# Costantini's
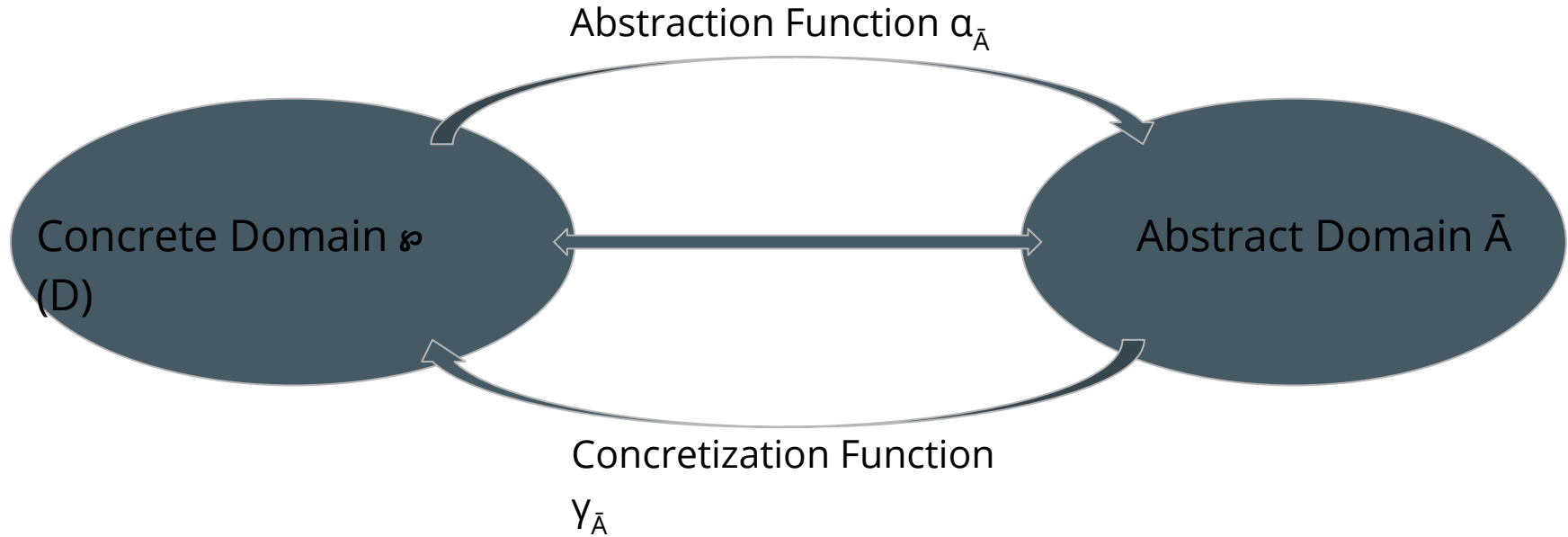# "Static Analysis of String Values"
# - A Summary

Koby Picker
Christian M. Maldonado

# Overview

- What is Abstract Interpretation?
- The Concrete Domain
- The Abstract Domains
  a. Character Inclusion
  b. Prefix and Suffix
  c. Bricks
  d. String Graphs

# Abstract Interpretation

Abstraction Function $\alpha_{\bar{A}}$

Concrete Domain ℘ (D)

Abstract Domain $\bar{A}$

Concretization Function $\gamma_{\bar{A}}$

# Concrete Domain

- Given an alphabet K, a finite set of characters ...

