



LUND
UNIVERSITY

EDAP15: Program Analysis

DATA FLOW ANALYSIS 4: GOING INTERPROCEDURAL

Christoph Reichenbach

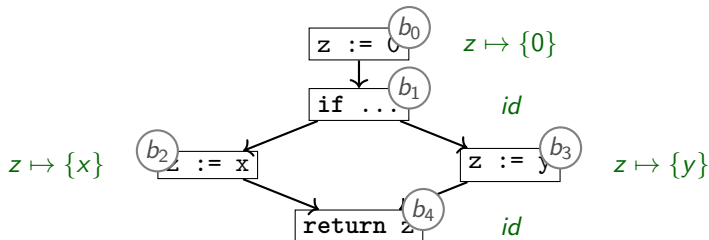


Welcome back!

- ▶ Homework #1 extended to Tuesday, noon (12:00)
- ▶ Homework #2 released
- ▶ Lectures this week:
 - ▶ Interprocedural Data Flow
 - ▶ Analysing Object-Oriented Programs

Summarising Procedures (Reaching Values)

$f(x, y) =$



► **Compose transfer functions:**

- $trans_{b_0} \circ trans_{b_1} = [z \mapsto 0]$
- $trans_{b_0} \circ trans_{b_1} \circ trans_{b_2} = [z \mapsto \{x\}]$
- $trans_{b_0} \circ trans_{b_1} \circ trans_{b_3} = [z \mapsto \{y\}]$
- $trans_{b_0} \circ trans_{b_1} \circ (trans_{b_2} \sqcup trans_{b_3}) = [z \mapsto \{x, y\}]$
- $trans_{b_0} \circ trans_{b_1} \circ (trans_{b_2} \sqcup trans_{b_3}) \circ trans_{b_4} = [z \mapsto \{x, y\}]$

Procedure Summaries vs Recursion

f calls g calls h calls f

- ▶ Requires additional analysis to identify who calls whom
- ▶ Compute summaries of mutually recursive functions together
- ▶ Recursive call edges analogous to loops

Procedure Summaries

- ▶ Composing transfer functions yields a combined transfer function for $f()$:

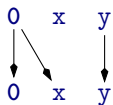
$$trans_f = [\mathbf{return} \mapsto \{x, y\}]$$

- ▶ Use $trans_f$ as transfer function for $f()$, discard f 's body
- ▶ **Opportunities:**
 - ▶ Can yield compact subroutine descriptions
 - ▶ Can speed up call site analysis dramatically
- ▶ **Challenges:**
 - ▶ More complex to implement
 - ▶ Recursion remains challenging
- ▶ **Limitations:**
 - ▶ Requires suitable representation for summary
 - ▶ Requires mechanism for abstracting and applying summary
 - ▶ Worst cases:
 - ▶ $trans_f$ is symbolic expression as complex as f itself

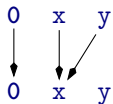
Representation Relations

Example procedure summary representation:

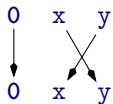
```
x := null;  
y := y;
```



```
if x != y {  
  x := y;  
}  
y := 1;
```



```
{ t := x  
  x := y  
  y := t }
```



'May be null' analysis

- ▶ $c \rightarrow d$:
if $P(c) \in \mathbf{in}_b$ then $P(d) \in \mathbf{out}_b$
- ▶ Representation Relations relate \mathbf{in}_b and \mathbf{out}_b variables \mathcal{V}
- ▶ $R \subseteq (\mathcal{V} \cup \{\mathbf{0}\}) \times (\mathcal{V} \cup \{\mathbf{0}\})$
- ▶ if $\langle \mathbf{0}, X \rangle \in R$:
 X always 'may be null' in \mathbf{out}_b
- ▶ if $\langle Y, X \rangle \in R$:
If Y 'may be null' in \mathbf{in}_b :
 $\Rightarrow X$ 'may be null' in \mathbf{out}_b

Representation Relations and Distributivity

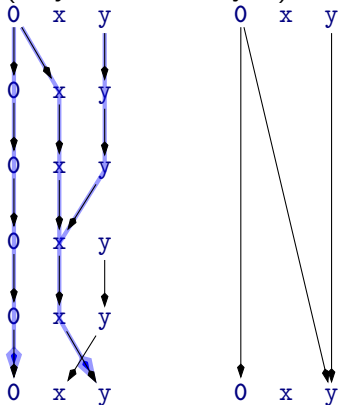
Composing Representation Relations

Representation Relations (*may be null analysis*):

```
x := null;  
y := y;
```

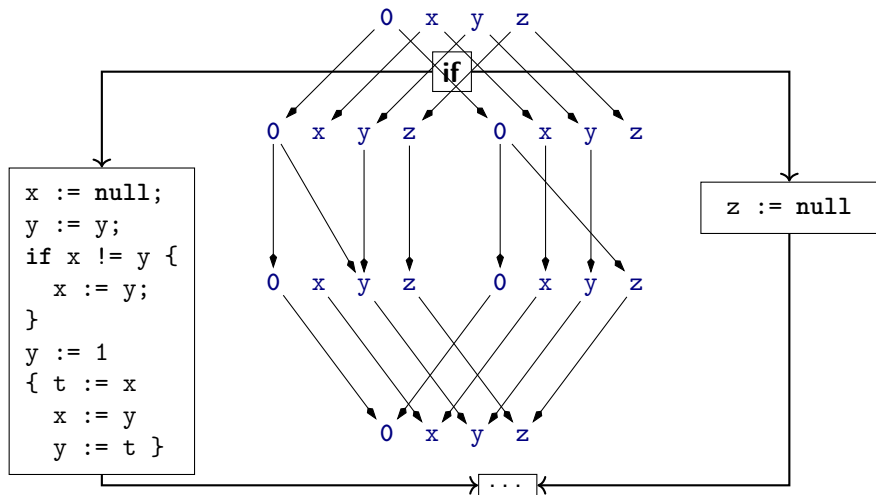
```
if x != y {  
  x := y;  
}  
y := 1;
```

```
{ t := x;  
  x := y;  
  y := t; }
```

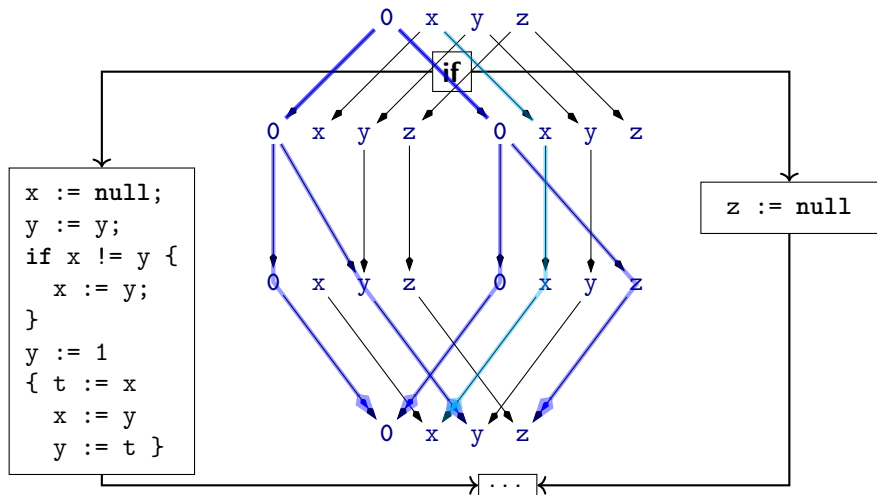


Composed representation relations are again representation relations

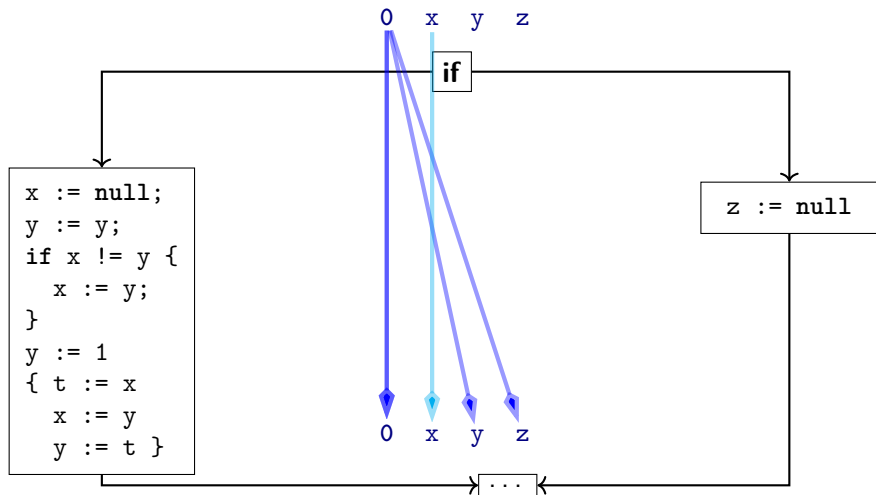
Joining Control-Flow Paths



Joining Control-Flow Paths



Joining Control-Flow Paths



Logical "Or"

