

Abstract Interpretation IV

Semantics and Application to Program Verification

Antoine Miné

École normale supérieure, Paris
year 2013–2014

Course 12
21 May 2014

Selected advanced topics:

- Disjunctive abstract domains
- Abstracting arrays
- Inter-procedural analyses
- Backward analyses

Practical session: help with the project

Disjunctive domains

Motivation

Remark: most domains abstract **convex sets** (conjunctions of constraints)

$\Rightarrow \cup^\#$ causes a loss of precision!

The need for non-convex invariants

```

X ← rand(10, 20);
Y ← rand(0, 1);
if Y > 0 then X ← -X;
• Z ← 100/X
  
```

Concrete semantics:

At •, $X \in [-20, -10] \cup [10, 20]$

\Rightarrow there is no division by zero

Abstract analysis:

Convex analyses (intervals, polyhedra) will find $X \in [-20, 20]$

(with intervals, $[-20, -10] \cup^\# [10, 20] = [-20, 20]$)

\Rightarrow possible division by zero

(false alarm)

Disjunctive domains

Principle:

generic constructions to lift any numeric abstract domain to a domain able to represent disjunctions exactly

Example constructions:

- powerset completion
unordered “soup” of abstract elements
- state partitioning
abstract elements keyed to selected subsets of environments
- decision tree abstract domains
efficient representation of state partitioning
- path-sensitive analyses
partition with respect to the history of execution

each construction has its strength and weakness
they can be combined during an analysis to exploit the best in each

Powerset completion

Given: $(\mathcal{E}^\#, \sqsubseteq, \gamma, \cup^\#, \cap^\#, \nabla, S^\# \llbracket \text{stat} \rrbracket)$

abstract domain $\mathcal{E}^\#$

ordered by \sqsubseteq , which also acts as a sound abstraction of \subseteq (i.e., $\sqsubseteq^\# = \sqsubseteq$)

with concretization $\gamma : \mathcal{E}^\# \rightarrow \mathcal{P}(\mathcal{E})$

sound abstractions $\cup^\#, \cap^\#, S^\# \llbracket \text{stat} \rrbracket$ of $\cup, \cap, S \llbracket \text{stat} \rrbracket$, and a widening ∇

Construct: $(\hat{\mathcal{E}}^\#, \hat{\sqsubseteq}, \hat{\gamma}, \hat{\cup}^\#, \hat{\cap}^\#, \hat{\nabla}, \hat{S}^\# \llbracket \text{stat} \rrbracket)$

- $\hat{\mathcal{E}}^\# \stackrel{\text{def}}{=} \mathcal{P}_{\text{finite}}(\mathcal{E}^\#)$ (finite sets of abstract elements)
- $\hat{\gamma}(A^\#) \stackrel{\text{def}}{=} \cup \{ \gamma(X^\#) \mid X^\# \in A^\# \}$ (join of concretizations)

Example:

using the interval domain for $\mathcal{E}^\#$

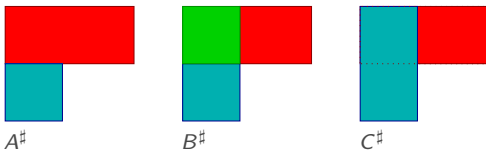
$$\hat{\gamma}\{[-10, -5], [2, 4], [0, 0], [2, 3]\} = [-10, -5] \cup \{0\} \cup [2, 4]$$

Ordering

Issue: how can we compare two elements of $\hat{\mathcal{E}}^\#$?

- $\hat{\gamma}$ is generally not injective
there is no canonical representation for $\hat{\gamma}(A^\#)$
- testing $\hat{\gamma}(A^\#) = \hat{\gamma}(B^\#)$ or $\hat{\gamma}(A^\#) \subseteq \hat{\gamma}(B^\#)$ is difficult

Example: powerset completion of the interval domain



$$A^\# = \{\{0\} \times \{0\}, [0, 1] \times \{1\}\}$$

$$B^\# = \{\{0\} \times \{0\}, \{0\} \times \{1\}, \{1\} \times \{1\}\}$$

$$C^\# = \{\{0\} \times [0, 1], [0, 1] \times \{1\}\}$$

$$\hat{\gamma}(A^\#) = \hat{\gamma}(B^\#) = \hat{\gamma}(C^\#)$$

$B^\#$ is more costly to represent: it requires three abstract elements instead of two
 $C^\#$ is a covering and not a partition ($\text{red} \cap \text{blue} = \{0\} \times \{1\} \neq \emptyset$)

Ordering (cont.)

Solution: sound approximation of \subseteq

$A^\# \hat{\sqsubseteq} B^\# \stackrel{\text{def}}{\iff} \forall X^\# \in A^\#: \exists Y^\# \in B^\#: X^\# \sqsubseteq Y^\#$ (Hoare powerdomain order)

- $\hat{\sqsubseteq}$ is a partial order (when \sqsubseteq is)
- $\hat{\sqsubseteq}$ is a sound approximation of \subseteq (when \sqsubseteq is)
 $(A^\# \hat{\sqsubseteq} B^\# \implies \hat{\gamma}(A^\#) \subseteq \hat{\gamma}(B^\#))$
- testing $\hat{\sqsubseteq}$ reduces to testing \sqsubseteq finitely many times

Example: powerset completion of the interval domain



$A^\#$



$B^\#$



$C^\#$

$$\hat{\gamma}(A^\#) = \hat{\gamma}(B^\#) = \hat{\gamma}(C^\#)$$

$$B^\# \hat{\sqsubseteq} A^\# \hat{\sqsubseteq} C^\#$$

Abstract operations

Abstract operators

- $\hat{S}^\# \llbracket stat \rrbracket A^\# \stackrel{\text{def}}{=} \{ S^\# \llbracket stat \rrbracket X^\# \mid X^\# \in A^\# \}$
 apply *stat* on each abstract element independently
- $A^\# \hat{\cup}^\# B^\# \stackrel{\text{def}}{=} A^\# \cup B^\#$
 keep elements from both arguments without applying any abstract operation
 $\hat{\cup}^\#$ is **exact**
- $A^\# \hat{\cap}^\# B^\# \stackrel{\text{def}}{=} \{ X^\# \cap^\# Y^\# \mid X^\# \in A^\#, Y^\# \in B^\# \}$
 $\hat{\cap}^\#$ is **exact** if $\cap^\#$ is (as \cup and \cap are distributive)

Galois connection:

in general, there is **no abstraction function** $\hat{\alpha}$ corresponding to $\hat{\gamma}$

Example: powerset completion $\hat{\mathcal{E}}^\#$ of the interval domain $\mathcal{E}^\#$

given the disc $S \stackrel{\text{def}}{=} \{ (x, y) \mid x^2 + y^2 \leq 1 \}$

$\alpha(S) = [-1, 1] \times [-1, 1]$ (optimal interval abstraction)

but there is no best abstraction in $\hat{\mathcal{E}}^\#$



$\alpha(S)$



not $\hat{\alpha}(S)$

Dynamic approximation

Issue: the size $|A^\#|$ of elements $A^\# \in \hat{\mathcal{E}}^\#$ is unbounded
 (every application of $\hat{\cup}^\#$ adds some more elements)
 \implies efficiency and convergence problems

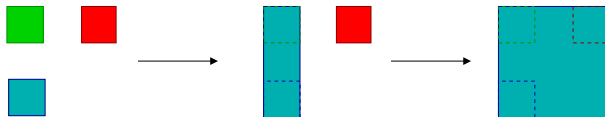
Solution: to reduce the size of elements

- redundancy removal

$\text{simplify}(A^\#) \stackrel{\text{def}}{=} \{X^\# \in A^\# \mid \forall Y^\# \neq X^\# \in A^\#: X^\# \not\subseteq Y^\#\}$
 (no loss of precision: $\hat{\gamma}(\text{simplify}(A^\#)) = \hat{\gamma}(A^\#)$)

- collapse: join elements in $\mathcal{E}^\#$

$\text{collapse}(A^\#) \stackrel{\text{def}}{=} \{\cup^\# \{X^\# \in A^\#\}\}$



(large loss of precision, but very effective: $|\text{collapse}(A^\#)| = 1$)

- partial collapse: limit $|A^\#|$ to a fixed size k by $\cup^\#$
 (but how to choose which elements to merge? no easy solution!)

