

Static Program Analysis

Part 11 – abstract interpretation

<http://cs.au.dk/~amoeller/spa/>

Anders Møller
Computer Science, Aarhus University

Abstract interpretation

Abstract interpretation provides a solid mathematical foundation for reasoning about static program analyses

- Is my analysis sound? (Does it safely approximate the actual program behavior?)
- Is it as precise as possible for the currently used analysis lattice? If not, where can precision losses arise? Which precision losses can be avoided (without sacrificing soundness)?

Answering such questions requires a precise definition of the semantics of the programming language, and precise definitions of the analysis abstractions in terms of the semantics

Agenda

- **Collecting semantics**
- Abstraction and concretization
- Soundness
- Optimality
- Completeness
- Trace semantics

Sign analysis, recap

$$\textit{Sign} = \begin{array}{c} \top \\ \swarrow \quad \searrow \\ + \quad - \quad 0 \\ \swarrow \quad \searrow \\ \perp \end{array}$$

$$\textit{States} = \textit{Vars} \rightarrow \textit{Sign}$$

$$\textit{States}^n$$

$$\llbracket v_1 \rrbracket = af_{v_1}(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket)$$

$$\llbracket v_2 \rrbracket = af_{v_2}(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket)$$

$$\vdots$$

$$\llbracket v_n \rrbracket = af_{v_n}(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket)$$

$$af(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket) = (af_{v_1}(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket), \dots, af_{v_n}(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket))$$

$af: \textit{States}^n \rightarrow \textit{States}^n$ is the analysis representation of the given program

$$\llbracket P \rrbracket = lfp(af)$$

Program semantics as constraint systems

$$\textit{ConcreteStates} = \textit{Vars} \hookrightarrow \mathbb{Z}$$

$$\llbracket v \rrbracket \subseteq \textit{ConcreteStates}$$

This is called a *reachable states collecting semantics*

The semantics of expressions

$$ceval: ConcreteStates \times Exp \rightarrow \mathcal{P}(\mathbb{Z})$$

$$ceval(\rho, X) = \{\rho(X)\}$$

$$ceval(\rho, I) = \{I\}$$

$$ceval(\rho, \mathbf{input}) = \mathbb{Z}$$

$$ceval(\rho, E_1 \mathbf{op} E_2) = \{v_1 \mathbf{op} v_2 \mid v_1 \in ceval(\rho, E_1) \wedge v_2 \in ceval(\rho, E_2)\}$$

$$ceval(R, E) = \bigcup_{\rho \in R} ceval(\rho, E)$$

Successors and joins

$$csucc: ConcreteStates \times Nodes \rightarrow \mathcal{P}(Nodes)$$

$$csucc(R, v) = \bigcup_{\rho \in R} csucc(\rho, v)$$

$$CJOIN(v) =$$

$$\{\rho \in ConcreteStates \mid \exists w \in Nodes: \rho \in \llbracket w \rrbracket \wedge v \in csucc(\rho, w)\}$$

Semantics of statements

$$\{\{X=E\}\} = \{\rho[X \mapsto z] \mid \rho \in CJOIN(v) \wedge z \in ceval(\rho, E)\}$$

$$\{\{\text{var } X_1, \dots, X_n\}\} = \{\rho[X_1 \mapsto z_1, \dots, X_n \mapsto z_n] \mid \rho \in CJOIN(v) \wedge z_1 \in \mathbb{Z} \wedge \dots \wedge z_n \in \mathbb{Z}\}$$

$$\{\{entry\}\} = \{\{\}\}$$

$$\{\{v\}\} = CJOIN(v)$$

The resulting constraint system

$$\begin{aligned}\llbracket v_1 \rrbracket &= cf_{v_1}(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket) \\ \llbracket v_2 \rrbracket &= cf_{v_2}(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket) \\ &\vdots \\ \llbracket v_n \rrbracket &= cf_{v_n}(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket)\end{aligned}$$

$$\begin{aligned}cf &: (\mathcal{P}(\textit{ConcreteStates}))^n \rightarrow (\mathcal{P}(\textit{ConcreteStates}))^n \\ cf(x_1, \dots, x_n) &= (cf_{v_1}(x_1, \dots, x_n), \dots, cf_{v_n}(x_1, \dots, x_n))\end{aligned}$$

is the semantic representation of the given program

$$x = cf(x)$$

$$\llbracket P \rrbracket = lfp(cf)$$

Example

```
var x;
x = 0;
while (input) {
  x = x + 2;
}
```

	solution 1	solution 2
$\{\{entry\}\}$	$\{\{\}\}$	$\{\{\}\}$
$\{\{var\ x\}\}$	$\{[x \mapsto z] \mid z \in \mathbb{Z}\}$	$\{[x \mapsto z] \mid z \in \mathbb{Z}\}$
$\{\{x = 0\}\}$	$\{[x \mapsto 0]\}$	$\{[x \mapsto 0]\}$
$\{\{input\}\}$	$\{[x \mapsto z] \mid z \in \{0, 2, 4, \dots\}\}$	$\{[x \mapsto z] \mid z \in \mathbb{Z}\}$
$\{\{x = x + 2\}\}$	$\{[x \mapsto z] \mid z \in \{2, 4, \dots\}\}$	$\{[x \mapsto z] \mid z \in \mathbb{Z}\}$
$\{\{exit\}\}$	$\{[x \mapsto z] \mid z \in \{0, 2, 4, \dots\}\}$	$\{[x \mapsto z] \mid z \in \mathbb{Z}\}$



the least solution

Kleene's fixed point theorem for complete join morphisms

$f : L \rightarrow L$ is a complete join morphism if $f(\bigsqcup A) = \bigsqcup_{a \in A} f(a)$ for every $A \subseteq L$

If f is a complete join morphism:

$$lfp(f) = \bigsqcup_{i \geq 0} f^i(\perp)$$

(even when L has infinite height!)

