

# Lecture Notes on Abstract Interpretation

André Platzer

Carnegie Mellon University || Karlsruhe Institute of Technology  
Lecture 28

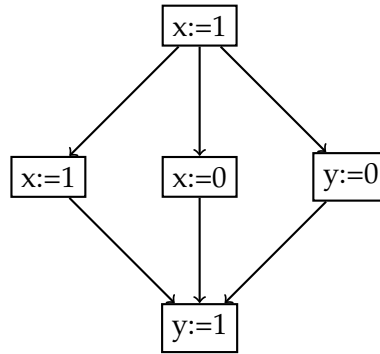
## 1 Introduction

Simple examples of abstract interpretation type ideas in more classical situations include sign abstraction of values into  $\{-, 0, +, ?\}$  or abstraction of values by remainders mod 4 [[WM95](#), Chapter 10]. The focus is on abstract interpretation uses and their connections with compilation and program analysis. This is a wide field and easily the topic of a whole semester. More information on abstract interpretation can be found in [[CC92](#), [CC77](#), [CC79](#)] and [[WM95](#), Chapter 10].

## 2 Abstract Interpretation by Example

Abstract interpretation generalizes the theory of monotone frameworks and dataflow analysis to a general principle of analyzing programs by defining an abstract semantics for it [[CC92](#), [CC77](#), [CC79](#), [WM95](#)]. In order to show the principle of abstract interpretation, without having to dig too much into the details, we consider an example where we abstractly interpret a program but still keep using monotone frameworks.

Suppose we want to check the property whether a variable  $x$  may be 0, which is a principle that can be useful for null pointer exception tests. As domain  $L$  for this we just choose the Boolean lattice  $\{true, false\}$ . The operator  $\sqcup$  is just logical disjunction ( $\vee$ ). The flow relation is the forward control flow. Initialization is *false*, because we now assume that pointers cannot be null unless they are assigned null. Transfer functions at the nodes make sense to choose from the constant functions *true*, *false* and the identity function *id*.



The transfer function for  $x := 1$  will be *true*, the one for  $x := 0$  will be *false* and the transfer function for  $y := \dots$  will be *id*. By fixed-point iteration on the above example we find that  $x = 1$  is possible after the program terminates. For a must analysis, instead, we would get that  $x = 1$  is not necessarily true after the program terminates.

For multiple variables, we can choose a Cartesian product  $\{true, false\}^n$  of the Boolean lattice and use projections to coordinates as further transfer functions for copying the value for  $y$  over to  $x$  at a move  $x := y$ .

Another example is an abstract interpretation that performs general analysis for constant propagation. The property space has the form  $\{x = \perp, x = ?\} \cup \{x = v : v \in \mathbb{Z}\}$ , where  $\perp$  means is the bottom of the semilattice for undefined,  $x = ?$  means that  $x$  has different possible nondeterministic values and  $x = v$  for a number  $v$  means that we can be certain that  $x$  will always have value  $v$  at this program point. Let's look at an example. We initialize with no information ( $\perp$ ) at all points, except the program init block, where we start with a nondeterministic initial value  $i = ?$ :

```

{ i=?, j=?, k=? }
i = 5; j = 0; k = 0;
{ i=⊥, j=⊥, k=⊥ }
while ( j <= i ) {
  { i=⊥, j=⊥, k=⊥ }
  i = i + 2; k = k + j; j = j + 1
  { i=⊥, j=⊥, k=⊥ }
  i = i - 2
  { i=⊥, j=⊥, k=⊥ }
}
{ i=⊥, j=⊥, k=⊥ }

```

Now we can execute the first line in the abstract semantics and then enter the loop in the abstract semantics and execute the loop body once

```

{ i=?, j=?, k=? }
i = 5; j = 0; k = 0;
{ i=5, j=0, k=0 }
while ( j <= i ) {
  { i=5, j=0, k=0 }
}

```

```

    i = i + 2; k = k + j; j = j + 1
    {i=7,j=1,k=0}
    i = i - 2
    {i=5,j=1,k=0}
  }
  {i=⊥,j=⊥,k=⊥}

