
Static Program Analysis using Abstract Interpretation

Introduction

Static Program Analysis

Static program analysis consists of **automatically discovering** properties of a program that hold for **all possible execution paths** of the program.

Static program analysis is **not**

- **Testing: manually checking** a property for **some** execution paths
- **Model checking: automatically checking** a property for **all** execution paths

Program Analysis for what?

- Optimizing compilers
- Semantic preprocessing:
 - Model checking
 - Automated test generation
- **Program verification**

Program Verification

- Check that every operation of a program will never cause an error (division by zero, buffer overrun, deadlock, etc.)
- Example:

```
int a[1000];  
for (i = 0; i < 1000; i++) {  
safe operation → a[i] = ... ; // 0 ≤ i ≤ 999  
}
```

```
buffer overrun → a[i] = ... ; // i = 1000;
```

Incompleteness of Program Analysis

- Discovering a sufficient set of properties for checking **every** operation of a program is an **undecidable** problem!
- Every non trivial behavioral property has (at least) NP complexity
- **False positives**: operations that are safe in reality but which cannot be decided safe or unsafe from the properties inferred by static analysis.

Precision versus Efficiency

Precision: number of program operations that can be decided safe or unsafe by an analyzer.

- Precision and computational complexity are strongly related
- Tradeoff precision/efficiency: limit in the average precision and scalability of a given analyzer
- Greater precision and scalability is achieved through **specialization**

Soundness

- What guarantees the soundness of the analyzer results?
- In dataflow analysis and type inference the soundness proof of the resolution algorithm is independent from the analysis specification
- An independent soundness proof precludes the use of test-and-try techniques
- **Need for analyzers correct by construction**

Abstract Interpretation

- A general methodology for designing static program analyzers that are:
 - Correct by construction
 - Generic
 - Easy to fine-tune
- Scalability is difficult to achieve but the payoff is worth the effort!

Approximation

The core idea of Abstract Interpretation is the formalization of the notion of **approximation**

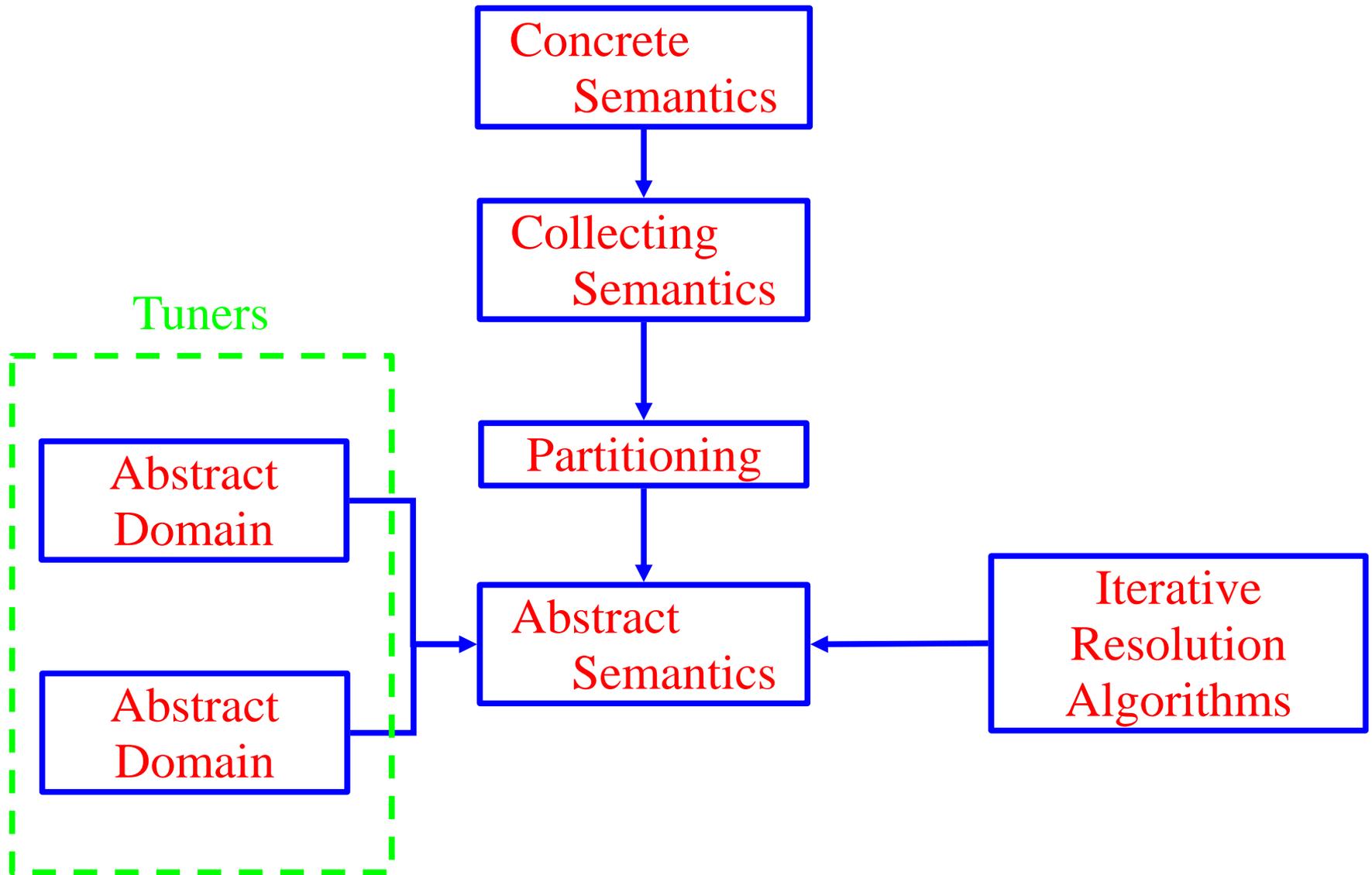
- An approximation of memory configurations is first defined
- Then the approximation of all atomic operations
- The approximation is automatically lifted to the whole program structure

Overview of Abstract Interpretation

- Start with a formal specification of the program semantics (the concrete semantics)
- Construct abstract semantic equations w.r.t. a parametric approximation scheme
- Use general fixpoints algorithms to solve the abstract semantic equations
- Try-and-test various instantiations of the approximation scheme in order to find the best fit

The Methodology of Abstract Interpretation

Methodology



Lattices and Fixpoints

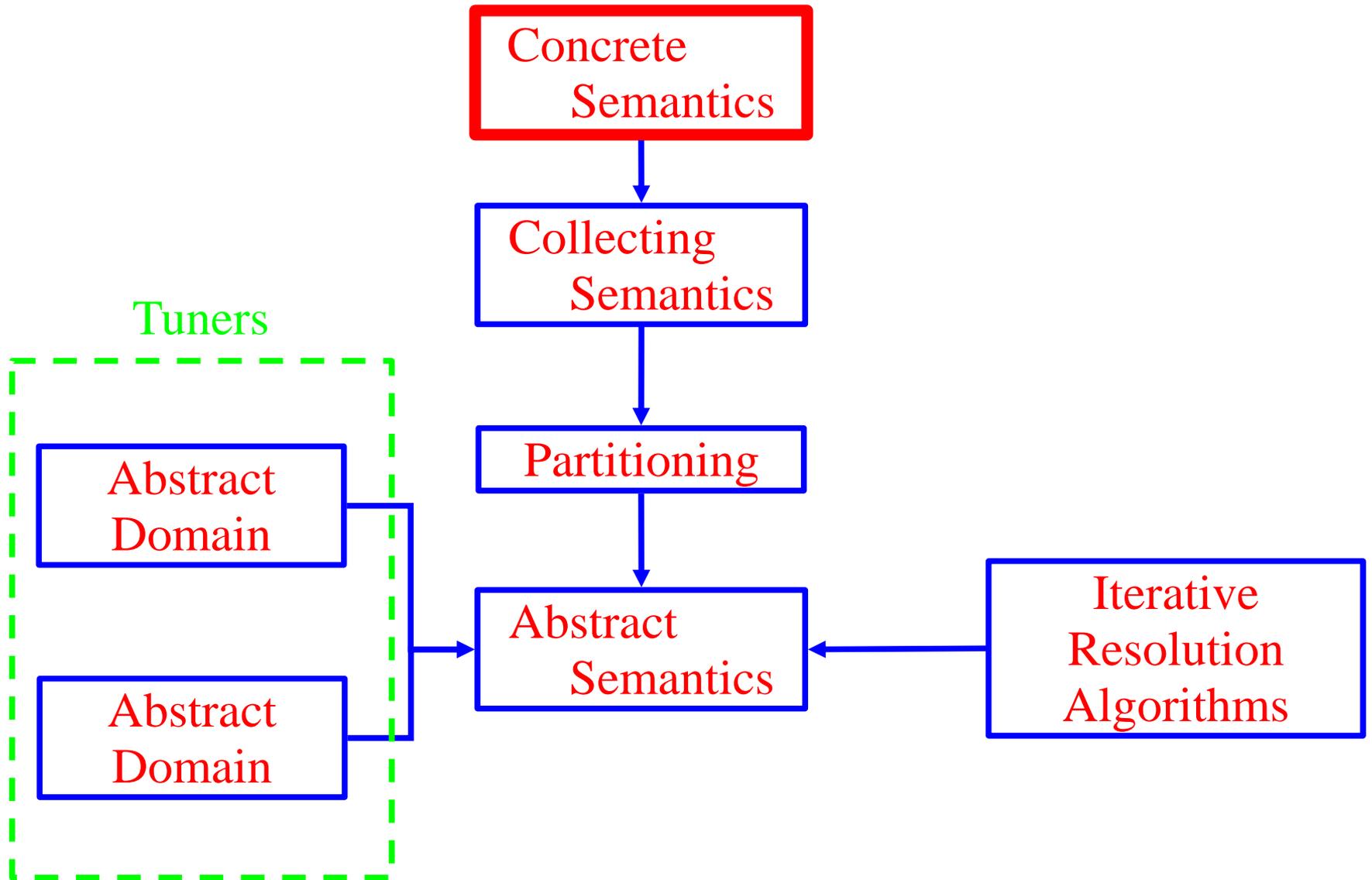
- A **lattice** $(L, \sqsubseteq, \perp, \sqcup, \top, \sqcap)$ is a partially ordered set (L, \sqsubseteq) with:
 - Least upper bounds (\sqcup) and greatest lower bounds (\sqcap) operators
 - A least element “bottom”: \perp
 - A greatest element “top”: \top
- L is **complete** if all least upper bounds exist
- A fixpoint X of $F: L \rightarrow L$ satisfies $F(X) = X$
- We denote by $\text{lfp } F$ the least fixpoint if it exists