(Proof: Exercise 11.8)

$(\mathcal{P}(\text{ConcreteStates}))^n$ is a complete lattice (a product of a powerset lattice)

cf is a complete join morphism (Exercise 11.9)

Tarski's fixed-point theorem

In a complete lattice L , every monotone function $f: L \rightarrow L$ has a unique least fixed point given by $\bigcap \{x \in L \mid f(x) \sqsubseteq x\}$.

(Proof in Chapter 6)

cf is monotone (Exercise 11.3)

Semantics vs. analysis

```
var a,b,c;  
a = 42;  
b = 87;  
if (input) {  
    c = a + b;  
} else {  
    c = a - b;  
}
```

$$\llbracket \mathbf{b} = 87 \rrbracket = \{[a \mapsto 42, b \mapsto 87, c \mapsto z] \mid z \in \mathbb{Z}\}$$

$$\llbracket \mathbf{c} = \mathbf{a} - \mathbf{b} \rrbracket = \{[a \mapsto 42, b \mapsto 87, c \mapsto -45]\}$$

$$\llbracket \mathit{exit} \rrbracket = \{[a \mapsto 42, b \mapsto 87, c \mapsto 129], [a \mapsto 42, b \mapsto 87, c \mapsto -45]\}$$

$$\llbracket \mathbf{b} = 87 \rrbracket = [a \mapsto +, b \mapsto +, c \mapsto \top]$$

$$\llbracket \mathbf{c} = \mathbf{a} - \mathbf{b} \rrbracket = [a \mapsto +, b \mapsto +, c \mapsto \top]$$

$$\llbracket \mathit{exit} \rrbracket = [a \mapsto +, b \mapsto +, c \mapsto \top]$$

Agenda

- Collecting semantics
- **Abstraction and concretization**
- Soundness
- Optimality
- Completeness
- Trace semantics

Abstraction functions for sign analysis

$$\alpha_a: \mathcal{P}(\mathbb{Z}) \rightarrow \text{Sign}$$

$$\alpha_b: \mathcal{P}(\text{ConcreteStates}) \rightarrow \text{States}$$

$$\alpha_c: (\mathcal{P}(\text{ConcreteStates}))^n \rightarrow \text{States}^n$$

$$\alpha_a(D) = \begin{cases} \perp & \text{if } D \text{ is empty} \\ + & \text{if } D \text{ is nonempty and contains only positive integers} \\ - & \text{if } D \text{ is nonempty and contains only negative integers} \\ \mathbf{0} & \text{if } D \text{ is nonempty and contains only the integer 0} \\ \top & \text{otherwise} \end{cases}$$

for any $D \in \mathcal{P}(\mathbb{Z})$

$$\alpha_b(R) = \sigma \text{ where } \sigma(X) = \alpha_a(\{\rho(X) \mid \rho \in R\})$$

for any $R \subseteq \text{ConcreteStates}$ and $X \in \text{Vars}$

$$\alpha_c(R_1, \dots, R_n) = (\alpha_b(R_1), \dots, \alpha_b(R_n))$$

for any $R_1, \dots, R_n \subseteq \text{ConcreteStates}$

Concretization functions for sign analysis

$$\gamma_a: \text{Sign} \rightarrow \mathcal{P}(\mathbb{Z})$$

$$\gamma_b: \text{States} \rightarrow \mathcal{P}(\text{ConcreteStates})$$

$$\gamma_c: \text{States}^n \rightarrow (\mathcal{P}(\text{ConcreteStates}))^n$$

$$\gamma_a(s) = \begin{cases} \emptyset & \text{if } s = \perp \\ \{1, 2, 3, \dots\} & \text{if } s = + \\ \{-1, -2, -3, \dots\} & \text{if } s = - \\ \{0\} & \text{if } s = \mathbf{0} \\ \mathbb{Z} & \text{if } s = \top \end{cases}$$

for any $s \in \text{Sign}$

$$\gamma_b(\sigma) = \{\rho \in \text{ConcreteStates} \mid \rho(X) \in \gamma_a(\sigma(X)) \text{ for all } X \in \text{Vars}\}$$

for any $\sigma \in \text{States}$

$$\gamma_c(\sigma_1, \dots, \sigma_n) = (\gamma_b(\sigma_1), \dots, \gamma_b(\sigma_n))$$

for any $(\sigma_1, \dots, \sigma_n) \in \text{States}^n$

Monotonicity of abstraction and concretization functions

Concretization functions are, like abstraction functions, naturally monotone.