Widening

Issue: for loops, abstract iterations $(A_n^\#)_{n \in \mathbb{N}}$ may not converge

- the size of $A_n^\#$ may grow arbitrarily large
- even if $|A_n^\#|$ is stable, some elements in $A_n^\#$ may not converge (if $\mathcal{E}^\#$ has infinite increasing sequences)

\implies we need a **widening** ∇

Widenings for powerset domains are **difficult to design**

Example widening: collapse after a fixed number N of iterations

$$A_{n+1}^\# \stackrel{\text{def}}{=} \begin{cases} A_n^\# \hat{\cup}^\# B_{n+1}^\# & \text{if } n < N \\ \text{collapse}(A_n^\#) \nabla \text{collapse}(B_{n+1}^\#) & \text{otherwise} \end{cases}$$

(this is very naïve, see Bagnara et al. STTT06 for more interesting widenings)

State partitioning

Principle:

- partition *a priori* \mathcal{E} into finitely many sets
- abstract each partition separately in \mathcal{E}^\sharp

Abstract domain:

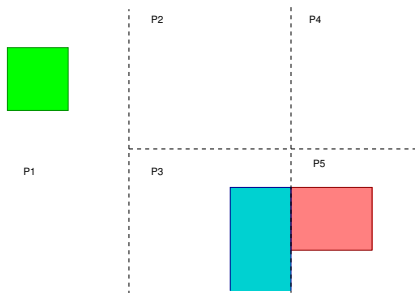
Given $P^\sharp \subseteq \mathcal{E}^\sharp$ such that:

- P^\sharp is finite
- $\bigcup \{ \gamma(X^\sharp) \mid X^\sharp \in P^\sharp \} = \mathcal{E}$
for generally, we have a covering, not a partitioning of \mathcal{E}
i.e., we can have $X^\sharp \neq Y^\sharp \in P^\sharp$ with $\gamma(X^\sharp) \cap \gamma(Y^\sharp) \neq \emptyset$

Then $\tilde{\mathcal{E}}^\sharp \stackrel{\text{def}}{=} P^\sharp \rightarrow \mathcal{E}^\sharp$

(representable in memory, as P^\sharp is finite)

Ordering



Example: $\mathcal{E}^\#$ is the interval domain

$P^\# = \{P_1, P_2, P_3, P_4, P_5\}$ where

$$P_1 = [-\infty, 0] \times [-\infty, +\infty]$$

$$P_2 = [0, 10] \times [0, +\infty]$$

$$P_3 = [0, 10] \times [-\infty, 0]$$

$$P_4 = [10, +\infty] \times [0, +\infty]$$

$$P_5 = [10, +\infty] \times [-\infty, 0]$$

$$X^\# = [P_1 \mapsto [-6, -5] \times [5, 6], P_2 \mapsto \perp, \\ P_3 \mapsto [9, 10] \times [-\infty, -1], P_4 \mapsto \perp, \\ P_5 \mapsto [10, 12] \times [-3, -1]]$$

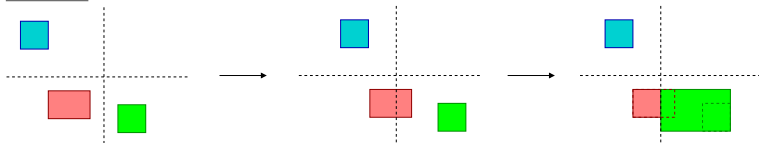
- $\tilde{\mathcal{E}}^\# \stackrel{\text{def}}{=} P^\# \rightarrow \mathcal{E}^\#$
- $\tilde{\gamma}(A^\#) \stackrel{\text{def}}{=} \bigcup \{ \gamma(A^\#(X^\#)) \cap \gamma(X^\#) \mid X^\# \in P^\# \}$
- $A^\# \sqsubseteq B^\# \stackrel{\text{def}}{\iff} \forall X^\# \in P^\#. A^\#(X^\#) \sqsubseteq B^\#(X^\#)$ (point-wise order)
- $\tilde{\alpha}(S) \stackrel{\text{def}}{=} \lambda X^\# \in P^\#. \alpha(S \cap \gamma(X^\#))$
(if $\mathcal{E}^\#$ enjoys a **Galois connection**, so does $\tilde{\mathcal{E}}^\#$)

Abstract operators

Abstract operators: point-wise extension from $\mathcal{E}^\#$ to $P^\# \rightarrow \mathcal{E}^\#$

- $A \tilde{\cup}^\# B \stackrel{\text{def}}{=} \lambda X^\# \in P^\#. A(X^\#) \cup^\# B(X^\#)$
- $A \tilde{\cap}^\# B \stackrel{\text{def}}{=} \lambda X^\# \in P^\#. A(X^\#) \cap^\# B(X^\#)$
- $A \tilde{\vee}^\# B \stackrel{\text{def}}{=} \lambda X^\# \in P^\#. A(X^\#) \vee^\# B(X^\#)$
- $\tilde{S}^\# \llbracket e \leq 0? \rrbracket A^\# \stackrel{\text{def}}{=} \lambda X^\# \in P^\#. S^\# \llbracket e \leq 0? \rrbracket A^\#(X^\#)$
- $\tilde{S}^\# \llbracket V \leftarrow e \rrbracket A^\#$ is **more complex**
as $S^\# \llbracket V \leftarrow e \rrbracket A^\#(X^\#)$ may escape $X^\#$

example: $X \leftarrow X + 2$



$$\tilde{S}^\# \llbracket V \leftarrow e \rrbracket A^\# \stackrel{\text{def}}{=} \lambda X^\#. \cup^\# \{ X^\# \cap^\# S^\# \llbracket V \leftarrow e \rrbracket A(Y^\#) \mid Y^\# \in P^\# \}$$

Example analysis

Example

```

X ← rand(10, 20);
Y ← rand(0, 1);
if Y > 0 then X ← -X;
• Z ← 100/X

```

Analysis:

- $\mathcal{E}^\#$ is the interval domain
- partition with respect to the sign of X
 $P^\# \stackrel{\text{def}}{=} \{X^+, X^-\}$ where
 $X^+ \stackrel{\text{def}}{=} [0, +\infty] \times \mathbb{Z} \times \mathbb{Z}$ and $X^- \stackrel{\text{def}}{=} [-\infty, 0] \times \mathbb{Z} \times \mathbb{Z}$
- at • we find:
 $X^+ \mapsto [X \in [10, 20], Y \mapsto [0, 0], Z \mapsto [0, 0]]$
 $X^- \mapsto [X \in [-20, -10], Y \mapsto [1, 1], Z \mapsto [0, 0]]$
 \implies no division by zero

Binary decision trees

Principle: data-structure to compactly represent partitions

Example: boolean partitions

- assume that variables have a type: $\mathbb{V} \stackrel{\text{def}}{=} \mathbb{V}_b \cup \mathbb{V}_n$
 - each $V \in \mathbb{V}_b$ has value in $\{0, 1\}$ (boolean variable)
 - each $V \in \mathbb{V}_n$ has value in \mathbb{Z} (numeric variable)

- $\mathcal{E} \simeq \{0, 1\}^{|\mathbb{V}_b|} \times \mathbb{Z}^{|\mathbb{V}_n|}$

$$P^\# \stackrel{\text{def}}{=} \{ \langle b_1, \dots, b_{|\mathbb{V}_b|} \rangle \times \mathbb{Z}^{|\mathbb{V}_n|} \mid b_1, \dots, b_{|\mathbb{V}_b|} \in \{0, 1\} \}$$

a partition corresponds to a precise valuation of all the boolean variables and no information on the numeric variables

- assume that $\mathcal{E}_n^\#$ abstracts $\mathcal{P}(\mathbb{V}_n \rightarrow \mathbb{Z})$ (numeric domain)

the boolean partitioning domain based on $\mathcal{E}_n^\#$ is:

$$\tilde{\mathcal{E}}^\# \stackrel{\text{def}}{=} \{0, 1\}^{|\mathbb{V}_b|} \rightarrow \mathcal{E}_n^\#$$

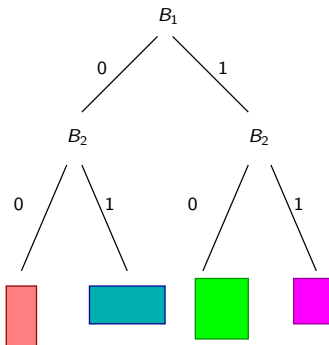
Binary decision trees (cont.)