- Strings = Sequence of characters (potentially infinite)

$$S = K^*, \text{ where } A^* \text{ is an ordered sequence of elements in } A$$

$$A^* = \{a_1 \ldots a_n : \forall i \in [1 \ldots n] : a_i \in A\}$$

# Concrete Semantics

**Table 2.** The concrete semantics, where $\top_B$ represents that the condition could be evaluated to **true** or **false** depending on the string in $S_1$ we are considering

$$\mathbb{S}[\![\text{new String(str)}]\!]() = \{\text{str}\}$$

$$\mathbb{S}[\![\text{concat}]\!](S_1, S_2) = \{s_1 s_2 : s_1 \in S_1 \land s_2 \in S_2\}$$

$$\mathbb{S}[\![\text{readLine}]\!]() = S$$

$$\mathbb{S}[\![\text{substring}_b^e]\!](S_1) = \{c_b..c_e : c_1..c_n \in S_1 \land n \geq e \land b \leq e\}$$

$$\mathbb{B}[\![\text{contains}_c]\!](S_1) = \begin{cases} \text{true} & \text{if } \forall s \in S_1 : c \in char(s) \\ \text{false} & \text{if } \forall s \in S_1 : c \notin char(s) \\ \top_B & \text{otherwise} \end{cases}$$

# Overview

- Abstract Interpretation
- The Concrete Domain
- The Abstract Domains
  a. Character Inclusion
  b. Prefix and Suffix
  c. Bricks
  d. String Graphs

# Character Inclusion - CI

CI consists of *certainly* contained characters and *maybe contained* characters

$$CI = \{(C, MC): C, MC \in \wp(K) \land C \subseteq MC\}$$
$$\cup \perp_{CI}$$

*CI* is partially ordered

We can define the *least upper bound* and *greatest lower bound*

# Semantics of $CI$

**Table 3.** The abstract semantics of $\overline{CI}$

$$\overline{\mathbb{S}_{CI}}[\![\text{new String(str)}]\!]() = (char(\text{str}), char(\text{str}))$$

$$\overline{\mathbb{S}_{CI}}[\![\text{concat}]\!]((\overline{C_1}, \overline{MC_1}), (\overline{C_2}, \overline{MC_2})) = (\overline{C_1} \cup \overline{C_2}, \overline{MC_1} \cup \overline{MC_2})$$

$$\overline{\mathbb{S}_{CI}}[\![\text{readLine}]\!]() = (\emptyset, \mathsf{K})$$

$$\overline{\mathbb{S}_{CI}}[\![\text{substring}_b^e]\!]((\overline{C_1}, \overline{MC_1})) = (\emptyset, \overline{MC_1})$$

$$\overline{\mathbb{B}_{CI}}[\![\text{contains}_c]\!]((\overline{C_1}, \overline{MC_1})) = \begin{cases} \text{true} & \text{if } c \in \overline{C_1} \\ \text{false} & \text{if } c \notin \overline{MC_1} \\ \top_\mathsf{B} & \text{otherwise} \end{cases}$$

# Prefix - $\mathcal{PR}$

- String = Sequence of characters which *begins* with a certain sequence of characters and ends with any string ($\epsilon$ included).

- Partial order:

$$\overline{\mathcal{PR}} = \mathsf{K}^* \cup \perp_{\overline{\mathcal{PR}}}$$

$$\overline{\mathsf{S}} \leq_{\overline{\mathcal{PR}}} \overline{\mathsf{T}} \Leftrightarrow \overline{\mathsf{S}} = \perp_{\overline{\mathcal{PR}}} \lor (\forall i \in [0, len(\overline{\mathsf{T}}) - 1] : len(\overline{\mathsf{T}}) \leq len(\overline{\mathsf{S}}) \land \overline{\mathsf{T}}[i] = \overline{\mathsf{S}}[i])$$

An abstract string S is smaller than T if T is a prefix of S or if S is the bottom of the domain

# Prefix (Cont.)

Least Upper Bound:

$\sqcup_{PR}$ (S$_1$, S$_2$)= Longest common prefix between two strings.