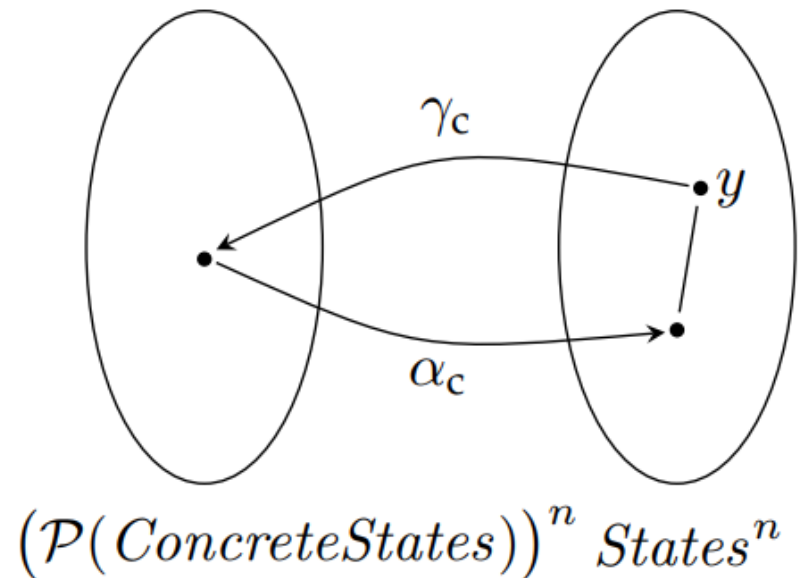
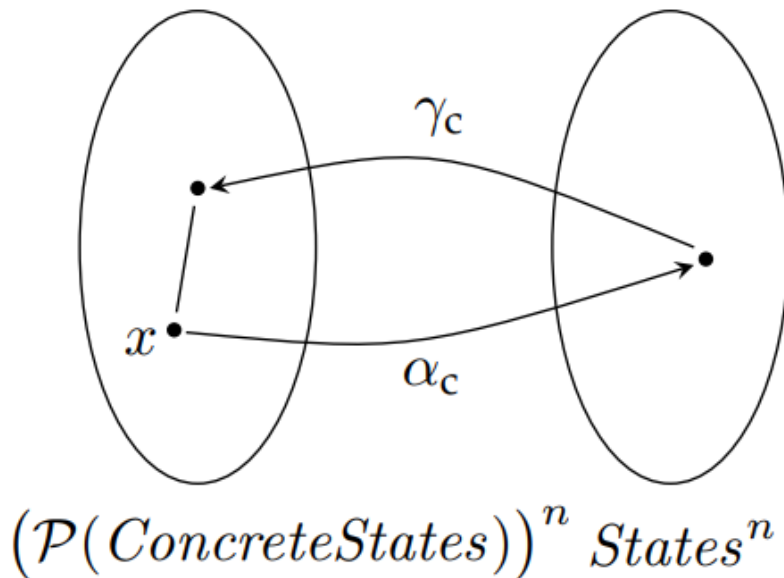
(A larger set of concrete values should correspond to a larger abstract state, and conversely)

Galois connections

The pair of monotone functions, α and γ , is called a *Galois connection* if

$\gamma \circ \alpha$ is extensive

$\alpha \circ \gamma$ is reductive



all three pairs of abstraction and concretization functions (α_a, γ_a) , (α_b, γ_b) , and (α_c, γ_c) from the sign analysis example are Galois connections

(Exercise 11.13)

Galois connections

For Galois connections, the concretization function uniquely determines the abstraction function and vice versa:

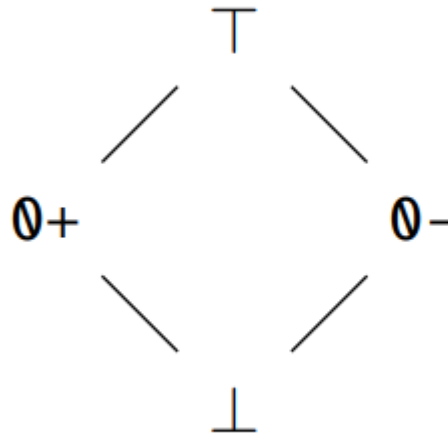
$$\gamma(y) = \bigsqcup \{x \in L_1 \mid \alpha(x) \sqsubseteq y\}$$

$$\alpha(x) = \bigsqcap \{y \in L_2 \mid x \sqsubseteq \gamma(y)\}$$

(Proof: Exercise 11.20)

Galois connections

For this lattice, given the “obvious” concretization function,
is there an abstraction function such that the concretization function
and the abstraction function form a Galois connection?



how should we define $\alpha_a(\{0\})$? (Exercise 11.22)

Representation functions

$$\beta: \mathbb{Z} \rightarrow \textit{Sign}$$

$$\beta(d) = \begin{cases} + & \text{if } d > 0 \\ - & \text{if } d < 0 \\ \mathbf{0} & \text{if } d = 0 \end{cases}$$