```

With those abstract values, we will repeat the loop, but we have to merge the previous information  $\{i=5,j=0,k=0\}$  from before the loop with the current information  $\{i=5,j=1,k=0\}$  from the end of the loop body and find a joint representation in the property space lattice by the  $\sqcup$  operator, giving  $\{i=5,j=?,k=0\}$  to keep the common  $i=5,k=0$  but nondeterministically overapproximate  $j$  with its multiple possible values. Then we execute the loop body again

```

  {i=?,j=?,k=?}
  i = 5; j = 0; k = 0;
  {i=5,j=0,k=0}
  while (j <= i) {
    {i=5,j=?,k=0}
    i = i + 2; k = k + j; j = j + 1
    {i=7,j=?,k=?}
    i = i - 2
    {i=5,j=?,k=?}
  }
  {i=⊥,j=⊥,k=⊥}

```

Again, merging the property values by the  $\sqcup$  operator and executing the loop body gives

```

  {i=?,j=?,k=?}
  i = 5; j = 0; k = 0;
  {i=5,j=0,k=0}
  while (j <= i) {
    {i=5,j=?,k=?}
    i = i + 2; k = k + j; j = j + 1
    {i=7,j=?,k=?}
    i = i - 2
    {i=5,j=?,k=?}
  }
  {i=5,j=?,k=?}

```

Here the property value at the loop entry didn't change, so we can propagate to the loop exit and the analysis terminates. Now we know, as good as our abstract semantics could represent, what values the variables can have at the various program points.

### 3 Abstract Interpretation by Example

Consider the following simple program

```

0
1 x = 1
2
3 while (x < 1000) {
4
5     x = x + 1
6
7 }
8
9 y = x

```

A run in the concrete semantics of the above program would start with the concrete state  $x = \perp, y = \perp$  where the initial value of  $x, y$  in line 0 is unknown. The program would do 999 iterations through the loop after which it terminates with the state  $y = x = 1000$ . Concrete execution just does not help much for static analysis of programs in general, because we won't know the dynamic data until runtime.

Instead, let us consider an abstract run in an abstract semantics where variables take on intervals as values (due to Cousot and Cousot [CC77]):

$$L = \{[a, b] : a, b \in \mathbb{N} \cup \{+\infty, -\infty\}\}$$

To unify notation, we write  $[-\infty, 5]$  for the left-open interval  $(-\infty, 5]$  here. Now a run of the above program in the interval abstract domain gives after 1 iteration

```

0 {x = [-∞, ∞], y = [-∞, ∞]}
1 x = 1
2 {x = [1, 1], y = [-∞, ∞]}
3 while (x < 1000) {
4     {x = [1, 1], y = [-∞, ∞]}
5     x = x + 1
6     {x = [2, 2], y = [-∞, ∞]}
7 }
8
9 y = x

```

and after 2 iterations where the information from lines 2 and 6 is merged to line 4

```

0 {x = [-∞, ∞], y = [-∞, ∞]}
1 x = 1
2 {x = [1, 1], y = [-∞, ∞]}
3 while (x < 1000) {
4     {x = [1, 2], y = [-∞, ∞]}
5     x = x + 1

```

```

6      {x = [2, 3], y = [-∞, ∞]}
7  }
8
9  y = x

```

and after 3 iterations where the information from lines 2 and 6 is merged to line 4

```

0  {x = [-∞, ∞], y = [-∞, ∞]}
1  x = 1
2  {x = [1, 1], y = [-∞, ∞]}
3  while (x < 1000) {
4      {x = [1, 3], y = [-∞, ∞]}
5      x = x + 1
6      {x = [2, 4], y = [-∞, ∞]}
7  }
8
9  y = x

```

We could keep on iterating, but this still takes an awfully large number of iterations to figure out, since the loop count is 1000. If the bound is not computable statically, we do not even know how often to iterate.