Dataflow via Graph Reachability

$$n = \langle b, v \rangle$$

- ▶ Assume binary lattice $(\{\top, \perp\}, \sqsubseteq, \sqcap, \sqcup)$
 - ▶ $\top \sqcup y = \top = x \sqcup \top$ and $\perp \sqcup \perp = \perp$
 - ▶ Typical for 'May' analysis ($P(x) = 'x \text{ may be null}'$)
- ▶ Encode Dataflow problem as *Graph-Reachability*
- ▶ Graph nodes $n = \langle b, v \rangle$
 - ▶ b : CFG node
 - ▶ v : Variable or $\mathbf{0}$
 - ▶ $\mathbf{0}$: $\langle b_1, \mathbf{0} \rangle \rightarrow \langle b_2, y \rangle$: $P(y)$ at b_2 holds always
 - ▶ Variable: $\langle b_1, x \rangle \rightarrow \langle b_2, y \rangle$: $P(x)$ at $b_1 \implies P(y)$ at b_2

Dataflow via Graph Reachability

$$n = \langle b, v \rangle$$

- ▶ Assume binary lattice $(\{\top, \perp\}, \sqsubseteq, \sqcap, \sqcup)$
 - ▶ $\top \sqcup y = \top = x \sqcup \top$ and $\perp \sqcup \perp = \perp$
 - ▶ Typical for 'May' analysis ($P(x) = \text{'x may be null'}$)
 - ▶ Equivalently for 'Must' analysis:
'x must be null' = not ('x may be non-null')
- ▶ Encode Dataflow problem as *Graph-Reachability*
- ▶ Graph nodes $n = \langle b, v \rangle$
 - ▶ b : CFG node
 - ▶ v : Variable or $\mathbf{0}$
 - ▶ $\mathbf{0}$: $\langle b_1, \mathbf{0} \rangle \rightarrow \langle b_2, y \rangle$: $P(y)$ at b_2 holds always
 - ▶ Variable: $\langle b_1, x \rangle \rightarrow \langle b_2, y \rangle$: $P(x)$ at $b_1 \implies P(y)$ at b_2

A Dataflow Worklist Algorithm: IFDS

- ▶ Call-site sensitive interprocedural data flow algorithm
- ▶ IFDS = (Interprocedural **F**inite **D**istributive **S**ubset problems)
- ▶ 'Exploded Supergraph': $G^\# = (N^\#, E^\#)$
 - ▶ $N^\# = N_{CFG} \times (\mathcal{V} \cup \{0\})$
 - ▶ Plus parameter/return call edges
- ▶ Property-of-interest holds if reachable from $\langle b_{main}^s, \mathbf{0} \rangle$
 - ▶ b_{main}^s is CFG *ENTER* node of main entry point
- ▶ **Key ideas:**
 - ▶ Worklist-based
 - ▶ Construct Representation Relations on demand
 - ▶ Construct 'Exploded Supergraph'
 - ▶ CFG of all functions $\times \mathcal{V} \cup \{0\}$

IFDS Datastructures

Instead of $\langle\langle b_0, v_0 \rangle, \langle b_3, v_0 \rangle\rangle$ we also write:

$$\langle b_0, v_0 \rangle \rightarrow \langle b_3, v_0 \rangle$$

WORKLIST edge

$\langle b_0, v_0 \rangle \dashrightarrow \langle b_3, v_0 \rangle$



PATHEDGE edge

All WORKLIST edges are also PATHEDGE edges

Result of our analysis

$N^\#$ -edge



SUMMARYINST

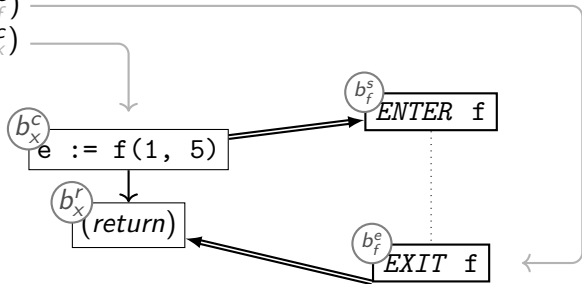
Generated from summary nodes

Otherwise equivalent to $N^\#$ -edges

IFDS Strategy

- ▶ Algorithm distinguishes between three types of nodes:

- ▶ Exit nodes (b_f^e)
- ▶ Call nodes (b_x^c)
- ▶ Other nodes



On-demand processing

```
Procedure propagate( $n_1 \rightarrow n_2$ ):  
begin  
  if  $n_1 \rightarrow n_2 \in \text{PATHEDGE}$  then  
    return  
   $\text{PATHEDGE} := \text{PATHEDGE} \cup \{n_1 \rightarrow n_2\}$   
   $\text{WORKLIST} := \text{WORKLIST} \cup \{n_1 \rightarrow n_2\}$   
end
```