Representation: for $\tilde{\mathcal{E}}^\# \stackrel{\text{def}}{=} \{0, 1\}^{|\mathbb{V}_b|} \rightarrow \mathcal{E}_n^\#$

binary trees:

- **nodes** are labelled with **boolean variables** $B_i \in \mathbb{V}_b$
- two children: $B_i = 0$ and $B_i = 1$
- **leaves** are **abstract elements in $\mathcal{E}_n^\#$**

(abstraction of $\mathcal{P}(\mathbb{V}_n \rightarrow \mathbb{Z})$)

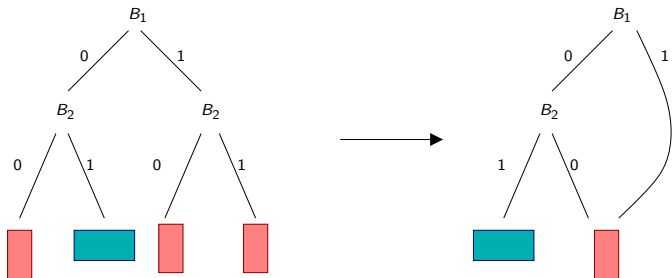


Reduced binary decision trees

Optimization: similar to **Reduced Ordered Binary Decision Diagrams**

- **merge** identical sub-trees (memory sharing)
- **remove** nodes if both children are identical

\Rightarrow we get a directed acyclic graphs



if $\gamma_n : \mathcal{E}_n^\# \rightarrow \mathbb{Z}^{|\mathbb{V}_n|}$ is injective and we use memoization
 then $\tilde{\gamma}(A^\#) = \tilde{\gamma}(B^\#) \iff A^\#$ and $B^\#$ are physically equal
 (i.e., $=$ in OCaml, which is faster to test than structural equality =)

Abstract operations

- **numeric operations:** performed independently on each leaf
(e.g., $\tilde{S}^\# \llbracket V \leftarrow e \rrbracket$ reverts to applying $S^\# \llbracket V \leftarrow e \rrbracket$ on each leaf)
- **boolean operations:** manipulate trees
 - $\tilde{S}^\# \llbracket B_i \leftarrow \mathbf{rand}(0,1) \rrbracket$: merge B_i 's subtrees recursively
 - $\tilde{S}^\# \llbracket B_i = 0? \rrbracket$: set all $B_i = 1$ branches to \perp
 - ...
- **binary operations:** $\tilde{U}^\#, \tilde{\cap}^\#, \tilde{\vee}, \tilde{\sqsubseteq}$
 - first, unify tree structures (unshare trees and add missing nodes)
 - then, apply the operation pair-wise on leaves
- optimization needs to be performed again after each operation
(ensures that abstract elements do not grow too large)

Example analysis

Example

```
X ← rand(0, 100);  
if X = 0 then B ← 0 else B ← 1;  
...  
• if B = 1 then • Y ← 100/X
```

Analysis: using the interval domain for \mathcal{E}_n^\sharp

at •, we can infer the invariant:

$$(B = 0 \implies X = 0) \wedge (B = 1 \implies X \in [1, 100])$$

at •, we deduce that $B = 1 \wedge X \in [1, 100]$

\implies there is no division by zero

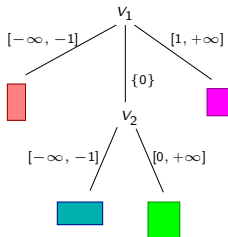
More tree-based partitioning

Other tree-based partitioning data-structure

we can extend partition trees in many ways

- allow n -array nodes
and partition wrt. abstract values

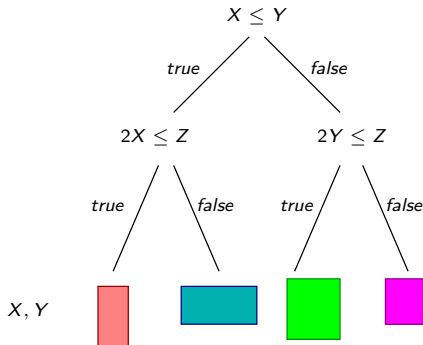
Example: partitioning integer variables in the interval domain



More tree-based partitioning

- partitioning with respect to predicates

Example: linear relations over $\mathbb{V} \stackrel{\text{def}}{=} \{X, Y, Z\}$



the same variables may appear in predicates and in the leaves

$\implies S^\sharp \llbracket \text{stat} \rrbracket$ must generally update both the nodes and the leaves

the set of node predicates may be fixed before the analysis
or chosen dynamically during the analysis

Path sensitivity

Principle: partition wrt. the **history of computation**

- keep different abstract elements for different execution **paths**
(i.e., different branches taken, different loop iterations)
- **avoid** merging with \cup^\sharp elements at control-flow **joins**
(at the end of **if** \dots **then** \dots **else**, or at loop head)

Intuition: as a program transformation

```
X ← rand(-50, 50);
if X ≥ 0 then
    Y ← X + 10
else
    Y ← X - 10;
assert Y ≠ 0
```

→

```
X ← rand(-50, 50);
if X ≥ 0 then
    Y ← X + 10;
    assert Y ≠ 0
else
    Y ← X - 10;
    assert Y ≠ 0
```

the **assert** is tested in the context of each branch
instead of after the control-flow join

the interval domain can prove the right assertion, but not the left one

Abstract domain

Formalization: limited hre to **if** \dots **then** \dots **else**

- \mathcal{L} denote **syntactic labels** of **if** \dots **then** \dots **else** instructions
- **history abstraction** $\mathbb{H} \stackrel{\text{def}}{=} \mathcal{L} \rightarrow \{\text{true}, \text{false}, \perp\}$

$H \in \mathbb{H}$ indicates the outcome of the last time we executed each test:

- $H(\ell) = \text{true}$: we took the **then** branch
- $H(\ell) = \text{false}$: we took the **else** branch
- $H(\ell) = \perp$: we never executed the test

Notes:

\mathbb{H} can remember the outcome of several successive tests

$\ell_1 : \text{if } \dots \text{ then } \dots \text{ else}; \ell_2 : \text{if } \dots \text{ then } \dots \text{ else}$

for tests in loops, H remembers only the last outcome

while \dots **do** $\ell : \text{if } \dots \text{ then } \dots \text{ else}$

we could extend \mathbb{H} to longer histories with $\mathbb{H} = (\mathcal{L} \rightarrow \{\text{true}, \text{false}, \perp\})^*$

we could extend \mathbb{H} to track loop iterations with $\mathbb{H} = \mathcal{L} \rightarrow \mathbb{N}$

- $\mathcal{E}^\# \stackrel{\text{def}}{=} \mathbb{H} \rightarrow \mathcal{E}^\#$

use a different abstract element for each abstract history

Abstract operators

- $\mathcal{E}^\# \stackrel{\text{def}}{=} \mathbb{H} \rightarrow \mathcal{E}^\#$
- $\gamma(A^\#) = \bigcup \{ \gamma(A^\#(H)) \mid H \in \mathbb{H} \}$
- $\underline{\llcorner}, \breve{\cup}^\#, \breve{\cap}^\#, \breve{\vee}$ are **point-wise**
- $\check{S}^\# \llbracket V \leftarrow e \rrbracket$ and $\check{S}^\# \llbracket e \leq 0? \rrbracket$ are **point-wise**
- $\check{S}^\# \llbracket \ell : \text{if } c \text{ then } s_1 \text{ else } s_2 \rrbracket A^\#$ is more complex
 - we **merge** all information about ℓ

$$C^\# = \lambda H. A^\#(H[\ell \mapsto \text{true}]) \cup^\# A^\#(H[\ell \mapsto \text{false}]) \cup^\# A^\#(H[\ell \mapsto \perp])$$
 - we compute the **then branch**, where $H(\ell) = \text{true}$

$$T'^\# = \check{S}^\# \llbracket s_1 \rrbracket (\check{S}^\# \llbracket c? \rrbracket T^\#)$$
 where

$$T^\# = \lambda H. C^\#(H) \text{ if } H(\ell) = \text{true}, \perp \text{ otherwise}$$
 - we compute the **else branch**, where $H(\ell) = \text{false}$

$$F'^\# = \check{S}^\# \llbracket s_2 \rrbracket (\check{S}^\# \llbracket \neg c? \rrbracket F^\#)$$
 where

$$F^\# = \lambda H. C^\#(H) \text{ if } H(\ell) = \text{false}, \perp \text{ otherwise}$$
 - we **join** both branches: $T'^\# \breve{\cup}^\# F'^\#$
 the join is exact as $\forall H \in \mathbb{H}$: either $T'^\#(H) = \perp$ or $F'^\#(H) = \perp$

\implies we get a semantic by induction on the syntax of the original program

Complex example

Linear interpolation

```

X ← rand(TX[0], TX[N]);
I ← 0;
while I < N && X > TX[I + 1] do I ← I + 1;
Y ← TY[I] + (X - TX[I]) × TS[I]

```

Concrete semantics: table-based interpolation based on the value of X

- look-up index I in the interpolation table: $TX[I] \leq X \leq TX[I + 1]$
- interpolate from value $TY[I]$ when $X = TX[I]$ with slope $TS[I]$

