"a"} \cdot \mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^* \wedge \\ x &\notin [0..9]^* \cdot \text{"a"} \cdot \mathcal{A}^* \end{aligned}$$

Word-based approach with incomplete procedures

$$\frac{x \in R_1 \quad x \notin R_2 \quad \mathcal{L}(R_1) \subseteq \mathcal{L}(R_2)}{\text{UNSAT}}$$

Use fast, incomplete procedure to verify $\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2)$

Notation: $\mathcal{L}(R)$ denotes the language generated by regex R

Proving $\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2)$

$$(1) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(R)} \quad (2) \frac{}{\mathcal{L}(\epsilon) \subseteq \mathcal{L}(R^*)} \quad (3) \frac{\text{for all } x \in \mathcal{L}(R), |x| = 1}{\mathcal{L}(R) \subseteq \mathcal{L}(\mathcal{A})}$$

$$(4) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(\mathcal{A}^*)} \quad (5) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(R^*)} \quad (6) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2)}{\mathcal{L}(R_1^*) \subseteq \mathcal{L}(R_2^*)}$$

$$(7) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2) \quad \mathcal{L}(R_2) \subseteq \mathcal{L}(R_3)}{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_3)}$$

$$(8) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(S_1) \quad \mathcal{L}(R_2) \subseteq \mathcal{L}(S_2)}{\mathcal{L}(R_1 \cdot R_2) \subseteq \mathcal{L}(S_1 \cdot S_2)} \quad (9) \frac{c_3 \preceq c_1 \quad c_2 \preceq c_4}{\mathcal{L}([c_1..c_2]) \subseteq \mathcal{L}([c_3..c_4])}$$

$c \preceq d$ iff c equals d or precedes d lexicographically ($c, d \in \mathcal{A}$)

Proving $\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2)$

$$(1) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(R)} \quad (2) \frac{}{\mathcal{L}(\epsilon) \subseteq \mathcal{L}(R^*)} \quad (3) \frac{\text{for all } x \in \mathcal{L}(R), |x| = 1}{\mathcal{L}(R) \subseteq \mathcal{L}(\mathcal{A})}$$

$$(4) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(\mathcal{A}^*)} \quad (5) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(R^*)} \quad (6) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2)}{\mathcal{L}(R_1^*) \subseteq \mathcal{L}(R_2^*)}$$

$$(7) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2) \quad \mathcal{L}(R_2) \subseteq \mathcal{L}(R_3)}{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_3)}$$

$$(8) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(S_1) \quad \mathcal{L}(R_2) \subseteq \mathcal{L}(S_2)}{\mathcal{L}(R_1 \cdot R_2) \subseteq \mathcal{L}(S_1 \cdot S_2)} \quad (9) \frac{c_3 \preceq c_1 \quad c_2 \preceq c_4}{\mathcal{L}([c_1..c_2]) \subseteq \mathcal{L}([c_3..c_4])}$$

$c \preceq d$ iff c equals d or precedes d lexicographically ($c, d \in \mathcal{A}$)

Proving $\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2)$

$$(1) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(R)} \quad (2) \frac{}{\mathcal{L}(\epsilon) \subseteq \mathcal{L}(R^*)} \quad (3) \frac{\text{for all } x \in \mathcal{L}(R), |x| = 1}{\mathcal{L}(R) \subseteq \mathcal{L}(\mathcal{A})}$$

$$(4) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(\mathcal{A}^*)} \quad (5) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(R^*)} \quad (6) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2)}{\mathcal{L}(R_1^*) \subseteq \mathcal{L}(R_2^*)}$$

$$(7) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2) \quad \mathcal{L}(R_2) \subseteq \mathcal{L}(R_3)}{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_3)}$$

$$(8) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(S_1) \quad \mathcal{L}(R_2) \subseteq \mathcal{L}(S_2)}{\mathcal{L}(R_1 \cdot R_2) \subseteq \mathcal{L}(S_1 \cdot S_2)}$$

$$(9) \frac{c_3 \preceq c_1 \quad c_2 \preceq c_4}{\mathcal{L}([c_1..c_2]) \subseteq \mathcal{L}([c_3..c_4])}$$

$c \preceq d$ iff c equals d or precedes d lexicographically ($c, d \in \mathcal{A}$)

Proving $\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2)$

$$(1) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(R)} \quad (2) \frac{}{\mathcal{L}(\epsilon) \subseteq \mathcal{L}(R^*)} \quad (3) \frac{\text{for all } x \in \mathcal{L}(R), |x| = 1}{\mathcal{L}(R) \subseteq \mathcal{L}(\mathcal{A})}$$

$$(4) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(\mathcal{A}^*)} \quad (5) \frac{}{\mathcal{L}(R) \subseteq \mathcal{L}(R^*)} \quad (6) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2)}{\mathcal{L}(R_1^*) \subseteq \mathcal{L}(R_2^*)}$$

$$(7) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_2) \quad \mathcal{L}(R_2) \subseteq \mathcal{L}(R_3)}{\mathcal{L}(R_1) \subseteq \mathcal{L}(R_3)}$$

$$(8) \frac{\mathcal{L}(R_1) \subseteq \mathcal{L}(S_1) \quad \mathcal{L}(R_2) \subseteq \mathcal{L}(S_2)}{\mathcal{L}(R_1 \cdot R_2) \subseteq \mathcal{L}(S_1 \cdot S_2)} \quad (9) \frac{c_3 \preceq c_1 \quad c_2 \preceq c_4}{\mathcal{L}([c_1..c_2]) \subseteq \mathcal{L}([c_3..c_4])}$$