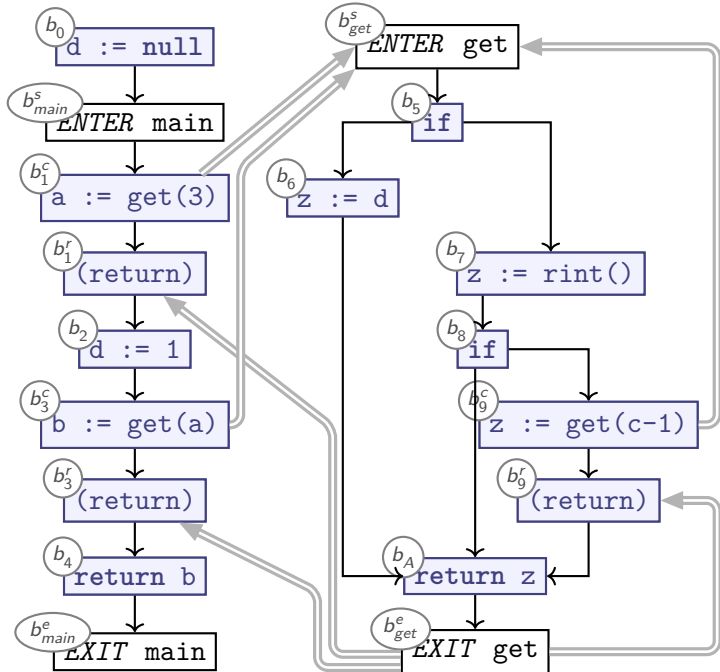
Running Example

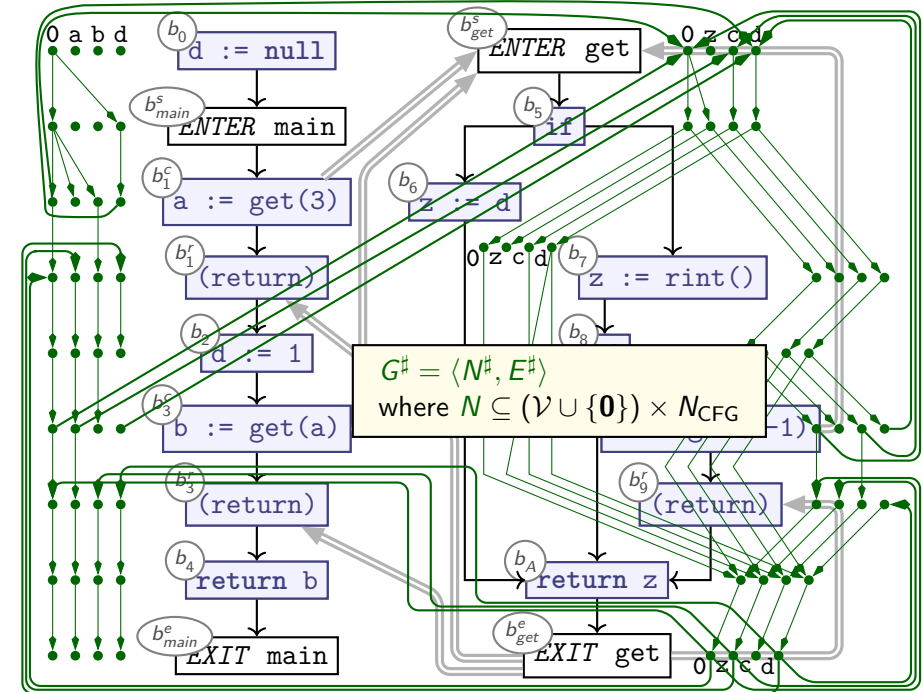
Teal-0: *main()*

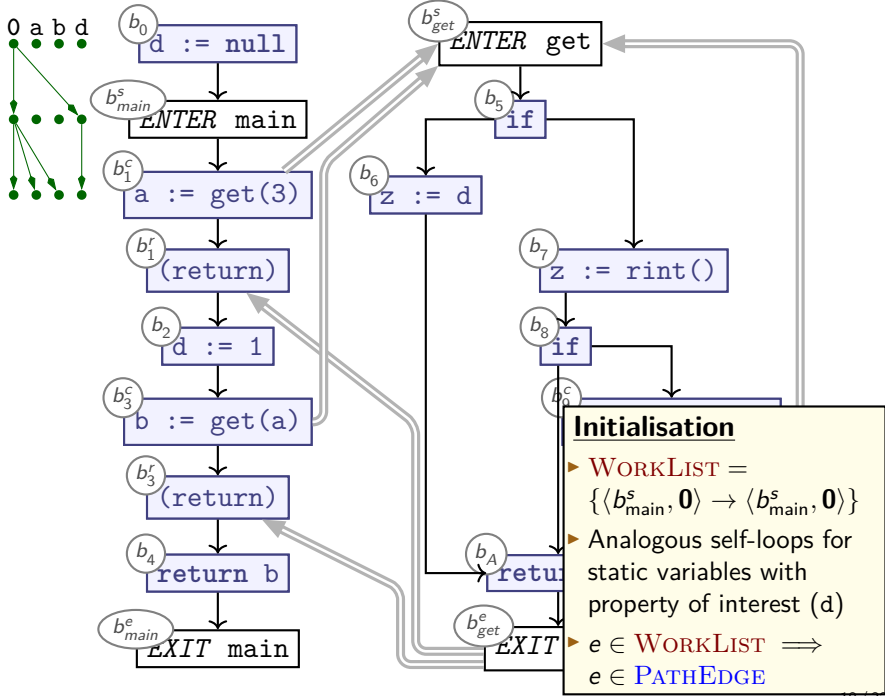
```
var default := null;
fun main() = {
  var a := get(3);
  default := 1;
  var b := get(3);
  return b;
}
```

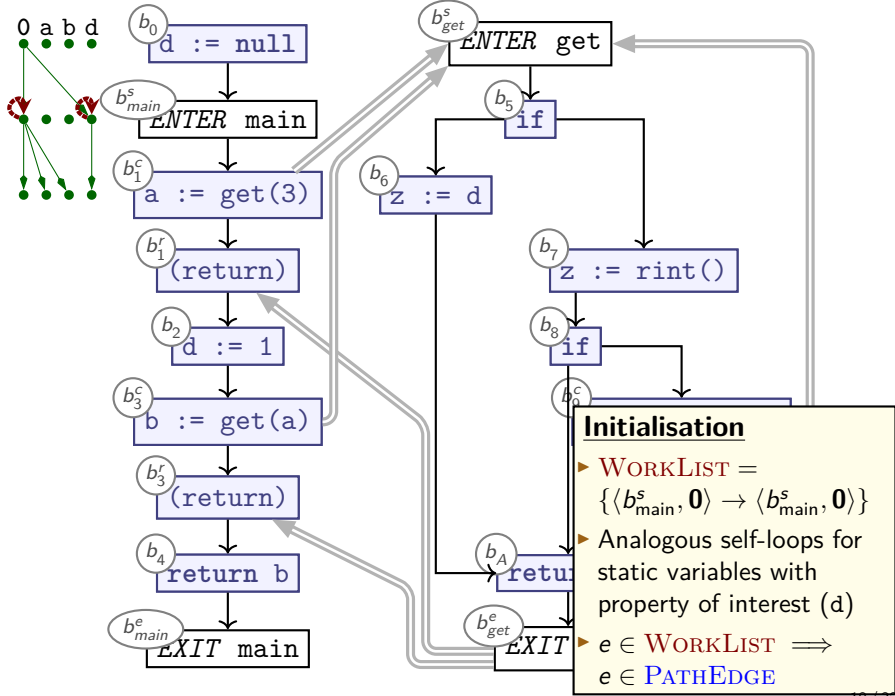
Teal-0: *get()*

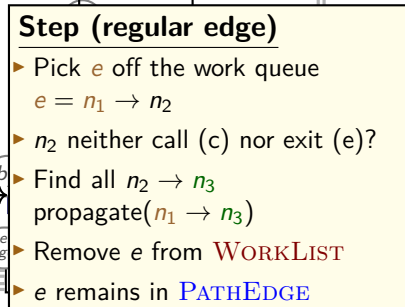
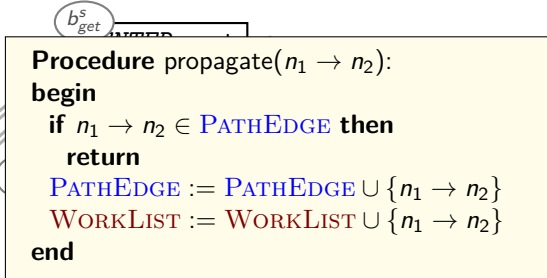
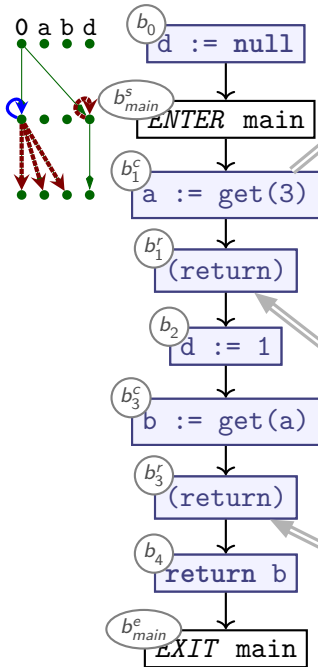
```
fun get(c) = {
  if c == 0 {
    z := default;
  } else {
    z := read_int();
    if z < 0 {
      z := get(c - 1);
    }
  }
  return z;
}
```