Analysis: in the interval domain

- without partitioning:

$$Y \in [\min TY, \max TY] + (X - [\min TX, \max TX]) \times [\min TS, \max TS]$$
- partitioning with respect to the **number of loop iterations**:

$$Y \in \bigcup_{I \in [0, N]} TY[I] + ([0, TX[I + 1] - TX[I]) \times TS[I]$$

(**more precise** as it keeps the relation between table indices)

Abstracting arrays

Example

Example: increasing subsequence

```

 $p[0] \leftarrow 0; B[0] \leftarrow A[0];$ 
 $i \leftarrow 1; k \leftarrow 1;$ 
while  $i < N$  do
    if  $A[i] > B[k - 1]$  then
         $B[k] \leftarrow A[i];$ 
         $p[k] \leftarrow i;$ 
         $k \leftarrow k + 1;$ 
     $i \leftarrow i + 1$ 

```

Given an array $A[0], \dots, A[N - 1]$
 the program computes an increasing sub-array $B[0], \dots, B[k - 1]$
 and the index sequence $p[0], \dots, p[k - 1]$

Invariants:

$1 \leq k \leq i \leq N$	$\forall x < k: B[x] = A[p[x]]$
$\forall x: 0 \leq p[x] < N$	$\forall x < k - 1: B[x + 1] > B[x]$

Overview

- Syntax and concrete semantics
- **Non-relational** abstract semantics
e.g., $\forall i: A[i] \leq \text{constant}$
 - application to interval analysis
- **Relational** (uniform) abstract semantics
e.g., $\forall i: A[i] \leq V$
 - expand and fold operations
 - application to polyhedral analysis
- **Non-uniform** abstraction
e.g., $\forall i: A[i] \leq i$

Syntax extension

Modified expressions and statements

$expr$	$::=$	V	$(\text{scalar access}, V \in \mathbb{V})$
		$A[expr]$	$(\text{array access}, A \in \mathbb{A})$
		\dots	
$stat$	$::=$	$V \leftarrow expr$	$(\text{scalar update}, V \in \mathbb{V})$
		$A[expr] \leftarrow expr$	$(\text{array update}, A \in \mathbb{A})$
		\dots	

Our language now has two ways to access the memory

- \mathbb{V} : scalar integer variables (as before)
- \mathbb{A} : **arrays** of integer values (new)
 - arrays are indexed by **positive integers**
 - arrays are **unbounded** (to simplify, we ignore overflows)

\implies an array A is similar to a map $A : \mathbb{N} \rightarrow \mathbb{Z}$

Concrete semantics

Concrete environments: $\mathcal{E} \stackrel{\text{def}}{=} \mathcal{P}((\mathbb{V} \cup (\mathbb{A} \times \mathbb{N})) \rightarrow \mathbb{Z})$

$\rho \in \mathcal{E}$ assigns an integer value to “memory cells” as follows:

- $\rho(V)$ for every scalar variable $V \in \mathbb{V}$
- $\rho(A, i)$ for every array position $A \in \mathbb{A}$, $i \geq 0$

Concrete semantics:

$$E[V] \rho \stackrel{\text{def}}{=} \{\rho(V)\}$$

$$E[A[e]] \rho \stackrel{\text{def}}{=} \{\rho(A, i) \mid i \in E[e] \rho\}$$

$$S[V \leftarrow e] R \stackrel{\text{def}}{=} \{\rho[V \mapsto v] \mid \rho \in R, v \in E[e] \rho\}$$

$$S[A[f] \leftarrow e] R \stackrel{\text{def}}{=} \{\rho[(A, i) \mapsto v] \mid \rho \in R, v \in E[e] \rho, i \in E[f] \rho, i \geq 0\}$$

...

Summarization abstraction

Goal: reuse existing numeric abstract domains

issue: numeric domains only abstract subsets of \mathbb{Z}^n , for finite n

solution: reduce \mathcal{E} to maps on **finite** set of **abstract variables**

Abstract variables: $\mathbb{V}^\# \stackrel{\text{def}}{=} \mathbb{V} \cup \mathbb{A}$

- scalar variables in \mathbb{V} are exactly represented in $\mathbb{V}^\#$
- the contents of an **array** $A \in \mathbb{A}$ is abstracted with a single **summary variable** A (modeling the contents of the whole array)
- $\mathbb{V}^\#$ is **finite**

Summarization Galois Connection:

$$(\mathcal{P}(\mathcal{E}), \subseteq) \xleftrightarrow[\alpha_s]{\gamma_s} (\mathcal{P}(\mathbb{V}^\# \rightarrow \mathbb{Z}), \subseteq)$$

- $\alpha_s(R) \stackrel{\text{def}}{=} \{ [\mathbb{V} \mapsto \rho(\mathbb{V}), \mathbb{A} \mapsto \rho(\mathbb{A}, \iota(\mathbb{A}))] \mid \rho \in R, \iota \in \mathbb{A} \rightarrow \mathbb{N} \}$
(folds all array elements (\mathbb{A}, i) into the abstract variable \mathbb{A})
- $\gamma_s(S) \stackrel{\text{def}}{=} \{ \rho \mid \forall \iota \in \mathbb{A} \rightarrow \mathbb{N}: [\mathbb{V} \mapsto \rho(\mathbb{V}), \mathbb{A} \mapsto \rho(\mathbb{A}, \iota(\mathbb{A}))] \in S \}$
(indeed, $\gamma_s(S) = \{ \rho \mid \alpha_s(\{\rho\}) \subseteq S \} = \cup \{ R \mid \alpha_s(R) \subseteq S \}$)

Non-relational abstraction

Reminder: Interval abstraction

- $\mathcal{P}(\mathbb{V}^\# \rightarrow \mathbb{Z})$ is abstracted into $\mathbb{V}^\# \rightarrow \mathcal{P}(\mathbb{Z})$ (Cartesian abstraction)
- $\mathcal{P}(\mathbb{Z})$ is abstracted as an interval in \mathbb{I}

(Note: the Cartesian and summarization abstractions commute)

Abstract semantics: in $\mathcal{E}^\# \stackrel{\text{def}}{=} \mathbb{V}^\# \rightarrow \mathbb{I}$

- $E^\# \llbracket V \rrbracket X^\# \stackrel{\text{def}}{=} X^\#(V)$
 $E^\# \llbracket A[e] \rrbracket X^\# \stackrel{\text{def}}{=} X^\#(A)$ (e is ignored)
- $S^\# \llbracket V \leftarrow e \rrbracket X^\# \stackrel{\text{def}}{=} X^\# [V \mapsto E^\# \llbracket e \rrbracket X^\#]$
 $S^\# \llbracket A[f] \leftarrow e \rrbracket X^\# \stackrel{\text{def}}{=} X^\# [A \mapsto X^\#(A) \cup^\# E^\# \llbracket e \rrbracket X^\#]$
 (f is ignored, we perform a **weak update** that accumulates values)
- assuming $X^\#(V) = X^\#(A) = [a, b]$:
 $S^\# \llbracket V \leq c \rrbracket X^\# \stackrel{\text{def}}{=} X^\# [V \mapsto [a, \min(b, c)]]$ if $a \leq c$, \perp otherwise
 $S^\# \llbracket A[e] \leq c \rrbracket X^\# \stackrel{\text{def}}{=} X^\#$ if $a \leq c$, \perp otherwise
 (we test for satisfiability but **do not** refine $X^\#(A)$; the case $A[e] \leq A[f]$ is similar)
- other operations are unchanged, including $\cap^\#, \cup^\#, \dots$

Interval analysis example

Example: increasing subsequence

```
 $p[0] \leftarrow 0; B[0] \leftarrow A[0];$   
 $i \leftarrow 1; k \leftarrow 1;$   
while  $i < N$  do  
    if  $A[i] > B[k - 1]$  then  
         $B[k] \leftarrow A[i];$   
         $p[k] \leftarrow i;$   
         $k \leftarrow k + 1;$   
     $i \leftarrow i + 1$ 
```

Analysis result:

Assuming that $N \in [N_\ell, N_h]$, $\forall x: A[x] \in [A_\ell, A_h]$, we get:

- $\forall x: p[x] \in [0, N_h - 1]$
- $\forall x: B[x] \in [\min(0, A_\ell), \max(0, A_h)]$

Variable duplication and fold

Reminders: adding and removing regular variables

$$S[\mathbf{add} \ V] R \stackrel{\text{def}}{=} \{ \rho[V \mapsto v] \mid \rho \in R, v \in \mathbb{Z} \}$$

$$S[\mathbf{del} \ V] R \stackrel{\text{def}}{=} \{ \rho|_{\text{dom}(\rho) \setminus \{V\}} \mid \rho \in R \}$$