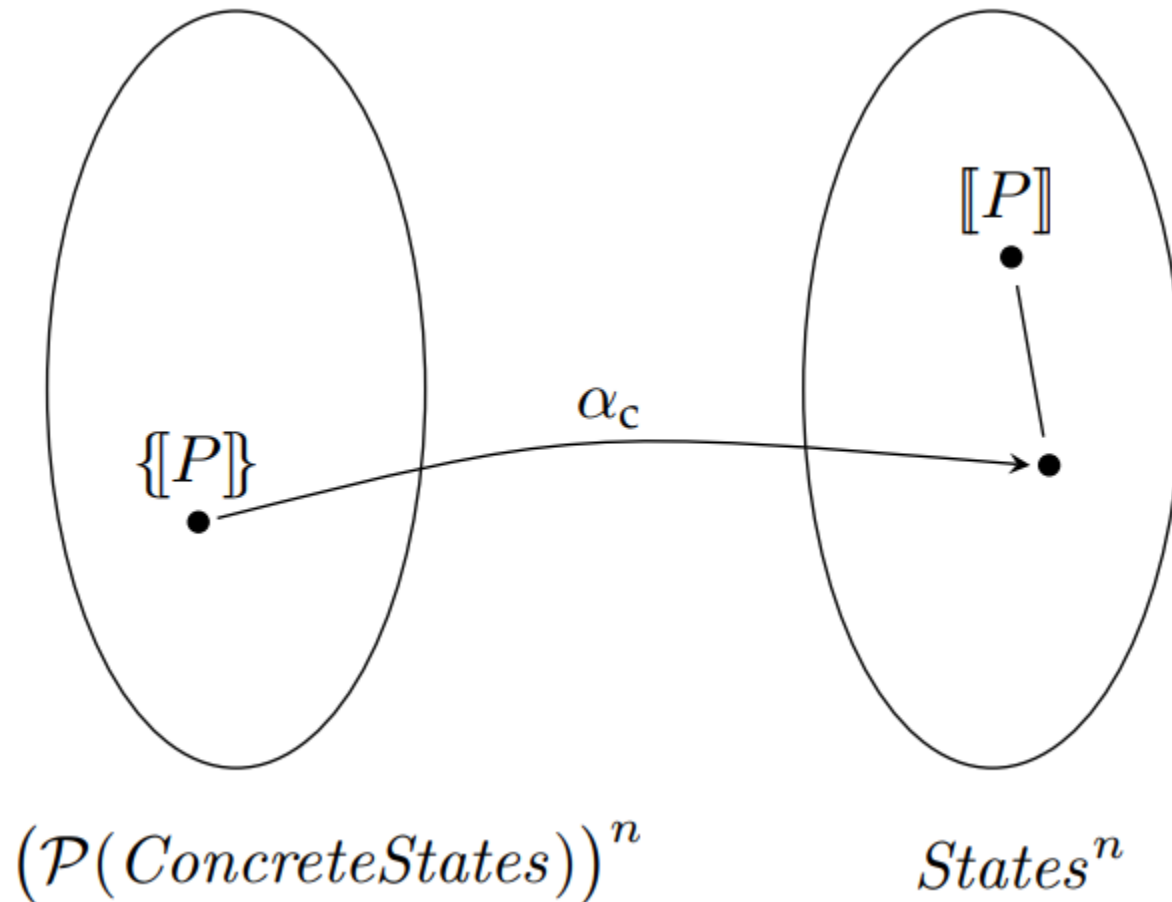
$$\alpha_a(D) = \bigsqcup \{\beta(d) \mid d \in D\}$$

Agenda

- Collecting semantics
- Abstraction and concretization
- **Soundness**
- Optimality
- Completeness
- Trace semantics