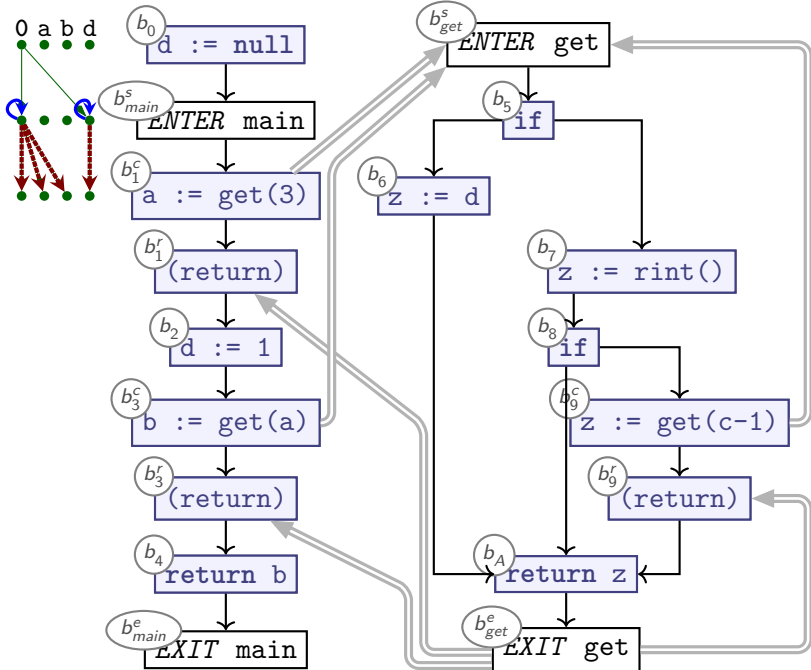


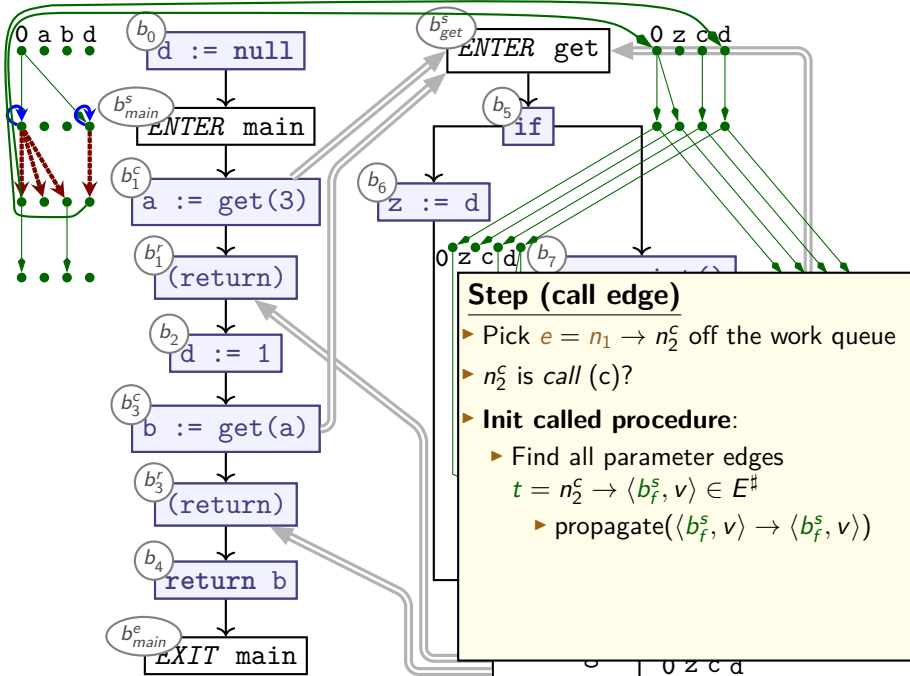


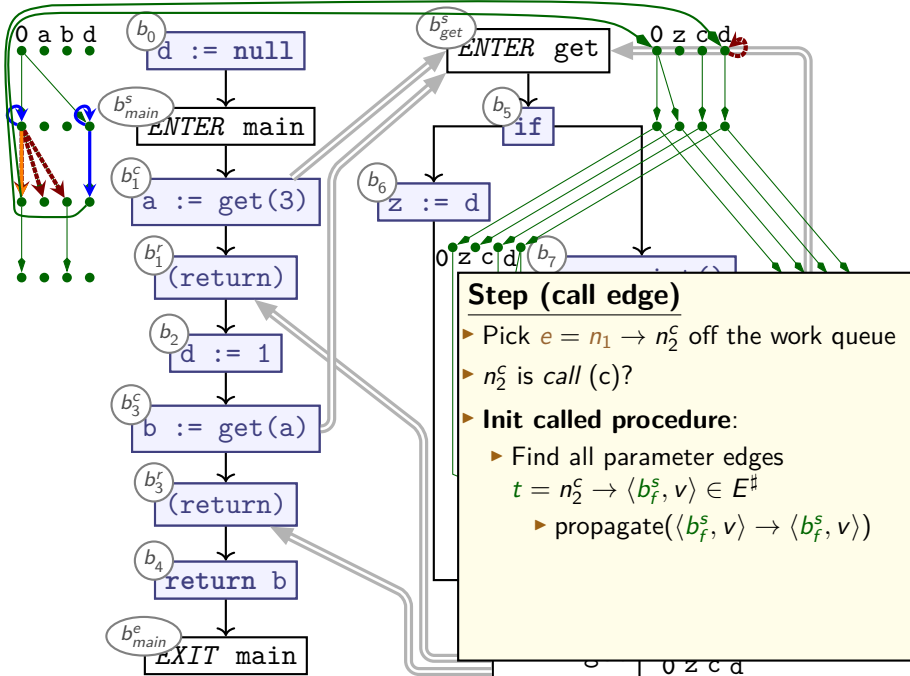


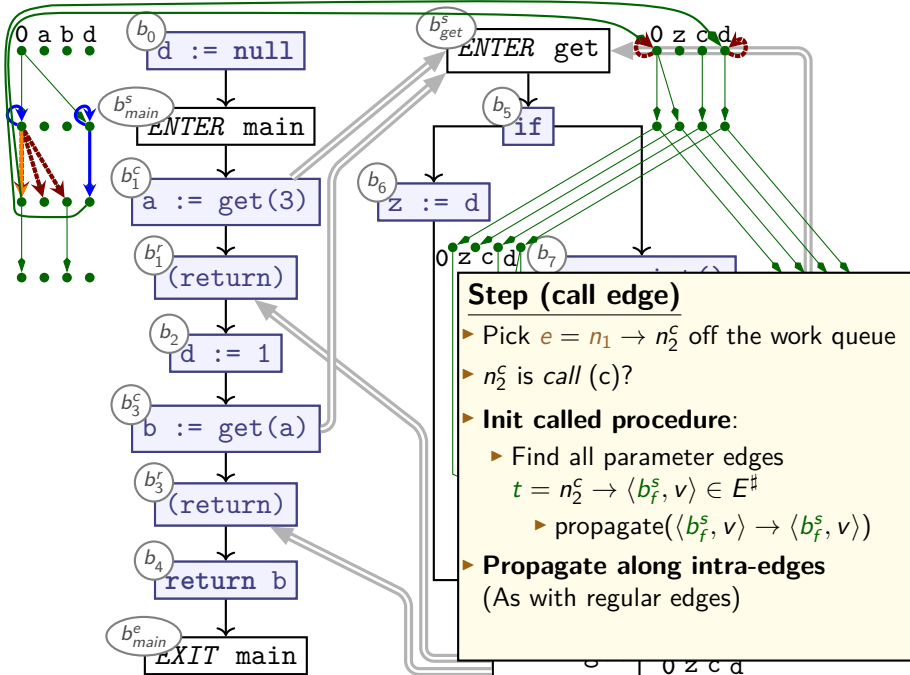