Expanding and folding: model dynamic summarization

$$S[\mathbf{expand} \ V \rightarrow V'] R \stackrel{\text{def}}{=} \{ \rho[V' \mapsto v] \mid \rho \in R \wedge \rho[V \mapsto v] \in R \}$$

$$S[\mathbf{fold} \ V \leftarrow V'] R \stackrel{\text{def}}{=} \{ \rho \mid \exists v: \rho[V' \mapsto v] \in R \vee \rho[V' \mapsto \rho(V), V \mapsto v] \in R \}$$

- **expand** **duplicates** a variable and its constraints
($1 \leq V \leq X \implies 1 \leq V \leq X \wedge 1 \leq V' \leq X$; but $V = V'$ does not hold!)
- **fold** **summarizes** V and V' into V
($1 \leq V \leq X \wedge 2 \leq V' \leq Y \implies 1 \leq V \leq X \vee 2 \leq V \leq Y$)
- **fold** is an **abstraction**, **expand** is its associated **concretization**:

$$\mathcal{P}(\mathbb{V} \rightarrow \mathbb{Z}) \xLeftrightarrow[S[\mathbf{expand} \ V \rightarrow V']]{S[\mathbf{fold} \ V \leftarrow V']} \mathcal{P}((\mathbb{V} \setminus \{V'\}) \rightarrow \mathbb{Z})$$

(we have a Galois insertion)

Relational expand and join

Polyhedral abstraction:

- **expand** can be exactly modeled by **copying constraints**:

$$S^\# \llbracket \mathbf{expand} \ V_a \rightarrow V_b \rrbracket \{ \sum_i \alpha_{ij} V_i \geq \beta_j \} \stackrel{\text{def}}{=} \\ \{ \sum_i \alpha_{ij} V_i \geq \beta_j \} \cup \{ \sum_{i \neq a} \alpha_{ij} V_i + \alpha_{aj} V_b \geq \beta_j \}$$

- **join** can be approximated using a **weak copy**:

$$S^\# \llbracket \mathbf{fold} \ V \leftarrow V' \rrbracket X^\# \stackrel{\text{def}}{=} S^\# \llbracket \mathbf{del} \ V' \rrbracket (X^\# \cup^\# S^\# \llbracket V \leftarrow V' \rrbracket X^\#)$$

(assignment that keeps new and old values, instead of replacing old by new)

example: $0 \leq V \leq 3 \wedge 10 \leq V' \leq 13 \implies 0 \leq V \leq 13$
which over-approximates $0 \leq V \leq 3 \vee 10 \leq V \leq 13$

- $S^\# \llbracket \mathbf{add} \ V \rrbracket$ keeps the constraint set unchanged
- $S^\# \llbracket \mathbf{del} \ V \rrbracket$ projects out V

Relational array abstraction

Goal: abstract $\mathcal{P}(\mathcal{E})$ using polyhedra over $\mathbb{V}^\# \stackrel{\text{def}}{=} \mathbb{V} \cup \mathbb{A}$

Principle: use temporary variables, join and expand

Abstract assignment: $S^\# \llbracket A[f] \leftarrow e \rrbracket X^\#$

- replace each array expression $A[expr]$ in e with a fresh copy of A
we get a new expression e' and environment $X_1^\#$

e.g., replace $B[expr]$ in $X^\#$, with B' in $X_1^\# \stackrel{\text{def}}{=} S^\# \llbracket \text{expand } B \rightarrow B' \rrbracket X^\#$

- create a new copy A' of A to hold the result

$X_2^\# \stackrel{\text{def}}{=} S^\# \llbracket \text{expand } A \rightarrow A' \rrbracket X_1^\#$

- assign e' into A'

$X_3^\# \stackrel{\text{def}}{=} S^\# \llbracket A' \leftarrow e' \rrbracket X_2^\#$

- fold A' back into A

$X_4^\# \stackrel{\text{def}}{=} S^\# \llbracket \text{fold } A \leftarrow A' \rrbracket X_3^\#$

- remove all fresh copies of arrays:

$S^\# \llbracket \text{del } B' \rrbracket X_4^\#$

The cases for $S^\# \llbracket V \leftarrow e \rrbracket$ and $S^\# \llbracket c? \rrbracket$ are similar, and a bit simpler

Polyhedral analysis example

Example: increasing subsequence

```

 $p[0] \leftarrow 0; B[0] \leftarrow A[0];$ 
 $i \leftarrow 1; k \leftarrow 1;$ 
while  $i < N$  do
    if  $A[i] > B[k - 1]$  then
         $B[k] \leftarrow A[i];$ 
         $p[k] \leftarrow i;$ 
         $k \leftarrow k + 1;$ 
     $i \leftarrow i + 1$ 

```

Analysis result:

Assuming that $\forall x: A[x] \in [A_\ell, A_h]$, we get:

- $\forall x: 0 \leq p[x] < N$
(which is stronger than $\forall k: 0 \leq p[k] < N_h$)
- $\forall x: B[x] \in [\min(0, A_\ell), \max(0, A_h)]$
($B \leq A$ would mean $\forall i, j: B[i] \leq A[j]$, which does not hold)

Beyond uniform abstractions

The summarization $\alpha_s : \mathcal{P}(\mathcal{E}) \rightarrow \mathcal{P}(\mathbb{V}^\# \rightarrow \mathbb{Z})$ is **uniform**:
it forgets relations between array element **indices** and element **values**

Non-uniform abstraction example: **array segmentation**

Initialization loop

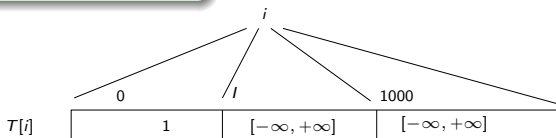
```

I ← 0;
while • I < 1000 do
  T[I] ← 1;
  I ← I + 1
  
```

we wish to analyze the loop without unrolling
at • we need to express the loop invariant:

$$\forall i < I: T[i] = 1$$

\Rightarrow at loop exit, T is initialized until 1000



abstract domain: partition the array contents into uniform segments
 segments have constant or symbolic bounds (0, I, 1000, ...)
 segments have a contents in an abstract domain (intervals, ...)

Inter-procedural analyses

Overview

- Syntax and concrete semantics
- Semantic **inlining**
simple and precise
but not efficient and may not terminate
- **Call-site** and call-stack abstraction
terminates even for recursive programs
parametric cost-precision trade-off
- **Tabulated** abstraction
optimal reuse of analysis partial results
- **Summary**-based abstraction
modular bottom-up analysis
leverage relational domains

in general, these different abstractions give incomparable results
(there is no clear winner)

Procedures

Syntax:

- \mathcal{F} finite set of procedure names
- $body : \mathcal{F} \rightarrow stat$: procedure bodies
- $main \in stat$: entry point body
- \mathbb{V}_G : set of **global** variables
- \mathbb{V}_f : set of **local** variables for procedure $f \in \mathcal{F}$
 procedure f can only access $\mathbb{V}_f \cup \mathbb{V}_G$
 $main$ has no local variable and can only access \mathbb{V}_G
- $stat ::= f(expr_1, \dots, expr_{|\mathbb{V}_f|}) \mid \dots$

procedure call, $f \in \mathcal{F}$, setting all its local variables

local variables double as **procedure arguments**

no special mechanism to return a value (a global variable can be used)

Concrete environments

Notes:

- when f calls g , we must **remember** the value of f 's locals \mathbb{V}_f in the semantics of g and **restore** them when returning
- **several copies** of each $V \in \mathbb{V}_f$ may exist at a given time (due to recursive calls, cycles in the call graph)

\implies concrete environments use **per-variable stacks**

Stacks: $\mathcal{S} \stackrel{\text{def}}{=} \mathbb{Z}^*$ (finite sequences of integers)

- **push**(v, s) $\stackrel{\text{def}}{=} v \cdot s$ ($v, v' \in \mathbb{Z}, s, s' \in \mathcal{S}$)
- **pop**(s) $\stackrel{\text{def}}{=} s'$ when $\exists v: s = v \cdot s'$, undefined otherwise
- **peek**(s) $\stackrel{\text{def}}{=} v$ when $\exists s': s = v \cdot s'$, undefined otherwise
- **set**(v, s) $\stackrel{\text{def}}{=} v \cdot s'$ when $\exists v': s = v' \cdot s'$, undefined otherwise