Greatest Lower Bound:

$$\sqcap_{\overline{\mathcal{PR}}}(\overline{S}_1, \overline{S}_2) = \begin{cases} \overline{S}_1 & \text{if } \overline{S}_1 \leq_{\overline{\mathcal{PR}}} \overline{S}_2 \\ \overline{S}_2 & \text{if } \overline{S}_2 \leq_{\overline{\mathcal{PR}}} \overline{S}_1 \\ \perp_{\overline{\mathcal{PR}}} \text{ otherwise} \end{cases}$$

# Semantics of $\mathcal{PR}$

**Table 4.** The abstract semantics of $\overline{\mathcal{PR}}$

$\overline{\mathbb{S}_{\mathcal{PR}}}[\![\text{new String(str)}]\!]() = \text{str}$

$\overline{\mathbb{S}_{\mathcal{PR}}}[\![\text{concat}]\!](\overline{\mathsf{p}}_1, \overline{\mathsf{p}}_2) = \overline{\mathsf{p}}_1$

$\overline{\mathbb{S}_{\mathcal{PR}}}[\![\text{readLine}]\!]() = \epsilon$

$\overline{\mathbb{S}_{\mathcal{PR}}}[\![\text{substring}_\mathsf{b}^\mathsf{e}]\!](\overline{\mathsf{p}}) = \begin{cases} \overline{\mathsf{p}}[\mathsf{b} \cdots \mathsf{e} - 1] & \text{if } \mathsf{e} \leq len(\overline{\mathsf{p}}) \\ \overline{\mathsf{p}}[\mathsf{b} \cdots len(\overline{\mathsf{p}}) - 1] & \text{if } \mathsf{e} > len(\overline{\mathsf{p}}) \wedge \mathsf{b} < len(\overline{\mathsf{p}}) \\ \epsilon & \text{otherwise} \end{cases}$

$\overline{\mathbb{B}_{\mathcal{PR}}}[\![\text{contains}_c]\!](\overline{\mathsf{p}}) = \begin{cases} \text{true if } \mathsf{c} \in char(\overline{\mathsf{p}}) \\ \top_\mathsf{B} \quad \text{otherwise} \end{cases}$

# Suffix - $\mathcal{SF}$

- String = Sequence of characters which *ends* with a certain sequence of characters.

$$\overline{\mathcal{SU}} = \mathsf{K}^* \cup \perp_{\overline{\mathcal{SU}}}$$

- The Suffix abstract domain is nearly analogous to the Prefix abstraction
- Partial Order:

$$\overline{\mathsf{S}} \leq_{\overline{\mathcal{SU}}} \overline{\mathsf{T}} \Leftrightarrow \overline{\mathsf{S}} = \perp_{\overline{\mathcal{SU}}} \vee (\forall i \in [0, len(\overline{\mathsf{T}}) - 1] : len(\overline{\mathsf{T}}) \leq len(\overline{\mathsf{S}}) \wedge$$
$$\overline{\mathsf{T}}[i] = \overline{\mathsf{S}}[i + len(\overline{\mathsf{S}}) - len(\overline{\mathsf{T}})])$$

# Suffix (Cont.)

Least Upper Bound:

$\sqcup_{SU}$ (S$_1$, S$_2$)= Longest common suffix between two strings.

Greatest Lower Bound:

$\sqcap_{SU}$ (S$_1$, S$_2$) =     Smallest suffix     if they are comparable

$\perp_{SU}$                              if they are not comparable

# Semantics of $\mathcal{SU}$

$$\overline{\mathbb{S}_{\mathcal{SU}}}[\![\text{new String(str)}]\!]() = \text{str}$$
$$\overline{\mathbb{S}_{\mathcal{SU}}}[\![\text{concat}]\!](\overline{s}_1, \overline{s}_2) = \overline{s}_2$$
$$\overline{\mathbb{S}_{\mathcal{SU}}}[\![\text{readLine}]\!]() = \epsilon$$
$$\overline{\mathbb{S}_{\mathcal{SU}}}[\![\text{substring}_b^e]\!](\overline{s}) = \epsilon$$
$$\overline{\mathbb{B}_{\mathcal{SU}}}[\![\text{contains}_c]\!](\overline{s}) =$$
$$= \begin{cases} \text{true if } c \in \mathit{char}(\overline{s}) \\ \top_{\mathbb{B}} \quad \text{otherwise} \end{cases}$$