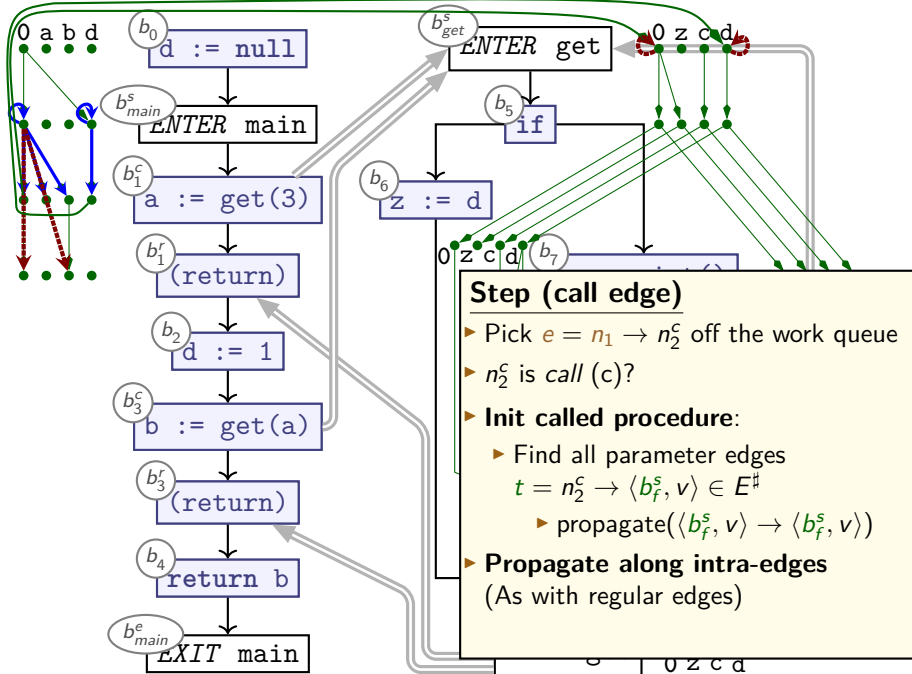


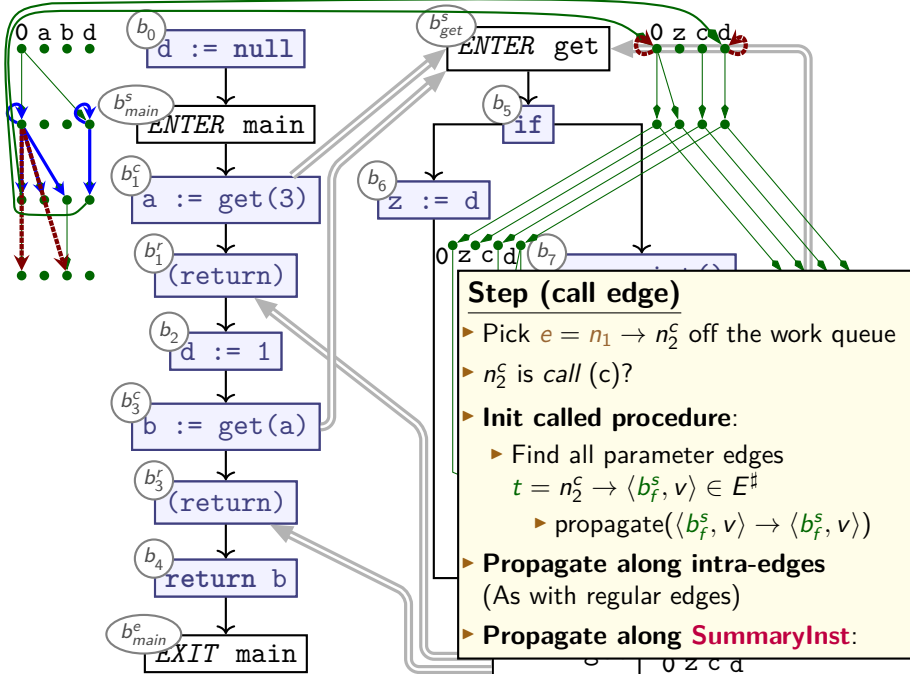


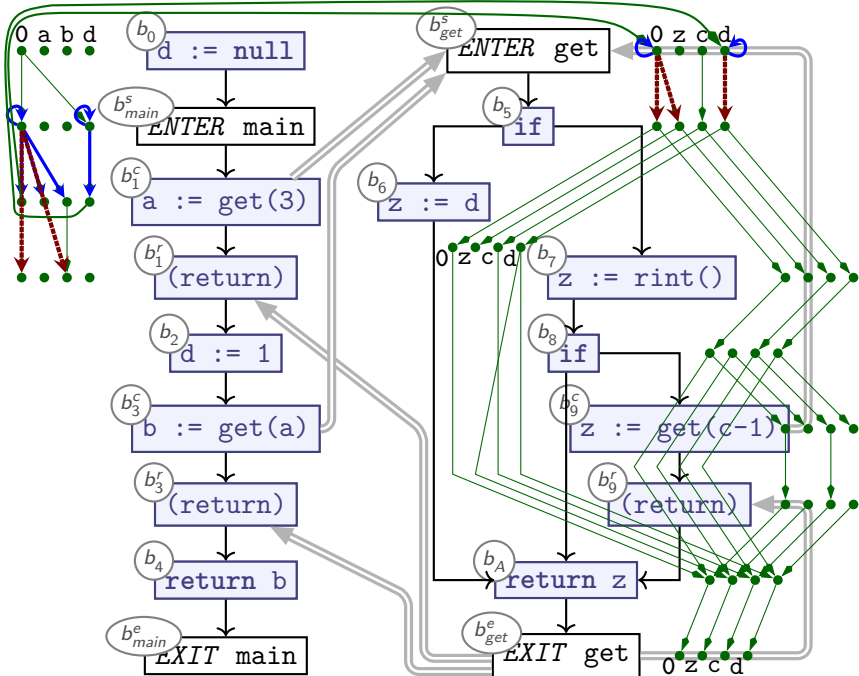


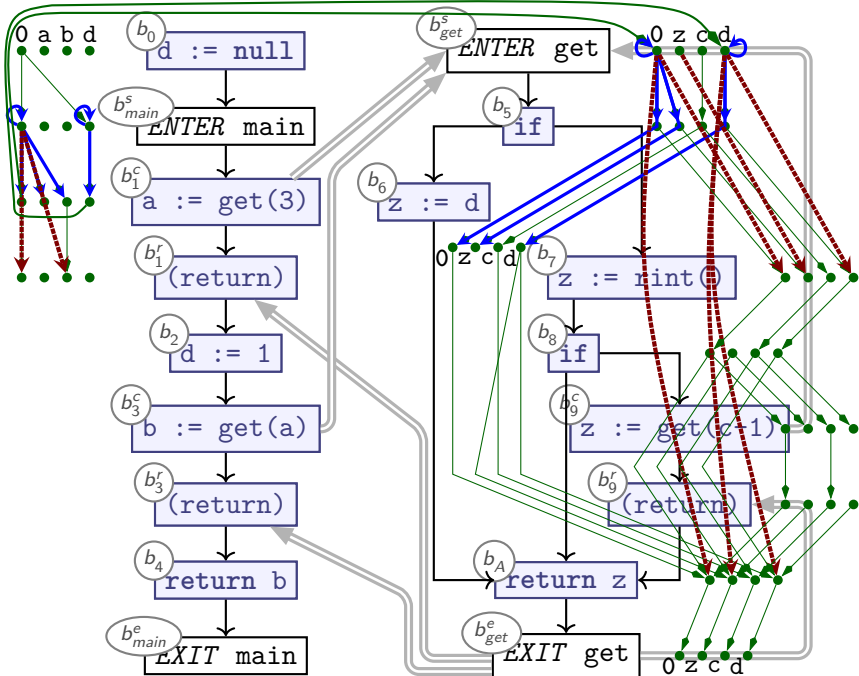
Step (call edge)