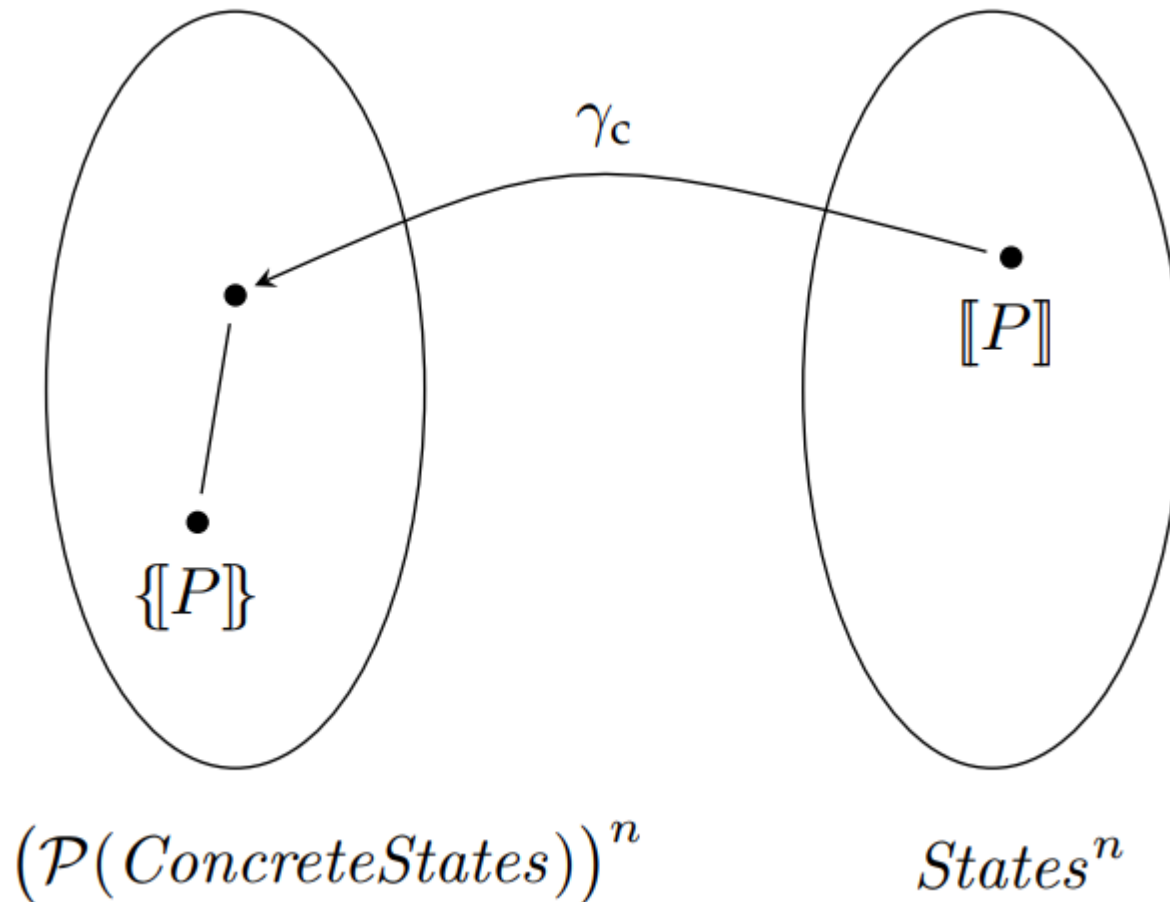
Soundness

$$\alpha(\{[P]\}) \subseteq [P]$$



Soundness

$$\{P\} \subseteq \gamma(\llbracket P \rrbracket)$$



Sound abstractions

$$\alpha_a(\text{ceval}(R, E)) \sqsubseteq \text{eval}(\alpha_b(R), E)$$

$$\text{csucc}(R, v) \subseteq \text{succ}(v) \text{ for any } R \subseteq \text{ConcreteStates}$$

$$\begin{aligned} \alpha_b(CJOIN(v)) &\sqsubseteq JOIN(v) \\ &\text{if } \alpha_b(\llbracket w \rrbracket) \sqsubseteq \llbracket w \rrbracket \text{ for all } w \in \text{Nodes}. \end{aligned}$$

Sound abstractions

if v represents an assignment statement $X = E$:

$$\begin{aligned} cf_v(\{[v_1]\}, \dots, \{[v_n]\}) &= \{\rho[X \mapsto z] \mid \rho \in CJOIN(v) \wedge z \in ceval(\rho, E)\} \\ af_v([v_1], \dots, [v_n]) &= \sigma[X \mapsto eval(\sigma, E)] \text{ where } \sigma = JOIN(v) \end{aligned}$$

$$\alpha_b(cf_v(R_1, \dots, R_n)) \sqsubseteq af_v(\alpha_b(R_1), \dots, \alpha_b(R_n))$$

(Exercise 11.33)