$c \preceq d$ iff c equals d or precedes d lexicographically ($c, d \in \mathcal{A}$)

Exercise

Using the proof rules above, prove that

$$\mathcal{L}([0..1]^* \cdot \mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^*) \subseteq \mathcal{L}([0..9]^* \cdot \mathcal{A}^*)$$

$$\begin{array}{c} \frac{0 \leq 0 \quad 1 \leq 9}{\mathcal{L}([0..1]) \subseteq \mathcal{L}([0..9])} \quad (9) \\ \frac{\mathcal{L}([0..1]) \subseteq \mathcal{L}([0..9])}{\mathcal{L}([0..1]^*) \subseteq \mathcal{L}([0..9]^*)} \quad (6) \quad \frac{}{\mathcal{L}(\mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^*) \subseteq \mathcal{L}(\mathcal{A}^*)} \quad (4) \\ \hline \frac{\mathcal{L}([0..1]^*) \subseteq \mathcal{L}([0..9]^*) \quad \mathcal{L}(\mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^*) \subseteq \mathcal{L}(\mathcal{A}^*)}{\mathcal{L}([0..1]^* \cdot \mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^*) \subseteq \mathcal{L}([0..9]^* \cdot \mathcal{A}^*)} \quad (8) \end{array}$$

Exercise

Using the proof rules above, prove that

$$\mathcal{L}([0..1]^* \cdot \mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^*) \subseteq \mathcal{L}([0..9]^* \cdot \mathcal{A}^*)$$

$$\frac{\frac{0 \leq 0 \quad 1 \leq 9}{\mathcal{L}([0..1]) \subseteq \mathcal{L}([0..9])} \quad (9) \quad \frac{\mathcal{L}([0..1]^*) \subseteq \mathcal{L}([0..9]^*)}{\mathcal{L}(\mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^*) \subseteq \mathcal{L}(\mathcal{A}^*)} \quad (4)}{\mathcal{L}([0..1]^* \cdot \mathcal{A}^* \cdot \text{"b"} \cdot \mathcal{A}^*) \subseteq \mathcal{L}([0..9]^* \cdot \mathcal{A}^*)} \quad (8)$$

More Information

Strings Papers

- “A DPLL(T) Theory Solver for a Theory of Strings and Regular Expressions” by Tianyi Liang, Andrew Reynolds, Cesare Tinelli, Clark Barrett, and Morgan Deters. In Proceedings of the 26th International Conference on Computer Aided Verification (CAV ’14), (Armin Biere and Roderick Bloem, eds.), July 2014, pp. 646-662. Vienna, Austria.
- “An Efficient SMT Solver for String Constraints” by Tianyi Liang, Andrew Reynolds, Nestan Tsiskaridze, Cesare Tinelli, Clark Barrett, and Morgan Deters. Formal Methods in System Design, vol. 48, no. 3, June 2016, pp. 206-234, Springer US.
- “Scaling up DPLL(T) String Solvers Using Context-Dependent Simplification” by Andrew Reynolds, Maverick Woo, Clark Barrett, David Brumley, Tianyi Liang, and Cesare Tinelli. In Proceedings of the 29th International Conference on Computer Aided Verification (CAV ’17), (Rupak Majumdar and Viktor Kuncak, eds.), July 2017, pp. 453-474. Heidelberg, Germany.
- “High-Level Abstractions for Simplifying Extended String Constraints in SMT” by Andrew Reynolds, Andres Nötzli, Clark Barrett, and Cesare Tinelli. In Proceedings of the 31st International Conference on Computer Aided Verification (CAV ’19), (Isil Dillig and Serdar Tasiran, eds.), July 2019, pp. 23-42. New York, New York.
- “Even Faster Conflicts and Lazier Reductions for String Solvers” by Andres Nötzli, Andrew Reynolds, Haniel Barbosa, Clark Barrett, and Cesare Tinelli. In Proceedings of the 34th International Conference on Computer Aided Verification (CAV ’22), (Sharon Shoham and Yakir Vizel, eds.), Aug. 2022, pp. 205-226. Haifa, Israel.

Amazon’s Zelkova Tool

- J. Backes et al., “Semantic-based Automated Reasoning for AWS Access Policies using SMT,” 2018 Formal Methods in Computer Aided Design (FMCAD), Austin, TX, 2018.