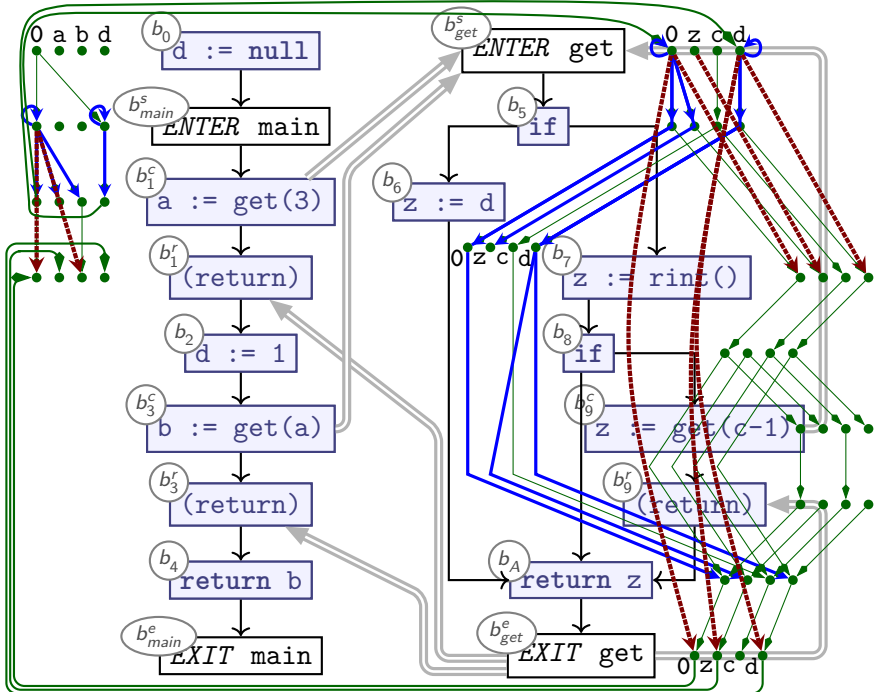
- ▶ Pick $e = n_1 \rightarrow n_2$ off the work queue
- ▶ n_2 is *call* (c)?
- ▶ **Init called procedure:**
 - ▶ Find all parameter edges $t = n_2^c \rightarrow \langle b_f^s, v \rangle \in E^\#$
 - ▶ propagate ($\langle b_f^s, v \rangle \rightarrow \langle b_f^s, v \rangle$)
- ▶ **Propagate along intra-edges**
(As with regular edges)

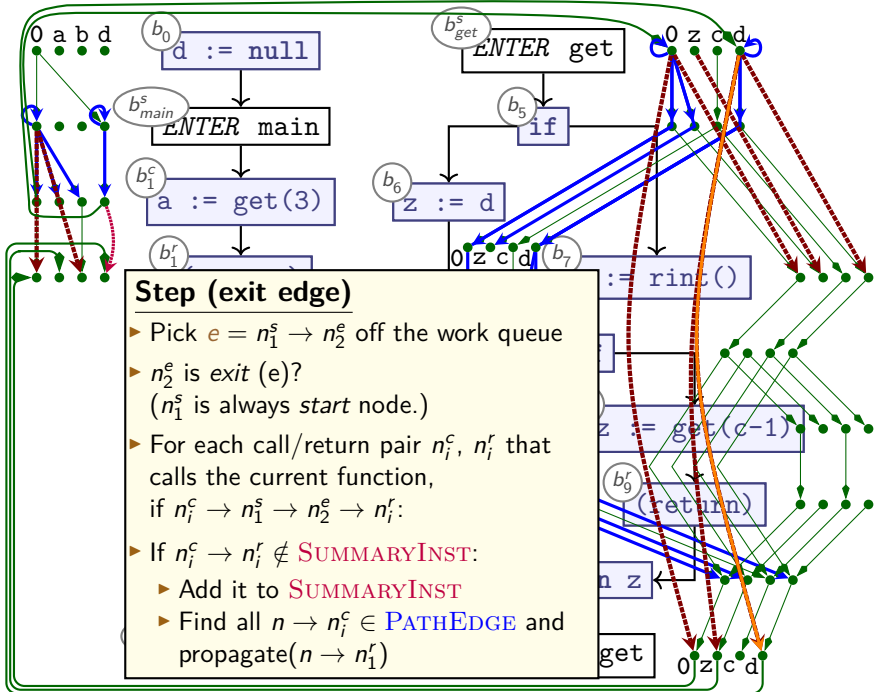






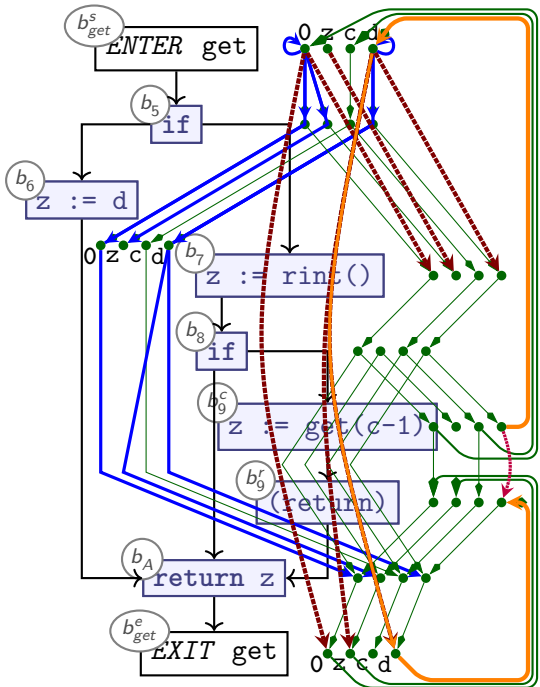
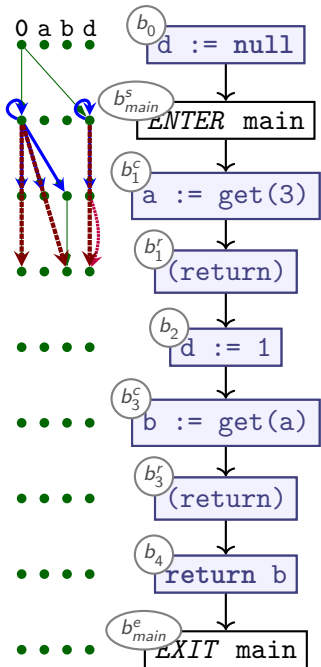