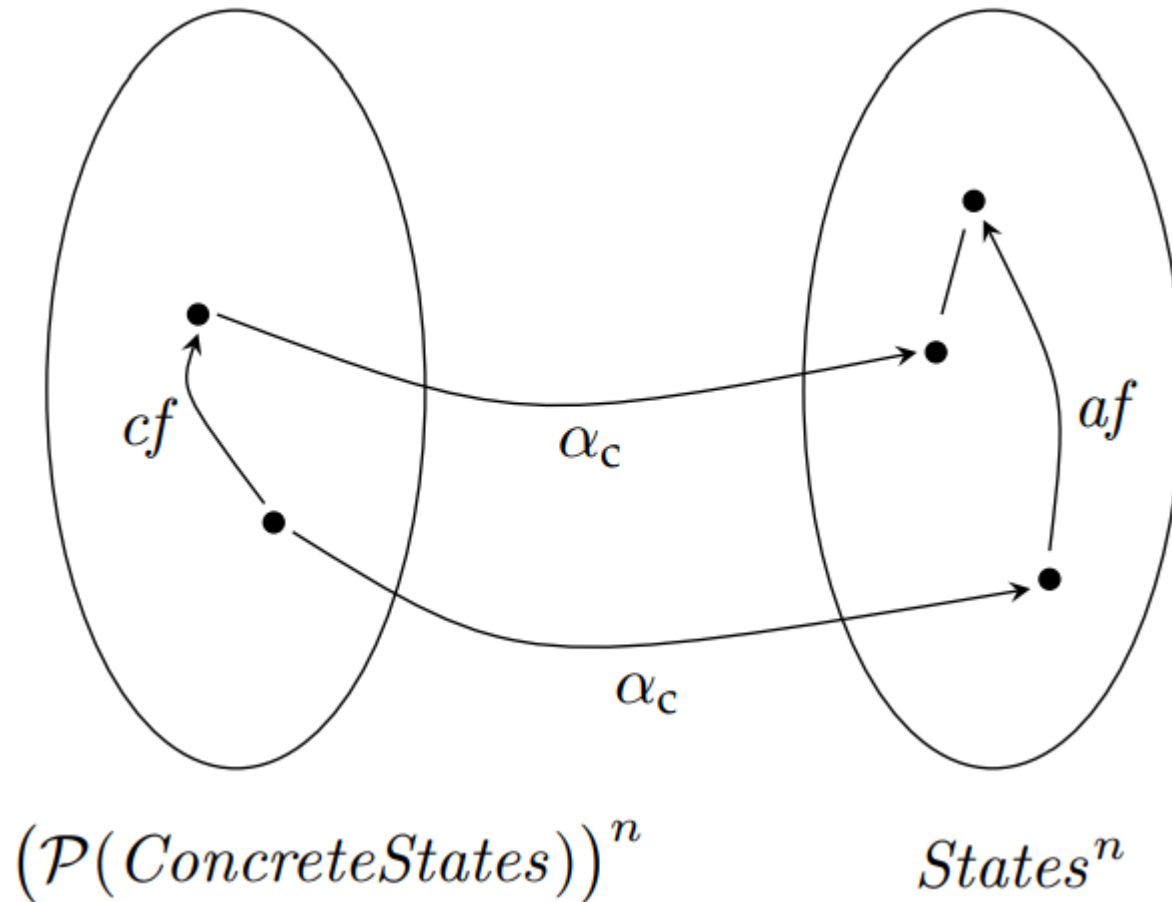
The two constraint systems

$$cf(\{v_1\}, \dots, \{v_n\}) = (cf_{v_1}(\{v_1\}, \dots, \{v_n\}), \dots, cf_{v_n}(\{v_1\}, \dots, \{v_n\}))$$

$$af(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket) = (af_{v_1}(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket), \dots, af_{v_n}(\llbracket v_1 \rrbracket, \dots, \llbracket v_n \rrbracket))$$

Sound abstractions

$$\alpha_c(cf(R_1, \dots, R_n)) \sqsubseteq af(\alpha_c(R_1, \dots, R_n))$$



Sound abstractions

$$\alpha \circ cf \sqsubseteq af \circ \alpha$$

$$cf \circ \gamma \sqsubseteq \gamma \circ af$$

Equivalent, if α and γ form a Galois connection

(Proof: Exercise 11.34)

The soundness theorem

If L_1 and L_2 are complete lattices with a concretization function $\gamma: L_2 \rightarrow L_1$, $cf: L_1 \rightarrow L_1$ and $af: L_2 \rightarrow L_2$ are monotone, and af is a sound abstraction of cf with respect to γ , i.e., $cf \circ \gamma \sqsubseteq \gamma \circ af$, then $lfp(cf) \sqsubseteq \gamma(lfp(af))$.

