Data-flow Analysis / Abstract Interpretation

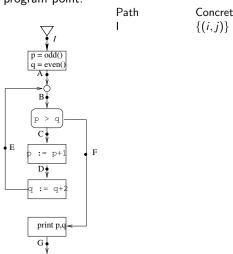
Deepak D'Souza and K. V. Raghavan

IISc

What is data-flow analysis

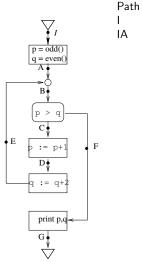
- "Computing 'safe' approximations to the set of values / behaviours arising dynamically at run time, statically or at compile time."
- Typically used by compiler writers to optimize running time of compiled code.
 - Constant propagation: Is the value of a variable constant at a particular program location.
 - Replace x := y + z by x := 17 during compilation.
- More recently, used for verifying properties of programs.

Collecting semantics of a program = set of (concrete) states occurring at each program point.



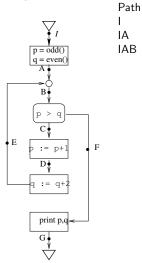
Concrete states $\{(i,j)\}$ (given)

Collecting semantics of a program = set of (concrete) states occurring at each program point.



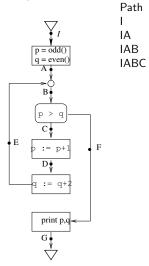
Concrete states $\{(i,j)\}$ (given) $\{(i,j)|i$ is odd, j is even $\}$

Collecting semantics of a program = set of (concrete) states occurring at each program point.



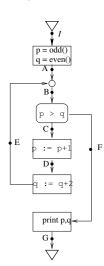
Concrete states $\{(i,j)\}$ (given) $\{(i,j)|i$ is odd, j is even $\}$ $\{(i,j)|i$ is odd, j is even $\}$

Collecting semantics of a program = set of (concrete) states occurring at each program point.



```
Concrete states \{(i,j)\} (given) \{(i,j)|i is odd, j is even\} \{(i,j)|i is odd, j is even\} \{(i,j)|i is odd, j is even, i>j\}
```

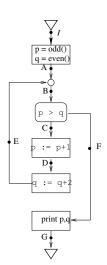
Collecting semantics of a program = set of (concrete) states occurring at each program point.



Path I IA IAB IABC IABCD

```
Concrete states  \{(i,j)\} \text{ (given)}   \{(i,j)|i \text{ is odd, } j \text{ is even} \}   \{(i,j)|i \text{ is odd, } j \text{ is even} \}   \{(i,j)|i \text{ is odd, } j \text{ is even, } i>j\}   \{(i,j)|i \text{ is even, } j \text{ is even, } i>j+1\}
```

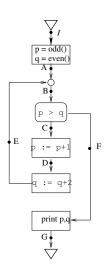
Collecting semantics of a program = set of (concrete) states occurring at each program point.



```
Path
I
IA
IAB
IABC
IABCD
IABCDE
```

```
Concrete states  \{(i,j)\} \text{ (given)}   \{(i,j)|i \text{ is odd, } j \text{ is even} \}   \{(i,j)|i \text{ is odd, } j \text{ is even} \}   \{(i,j)|i \text{ is odd, } j \text{ is even, } i>j\}   \{(i,j)|i \text{ is even, } j \text{ is even, } i>j+1\}   \{(i,j)|i \text{ is even, } j \text{ is even, } i\geq j\}
```

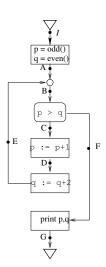
Collecting semantics of a program = set of (concrete) states occurring at each program point.



Path
I
IA
IAB
IABC
IABCD
IABCDE
IABCDEB

```
Concrete states \{(i,j)\} (given) \{(i,j)|i is odd, j is even\} \{(i,j)|i is odd, j is even\} \{(i,j)|i is odd, j is even, i>j\} \{(i,j)|i is even, j is even, i>j+1\} \{(i,j)|i is even, j is even, i\geq j \{(i,j)|i is even, j is even.
```

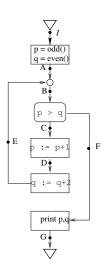
Collecting semantics of a program = set of (concrete) states occurring at each program point.