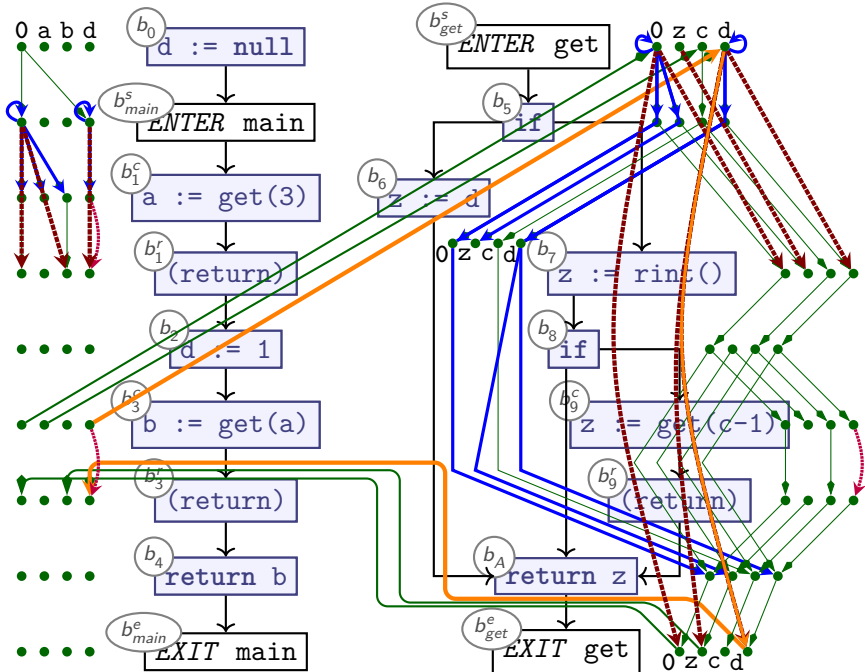


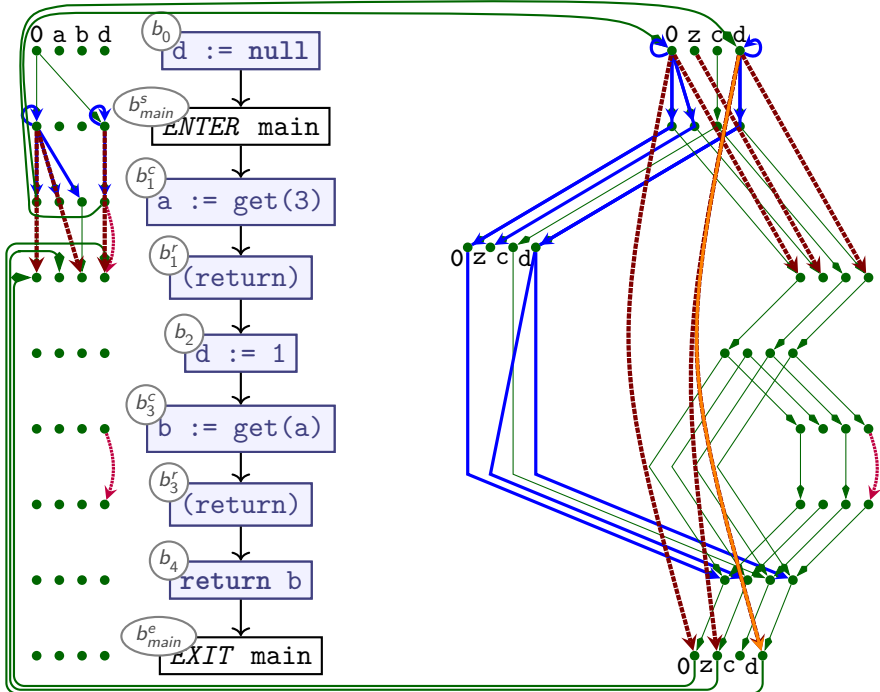


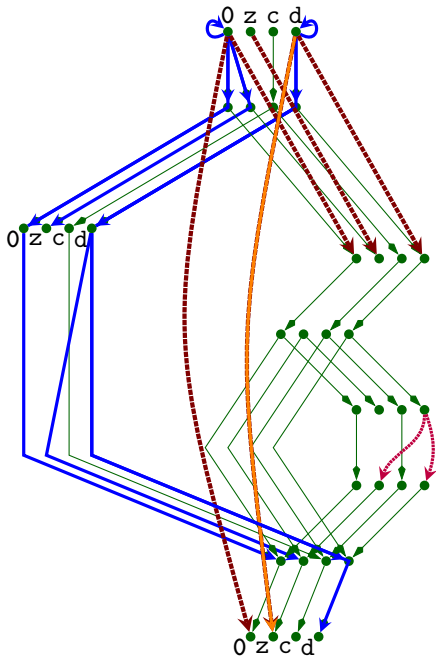
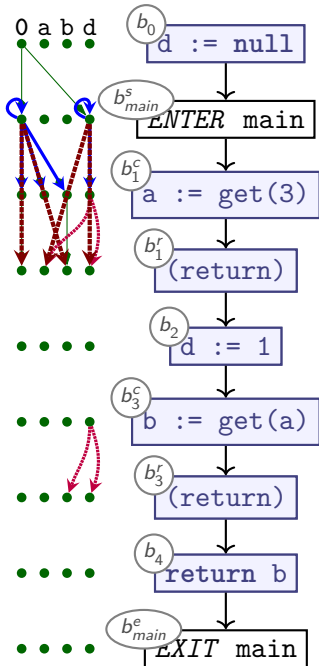
Step (exit edge)

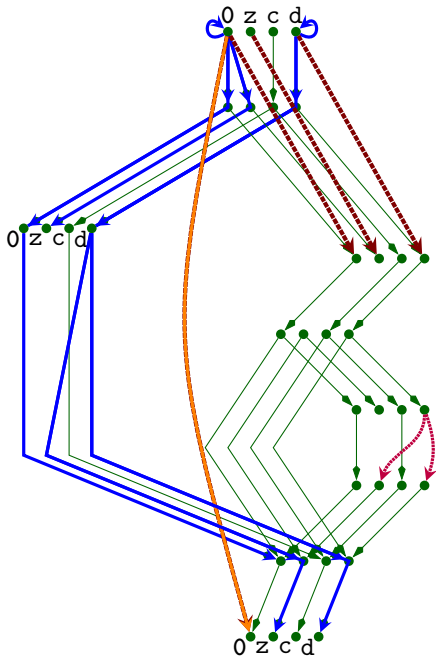
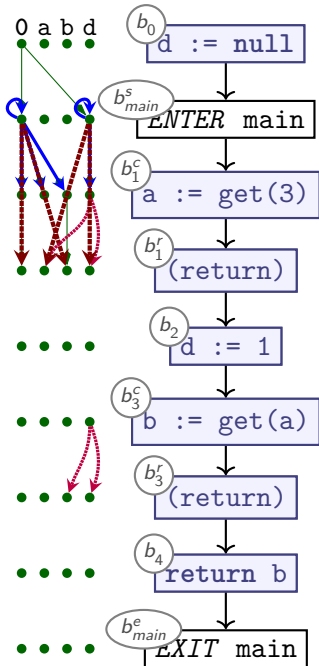
- ▶ Pick $e = n_1^s \rightarrow n_2^e$ off the work queue
- ▶ n_2^e is exit (e)?
(n_1^s is always *start* node.)
- ▶ For each call/return pair n_i^c, n_i^r that calls the current function, if $n_i^c \rightarrow n_1^s \rightarrow n_2^e \rightarrow n_i^r$:
- ▶ If $n_i^c \rightarrow n_i^r \notin \text{SUMMARYINST}$:
 - ▶ Add it to **SUMMARYINST**
 - ▶ Find all $n \rightarrow n_i^c \in \text{PATHEDGE}$ and propagate($n \rightarrow n_i^r$)

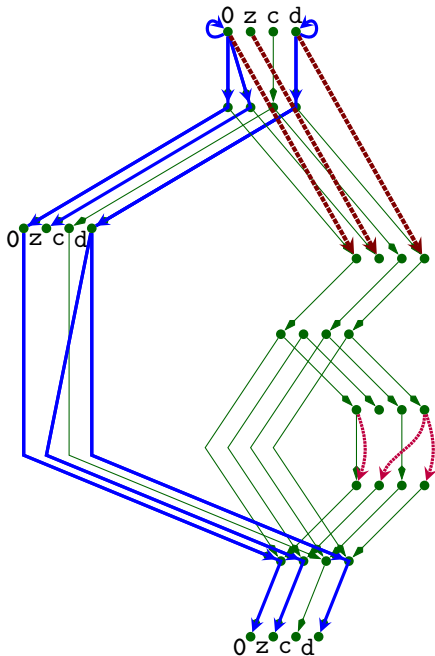
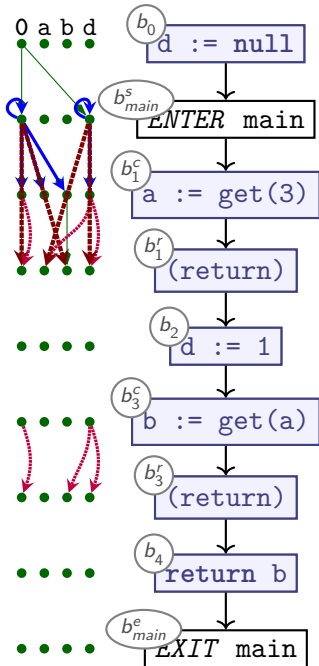