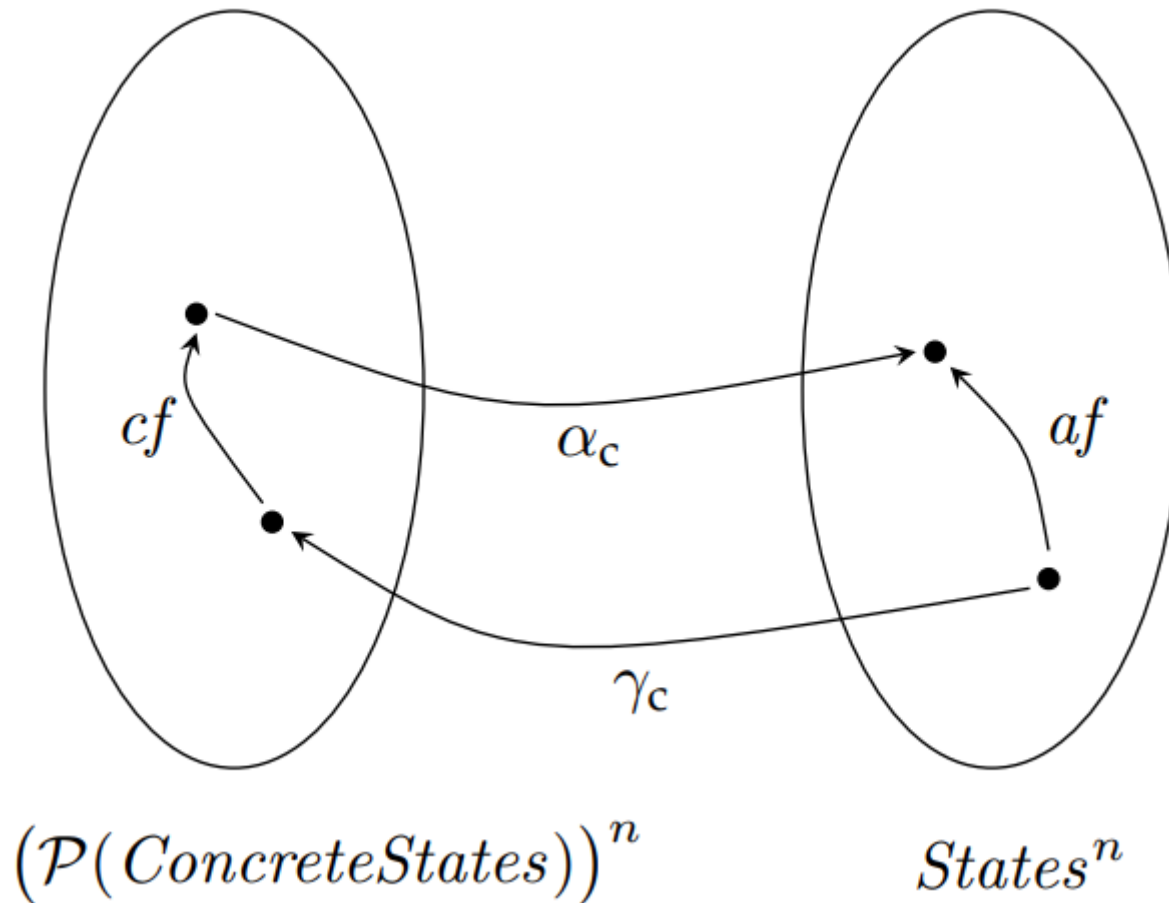
Agenda

- Collecting semantics
- Abstraction and concretization
- Soundness
- **Optimality**
- Completeness
- Trace semantics

Optimal abstractions

af is an optimal abstraction of cf if

$$af = \alpha \circ cf \circ \gamma$$



(compare with slide 28)

Optimal abstractions in sign analysis?

$\hat{*}$ is optimal:

$$s_1 \hat{*} s_2 = \alpha_a(\gamma_a(s_1) \cdot \gamma_a(s_2))$$

eval is not optimal:

$$\sigma(\mathbf{x}) = \top$$

$$eval(\sigma, \mathbf{x} - \mathbf{x}) = \top$$

$$\alpha_b(ceval(\gamma_b(\sigma), \mathbf{x} - \mathbf{x})) = \mathbf{0}$$

Even if we could make *eval* optimal, the analysis result is not always optimal:

```
x = input;
```

```
y = x;
```

```
z = x - y;
```

Agenda

- Collecting semantics
- Abstraction and concretization
- Soundness
- Optimality
- **Completeness**
- Trace semantics

Completeness

$$\alpha(\{P\}) \sqsupseteq \llbracket P \rrbracket$$

(compare with slide 23)

Sound *and* complete: $\alpha(\{P\}) = \llbracket P \rrbracket$

(Intuitively, the analysis result is the most precise possible for the currently used lattice)

Not the same as $\{P\} = \gamma(\llbracket P \rrbracket)$ (called “exact”)

(Intuitively, the analysis result exactly captures the semantics of the program)

Completeness in sign analysis?

$\hat{*}$ is complete:

$$\alpha_a(D_1) \hat{*} \alpha_a(D_2) \sqsubseteq \alpha_a(D_1 \cdot D_2)$$

$\hat{+}$ is *not* complete

$$\alpha_a(D_1) \hat{+} \alpha_a(D_2) \not\sqsubseteq \alpha_a(D_1 + D_2)$$

Sign analysis is sound and complete for *some* programs,
but not for all programs

Agenda

- Collecting semantics
- Abstraction and concretization
- Soundness
- Optimality
- Completeness
- **Trace semantics**

Limitations of the reachable states collecting semantics

The reachable states collecting semantics “collects” a set of concrete states for each program point

$$ConcreteStates = Vars \hookrightarrow \mathbb{Z}$$

$$\llbracket v \rrbracket \subseteq ConcreteStates$$

That is not always sufficient – a trivial example:

```
main(x) {  
    return x;  
}
```

For this program, the collecting semantics doesn’t allow us to express properties such as “in any execution, the return value is the same as the input value”

Trace semantics

$$Traces = (Nodes \times ConcreteStates)^*$$

$$ct_v: ConcreteStates \rightarrow \mathcal{P}(ConcreteStates)$$

$$ct_{X=E}(\rho) = \{\rho[X \mapsto z] \mid z \in ceval(\rho, E)\}$$

$$ct_v(\rho) = \{\rho\} \quad \text{for other kinds of nodes}$$

$$\langle\!\langle P \rangle\!\rangle \in \mathcal{P}(Traces)$$

$$(entry, []) \in \langle\!\langle P \rangle\!\rangle$$

$$\pi \cdot (v, \rho) \in \langle\!\langle P \rangle\!\rangle \wedge v' \in csucc(\rho, v) \wedge \rho' \in ct_{v'}(\rho) \implies \pi \cdot (v, \rho) \cdot (v', \rho') \in \langle\!\langle P \rangle\!\rangle$$

Example

```
main(x) {  
    return x;  
}
```

$$\langle\!\langle P \rangle\!\rangle = \{ (entry, [x \mapsto 0]) \cdot (return\ x, [x \mapsto 0]) \cdot (exit, [x \mapsto 0]), \\ (entry, [x \mapsto 1]) \cdot (return\ x, [x \mapsto 1]) \cdot (exit, [x \mapsto 1]), \\ \dots \}$$

Reachable states is an abstraction of traces

the relation between the reachable states collecting semantics and the trace semantics can be expressed as a Galois connection induced by an abstraction function

$$\alpha_t: \mathcal{P}(\text{Traces}) \rightarrow (\mathcal{P}(\text{ConcreteStates}))^n$$