Path
I
IA
IAB
IABC
IABCD
IABCDE
IABCDEB
IABCDEB

```
Concrete states  \{(i,j)\} \text{ (given)}   \{(i,j)|i \text{ is odd, } j \text{ is even} \}   \{(i,j)|i \text{ is odd, } j \text{ is even} \}   \{(i,j)|i \text{ is odd, } j \text{ is even, } i > j \}   \{(i,j)|i \text{ is even, } j \text{ is even, } i > j+1 \}   \{(i,j)|i \text{ is even, } j \text{ is even, } i \geq j \}   \{(i,j)|i \text{ is even, } j \text{ is even, } i \geq j \}   \{(i,j)|i \text{ is even, } j \text{ is even, } i > j \}
```

Collecting semantics of a program = set of (concrete) states occurring at each program point.



```
Path
                       Concrete states
                       \{(i,j)\} (given)
IΑ
                      \{(i, j)|i \text{ is odd}, j \text{ is even}\}
IAB
                      \{(i, j)|i \text{ is odd}, j \text{ is even}\}
IABC
                      \{(i, j)|i \text{ is odd}, j \text{ is even}, i > j\}
IABCD
                       \{(i,j)|i is even, j is even, i > j+1\}
IABCDE
                      \{(i, j)|i \text{ is even, } j \text{ is even, } i > j\}
IABCDEB
                       \{(i, j)|i \text{ is even, } j \text{ is even, } i > j\}
                      \{(i, j)|i \text{ is even, } j \text{ is even, } i > j\}
IABCDEBC
IABCDEBCD
                      \{(i, j)|i \text{ is odd, } j \text{ is even, } i > j + 1\}
```

Collecting semantics of a program = set of (concrete) states occurring at each program point.

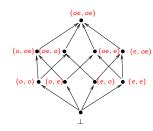
```
Path
                                                       Concrete states
                                                       \{(i, j)\}\ (given)
                                                       \{(i, j)|i \text{ is odd}, j \text{ is even}\}
                                IΑ
                                IAB
                                                       \{(i, j)|i \text{ is odd}, j \text{ is even}\}
      p = odd()
       q = even()
                                IABC
                                                       \{(i, j)|i \text{ is odd}, j \text{ is even}, i > j\}
                                IABCD
                                                       \{(i,j)|i is even, j is even, i > j+1\}
                                                       \{(i, j)|i \text{ is even, } j \text{ is even, } i \geq j\}
                                IABCDE
         Β •
                                IABCDEB
                                                       \{(i, j)|i is even, j is even, i > j\}
       p > q
                               IABCDEBC \{(i, j)|i \text{ is even, } j \text{ is even, } i > j\}
                               IABCDEBCD \{(i, j)|i \text{ is odd}, j \text{ is even}, i > j + 1\}
E
                             Therefore, collecting semantics:
        D_{\psi}
                                     \{(i, j)\}
         := q+2
                               A \{(i, j)|i \text{ odd}, j \text{ even}\}
                               B \{(i,j)|i \text{ odd}, j \text{ even}\} \cup \{(i,j)|i \text{ even}, j \text{ even}, i > j\}
                               C \{(i, j) | j \text{ even}, i > j \}
        print p,q-
                                D \{(i, j) | j \text{ even}, i > j + 1\}
                               E \{(i, j)|j \text{ even}, i > j\}
                                   \{(i,j)|i \text{ odd}, j \text{ even}, i < j\} \cup \{(i,j)|i \text{ even}, j \text{ even}, i = j\}
                               G \{(i,j)|i \text{ odd}, j \text{ even}, i < j\} \cup \{(i,j)|i \text{ even}, j \text{ even}, i = j\}
```