Fixpoint Theorems

- Knaster-Tarski theorem: If $F: L \rightarrow L$ is monotone and L is a complete lattice, the set of fixpoints of F is also a complete lattice.
- Kleene theorem: If $F: L \rightarrow L$ is monotone, L is a complete lattice and F preserves all least upper bounds then $\text{lfp } F$ is the limit of the sequence:

$$\begin{cases} F_0 & = \perp \\ F_{n+1} & = F(F_n) \end{cases}$$

Methodology



Concrete Semantics

Small-step operational semantics: (Σ, \rightarrow)

$s = \langle \text{program point}, \text{env} \rangle$

$s \rightarrow s'$

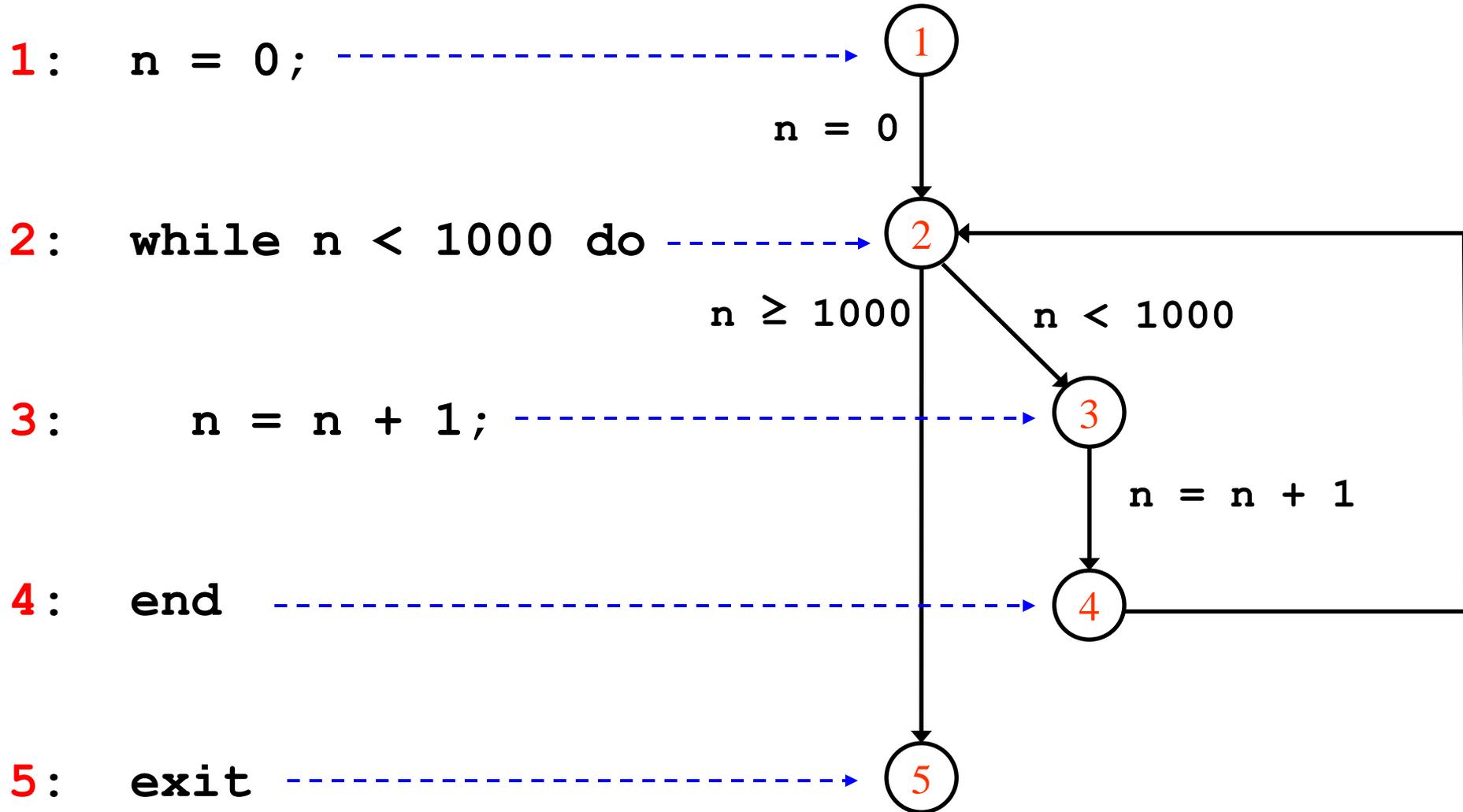
Example:

```
1:  n = 0;  
2:  while n < 1000 do  
3:    n = n + 1;  
4:  end  
5:  exit
```

$\langle 1, n \Rightarrow \Omega \rangle \rightarrow \langle 2, n \Rightarrow 0 \rangle \rightarrow \langle 3, n \Rightarrow 0 \rangle \rightarrow \langle 4, n \Rightarrow 1 \rangle$
 $\rightarrow \langle 2, n \Rightarrow 1 \rangle \rightarrow \dots \rightarrow \langle 5, n \Rightarrow 1000 \rangle$

Undefined value

Control Flow Graph



Transition Relation

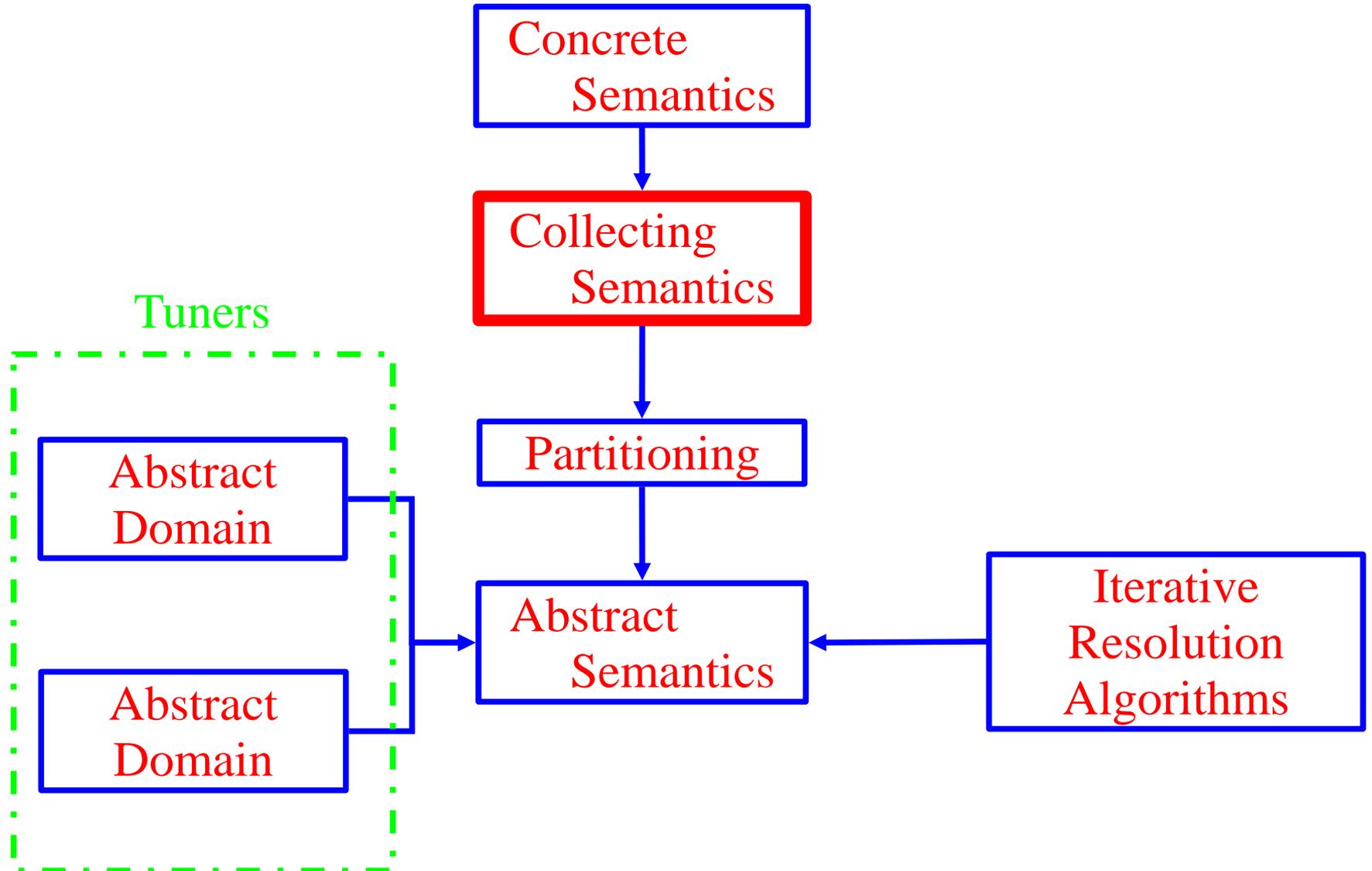
Control flow graph: $\textcircled{i} \xrightarrow{\text{op}} \textcircled{j}$