We can iterate until we reach a fixpoint. And we can also speed up convergence by jumping ahead in the lattice using a widening operator  $\nabla : L \times L \rightarrow L$  that combines information from two lattice elements to a joint overapproximation of the two. For intervals let us jump ahead to  $\pm\infty$  whenever our interval bounds are not inclusive:

$$[a, b] \nabla [a', b'] := \left[ \begin{cases} a & \text{if } a \leq a' \\ -\infty & \text{otherwise} \end{cases}, \begin{cases} b & \text{if } b' \leq b \\ +\infty & \text{otherwise} \end{cases} \right]$$

So in the 4th iteration, instead of doing a standard image iteration, let us widening for computing line 4 from the previous two values  $[1, 3] \nabla [1, 4] = [1, \infty]$ :

```

0  {x = [-∞, ∞], y = [-∞, ∞]}
1  x = 1
2  {x = [1, 1], y = [-∞, ∞]}
3  while (x < 1000) {
4      {x = [1, ∞], y = [-∞, ∞]}
5      x = x + 1
6      {x = [2, ∞], y = [-∞, ∞]}
7  }
8
9  y = x

```

In iteration 5, we obtain precise information by intersecting with the loop guards

```

0  {x = [-∞, ∞], y = [-∞, ∞]}
1  x = 1

```

```

2  {x = [1, 1], y = [-∞, ∞]}
3  while (x < 1000) {
4      {x = [1, 999], y = [-∞, ∞] since x = [1, ∞] ∩ [1, 999] = [1, 999]}
5      x = x + 1
6      {x = [2, 1000], y = [-∞, ∞]}
7  }
8  {x = [1000, 1000], y = [-∞, ∞] since x = [2, 1000] ∩ [1000, ∞] = [1000, 1000]}
9  y = x
10 {x = [1000, 1000], y = [1000, 1000]}

```

What we want the widening operator  $\nabla$  to satisfy is that it is like a union ( $\cup$ ) but could be a bigger element of the lattice:

$$x \leq x \nabla y \quad y \leq x \nabla y$$

We also want iterated uses of the widening operator to become a fixpoint eventually. That is

$$x_0 \nabla x_1 \nabla x_2 \nabla x_3 \nabla \dots$$

is a finite sequence, for any  $x_i \in L$ .

This seems very powerful and it is, as a framework for static program analysis. The particular abstract domain of intervals alone, however, is insufficient. A simple variation of the above example shows that the example is misleading and real programs more complicated:

```

0  {x = [-∞, ∞], y = [-∞, ∞]}
1  x = 1
2  {x = [1, 1], y = [-∞, ∞]}
3  y = 1
4  {x = [1, 1], y = [1, 1]}
5  while (x < 1000) {
6      {x = [1, 999], y = [1, ∞] since x = [1, ∞] ∩ [1, 999] = [1, 999]}
7      x = x + 1
8      {x = [2, 1000], y = [1, ∞]}
9      y = y + 1
10     {x = [2, 1000], y = [2, ∞]}
11 }
12 {x = [1000, 1000], y = [1, ∞] since x = [2, 1000] ∩ [1000, ∞] = [1000, 1000]}

```

This result is perfectly correct but rather useless as far as  $y$  is concerned, because it does not constrain the values of  $y$ , except for positivity, simply because  $y$  did not occur directly in the loop exit condition.

But the abstract interpretation framework still applies. Abstract domains that can handle the above example need correlations of variables, i.e., they need to capture variable correlations like  $0 \leq x - y \leq 1$ . Difference-bounds matrix [Min01] are a fast abstract domain for this purpose. General convex polyhedra can be useful too. This is possible but out of scope for this lecture. We only show the cheaper difference logic, where a

fast implementation are difference-bounds matrices. We adjoin an extra information of difference-bounds to the abstract domain  $L$ . As an optimization, we simply bootstrap from the converged values of  $x$  and  $y$  in our interval domain, since those would be found after some number of iterations anyhow. Better values are possible now, but not worse values. First, we need to figure out what the effect of the assignment  $x = x + 1$  will be on the abstract value  $l \leq x - y \leq u$  in the difference-bounds.

$$l \leq x - y \leq u \quad \xrightarrow{x:=x+1} \quad l + 1 \leq \underbrace{(x + 1)}_{x_{new}} - y \leq u + 1$$