An abstract interpretation

Components of an abstract interpretation:

- Set of abstract states D, forming a complete lattice.
- "Concretization" function $\gamma: D \to 2^{State}$, which associates a set of concrete states with each abstract state.
- Set of transfer functions, with signature $D \rightarrow D$, which interpret assignments and conditionals using abstract states:
 - a function $f_n: D \to D$ for each assignment node n
 - a function $f_{n,T}:D\to D$ and a function $f_{n,F}:D\to D$ for each conditional node n

Abstract lattice D



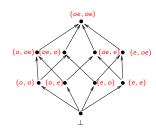
Transfer function for an assignment node n: p := p+q

$$f_n(s) = \begin{cases} \bot & \text{if } s \text{ is } \bot \\ (o, s[q]) & \text{if } s[p] \text{ is o and } s[q] \text{ is e,} \\ & \text{or } s[p] \text{ is e and } s[q] \text{ is o} \\ (e, s[q]) & \text{if both } s[p] \text{ and } s[q] \text{ are o} \\ & \text{or both } s[p] \text{ and } s[q] \text{ are e} \\ (oe, s[q]) & \text{otherwise} \end{cases}$$

ullet The concretization function γ

•
$$\gamma((oe, oe)) =$$

Abstract lattice D



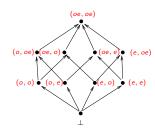
Transfer function for an assignment node n: p := p+q

$$f_n(s) = \begin{cases} \bot & \text{if } s \text{ is } \bot \\ (o, s[q]) & \text{if } s[p] \text{ is o and } s[q] \text{ is e,} \\ & \text{or } s[p] \text{ is e and } s[q] \text{ is o} \\ (e, s[q]) & \text{if both } s[p] \text{ and } s[q] \text{ are o} \\ & \text{or both } s[p] \text{ and } s[q] \text{ are e} \\ (oe, s[q]) & \text{otherwise} \end{cases}$$

ullet The concretization function γ

•
$$\gamma((oe, oe)) = State, \gamma(\bot) =$$

Abstract lattice D



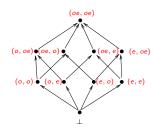
Transfer function for an assignment node n: p := p+q

$$f_n(s) = \begin{cases} \bot & \text{if } s \text{ is } \bot \\ (o, s[q]) & \text{if } s[p] \text{ is o and } s[q] \text{ is e,} \\ & \text{or } s[p] \text{ is e and } s[q] \text{ is o} \\ (e, s[q]) & \text{if both } s[p] \text{ and } s[q] \text{ are o} \\ & \text{or both } s[p] \text{ and } s[q] \text{ are e} \\ (oe, s[q]) & \text{otherwise} \end{cases}$$

ullet The concretization function γ

•
$$\gamma((oe, oe)) = State, \gamma(\bot) = \emptyset, \gamma((o, oe)) =$$

Abstract lattice D



Transfer function for an assignment node n: p := p+q

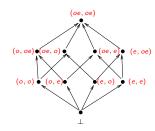
$$f_n(s) = \begin{cases} \bot & \text{if } s \text{ is } \bot \\ (o, s[q]) & \text{if } s[p] \text{ is o and } s[q] \text{ is e,} \\ & \text{or } s[p] \text{ is e and } s[q] \text{ is o} \\ (e, s[q]) & \text{if both } s[p] \text{ and } s[q] \text{ are o} \\ & \text{or both } s[p] \text{ and } s[q] \text{ are e} \\ (oe, s[q]) & \text{otherwise} \end{cases}$$

• The concretization function γ

•
$$\gamma((oe, oe)) = State, \ \gamma(\bot) = \emptyset, \ \gamma((o, oe)) = \{(m, n) \mid m \text{ is odd}\}$$