Operational semantics: $\langle \textcircled{i}, \varepsilon \rangle \rightarrow \langle \textcircled{j}, \llbracket \text{op} \rrbracket \varepsilon \rangle$

Semantics of op



Methodology



Collecting Semantics

The collecting semantics is the **set of observable behaviours** in the operational semantics. It is the starting point of any analysis design.

- The set of all descendants of the initial state
- The set of all descendants of the initial state that can reach a final state
- The set of all finite traces from the initial state
- The set of all finite and infinite traces from the initial state
- etc.

Which Collecting Semantics?

- Buffer overrun, division by zero, arithmetic overflows: **state properties**
- Deadlocks, un-initialized variables: **finite trace properties**
- Loop termination: **finite and infinite trace properties**

State properties

The set of descendants of the initial state s_0 :

$$S = \{s \mid s_0 \rightarrow \dots \rightarrow s\}$$

Theorem: $F : (\wp(\Sigma), \subseteq) \rightarrow (\wp(\Sigma), \subseteq)$

$$F(S) = \{s_0\} \cup \{s' \mid \exists s \in S: s \rightarrow s'\}$$

$$S = \text{lfp } F$$

Example

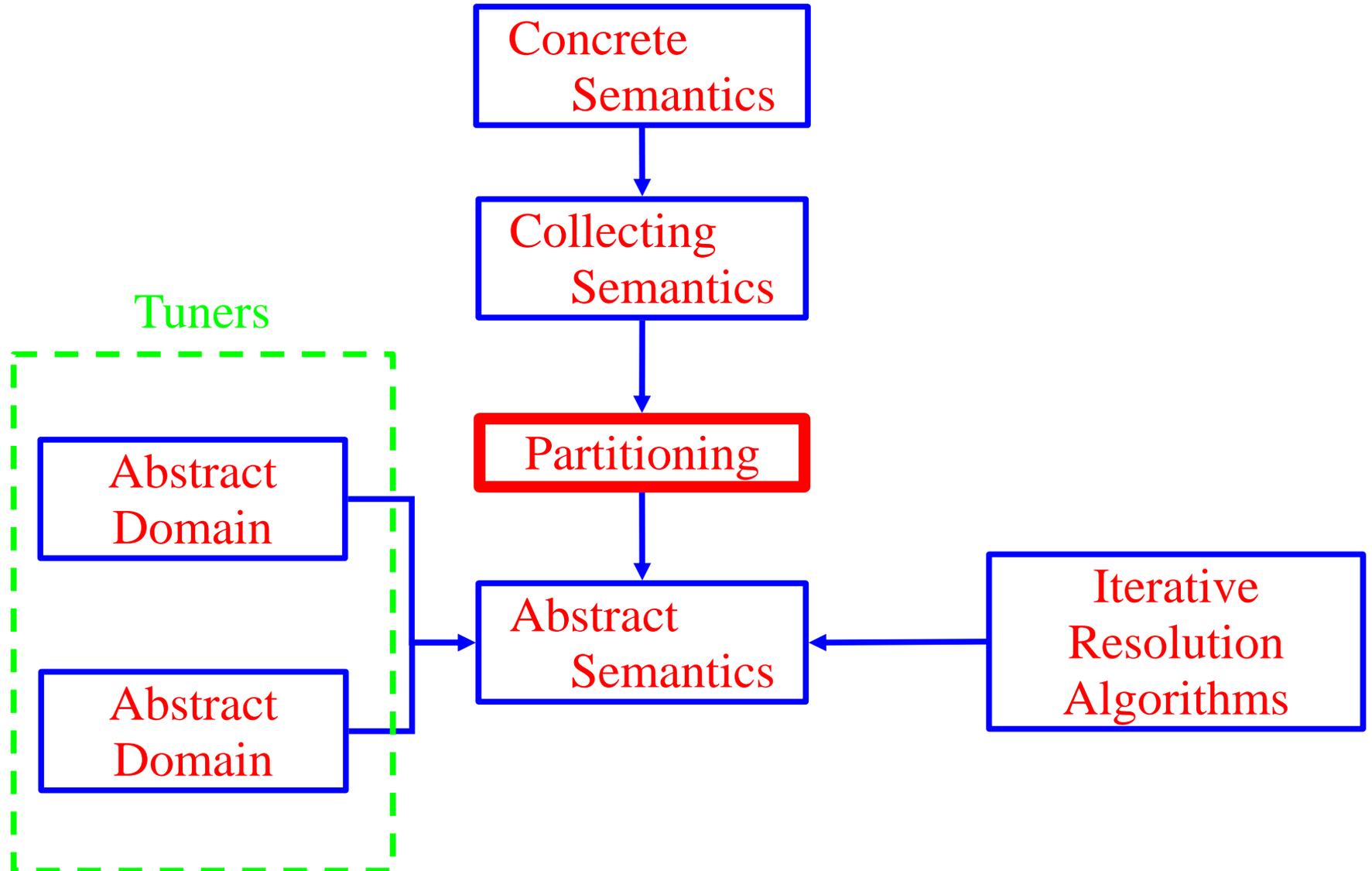
```
1:  n = 0;  
2:  while n < 1000 do  
3:    n = n + 1;  
4:  end  
5:  exit
```

$$S = \{ \langle 1, n \Rightarrow \Omega \rangle, \langle 2, n \Rightarrow 0 \rangle, \langle 3, n \Rightarrow 0 \rangle, \langle 4, n \Rightarrow 1 \rangle, \langle 2, n \Rightarrow 1 \rangle, \dots, \langle 5, n \Rightarrow 1000 \rangle \}$$

Computation

- $F_0 = \emptyset$
- $F_1 = \{ \langle 1, n \Rightarrow \Omega \rangle \}$
- $F_2 = \{ \langle 1, n \Rightarrow \Omega \rangle, \langle 2, n \Rightarrow 0 \rangle \}$
- $F_3 = \{ \langle 1, n \Rightarrow \Omega \rangle, \langle 2, n \Rightarrow 0 \rangle, \langle 3, n \Rightarrow 0 \rangle \}$
- $F_4 = \{ \langle 1, n \Rightarrow \Omega \rangle, \langle 2, n \Rightarrow 0 \rangle, \langle 3, n \Rightarrow 0 \rangle, \langle 4, n \Rightarrow 1 \rangle \}$
- ...