Similarly for the assignment  $y = y + 1$ :

$$l \leq x - y \leq u \quad \xrightarrow{y:=y+1} \quad l - 1 \leq x - \underbrace{(y + 1)}_{y_{new}} - y \leq u - 1$$

After the first iteration, we get

```

0  {x = [-∞, ∞], y = [-∞, ∞], ∞ ≤ x - y ≤ ∞}
1  x = 1
2  {x = [1, 1], y = [-∞, ∞], ∞ ≤ x - y ≤ ∞}
3  y = 1
4  {x = [1, 1], y = [1, 1], 0 ≤ x - y ≤ 0}
5  while (x < 1000) {
6      {x = [1, 999], y = [1, ∞], 0 ≤ x - y ≤ 0}
7      x = x + 1
8      {x = [2, 1000], y = [1, ∞], 1 ≤ x - y ≤ 1}
9      y = y + 1
10     {x = [2, 1000], y = [2, ∞], 0 ≤ x - y ≤ 0}
11 }
12 {x = [1000, 1000], y = [1000, 1000], 0 ≤ x - y ≤ 0}

```

At which the fixpoint is reached immediately. Note that line 4 uses that the abstract value  $x = [1, 1], y = [1, 1]$  in the interval domain is communicated to the best corresponding constraint expressible as difference bounds:  $0 \leq x - y \leq 0$ :

$$x = [a, b], y = [c, d] \rightsquigarrow a - d \leq x - y \leq b - c$$

Hence, it is important that the abstract domains “talk” to each other. Conversely, in line 12, the abstract state  $0 \leq x - y \leq 0$  in the difference bounds can “talk” to the interval domain and synchronize to the best constraint that follows from the difference bounds in combination with the known individual interval bounds as follows:

$$x = [a, b], y = [c, d], l \leq x - y \leq u \rightsquigarrow x = [\max(a, c + l), \min(b, d + u)]$$

since  $l \leq x - y$  implies  $x \geq y + l$ , yet  $y \geq c$ . Similarly  $x - y \leq u$  implies  $x \leq y + u$  with  $y \leq d$ .

When widening was too aggressive, a dual operator called narrowing  $\Delta : L \times L \rightarrow L$  can be used as well. It is supposed to be like an intersection ( $\cap$ ) but could be bigger:

$$x \cap y \leq x \Delta y$$

We also want iterated uses of the widening operator to become a fixedpoint eventually. That is

$$x_0 \Delta x_1 \Delta x_2 \Delta x_3 \Delta \dots$$

is a finite sequence, for any  $x_i \in L$ .

## Quiz

1. What would happen if you had initialized  $x = [-\infty, \infty]$  everywhere to express that you don’t know initially what value  $x$  would have? Would that be the same as initializing  $x = \perp$ ?



2. Can abstract interpretation with interval bounds be used to perform analysis for possible occurrences of divisions by zero?
3. Show how abstract interpretation with the interval bounds domain can be used to perform array bounds checking optimizations.
4. To convince yourself under which circumstance narrowing  $\Delta$  may become necessary after widening, consider the example

```
0
1 x = 1
2
3 while (x < 1000) {
4
5     x = x + 1
6
7     if (x > 20) break;
8
9 }
```

5. Define a narrowing operator  $\Delta$  for the above case and show how to use it successfully.

## References

- [CC77] Patrick Cousot and Radhia Cousot. Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *POPL*, pages 238–252, 1977.
- [CC79] Patrick Cousot and Radhia Cousot. Systematic design of program analysis frameworks. In *POPL*, pages 269–282, 1979.
- [CC92] Patrick Cousot and Radhia Cousot. Abstract interpretation and application to logic programs. *J. Log. Program.*, 13(2&3):103–179, 1992.
- [Min01] Antoine Miné. A new numerical abstract domain based on difference-bound matrices. In Olivier Danvy and Andrzej Filinski, editors, *PADO*, volume 2053, pages 155–172. Springer, 2001.
- [WM95] Reinhard Wilhelm and Dieter Maurer. *Compiler Design*. Addison-Wesley, 1995. [doi:10.1007/978-3-642-59081-8](https://doi.org/10.1007/978-3-642-59081-8).