•
$$\gamma((o,e)) =$$

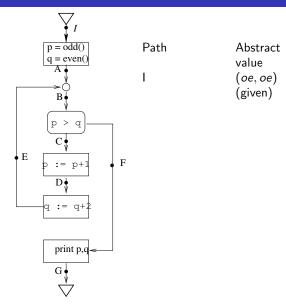
Abstract lattice D

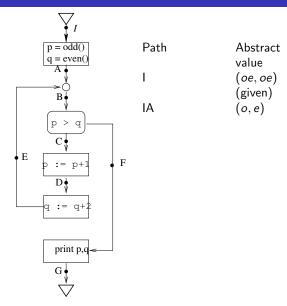


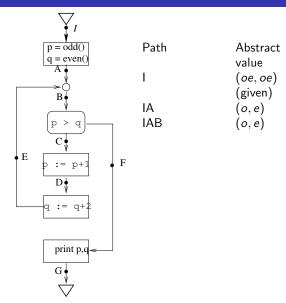
Transfer function for an assignment node n: p := p+q

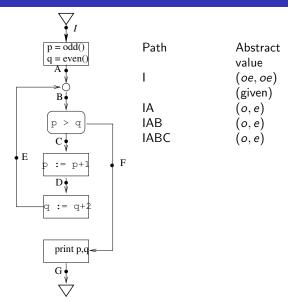
$$f_n(s) = \begin{cases} \bot & \text{if } s \text{ is } \bot \\ (o, s[q]) & \text{if } s[p] \text{ is o and } s[q] \text{ is e,} \\ & \text{or } s[p] \text{ is e and } s[q] \text{ is o} \\ (e, s[q]) & \text{if both } s[p] \text{ and } s[q] \text{ are o} \\ & \text{or both } s[p] \text{ and } s[q] \text{ are e} \\ (oe, s[q]) & \text{otherwise} \end{cases}$$

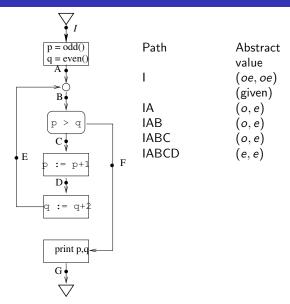
- The concretization function γ
 - $\gamma((oe, oe)) = State, \ \gamma(\bot) = \emptyset, \ \gamma((o, oe)) = \{(m, n) \mid m \text{ is odd}\}$
 - $\gamma((o,e)) = \{(m,n) \mid m \text{ is odd and } n \text{ is even}\}, \ldots$