Environments: $\mathcal{E} \stackrel{\text{def}}{=} (\cup_{f \in \mathcal{F}} \mathbb{V}_f \cup \mathbb{V}_G) \rightarrow \mathcal{S}$

for \mathbb{V}_G , stacks are not necessary but simplify the presentation

traditionally, there is a single global stack for all local variables

using per-variable stacks instead will make the analysis presentation simpler

Concrete semantics

Concrete semantics: on $\mathcal{E} \stackrel{\text{def}}{=} (\cup_{f \in \mathcal{F}} \mathbb{V}_f \cup \mathbb{V}_G) \rightarrow \mathcal{S}$

variable read and update only consider the **top of the stack**

procedure calls **push** and **pop** local variables

- $E[V] \rho \stackrel{\text{def}}{=} \text{peek}(\rho(V))$
- $S[V \leftarrow e] R \stackrel{\text{def}}{=} \{ \rho[V \mapsto \text{set}(x, \rho(V))] \mid \rho \in R, x \in E[e] \rho \}$
- $S[f(e_{V_1}, \dots, e_{V_n})] R = R_3$, where:
 - $R_1 \stackrel{\text{def}}{=} \{ \rho[\forall V \in \mathbb{V}_f: V \mapsto \text{push}(x_V, \rho(V))] \mid \rho \in R, \forall V \in \mathbb{V}_f: x_V \in E[e_V] \rho \}$
(evaluate each argument e_V and push its value x_V on the stack $\rho(V)$)
 - $R_2 \stackrel{\text{def}}{=} S[\text{body}(f)] R_1$ (evaluate the procedure body)
 - $R_3 \stackrel{\text{def}}{=} \{ \rho[\forall V \in \mathbb{V}_f: V \mapsto \text{pop}(\rho(V))] \mid \rho \in R_2 \}$ (pop local variables)
- initial environment: $\rho_0 \stackrel{\text{def}}{=} \lambda V \in \mathbb{V}_G. 0$

(other statements are unchanged)

Semantic inlining

Naïve abstract procedure call: mimic the concrete semantics

- assign **abstract variables** to stack positions:

$$\mathbb{V}^\# \stackrel{\text{def}}{=} \mathbb{V}_G \cup (\cup_{f \in \mathcal{F}} \mathbb{V}_f \times \mathbb{N})$$

$\mathbb{V}^\#$ is **infinite**, but each abstract environment uses finitely many variables

- $\mathcal{E}_V^\#$ abstracts $\mathcal{P}(V \rightarrow \mathbb{Z})$, for any finite $V \subseteq \mathbb{V}^\#$

$V \in \mathbb{V}_f$ denotes $(V, 0)$ in $\mathbb{V}^\#$

push V : shift variables, replacing (V, i) with $(V, i + 1)$, then add $(V, 0)$

pop V : remove $(V, 0)$ and shift each (V, i) to $(V, i - 1)$

- $S^\# \llbracket f(e_1, \dots, e_n) \rrbracket X^\#$ is then reduced to:

$$X_1^\# = S^\# \llbracket \text{push } V_1; \dots; \text{push } V_n \rrbracket X^\# \quad (\text{add fresh variables for } \mathbb{V}_f)$$

$$X_2^\# = S^\# \llbracket V_1 \leftarrow e_1; \dots; V_n \leftarrow e_n \rrbracket X_1^\# \quad (\text{bind arguments to locals})$$

$$X_3^\# = S^\# \llbracket \text{body}(f) \rrbracket X_2^\# \quad (\text{execute the procedure body})$$

$$X_4^\# = S^\# \llbracket \text{pop } V_1; \dots; \text{pop } V_n \rrbracket X_3^\# \quad (\text{delete local variables})$$

Limitations:

- does not terminate in case of **unbounded recursivity**
- requires **many abstract variables** to represent the stacks
- procedures must be **re-analyzed for every call**
(full context-sensitivity: precise but costly)

Example

Example

main :

$R \leftarrow -1;$

$f(\text{rand}(5, 10));$

$f(80)$

$f(X)$:

$R \leftarrow 2 \times X;$

if $R > 100$ **then** $R \leftarrow 0$

Analysis using intervals

- after the first call to f , we get $R \in [10, 20]$
- after the second call to f , we get $R = 0$

Call-site abstraction

Abstracting stacks: into a **fixed, bounded** set $\mathbb{V}^\#$ of variables

- $\mathbb{V}^\# \stackrel{\text{def}}{=} \bigcup_{f \in \mathcal{F}} \{V, \hat{V} \mid V \in \mathbb{V}_f\} \cup \mathbb{V}_G$
 two copies of each local variable
 V abstracts the value at the top of the stack (current call)
 \hat{V} abstracts the rest of the stack
- $S^\#[\text{push } V] X^\# \stackrel{\text{def}}{=} X^\# \cup^\# S^\#[\hat{V} \leftarrow V] X^\#$
 $S^\#[\text{pop } V] X^\# \stackrel{\text{def}}{=} X^\# \cup^\# S^\#[V \leftarrow \hat{V}] X^\#$
 weak updates, similar to array manipulation
 no need to create and delete variables dynamically
- assignments and tests always access V , not \hat{V}
 \implies strong update (precise)

Note: when there is no recursivity, \hat{V} , **push** and **pop** can be omitted

Call-site abstraction

Principle: merge all the contexts in which each function is called

- we maintain two global maps $\mathcal{F} \rightarrow \mathcal{E}^\#$:

$C^\#(f)$: abstracts the environments when calling f

$R^\#(f)$: abstracts the environments when returning from f

(gather environments from all possible calls to f , disregarding the call sites)

- during the analysis, when encountering a call $S^\# \llbracket body(f) \rrbracket X^\#$:

we return $R^\#(f)$

but we also replace $C^\#$ with $C^\#[f \mapsto C^\#(f) \cup^\# X^\#]$

- $R^\#(f)$ is computed from $C^\#(f)$ as

$$R^\#(f) = S^\# \llbracket body(f) \rrbracket (C^\#(f))$$

Call-site abstraction

Fixpoint:

there may be **circular dependencies** between C^\sharp and R^\sharp

e.g., in $f(2); f(3)$, the input for $f(3)$ depends on the output from $f(2)$

\implies we compute a fixpoint for C^\sharp by iteration:

- initially, $\forall f: C^\sharp(f) = R^\sharp(f) = \perp$
- analyze *main*
- while $\exists f: C^\sharp(f)$ **not stable**
 - apply **widening** ∇ to the iterates of $C^\sharp(f)$
 - update** $R^\sharp(f) = S^\sharp \llbracket \text{body}(f) \rrbracket C^\sharp(f)$
 - analyze** *main* and all the procedures **again**
(this may modify some $C^\sharp(g)$)

\implies using ∇ , the analysis always terminates in finite time

we can be more efficient and avoid re-analyzing procedures when not needed

e.g., use a workset algorithm, track procedure dependencies, etc.

Example

Example

main :

$R \leftarrow -1;$

$f(\text{rand}(5, 10));$

$f(80)$

$f(X)$:

$R \leftarrow 2 \times X;$

if $R > 100$ **then** $R \leftarrow 0$

Analysis: using intervals (without widening as there is no dependency)

- first analysis of *main*: we get \perp (as $R^\sharp(f) = \perp$)
but $C^\sharp(f) = [R \mapsto [-1, -1], X \mapsto [5, 10]]$
- first analysis of *f*: $R^\sharp(f) = [R \mapsto [10, 20], X \mapsto [5, 10]]$
- second analysis of *main*: we get
 $C^\sharp(f) = [R \mapsto [-1, 20], X \mapsto [5, 80]]$
- second analysis of *f*: $R^\sharp(f) = [R \mapsto [0, 100], X \mapsto [5, 80]]$
- final analysis of *main*, we find $R \in [0, 100]$ at the program end
(less precise than $R = 0$ found by semantic inlining!)