Methodology



Partitioning

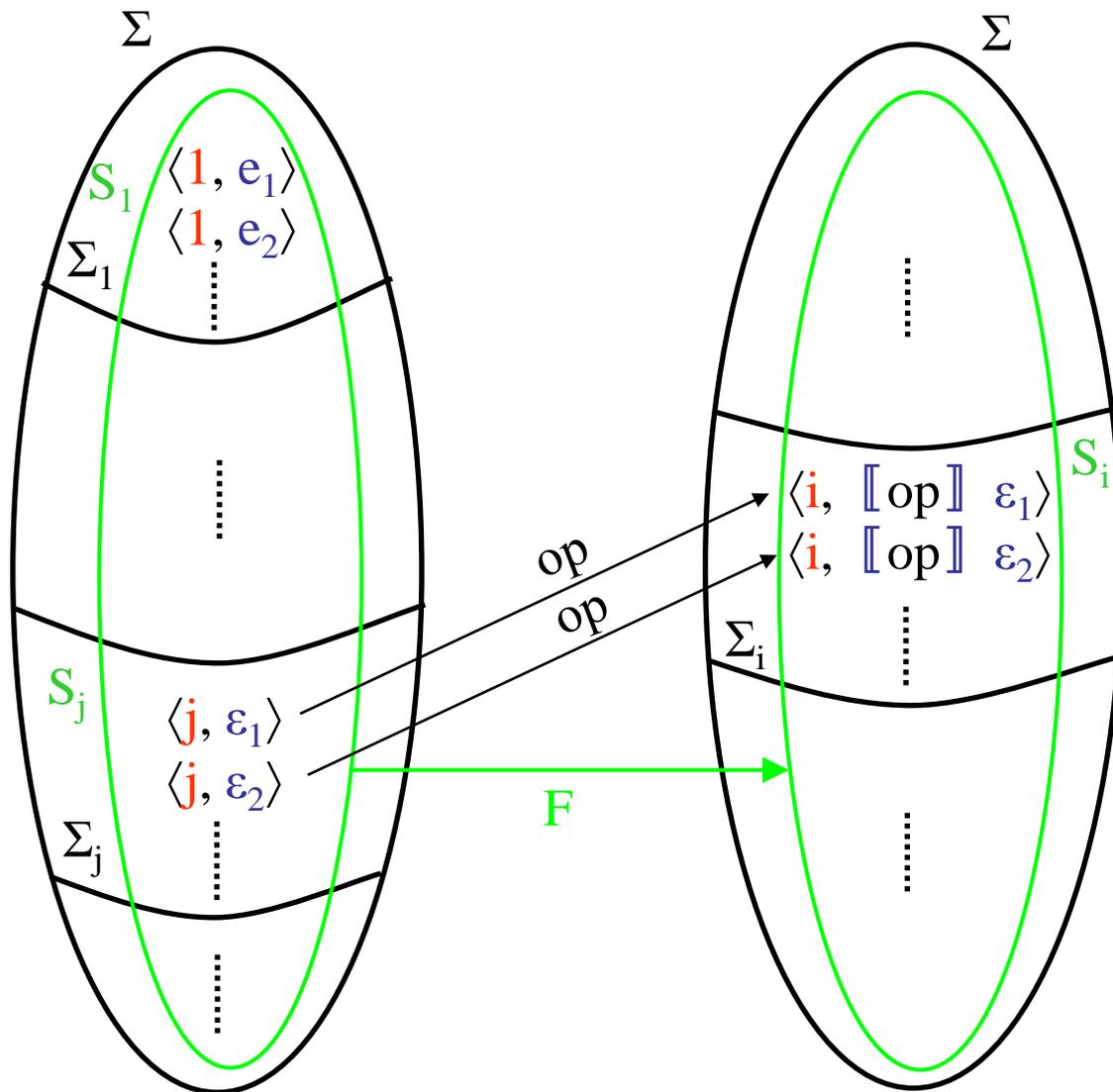
We partition the set S of states w.r.t. program points:

- $\Sigma = \Sigma_1 \oplus \Sigma_2 \oplus \dots \oplus \Sigma_n$
- $\Sigma_i = \{ \langle k, \epsilon \rangle \in \Sigma \mid k = i \}$
- $F(S_1, \dots, S_n)_0 = \{ s_0 \}$
- $F(S_1, \dots, S_n)_i = \{ s' \in S_i \mid \exists j \exists s \in S_j : s \rightarrow s' \}$

i.e.

$$F(S_1, \dots, S_n)_i = \{ \langle i, \llbracket \text{op} \rrbracket \epsilon \rangle \mid \textcircled{j} \xrightarrow{\text{op}} \textcircled{i} \in \text{CFG}(P) \}$$

Illustration



Semantic Equations

- Notation: E_i = set of environments at program point i
- System of semantic equations:

$$E_i = \bigcup \{ \llbracket \text{op} \rrbracket E_j \mid \textcircled{j} \xrightarrow{\text{op}} \textcircled{i} \in \text{CFG}(P) \}$$

- Solution of the system $S = \text{lfp } F$

Example

```
1:  n = 0;  
2:  while n < 1000 do  
3:    n = n + 1;  
4:  end  
5:  exit
```

$$E_1 = \{n \Rightarrow \Omega\}$$

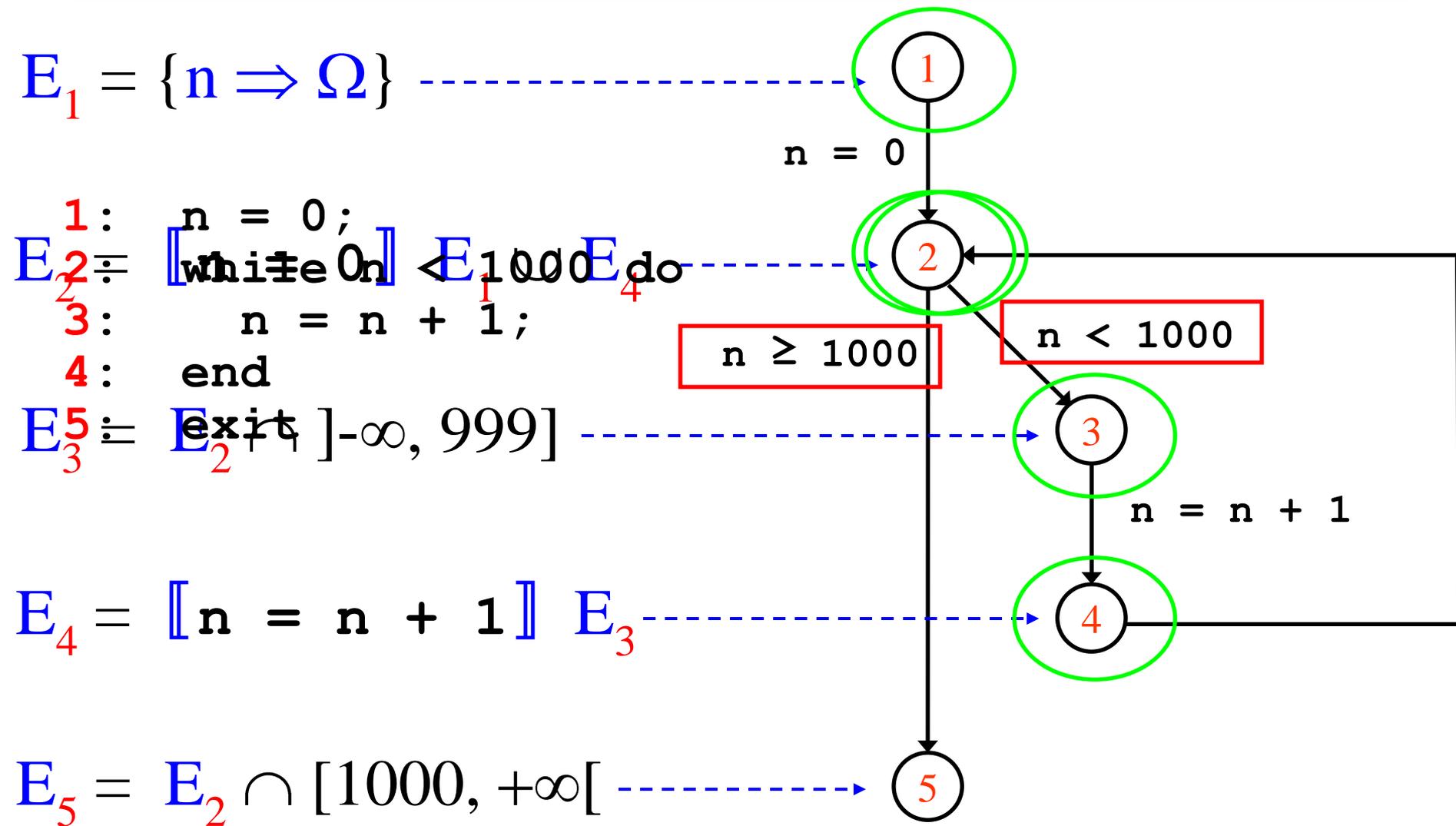
$$E_2 = \llbracket n = 0 \rrbracket E_1 \cup E_4$$