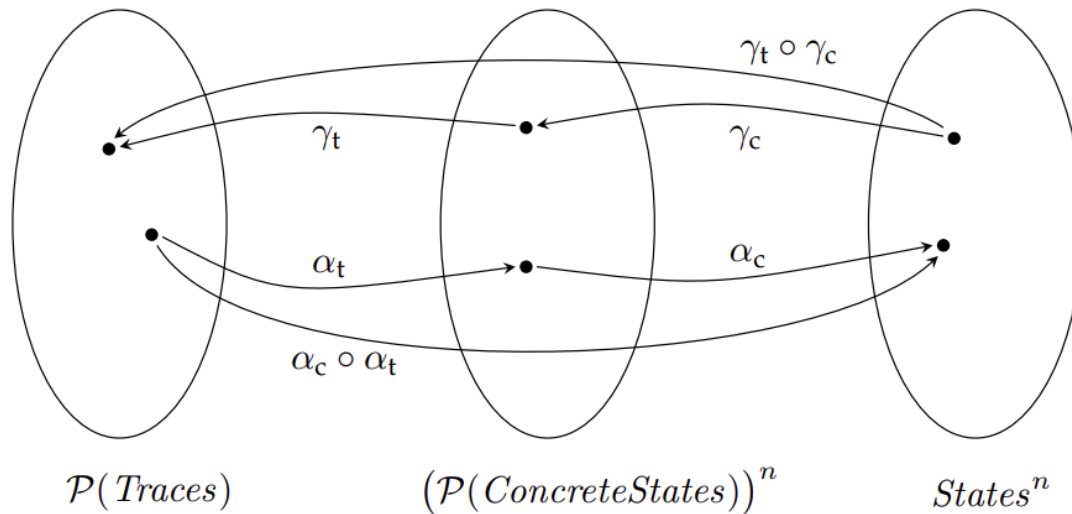
defined by

$$\alpha_t(T) = (R_1, \dots, R_n)$$

where $R_i = \{\rho \mid \dots \cdot (v_i, \rho) \cdot \dots \in T\}$ for each $i = 1, \dots, n$.

Composition of Galois connections

Exercise 11.57: Let $\alpha_1: L_1 \rightarrow L_2$, $\gamma_1: L_2 \rightarrow L_1$, $\alpha_2: L_2 \rightarrow L_3$, and $\gamma_2: L_3 \rightarrow L_2$. Assume both (α_1, γ_1) and (α_2, γ_2) are Galois connections. Prove that $(\alpha_2 \circ \alpha_1, \gamma_1 \circ \gamma_2)$ is then also a Galois connection.



Exercise 11.58: Prove that the reachable states collecting semantics is sound with respect to the trace semantics. (Even though the collecting semantics is not a computable analysis, we can still apply the notion of soundness and the proof techniques from Section 11.3.)

Soundness of reaching definitions analysis

Instrumented trace semantics:

$$Traces_{RD} = (Nodes \times ConcreteStates \times ReachingDefs)^*$$

$$ReachingDefs = Vars \hookrightarrow Nodes$$

$$(entry, [], []) \in \llbracket P \rrbracket_{RD}$$

$$\pi \cdot (v, \rho, \delta) \in \llbracket P \rrbracket_{RD} \wedge v' \in csucc(\rho, v) \wedge \rho' \in ct_{v'}(\rho) \implies \\ \pi \cdot (v, \rho, \delta) \cdot (v', \rho', \delta') \in \llbracket P \rrbracket_{RD}$$

$$\delta' = \begin{cases} \delta[X \mapsto v'] & \text{if } v' \text{ is an assignment } X = E \\ \delta & \text{otherwise} \end{cases}$$

Soundness of reaching definitions analysis

Abstraction:

$$\alpha_{RD}: \mathcal{P}(\text{Traces}_{RD}) \rightarrow (\mathcal{P}(\text{Nodes}))^n$$

$$\alpha_{RD}(T) = (R_1, \dots, R_n)$$

$$R_i = \{w \mid \dots (v_i, \rho, \delta) \dots \in T \text{ where } \delta(X) = w \text{ for some } \rho, \delta, X\}$$

Soundness:

$$\alpha_{RD}(\langle P \rangle_{RD}) \sqsubseteq \llbracket P \rrbracket_{RD}$$

Soundness proofs for other analyses

Exercise 11.61: Use the approach from Section 11.3 and an appropriate instrumented trace semantics to prove that the available expressions analysis from Section 5.5 is sound. (This is more tricky than Exercise 11.60, because available expressions analysis is a “must” analysis!)

Exercise 11.62: Use the approach from Section 11.3 and an appropriate instrumented trace semantics to prove that the live variables analysis from Section 5.4 is sound. As part of this, you need to specify an appropriate collecting semantics that formally captures what it means for a variable to be live (see the informal definition in Section 5.4). (This is more tricky than Exercise 11.60, because live variables analysis is a “backward” analysis!)

Conclusions

Abstract interpretation provides a solid mathematical foundation for reasoning about soundness and precision of static program analyses

We need

- the static analysis (the analysis lattices and constraint rules)
- the language semantics (a suitable collecting semantics)
- abstraction/concretization functions that specify the meaning of the elements in the analysis lattice in terms of the semantic lattice

... and then

- if each constituent of the analysis is a *sound abstraction* of its semantic counterpart, then the analysis is sound (according to the soundness theorem)
- if an abstraction is *optimal*, then it is as precise as possible (yet sound), relative to the choice of analysis lattice
- if the analysis is *sound and complete*, then the analysis result is as precise as possible (yet sound), relative to the choice of analysis lattice