Partial context-sensitivity

Variants: k -limiting, k is a constant

- **stack:**

assign a distinct variable for the k highest levels of V
 abstract the lower (unbounded) stack part with \hat{V}
 (more precise than keeping only the top of the stack separately)

- **context-sensitivity:**

each syntactic call has a unique **call-site** $\ell \in \mathcal{L}$
 a call stack is a sequence of nested call sites: $c \in \mathcal{L}^*$
 an **abstract call stack** remembers the last k call sites: $c^\# \in \mathcal{L}^k$
 the $C^\#$ and $R^\#$ maps now distinguish abstract call stacks
 $C^\#, R^\# : \mathcal{L}^k \rightarrow \mathcal{E}^\#$
 (more precise than a partitioning by function only)

larger k give more precision but less efficiency

Example: context-sensitivity

Example

main :

$R \leftarrow -1;$

$\ell_1 : f(\text{rand}(5, 10));$

$\ell_2 : f(80)$

$f(X) :$

$R \leftarrow 2 \times X;$

if $R > 100$ **then** $R \leftarrow 0$

Analysis: using intervals and $k = 1$

- $C^\sharp(\ell_1) = [R \mapsto [-1, 1], X \mapsto [5, 10]]$
 $\implies R^\sharp(\ell_1) = [R \mapsto [10, 20], X \mapsto [5, 10]]$
- $C^\sharp(\ell_2) = [R \mapsto [10, 20], X \mapsto [80, 80]]$
 $\implies R^\sharp(\ell_2) = [R \mapsto [0, 0], X \mapsto [80, 80]]$
- at the end of the analysis, we get $R = 0$
 (more precise than $R \in [0, 100]$ found without context-sensitivity)

Cardinal power

Principle:

the semantic of a function is $S[\![\text{body}(f)]\!] : \mathcal{P}(\mathcal{E}) \rightarrow \mathcal{P}(\mathcal{E})$

abstract it as an **abstract function** in $\mathcal{E}^\# \rightarrow \mathcal{E}^\#$

(we use a partial function as the image of most abstract elements is not useful)

Analysis: tabulated analysis

- use a global partial map $F^\# : \mathcal{F} \times \mathcal{E}^\# \rightarrow \mathcal{E}^\#$
- $F^\#$ is initially empty, and is filled **on-demand**
- when encountering $S^\#[\![\text{body}(f)]\!] X^\#$
 return $F^\#(f, X^\#)$ if defined
 else, compute $S^\#[\![\text{body}(f)]\!] X^\#$, store it in $F^\#(f, X^\#)$ and return it

Optimizations: trade precision for efficiency

- if $X^\# \sqsubseteq Y^\#$ and $F^\#(f, X^\#)$ is not defined, we can use $F^\#(f, Y^\#)$ instead
- if the size of $F^\#$ grows too large, use $F^\#(f, \top)$ instead
 (sound, and ensures that the analysis terminates in finite time)

Example

Example

main :

$R \leftarrow -1;$
 $f(\text{rand}(5, 10));$
 $f(80)$

$f(X)$:

$R \leftarrow 2 \times X;$
if $R > 100$ **then** $R \leftarrow 0$

Analysis using intervals

- $F^\# =$

$$[(f, [R \mapsto [-1, -1], X \mapsto [5, 10]]) \mapsto [R \mapsto [10, 20], X \mapsto [5, 10]],$$

$$(f, [R \mapsto [10, 20], X \mapsto [80, 80]]) \mapsto [R \mapsto [0, 0], X \mapsto [80, 80]]]$$
- at the end of the analysis, we get again $R = 0$

(here, the function partitioning gives the same result as the call-site partitioning)

Dynamic partitioning: complex example

Example: McCarthy's 91 function

main :

Mc(rand(0, $+\infty$))

Mc(*n*) :

if *n* > 100 **then** *r* \leftarrow *n* - 10

else *Mc*(*n* + 11); *Mc*(*r*)

- in the concrete, when terminating:
 $r = n - 10$ when $n > 101$, and $r = 91$ when $n \in [0, 101]$
- using a widening ∇ to choose tabulated abstract values $F^\#(f, X^\#)$
 we find:

$n \in [0, 72]$	\Rightarrow	$r = 91$
$n \in [73, 90]$	\Rightarrow	$r \in [91, 101]$
$n \in [91, 101]$	\Rightarrow	$r = 91$
$n \in [102, 111]$	\Rightarrow	$r \in [91, 101]$
$n \in [112, +\infty]$	\Rightarrow	$r \in [91, +\infty]$

(source: Bourdoncle, JFP 1992)

Summary-based analyses

Principle:

- abstract the **input-output relation** using a relational domain
- analyze each procedure out of context
no information about its possible arguments
- analyze a procedure given the analysis of the procedures it calls
bottom-up analysis, from leaf functions to *main*
 \implies completely **modular analysis**
 (for recursive calls, we still need to iterate the analysis of call cycles, with ∇)

Analysis:

- analyze f with abstract variables $\mathbb{V}_f^\# \stackrel{\text{def}}{=} \{ \mathbf{V}, \mathbf{V}' \mid V \in \mathbb{V}_G \cup \mathbb{V}_f \}$
 \mathbf{V}' denotes the current value of the variable
 \mathbf{V} denotes the value of the variable at the function entry
- at the beginning of the procedure, start with $\forall V \in \mathbb{V}_G \cup \mathbb{V}_f: V = V'$
 the analysis updates only V' , never V
 at the end of the procedure, the invariant gives an input-output relation
 it **summarizes** the effect of the procedure, store it as $T^\#(f)$
- $S^\#[\llbracket \text{body}(f) \rrbracket X^\#]$ can be computed using $T^\#(f)$ and variable substitution
 $S^\#[\llbracket \forall i: \text{del } V_i'' \rrbracket (X^\#[\llbracket \forall i: V_i'' / V_i \rrbracket] \cap^\# T^\#(f)[\llbracket \forall i: V_i'' / V_i \rrbracket])$

Example

Example

max(*a*, *b*) :

if *a* > *b* **then** *r* \leftarrow *a*;
else *r* \leftarrow *b*; *c* \leftarrow *c* + 1;

main :

x \leftarrow [0, 10]; *y* \leftarrow [0, 10];
c \leftarrow 0; *max*(*x*, *y*);
r \leftarrow *r* - *x*

Analysis using polyhedra

- the analysis of *max* gives:

$$r' \geq a \wedge r' \geq b \wedge c' \geq c \wedge c' \leq c + 1 \wedge a = a' \wedge b = b' \wedge x = x' \wedge y = y'$$

- at *main*'s call to *max*

before *max*: $c' = 0 \wedge x' \in [0, 10] \wedge y' \in [0, 10]$

applying the summary: $c' \in [0, 1] \wedge x' \in [0, 10] \wedge y' \in [0, 10] \wedge r' \geq x' \wedge r' \geq y'$

at the end of the program, $x \in [0, 10]$, $y \in [0, 10]$, $r \in [0, 10]$, $c \in [0, 1]$

the method requires a **relational domain** to infer interesting input-output relations
 it compensates for the lack of information about the entry point

Backward analysis

Forward versus backward analysis

Example

```
Y ← 0;  
while Y ≤ X do Y ← Y + 1
```

Forward analysis:

- given $X \in [-10, 10]$ at the **beginning** of the program
 $Y \in [0, 11]$ at the **end** of the program

Backward analysis:

- to have $Y \in [10, 20]$ at the **end** of the program
we must have $X \in [9, 19]$ at the **beginning** of the program

Concrete semantics: forward

$$\underline{S[\textit{stat}] : \mathcal{P}(\mathcal{E}) \rightarrow \mathcal{P}(\mathcal{E})}$$

$$S[\textit{skip}] R \stackrel{\text{def}}{=} R$$

$$S[s_1; s_2] R \stackrel{\text{def}}{=} S[s_2] (S[s_1] R)$$

$$S[V \leftarrow e] R \stackrel{\text{def}}{=} \{ \rho[V \mapsto v] \mid \rho \in R, v \in E[e] \rho \}$$

$$S[c?] R \stackrel{\text{def}}{=} \{ \rho \in R \mid \text{true} \in C[c] \rho \}$$

$$S[\textit{if } c \textit{ then } s_1 \textit{ else } s_2] R \stackrel{\text{def}}{=} S[s_1] (S[c?] R) \cup S[s_2] (S[\neg c?] R)$$

$$S[\textit{while } c \textit{ do } s] R \stackrel{\text{def}}{=} S[\neg c?] (\text{Ifp } \lambda I. R \cup S[s] (S[c?] I))$$

Concrete semantics: backward

$$\overleftarrow{S} \llbracket \text{stat} \rrbracket : \mathcal{P}(\mathcal{E}) \rightarrow \mathcal{P}(\mathcal{E})$$

$$\overleftarrow{S} \llbracket \text{skip} \rrbracket F$$

$$\stackrel{\text{def}}{=} F$$

$$\overleftarrow{S} \llbracket s_1; s_2 \rrbracket F$$

$$\stackrel{\text{def}}{=} \overleftarrow{S} \llbracket s_1 \rrbracket (\overleftarrow{S} \llbracket s_2 \rrbracket F)$$

$$\overleftarrow{S} \llbracket V \leftarrow e \rrbracket F$$

$$\stackrel{\text{def}}{=} \{ \rho \mid \exists v \in E \llbracket e \rrbracket \rho : \rho[V \mapsto v] \in F \}$$

$$\overleftarrow{S} \llbracket c? \rrbracket F$$

$$\stackrel{\text{def}}{=} \{ \rho \in F \mid \text{true} \in C \llbracket c \rrbracket \rho \}$$

$$\overleftarrow{S} \llbracket \text{if } c \text{ then } s_1 \text{ else } s_2 \rrbracket F$$

$$\stackrel{\text{def}}{=} \overleftarrow{S} \llbracket c? \rrbracket (\overleftarrow{S} \llbracket s_1 \rrbracket F) \cup \overleftarrow{S} \llbracket \neg c? \rrbracket (\overleftarrow{S} \llbracket s_2 \rrbracket F)$$

$$\overleftarrow{S} \llbracket \text{while } c \text{ do } s \rrbracket F$$

$$\stackrel{\text{def}}{=} \text{lfp } \lambda I. \overleftarrow{S} \llbracket \neg c? \rrbracket F \cup \overleftarrow{S} \llbracket c? \rrbracket (\overleftarrow{S} \llbracket s \rrbracket I)$$