$$E_3 = E_2 \cap]-\infty, 999]$$

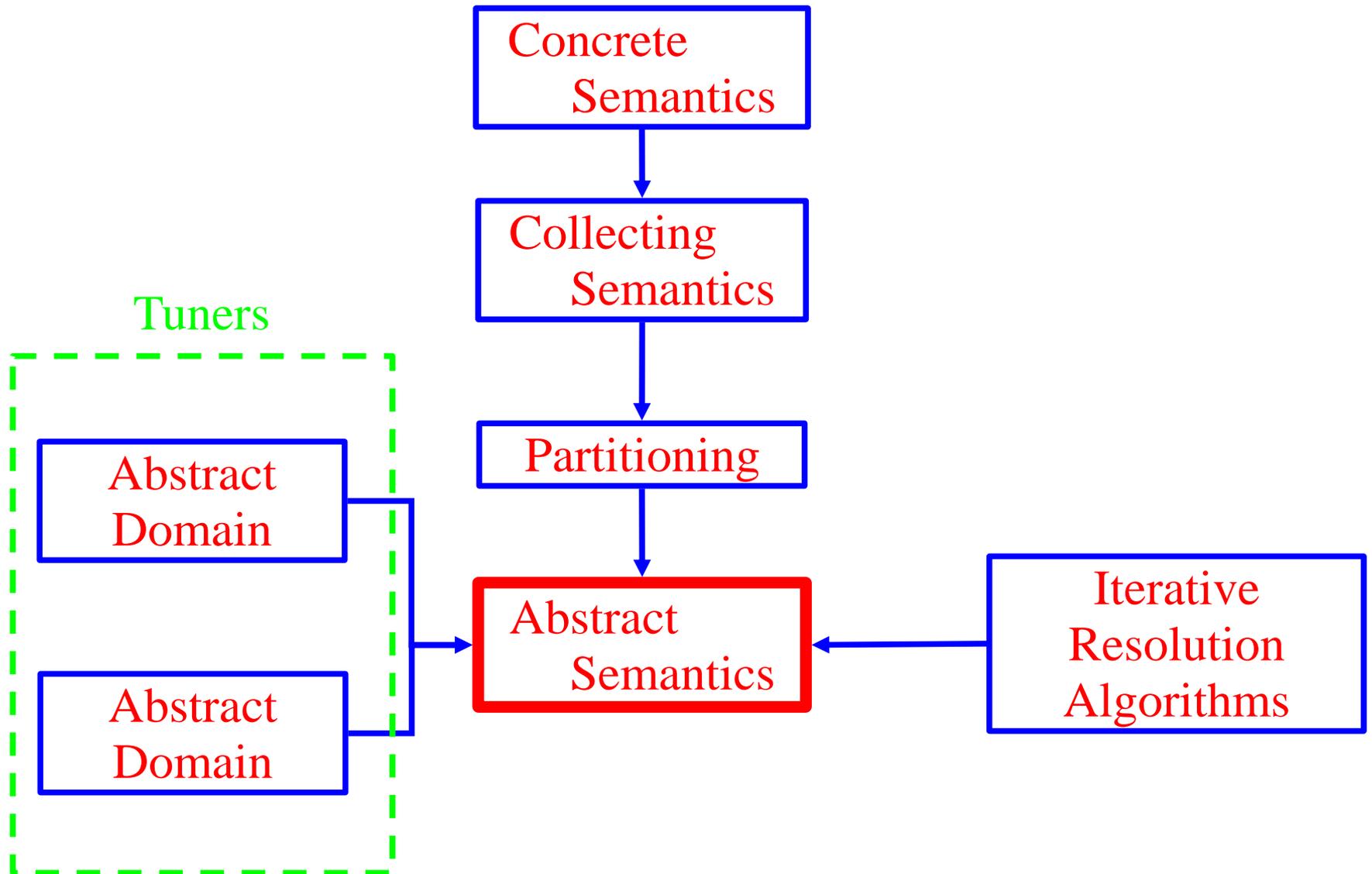
$$E_4 = \llbracket n = n + 1 \rrbracket E_3$$

$$E_5 = E_2 \cap [1000, +\infty[$$

Example



Methodology



Approximation

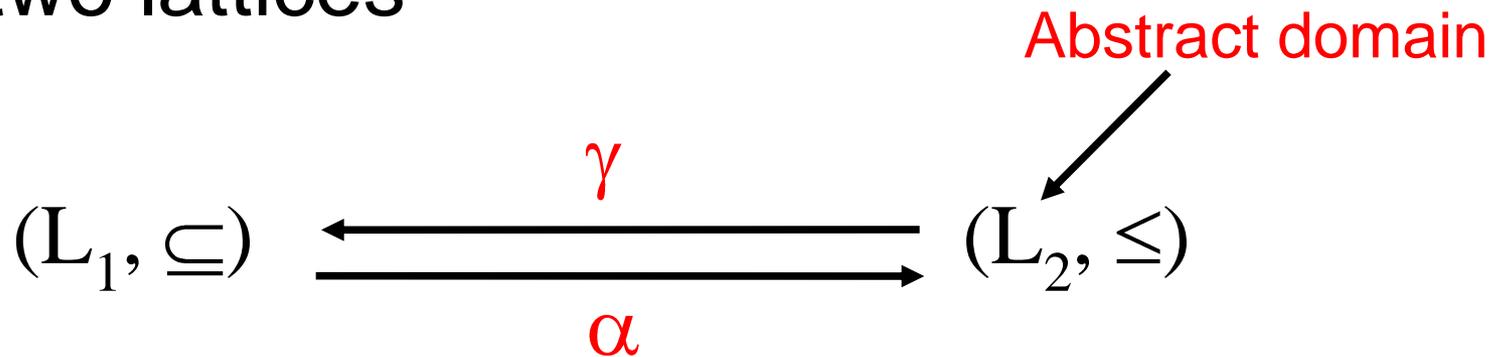
Problem: Compute a sound approximation $S^\#$
of S

$$S \subseteq S^\#$$

Solution: Galois connections

Galois Connection

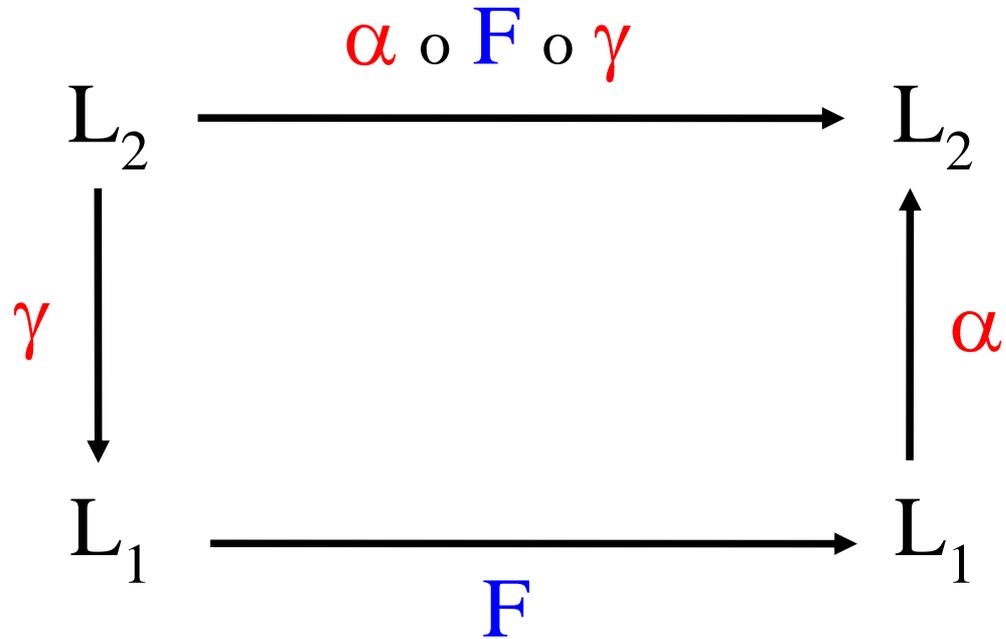
L_1, L_2 two lattices



- $\forall x \forall y : \alpha(x) \leq y \Leftrightarrow x \subseteq \gamma(y)$
- $\forall x \forall y : x \subseteq \gamma \circ \alpha(x) \ \& \ \alpha \circ \gamma(y) \leq y$

Fixpoint Approximation

Abstract
computation



Concrete
computation

Theorem:

$$\text{lfp } F \subseteq \gamma (\text{lfp } \alpha \circ F \circ \gamma)$$

Abstracting the Collecting Semantics

- Find a Galois connection:

$$(\wp(\Sigma), \subseteq) \begin{array}{c} \xleftarrow{\gamma} \\ \xrightarrow{\alpha} \end{array} (\Sigma^\#, \leq)$$

- Find a function: $\alpha \circ F \circ \gamma \leq F^\#$

Abstract Algebra

- Notation: \mathbf{E} the set of all environments
- Galois connection:

$$(\wp(\mathbf{E}), \subseteq) \begin{array}{c} \xleftarrow{\gamma} \\ \xrightarrow{\alpha} \end{array} (\mathbf{E}^\#, \leq)$$

- \cup, \cap approximated by $\cup^\#, \cap^\#$
- Semantics $\llbracket \mathbf{op} \rrbracket$ approximated by $\llbracket \mathbf{op} \rrbracket^\#$

$$\alpha \circ \llbracket \mathbf{op} \rrbracket \circ \gamma \subseteq \llbracket \mathbf{op} \rrbracket^\#$$