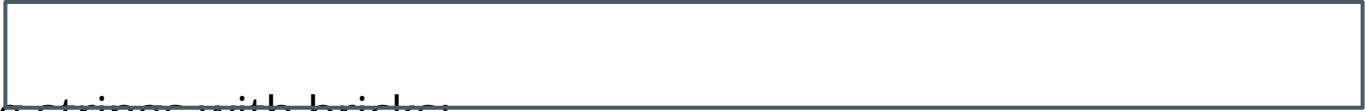(a) The abstract semantics of $\underline{\mathcal{SU}}$

# Bricks - *BR*

Significantly, Bricks capture both *inclusion* and *order*

An example brick:

Representing strings with bricks:

# Bricks – *BR* – Definition and partial order

**Definition: BR =** $B^*$ , *that is, the set of all finite sequences composed of bricks*

Partial order between single bricks:

# Bricks – *BR* – Least upper bound

**Definition: BR =** $B^*$ *, that is, the set of all finite sequences composed of bricks*

LUB between single bricks:

<br>

*The LUB is the union of each brick's set of strings, and the union of their repetition intervals.*

# Bricks - *BR* - Least upper bound example

$L_1 = [star, grape]^{0,1}[fruit]^{0,1}$      $L_2 = [grape]^{1,1} [tomato]^{0,1}$

$[star, grape]^{0,1}, [grape]^{1,1}$      $[fruit]^{0,1}, [tomato]^{0,1})$

$= [star, grape]^{0,1}[fruit, tomato]^{0,1}$

*Derives* ε, *"starfruit", "grapefruit", "grapetomato", "startomato", and each singleton string*

# Bricks – *BR* – Widening operator

The widening operator: Given $n = \max(len(L_1), len(L_2))$ , define constants $k_L$, $k_S$, $k_I$

$$\nabla_{\overline{\mathcal{BR}}}(\overline{L}_1, \overline{L}_2) = \begin{cases} \top_{\overline{\mathcal{BR}}} & \text{if } (\overline{L}_1 \not\leq_{\overline{\mathcal{BR}}} \overline{L}_2 \wedge \overline{L}_2 \not\leq_{\overline{\mathcal{BR}}} \overline{L}_1) \vee \\ & \qquad (\exists i \in [1,2] : len(\overline{L}_i) > k_L) \\ w(\overline{L}_1, \overline{L}_2) & \text{otherwise} \end{cases}$$

where $w(\overline{L}_1, \overline{L}_2) = [\overline{\mathcal{B}}_1^{\text{new}}(\overline{L}_1[1], \overline{L}_2[1]); \overline{\mathcal{B}}_2^{\text{new}}(\overline{L}_1[2], \overline{L}_2[2]); \ldots; \overline{\mathcal{B}}_n^{\text{new}}(\overline{L}_1[n], \overline{L}_2[n])]$,

$$\overline{\mathcal{B}}_i^{\text{new}}([\overline{S}_{1i}]^{m_{1i}, M_{1i}}, [\overline{S}_{2i}]^{m_{2i}, M_{2i}}) = \begin{cases} \top_{\overline{\mathcal{B}}} & \text{if } |\overline{S}_{1i} \cup \overline{S}_{2i}| > k_S \\ & \qquad \vee \overline{L}_1[i] = \top_{\overline{\mathcal{B}}} \vee \overline{L}_2[i] = \top_{\overline{\mathcal{B}}} \\ [\overline{S}_{1i} \cup \overline{S}_{2i}]^{(0,\infty)} & \text{if } (M - m) > k_I \\ [\overline{S}_{1i} \cup \overline{S}_{2i}]^{(m,M)} & \text{otherwise} \end{cases}$$