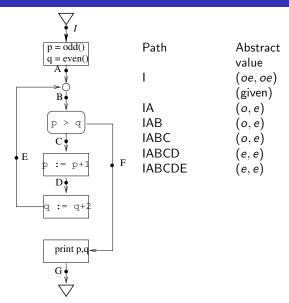


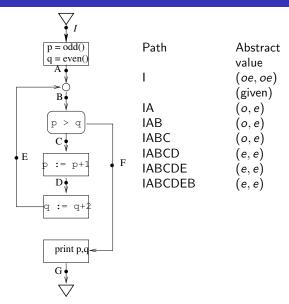


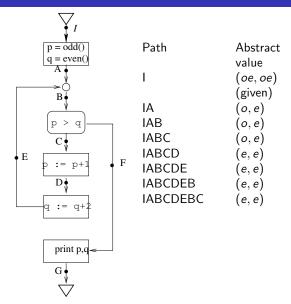


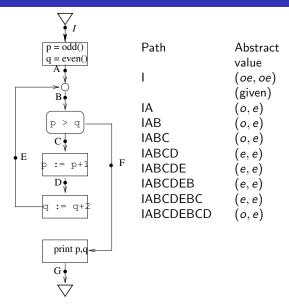


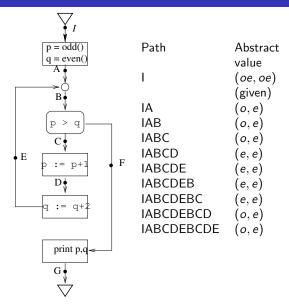


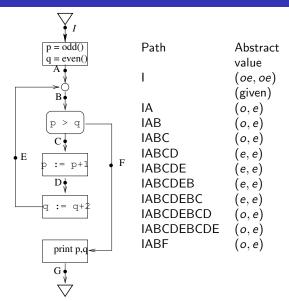


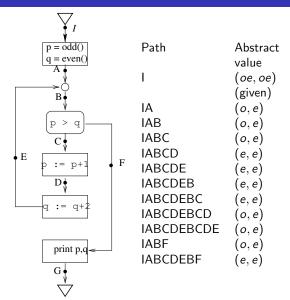


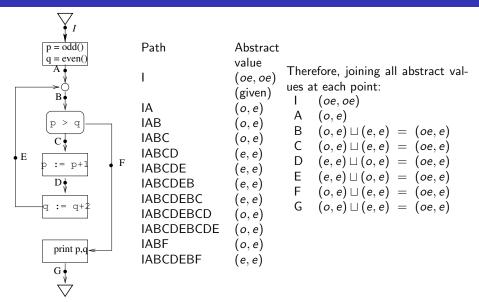


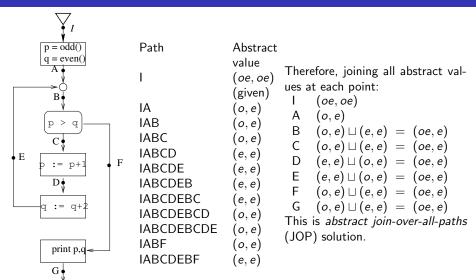












Comparison of abstract JOP states and collecting states

```
Abstract JOP: Collecting states: A (o,e) A \{(i,j)|i \text{ odd, } j \text{ even}\}
B (oe,e) B \{(i,j)|i \text{ odd, } j \text{ even}\} \cup \{(i,j)|i \text{ even, } i \geq j\}
C (oe,e) C \{(i,j)|j \text{ even, } i > j\}
D (oe,e) D \{(i,j)|j \text{ even, } i > j+1\}
E (oe,e) E \{(i,j)|j \text{ even, } i \geq j\}
F (oe,e) F \{(i,j)|i \text{ odd, } j \text{ even, } i < j\} \cup \{(i,j)|i \text{ even, } j \text{ even, } i = j\}
G (oe,e) G \{(i,j)|i \text{ odd, } j \text{ even, } i < j\} \cup \{(i,j)|i \text{ even, } j \text{ even, } i = j\}
```

Comparison of abstract JOP states and collecting states

```
Abstract JOP: Collecting states: A (o,e) A \{(i,j)|i \text{ odd, } j \text{ even}\}
B (oe,e) B \{(i,j)|i \text{ odd, } j \text{ even}\} \cup \{(i,j)|i \text{ even, } i \geq j\}
C (oe,e) C \{(i,j)|j \text{ even, } i > j\}
D (oe,e) D \{(i,j)|j \text{ even, } i > j+1\}
E (oe,e) E \{(i,j)|j \text{ even, } i \geq j\}
F (oe,e) F \{(i,j)|i \text{ odd, } j \text{ even, } i < j\} \cup \{(i,j)|i \text{ even, } j \text{ even, } i = j\}
G (oe,e) G \{(i,j)|i \text{ odd, } j \text{ even, } i < j\} \cup \{(i,j)|i \text{ even, } j \text{ even, } i = j\}
```

Note that at each point γ image of abstract solution is over-approximation of collecting states.

A given abstract interpretation is said to be *correct* if, for all abstract states $d_0 \in D$, for all programs P and for all program points p in P,

 γ image of join of all abstract states arising at p (i.e., abstract JOP solution at p), with d_0 as the initial abstract value at P's entry

collecting semantics at p, with $\gamma(d_0)$ as the initial set of concrete states at P's entry

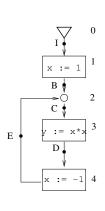
A given abstract interpretation is said to be *correct* if, for all abstract states $d_0 \in D$, for all programs P and for all program points p in P,

 γ image of join of all abstract states arising at p (i.e., abstract JOP solution at p), with d_0 as the initial abstract value at P's entry

collecting semantics at p, with $\gamma(d_0)$ as the initial set of concrete states at P's entry

We will study later certain sufficient conditions for a given abstract interpretation to be correct.