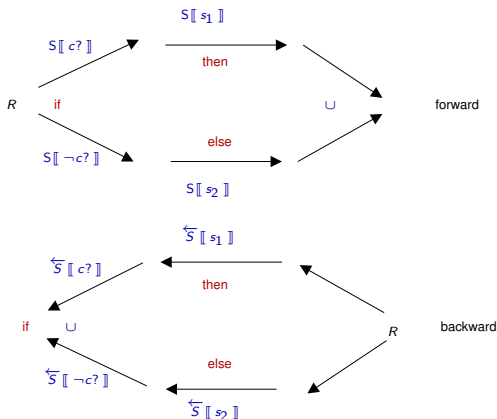
note:

statement order is inverted (s_2 before s_1 , s_1 before $c?$, etc.)

$\overleftarrow{S} \llbracket c? \rrbracket$ is unchanged

Concrete semantics: flow intuition

Intuition: information propagation for **if** \dots **then** \dots **else**



$$\begin{aligned}
 S \llbracket \text{if } c \text{ then } s_1 \text{ else } s_2 \rrbracket R &= S \llbracket s_1 \rrbracket (S \llbracket c? \rrbracket R) \cup S \llbracket s_2 \rrbracket (S \llbracket \neg c? \rrbracket R) \\
 \overleftarrow{S} \llbracket \text{if } c \text{ then } s_1 \text{ else } s_2 \rrbracket F &= \overleftarrow{S} \llbracket c? \rrbracket (\overleftarrow{S} \llbracket s_1 \rrbracket F) \cup \overleftarrow{S} \llbracket \neg c? \rrbracket (\overleftarrow{S} \llbracket s_2 \rrbracket F)
 \end{aligned}$$

Core property

Executions

- $S \llbracket \text{stat} \rrbracket R$
set of all possible states at the program end
when starting in a state in R
- $\overleftarrow{S} \llbracket \text{stat} \rrbracket F$
set of all the states at the program entry
such that at least one execution ends in a state in F

Correspondence: $\iota \in \overleftarrow{S} \llbracket \text{stat} \rrbracket \{\phi\} \iff \phi \in S \llbracket \text{stat} \rrbracket \{\iota\}$

Note: trace semantics and trace abstractions

the notion of “program execution” can be formalized as trace semantics:

$T \stackrel{\text{def}}{=} \text{Ifp } \lambda X. I \cup \{ \langle \rho_1, \dots, \rho_{n+1} \rangle \mid \langle \rho_1, \dots, \rho_n \rangle \in X \wedge \rho_n \rightarrow \rho_{n+1} \}$

$S \llbracket \cdot \rrbracket$ and $\overleftarrow{S} \llbracket \cdot \rrbracket$ are abstractions that only remember the end or beginning of traces

$S \llbracket \text{stat} \rrbracket \{\rho\} \simeq \{ \rho' \mid \exists \langle \rho_1, \dots, \rho_n \rangle \in X \wedge \rho_n \vdash \rho, \rho = \rho_1, \rho' = \rho_n \}$

$\overleftarrow{S} \llbracket \text{stat} \rrbracket \{\rho'\} \simeq \{ \rho \mid \exists \langle \rho_1, \dots, \rho_n \rangle \in X \wedge \rho_n \vdash \rho', \rho = \rho_1, \rho' = \rho_n \}$

Abstraction semantics

Goal: construct $\overleftarrow{S}^\# \llbracket stat \rrbracket$ that soundly approximates $\overleftarrow{S} \llbracket stat \rrbracket$

We can define, by induction:

$$\overleftarrow{S}^\# \llbracket \text{skip} \rrbracket F^\# \stackrel{\text{def}}{=} F^\#$$

$$\overleftarrow{S}^\# \llbracket s_1; s_2 \rrbracket F^\# \stackrel{\text{def}}{=} \overleftarrow{S}^\# \llbracket s_1 \rrbracket (\overleftarrow{S}^\# \llbracket s_2 \rrbracket F^\#)$$

$$\overleftarrow{S}^\# \llbracket c? \rrbracket F^\# \stackrel{\text{def}}{=} S^\# \llbracket c? \rrbracket F^\#$$

$$\overleftarrow{S}^\# \llbracket \text{if } c \text{ then } s_1 \text{ else } s_2 \rrbracket F^\# \stackrel{\text{def}}{=} \overleftarrow{S}^\# \llbracket c? \rrbracket (\overleftarrow{S}^\# \llbracket s_1 \rrbracket F^\#) \cup^\# \overleftarrow{S}^\# \llbracket \neg c? \rrbracket (\overleftarrow{S}^\# \llbracket s_2 \rrbracket F^\#)$$

$$\overleftarrow{S}^\# \llbracket \text{while } c \text{ do } s \rrbracket F^\# \stackrel{\text{def}}{=} \lim \lambda I^\#. I^\# \nabla (\overleftarrow{S}^\# \llbracket \neg c? \rrbracket F^\# \cup^\# \overleftarrow{S}^\# \llbracket c? \rrbracket (\overleftarrow{S}^\# \llbracket s \rrbracket I^\#))$$

Abstract operators:

- we can reuse $\cup^\#$, ∇ and $S^\# \llbracket c? \rrbracket$
- only $S^\# \llbracket V \leftarrow e \rrbracket$ needs to be defined on a per-domain basis

Backward assignment

Concrete assignment:

$$\overleftarrow{S} \llbracket V \leftarrow e \rrbracket F \stackrel{\text{def}}{=} \{ \rho \mid \exists v \in E \llbracket e \rrbracket \rho: \rho[V \mapsto v] \in F \}$$

Abstract assignment examples:

- affine assignment in **polyhedra**

$$\overleftarrow{S}^\# \llbracket V_j \leftarrow \sum_i \alpha_i V_i + \beta \rrbracket X^\#$$

\implies **substitute** V_j with $\sum_i \alpha_i V_i + \beta$ in each constraint

(similar to the computation of weakest preconditions $wlp(X \leftarrow e, P) = P[e/X]$)

- **intervals**

$$\overleftarrow{S}^\# \llbracket V \leftarrow V + \mathbf{rand}(a, b) \rrbracket X^\# = S^\# \llbracket V \leftarrow V - \mathbf{rand}(a, b) \rrbracket X^\#$$

using substitution is also possible but does not always give interval constraints
we then need to solve or approximate an optimization problem: $\min V, \max V$

- **fall-back** (e.g., non-affine assignments in polyhedra)

$$\overleftarrow{S}^\# \llbracket V \leftarrow e \rrbracket X^\# \stackrel{\text{def}}{=} S^\# \llbracket V \leftarrow [-\infty, +\infty] \rrbracket X^\#$$

(same fall-back operation as for forward assignment)

Backward-forward combination

Goal: given initial states I and final states F
 consider only executions that start in I and end in F

Application: analysis specialization to remove false alarms

Example

```

 $X \leftarrow \text{rand}(-100, 100);$ 
if  $X = 0$  then  $X \leftarrow 1;$ 
  •  $Y \leftarrow 100/X$ 
  
```

Analysis: using the interval domain

- a forward analysis finds $X \in [-100, 100]$ at •
 \implies **false alarm** for division by zero
- backward analysis from • **assuming** $X = 0$
 we find \perp at the program entry
 \implies no execution can trigger the division by zero
 (we have **removed** the false alarm)

more complex combinations exist, such as iterated forward and backward analyses

Necessary versus sufficient conditions

Example

```
Y ← 0; I ← 0;
while I ≤ X do Y ← Y + rand(1, 2); I ← I + 1
assert(Y ∈ [10, 30])
```

In case of non-determinism, $\overleftarrow{S} \llbracket \rrbracket F$ gives the initial states such that **at least one** execution terminates in F : it is a **necessary** conditions

We can also consider **sufficient conditions**
initial states such that **all** executions terminate in F

Examples: preconditions ensuring the assertion

(strongest) **necessary** precondition: $X \in [5, 30]$

(weakest) **sufficient** precondition: $X \in [10, 15]$

Note:

strongest necessary conditions can be over-approximated
weakest sufficient conditions must be **under-approximated**
⇒ leads to very different abstract operations