# Bricks – *BR* – Semantics

**Table 5.** The abstract semantics of $\overline{\mathcal{BR}}$

$$\overline{\mathbb{S}_{\mathcal{BR}}}[\![\text{new String(str)}]\!]() = [\{\text{str}\}]^{1,1}$$

$$\overline{\mathbb{S}_{\mathcal{BR}}}[\![\text{concat}]\!](\overline{b}_1, \overline{b}_2) = \overline{normBricks}(\overline{concatList}(\overline{b}_1, \overline{b}_2))$$

$$\overline{\mathbb{S}_{\mathcal{BR}}}[\![\text{readLine}]\!]() = \top_{\overline{\mathcal{BR}}}$$

$$\boxed{\overline{T}' = \{\overline{t}.\text{substring}(b, e) \forall \overline{t} \in \overline{T}\}}$$

$$\overline{\mathbb{S}_{\mathcal{BR}}}[\![\text{substring}_b^e]\!](\overline{b}) = \begin{cases} [\overline{T}']^{1,1} & \text{if } \overline{b}[0] = [\overline{T}]^{1,1} \wedge \forall \overline{t} \in \overline{T} : len(\overline{t}) \geq e \\ \top_{\overline{\mathcal{BR}}} & \text{otherwise} \end{cases}$$

$$\overline{\mathbb{B}_{\mathcal{BR}}}[\![\text{contains}_c]\!](\overline{b}) = \begin{cases} \text{true} & \text{if } \exists \overline{B} \in \overline{b} : \overline{B} = [\overline{T}]^{m,M} \wedge 1 \leq m \leq M \wedge (\forall \overline{t} \in \overline{T} : c \in char(\overline{t})) \\ \text{false} & \text{if } \forall [\overline{T}]^{m,M} \in \overline{b}, \forall \overline{t} \in \overline{T} : c \notin char(\overline{t}) \\ \top_B & \text{otherwise} \end{cases}$$

# Type Graphs

- A type graph T is triplet $(N, A_F, A_B)$ where $(N, A_F)$ is a rooted tree whose arcs in $A_F$ are called forward arcs, and $A_B$ is a restricted class of arcs, backward arcs, superimposed on $(N, A_F)$.
- Suitable for representing a set of terms
- A node $n \in N$ can can be in one of three classes:
  a. Simple
  b. Functor
  c. OR
- n/i denotes the i-th son of node n, and the set of sons of a node n is then denoted as {n/1,..., n/k}

# String Graphs - $\mathcal{SG}$



**Fig. 5.** An example of string graph

- Adaptation of a Type Graph to strings
- Differences:
  a. Simple nodes have labels from the set $\{max, \perp, \epsilon\} \cup K$
  b. The only functor is *concat*
- $\mathcal{SG} = \mathcal{NSG}$, where $\mathcal{NSG}$ is the set of all Normal String Graphs.
- $\perp_{\mathcal{SG}}$ = A string graph made by one bottom node
- $\mathsf{T}_{\mathcal{SG}}$ = A string graph made by only one node, a **max**-node
- Partial order:

$$\overline{\mathsf{T}}_1 \leq_{\overline{\mathcal{SG}}} \overline{\mathsf{T}}_2 \Leftrightarrow \overline{\mathsf{T}}_1 = \perp_{\overline{\mathcal{SG}}} \vee (\leq (\overline{n}_0, \overline{m}_0, \emptyset) : \overline{n}_0 = \overline{root}(\overline{\mathsf{T}}_1) \wedge \overline{m}_0 = \overline{root}(\overline{\mathsf{T}}_2))$$
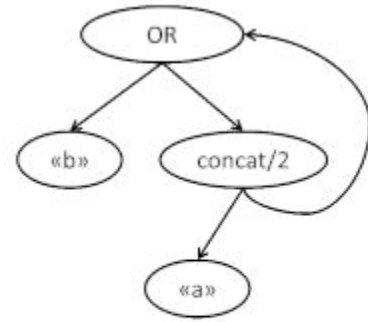
# String Graphs - $\mathcal{SG}$ (Cont.)