Another example program



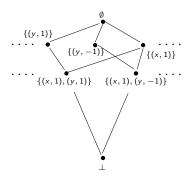
$$\begin{array}{ll} \text{Path} & \text{Characterization of concrete states} \\ \text{I} & \textit{true} \ (\text{given}) \\ \text{IB} & \text{x} = 1 \\ \text{IBC} & \text{x} = 1 \\ \text{IBCD} & \text{x} = 1 \land \text{y} = 1 \\ \text{IBCDE} & \text{x} = -1 \land \text{y} = 1 \\ \text{IBCDEC} & \text{x} = -1 \land \text{y} = 1 \\ \text{IBCDECD} & \text{x} = -1 \land \text{y} = 1 \\ \dots & \text{x} = -1 \land \text{y} = 1 \\ \end{array}$$

Therefore, collecting semantics:

Point Characterization of concrete states | true | B |
$$x = 1$$
 | C | $(x = 1) \lor (x = -1 \land y = 1)$ | D | $(y = 1) \land (x = -1 \lor x = 1)$ | E | $x = -1 \land y = 1$

Abstract interpretation for constant propagation

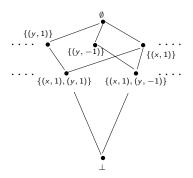
Abstract lattice D



• Concretization function: What is $\gamma(d)$?

Abstract interpretation for constant propagation

Abstract lattice D



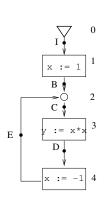
• Concretization function: What is $\gamma(d)$?

$$\begin{array}{cccc} \bot & \mapsto & \{\} \\ \emptyset & \mapsto & \mathit{State} \\ \{(x,c)\} & \mapsto & \{(c,j)|\ j \ \mathsf{is any value} \} \\ \{(x,c),(y,d)\} & \mapsto & \{(c,d)\} \end{array}$$

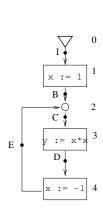
Abstract interpretation for constant propagation – contd.

Transfer function for assignment node n of the form x := exp.

$$f_n(P) = \bot$$
, if P is \bot
= $\{(y,c) \in P \mid y \text{ is any variable other than } x\} \cup \{(x,d)\}$,
if all variables in exp have constant values in P , and if
exp evaluates to d with these constant values
= $\{(y,c) \in P \mid y \text{ is any variable other than } x\}$, otherwise

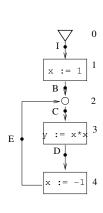


Path Abstract value at end of path I \emptyset

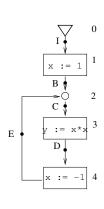




Abstract value at end of path \emptyset $\{(x,1)\}$

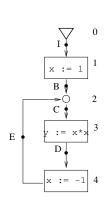


Path I IB IBC Abstract value at end of path \emptyset $\{(x,1)\}$ $\{(x,1)\}$

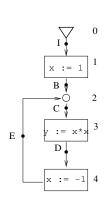




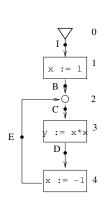
```
Abstract value at end of path \emptyset \{(x,1)\} \{(x,1)\} \{(x,1),(y,1)\}
```



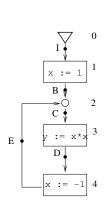
 $\begin{array}{lll} \text{Path} & \text{Abstract value at end of path} \\ \text{I} & \emptyset \\ \text{IB} & \{(x,1)\} \\ \text{IBC} & \{(x,1)\} \\ \text{IBCD} & \{(x,1),(y,1)\} \\ \text{IBCDE} & \{(x,-1),(y,1)\} \end{array}$