Abstract Semantic Equations

```
1:  n = 0;  
2:  while n < 1000 do  
3:    n = n + 1;  
4:  end;  
5:  exit;
```

$$E_1^\# = \alpha (\{n \Rightarrow \Omega\})$$

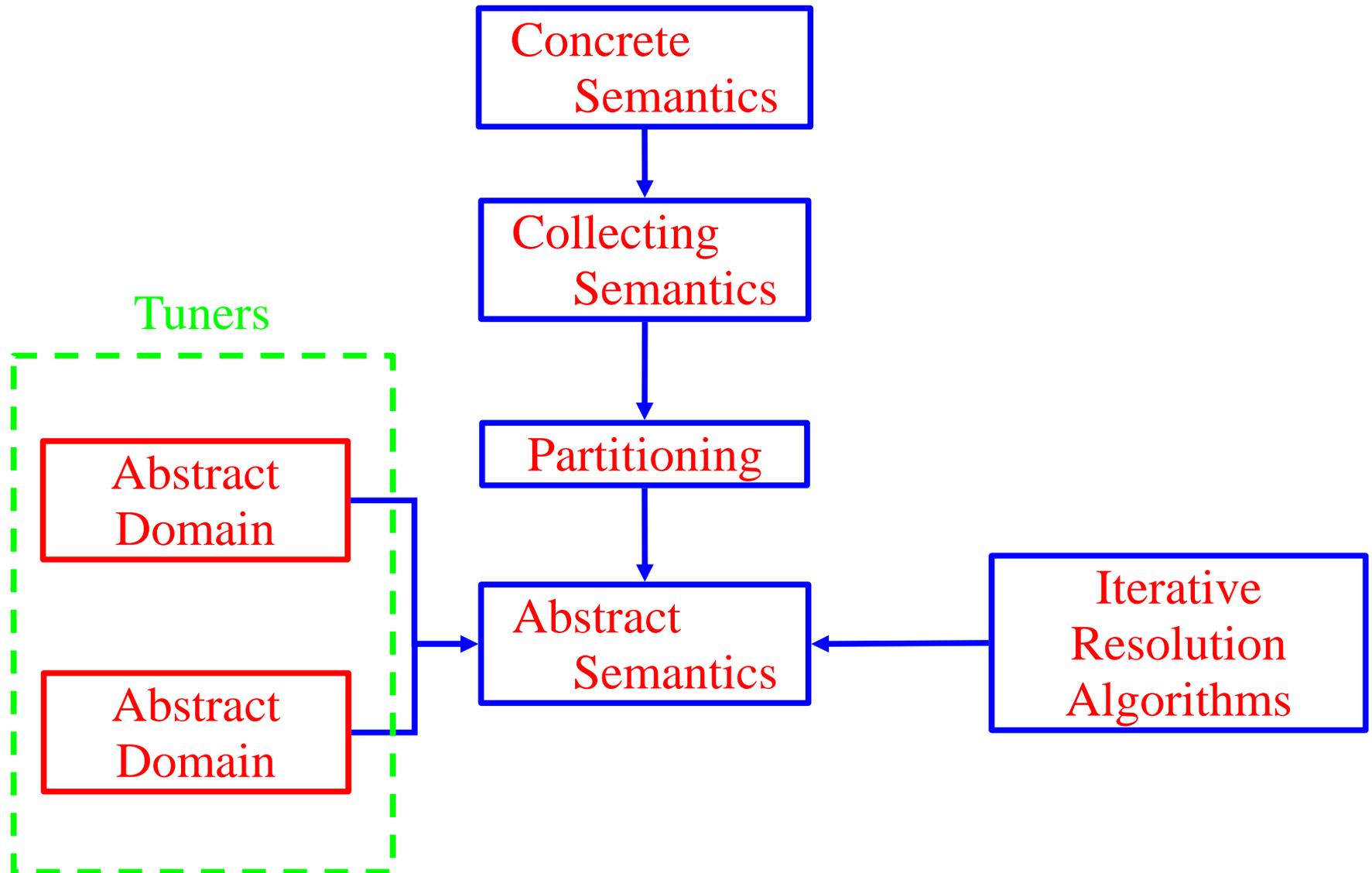
$$E_2^\# = \llbracket n = 0 \rrbracket \# E_1^\# \cup^\# E_4^\#$$

$$E_3^\# = E_2^\# \cap^\# \alpha (]-\infty, 999])$$

$$E_4^\# = \llbracket n = n + 1 \rrbracket \# E_3^\#$$

$$E_5^\# = E_2^\# \cap^\# \alpha ([1000, +\infty[)$$

Methodology



Abstract Domains

Various kinds of approximations:

- Signs (non relational)

$$x \Rightarrow +, y \Rightarrow -, \dots$$

- Intervals (nonrelational):

$$x \Rightarrow [3, 9], y \Rightarrow [-23, 4], \dots$$

- Polyhedra (relational):

$$x + y - 2z \leq 10, \dots$$

- Difference-bound matrices (weakly relational):

$$y - x \leq 5, z - y \leq 10, \dots$$

Example: intervals

```
1:  n = 0;  
2:  while n < 1000 do  
3:    n = n + 1;  
4:  end  
5:  exit
```

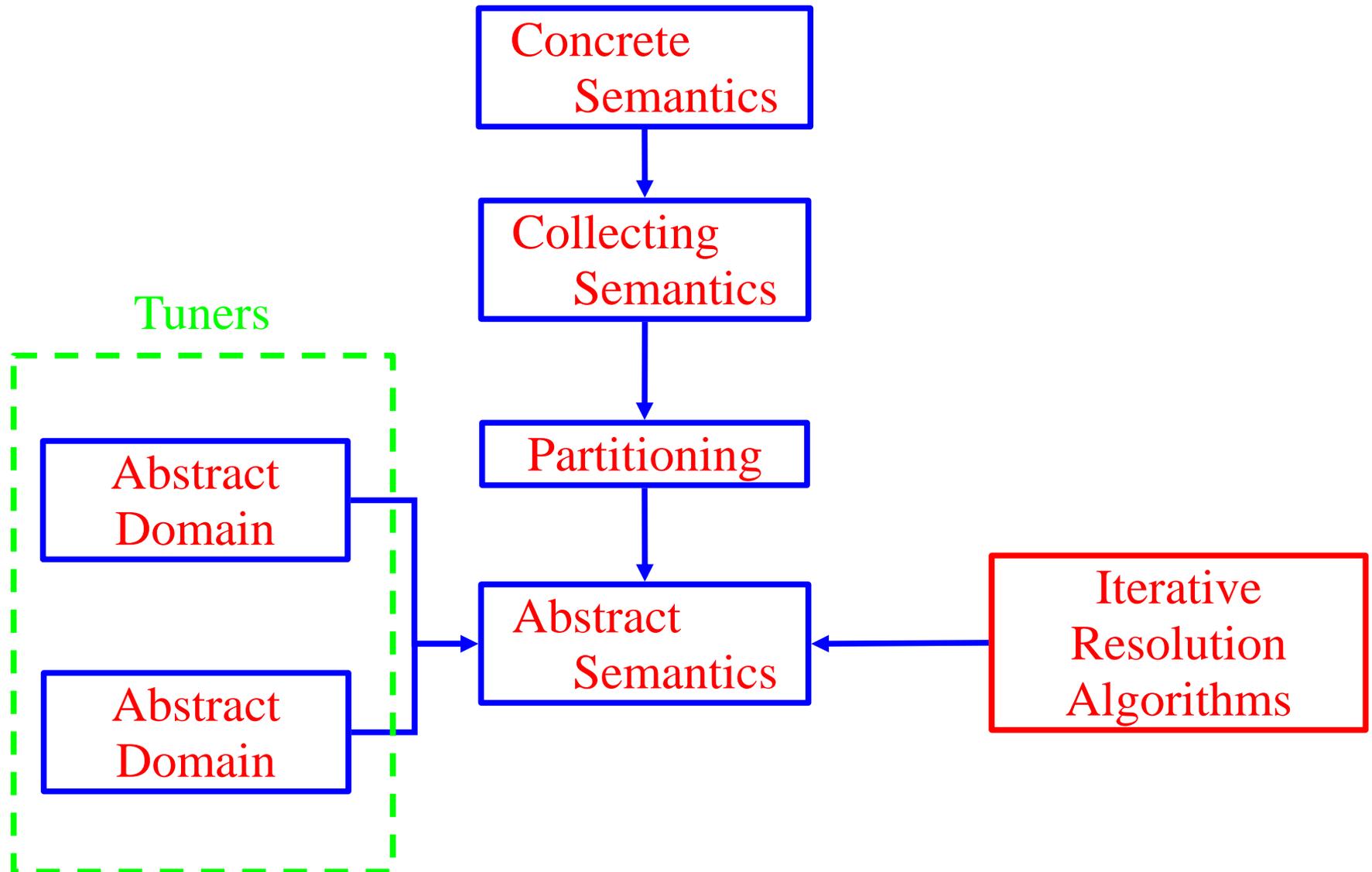
- Iteration 1: $E_2^\# = [0, 0]$
- Iteration 2: $E_2^\# = [0, 1]$
- Iteration 3: $E_2^\# = [0, 2]$
- Iteration 4: $E_2^\# = [0, 3]$
- ...