Least Upper Bound:

$$\sqcup_{SG} (T_1, T_2) = normStringGraph ( \text{ OR}(T_1, T_2) )$$

# Semantics of $\mathcal{SG}$

**Table 6.** The abstract semantics of $\overline{\mathcal{SG}}$

$$\overline{\mathbb{S}_{\mathcal{SG}}}[\![\mathbf{new\ String(str)}]\!]() = \mathsf{concat}/\mathsf{k}\{\mathbf{str}[\mathsf{i}] : \mathsf{i} \in [0, \mathsf{k}-1]\}$$

$$\overline{\mathbb{S}_{\mathcal{SG}}}[\![\mathbf{concat}]\!](\overline{\mathsf{t}}_1, \overline{\mathsf{t}}_2) = \overline{normStringGraph}(\mathsf{concat}/2\{\overline{\mathsf{t}}_1, \overline{\mathsf{t}}_2\})$$

$$\overline{\mathbb{S}_{\mathcal{SG}}}[\![\mathbf{readLine}]\!]() = \top_{\overline{\mathcal{SG}}}$$

$$\overline{\mathbb{S}_{\mathcal{SG}}}[\![\mathbf{substring}_\mathsf{b}^\mathsf{e}]\!](\overline{\mathsf{t}}) = \begin{cases} \overline{\mathbf{res}} & \text{if } \overline{root}(\overline{\mathsf{t}}) = \mathsf{concat}/\mathsf{k} \wedge \forall \mathsf{i} \in [0, \mathsf{e}-1] : \overline{lb}(\overline{root}(\overline{\mathsf{t}})/\mathsf{i}) \in \mathsf{K} \\ \top_{\overline{\mathcal{SG}}} & \text{otherwise} \end{cases}$$

$$\overline{\mathbb{B}_{\mathcal{SG}}}[\![\mathbf{contains}_\mathsf{c}]\!](\overline{\mathsf{t}}) = \begin{cases} \text{true} & \text{if } \exists \overline{\mathsf{m}} \in \overline{\mathsf{t}} : \overline{\mathsf{m}} = \mathsf{concat}/\mathsf{k} \wedge \mathsf{OR} \notin \overline{path}(\overline{root}, \overline{\mathsf{m}}) \wedge \\ & \qquad\qquad\qquad\qquad\qquad\qquad \exists \mathsf{i} : \overline{lb}(\overline{\mathsf{m}}/\mathsf{i}) = \mathsf{c} \\ \text{false} & \text{if } \nexists \overline{\mathsf{n}} \in \overline{\mathsf{t}} : \overline{lb}(\overline{\mathsf{n}}) = \mathsf{max} \vee \overline{lb}(\overline{\mathsf{n}}) = \mathsf{c} \\ \top_\mathsf{B} & \text{otherwise} \end{cases}$$

# Conclusion

- Two axes of precision in string value analyzers:
  - Character containment in a string
  - Position in the string
- Character inclusion ($\mathcal{CI}$)
  - Considers character containment
  - Discards the order
- Prefix ($\mathcal{PR}$) and Suffix ($\mathcal{SU}$)
  - Collect only partial information about character containment
  - Consider order only in the initial/final segment of the string

# Conclusion (Cont.)

- **Bricks ($\mathcal{BR}$)**
  - Considers character containment
  - Considers order inside the string
- **String Graph ($\mathcal{SG}$)**
  - Considers character containment
  - Considers order inside the string

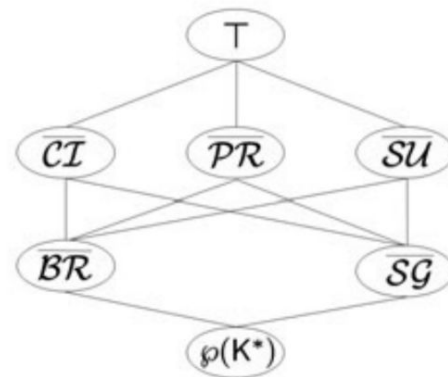So $\mathcal{BR}$ and $\mathcal{SG}$ seem to be the most precise.



**Fig. 7.** The hierarchy of abstract domains

# Reference