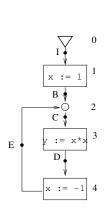
```
\begin{array}{ll} \text{Path} & \text{Abstract value at end of path} \\ \text{I} & \emptyset \\ \text{IBC} & \{(x,1)\} \\ \text{IBCD} & \{(x,1)\} \\ \text{IBCD} & \{(x,1),(y,1)\} \\ \text{IBCDE} & \{(x,-1),(y,1)\} \\ \text{IBCDEC} & \{(x,-1),(y,1)\} \end{array}
```



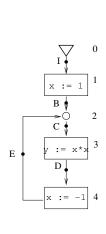
```
Path
            Abstract value at end of path
ΙB
            \{(x,1)\}
            \{(x,1)\}
IBC
            \{(x,1),(y,1)\}
IBCD
       \{(x,-1),(y,1)\}
IBCDE
IBCDEC \{(x, -1), (y, 1)\}
IBCDECD \{(x, -1), (y, 1)\}
```



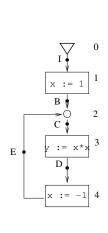
```
Path
            Abstract value at end of path
ΙB
            \{(x,1)\}
            \{(x,1)\}
IBC
IBCD
            \{(x,1),(y,1)\}
       \{(x,-1),(y,1)\}
IBCDE
IBCDEC \{(x, -1), (y, 1)\}
IBCDECD \{(x, -1), (y, 1)\}
            \{(x,-1),(y,1)\}
```



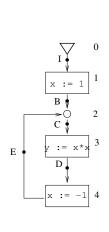
```
Path
            Abstract value at end of path
ΙB
            \{(x,1)\}
           \{(x,1)\}
IBC
IBCD \{(x,1),(y,1)\}
IBCDE \{(x, -1), (y, 1)\}
IBCDEC \{(x, -1), (y, 1)\}
IBCDECD \{(x, -1), (y, 1)\}
            \{(x,-1),(y,1)\}
Point
       Abstract JOP value
```



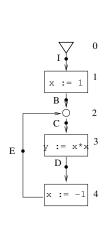
```
Path
            Abstract value at end of path
ΙB
            \{(x,1)\}
            \{(x,1)\}
IBC
IBCD \{(x,1),(y,1)\}
IBCDE \{(x, -1), (y, 1)\}
IBCDEC \{(x, -1), (y, 1)\}
IBCDECD \{(x, -1), (y, 1)\}
            \{(x,-1),(y,1)\}
Point
       Abstract JOP value
       \{(x,1)\}
```



```
Path
            Abstract value at end of path
ΙB
            \{(x,1)\}
            \{(x,1)\}
IBC
IBCD \{(x,1),(y,1)\}
IBCDE \{(x, -1), (y, 1)\}
IBCDEC \{(x, -1), (y, 1)\}
IBCDECD \{(x, -1), (y, 1)\}
            \{(x, -1), (y, 1)\}
Point
       Abstract JOP value
       \{(x,1)\}
```



```
Path
            Abstract value at end of path
ΙB
            \{(x,1)\}
            \{(x,1)\}
IBC
IBCD \{(x, 1), (y, 1)\}
IBCDE \{(x, -1), (y, 1)\}
IBCDEC \{(x, -1), (y, 1)\}
IBCDECD \{(x, -1), (y, 1)\}
            \{(x, -1), (y, 1)\}
Point
        Abstract JOP value
       \{(x,1)\}
        \{(y,1)\}
```



```
Path
            Abstract value at end of path
ΙB
            \{(x,1)\}
            \{(x,1)\}
IBC
IBCD \{(x, 1), (y, 1)\}
IBCDE \{(x, -1), (y, 1)\}
IBCDEC \{(x, -1), (y, 1)\}
IBCDECD \{(x, -1), (y, 1)\}
            \{(x, -1), (y, 1)\}
Point
       Abstract JOP value
       \{(x,1)\}
  \{(y,1)\}
       \{(x,-1),(y,1)\}
```

Correctness in previous example

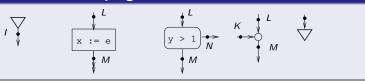
Verify that

- at points I, B and E $\gamma(\text{abstract JOP value}) = \text{collecting semantics}.$
- at points C and D $\gamma({\rm abstract\ JOP\ value}) \supset {\rm collecting\ semantics}.$
- the abstract transfer functions given are the best possible for the given lattice L. That is, imprecision is due to the lattice, not the transfer functions.