Problem

How to cope with lattices of infinite height?

Solution: automatic extrapolation operators

Methodology



Widening operator

Lattice (L, \leq) : $\nabla : L \times L \rightarrow L$

- Abstract union operator:

$$\forall x \forall y : x \leq x \nabla y \ \& \ y \leq x \nabla y$$

- Enforces convergence: $(x_n)_{n \geq 0}$

$$\begin{cases} y_0 = x_0 \\ y_{n+1} = y_n \nabla x_{n+1} \end{cases}$$

$(y_n)_{n \geq 0}$ is ultimately stationary

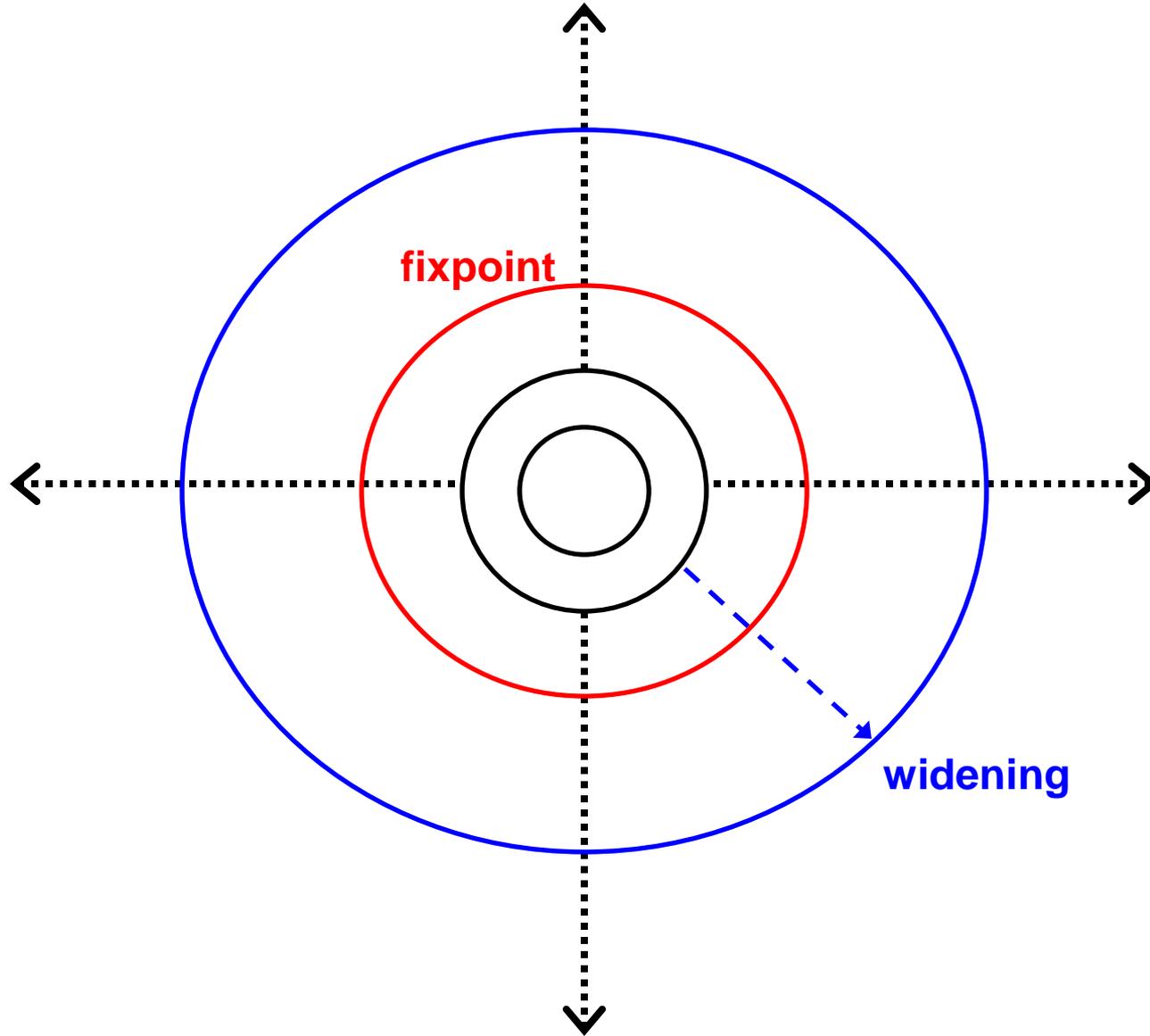
Widening of intervals

$$[a, b] \nabla [a', b']$$

- If $a \leq a'$ then a else $-\infty$
- If $b' \leq b$ then b else $+\infty$

↪ Open unstable bounds (jump over the fixpoint)

Widening and Fixpoint



Iteration with widening

```
1:  n = 0;  
2:  while n < 1000 do  
3:    n = n + 1;  
4:  end  
5:  exit
```

$$(E_2^\#)_{n+1} = (E_2^\#)_n \nabla \left(\llbracket n = 0 \rrbracket \# (E_1^\#)_n \cup^\# (E_4^\#)_n \right)$$

Iteration 1 (union): $E_2^\# = [0, 0]$

Iteration 2 (union): $E_2^\# = [0, 1]$

Iteration 3 (widening): $E_2^\# = [0, +\infty] \Rightarrow$ stable

Imprecision at loop exit

```
1:  n = 0;  
2:  while n < 1000 do  
3:    n = n + 1;  
4:  end  
5:  exit; t[n] = 0; // t has 1500 elements
```

False positive!!!



- $E_5^\# = [1000, +\infty[$

Narrowing operator

Lattice (L, \leq) : $\Delta : L \times L \rightarrow L$

- Abstract intersection operator:

$$\forall x \forall y : x \cap y \leq x \Delta y$$

- Enforces convergence: $(x_n)_{n \geq 0}$

$$\begin{cases} y_0 & = x_0 \\ y_{n+1} & = y_n \Delta x_{n+1} \end{cases}$$

$(y_n)_{n \geq 0}$ is ultimately stationary

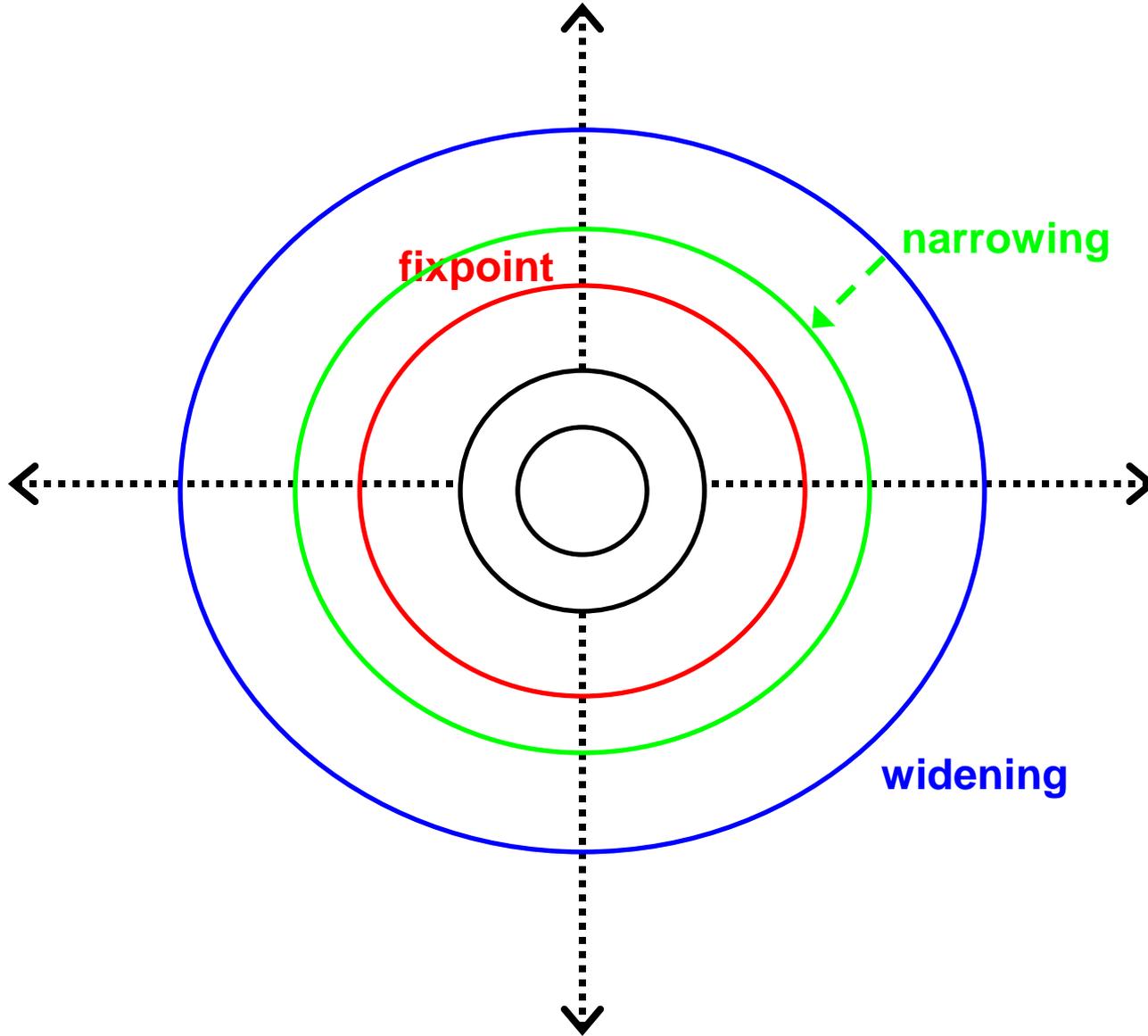
Narrowing of intervals

$$[a, b] \Delta [a', b']$$

- If $a = -\infty$ then a' else a
- If $b = +\infty$ then b' else b

↪ Refines open bounds

Narrowing and Fixpoint



Iteration with narrowing

```
1:  n = 0;  
2:  while n < 1000 do  
3:    n = n + 1;  
4:  end  
5:  t[n] = 0;
```

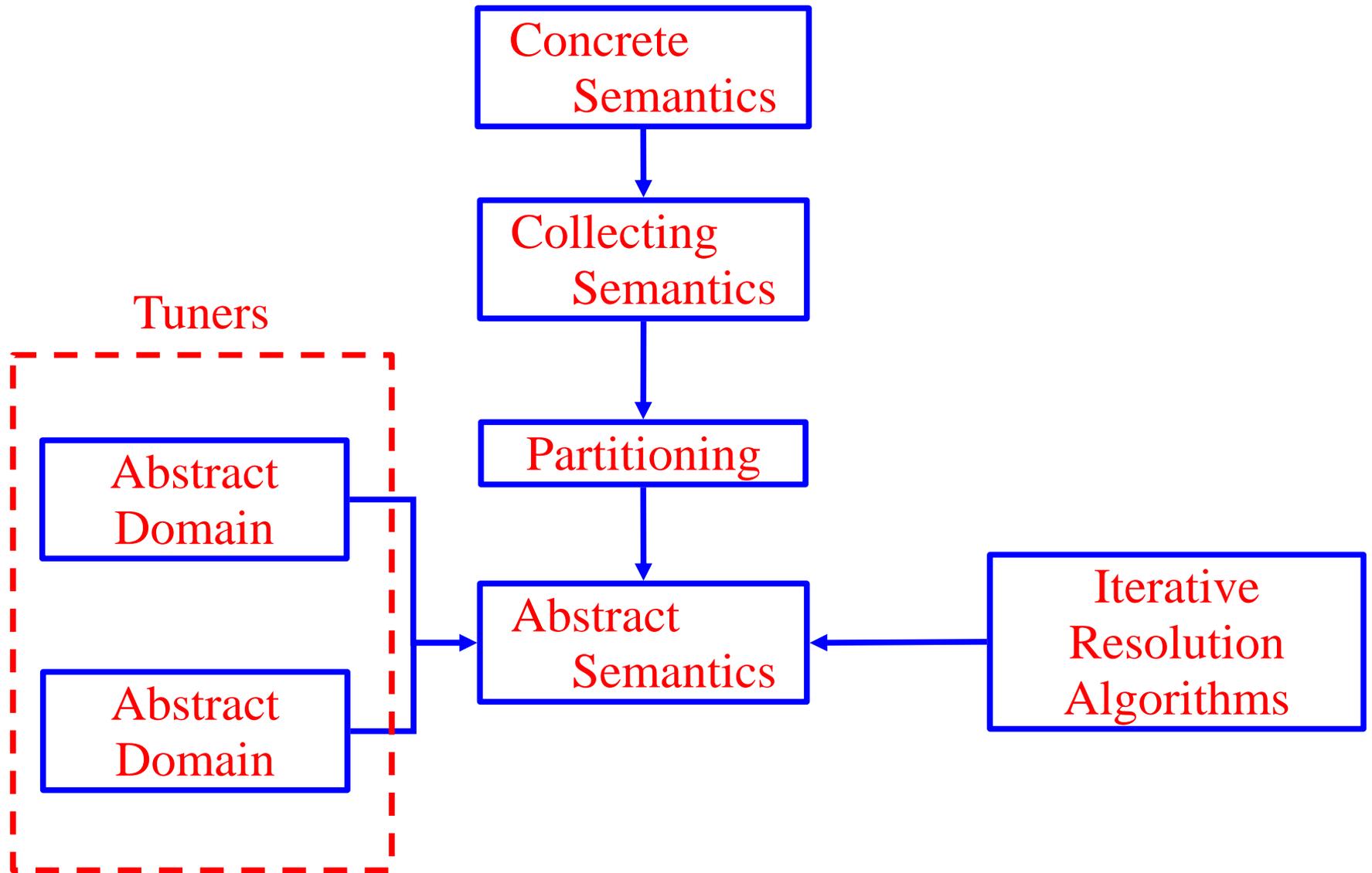
$$(E_2^\#)_{n+1} = (E_2^\#)_n \Delta \left(\llbracket n = 0 \rrbracket \# (E_1^\#)_n \cup^\# (E_4^\#)_n \right)$$

Beginning of iteration: $E_2^\# = [0, +\infty[$

Iteration 1: $E_2^\# = [0, 1000] \Rightarrow$ stable

Consequence: $E_5^\# = [1000, 1000]$

Methodology



Tuning the abstract domains

```
1:  n = 0;  
2:  k = 0;  
3:  while n < 1000 do  
4:    n = n + 1;  
5:    k = k + 1;  
6:  end  
7:  exit
```

- Intervals:

$$E_4^\# = \langle n \Rightarrow [0, 1000], k \Rightarrow [0, +\infty[\rangle$$

- Convex polyhedra:

$$E_4^\# = \langle 0 \leq n \leq 1000, 0 \leq k \leq 1000, n - k = 0 \rangle$$

Annotated Bibliography

References

- The historic paper:
 - **Patrick Cousot & Radhia Cousot**. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Conference Record of the Fourth Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 238—252, Los Angeles, California, 1977. ACM Press, New York, NY, USA.
- Accessible introductions to the theory:
 - **Patrick Cousot**. Semantic foundations of program analysis. In S.S. Muchnick and N.D. Jones, editors, *Program Flow Analysis: Theory and Applications*, Ch. 10, pages 303—342, Prentice-Hall, Inc., Englewood Cliffs, New Jersey, U.S.A., 1981.
 - **Patrick Cousot & Radhia Cousot**. Abstract interpretation and application to logic programs. *Journal of Logic Programming*, 13(2—3):103—179, 1992.
- Beyond Galois connections, a presentation of relaxed frameworks:
 - **Patrick Cousot & Radhia Cousot**. Abstract interpretation frameworks. *Journal of Logic and Computation*, 2(4):511—547, August 1992.
- A thorough description of a static analyzer with all the proofs (difficult to read):
 - **Patrick Cousot**. The Calculational Design of a Generic Abstract Interpreter. Course notes for the NATO International Summer School 1998 on Calculational System Design. Marktobendorf, Germany, 28 July—9 august 1998, organized by F.L. Bauer, M. Broy, E.W. Dijkstra, D. Gries and C.A.R. Hoare.

References

- The abstract domain of intervals:
 - **Patrick Cousot & Radhia Cousot**. Static Determination of Dynamic Properties of Programs. In B. Robinet, editor, *Proceedings of the second international symposium on Programming*, Paris, France, pages 106—130, april 13-15 1976, Dunod, Paris.
- The abstract domain of convex polyhedra:
 - **Patrick Cousot & Nicolas Halbwachs**. Automatic discovery of linear restraints among variables of a program. In *Conference Record of the Fifth Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 84—97, Tucson, Arizona, 1978. ACM Press, New York, NY, USA.
- Weakly relational abstract domains:
 - **Antoine Miné**. The Octagon Abstract Domain. In Analysis, Slicing and Transformation (part of Working Conference on Reverse Engineering), October 2001, IEEE, pages 310-319.
 - **Antoine Miné**. A New Numerical Abstract Domain Based on Difference-Bound Matrices. In Program As Data Objects II, May 2001, LNCS 2053, pages 155-172.
- Classical data flow analysis:
 - **Steven Muchnick**. *Advanced Compiler Design and Implementation*. Morgan Kaufmann, 1997.