Formal definition of control-flow graphs

Programs are finite directed graphs with following nodes (statements):

Nodes or statements in a program



• Expressions:

$$e := c | x | e + e | e - e | e * e.$$

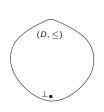
Boolean expressions:

$$be ::= tt \mid ff \mid e \leq e \mid e = e \mid \neg be \mid be \lor be \mid be \land be.$$

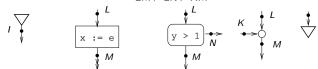
• Assume unique initial program point 1.

Formal definition of an abstract interpretation

- Complete join semi-lattice (D, \leq) , with a least element \perp .
- Concretization function $\gamma: D \to 2^{State}$
- $\bot \in D$ represents unreachability of the program point (i.e., $\gamma(\bot)$ should be equal to \emptyset). Also, $\gamma(\top)$ should be *State*.



• We require transfer functions f_{LM} , f_{LN} , f_{KM} for all scenarios below:



- We assume transfer functions are monotonic, and satisfy $f(\bot) = \bot$.
- For junction nodes, both transfer functions should be identity

What we want to compute for a given program

- Path in a program: Sequence of connected edges or program points, beginning at initial point I
- Transfer functions extend to paths in program:

$$f_{IBCD} = f_{CD} \circ f_{BC} \circ f_{IB}$$
.

where $(f_a \circ f_b)(x)$ is defined as $f_a(f_b(x))$.

- f_p is $\lambda d. \perp \Rightarrow$ path p is infeasible.
- Join over all paths (JOP) definition: For each program point N

$$d_N = \bigsqcup_{\text{paths } p \text{ from } I \text{ to } N} f_p(d_0).$$

where d_0 is a given initial abstract value at entry node.

Formalization of collecting semantics

- Let *Val* be the set of all *concrete* values; e.g., *Integer* ∪ *Boolean*.
- State is normally the domain $Var \rightarrow Val$. However, in general, it can be any semantic domain.
- Program semantics is given by the functions $nstate_{MN}: State \rightarrow 2^{State}$



• These induce the functions $nstate': 2^{State} \rightarrow 2^{State}$

$$\textit{nstate'}_{\textit{MN}}(\textit{S}_1 \in 2^{\textit{State}}) = \bigcup_{\textit{s}_1 \in \textit{S}_1} \textit{nstate}_{\textit{MN}}(\textit{s}_1)$$

Formalization of collecting semantics – contd.

- ullet Collecting semantics SS is a map $ProgramPoints
 ightarrow 2^{State}$
- It can be seen that at each program point N, the collecting semantics SS(N) is equal to:

$$\bigcup_{p \text{ path from } I \text{ to } N} nstate'_p(S_0).$$

where I is entry point of CFG, S_0 is the given initial set of states, and $nstate'_p$ is composition of nstate' functions of edges that constitute p.

Instantiating a generic lattice and transfer functions implementation to a given program

- Implement a generic lattice element class, such that when one constructs a lattice element, one can provide the variable names along with other internal details
 - For e.g., $\{(p, o), (q, e)\}$, or $\{(x, 1), (y, 2)\}$
- Implement just three generic transfer functions: one for assignment statements, one for *true* branch of conditionals, and one for *false* branch of conditionals
 - Each transfer function takes, in addition to the incoming data flow fact (i.e., lattice element), the actual assignment/condition node as an argument