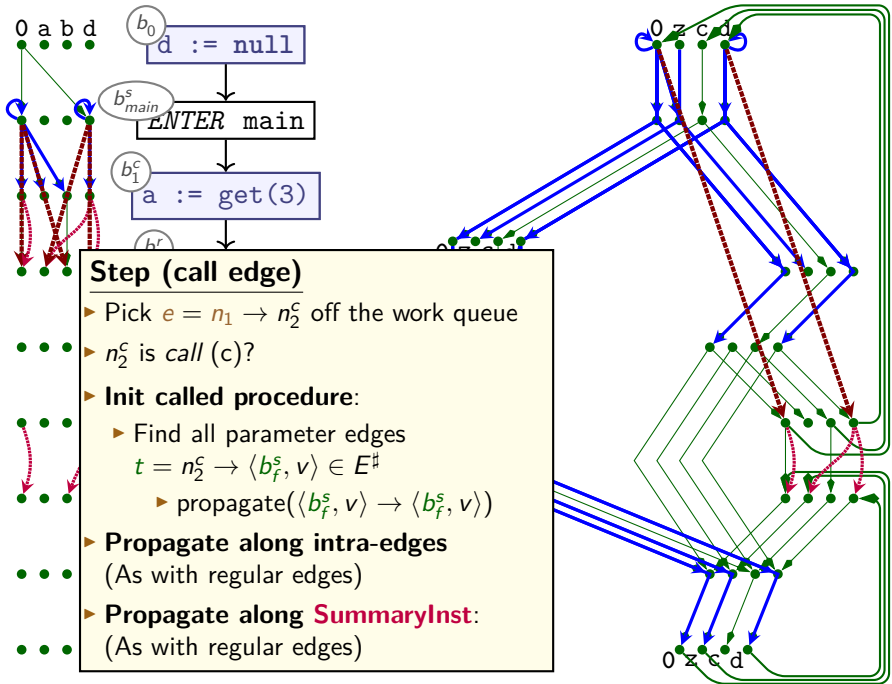


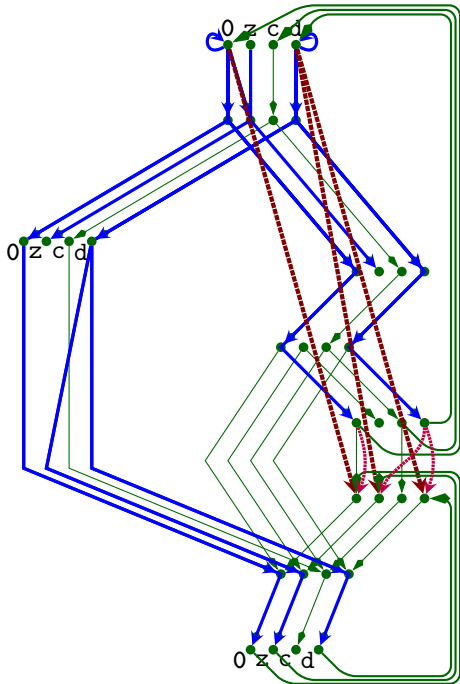
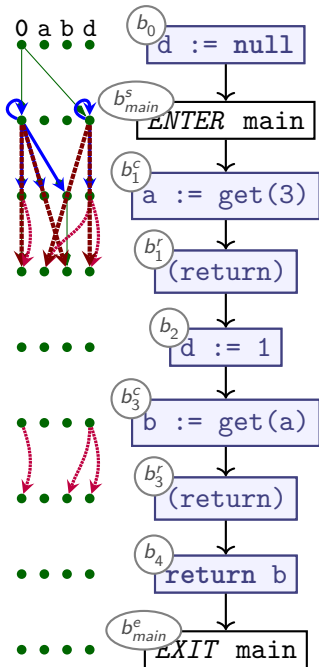


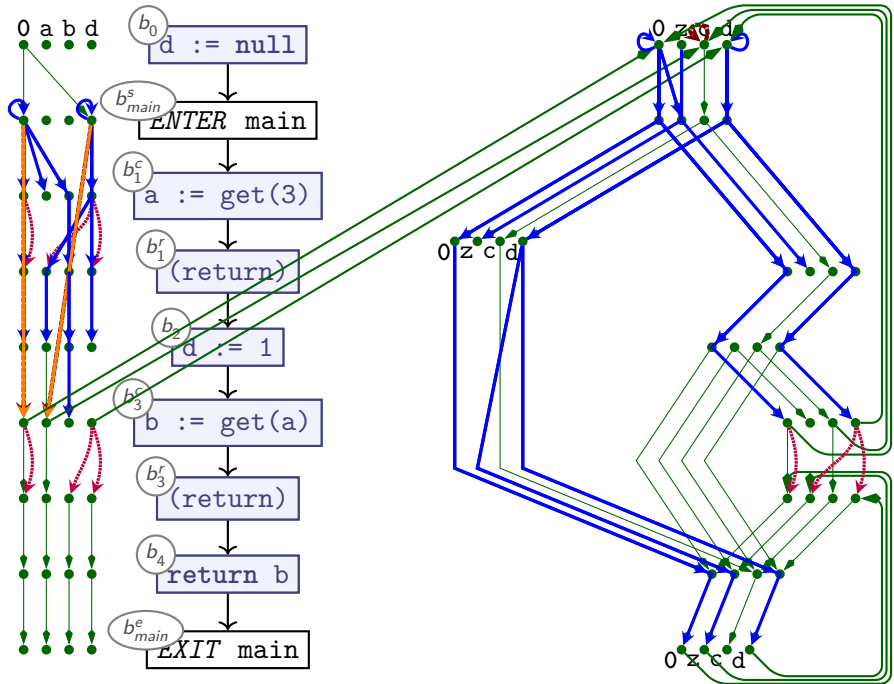


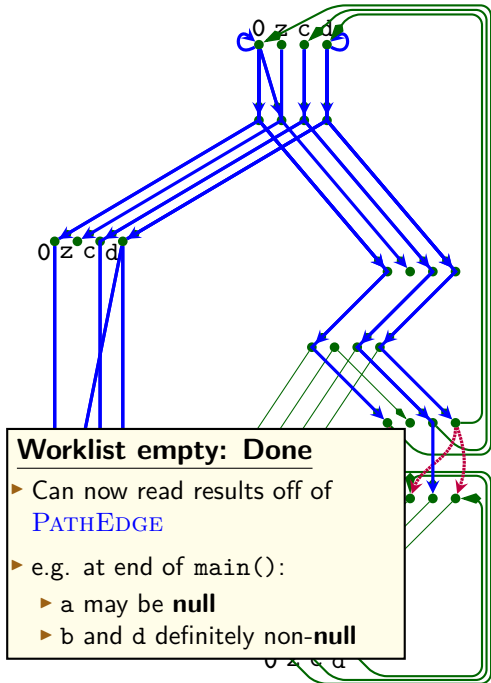
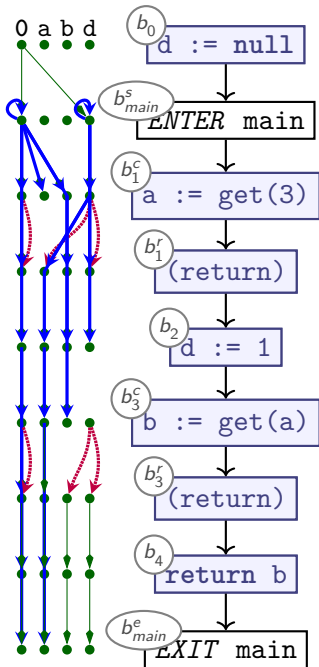












Worklist empty: Done

- ▶ Can now read results off of **PATHEDGE**
- ▶ e.g. at end of main():
 - ▶ a may be **null**
 - ▶ b and d definitely **non-null**

The IFDS Algorithm: Initialisation and Propagation)

Procedure Init():

begin

WORKLIST := **PATHEDGE** := \emptyset

propagate($\langle b_{\text{main}}^s, \mathbf{0} \rangle \rightarrow \langle b_{\text{main}}^s, \mathbf{0} \rangle$)

ForwardTabulate()

end

Procedure propagate($n_1 \rightarrow n_2$):

begin

if $n_1 \rightarrow n_2 \in \mathbf{PATHEDGE}$ **then**

return

$\mathbf{PATHEDGE} := \mathbf{PATHEDGE} \cup \{n_1 \rightarrow n_2\}$

$\mathbf{WORKLIST} := \mathbf{WORKLIST} \cup \{n_1 \rightarrow n_2\}$

end

IFDS: Forward Tabulation

Procedure ForwardTabulate():

begin

while $n_0 \rightarrow n_1 \in \text{WORKLIST}$ **do**

WorkList := **WorkList** $\setminus \{n_0 \rightarrow n_1\}$

$\langle b_0, v_0 \rangle = n_0$; $\langle b_1, v_1 \rangle = n_1$

if b_1 is neither *Call* nor *Exit* node **then**

foreach $n_1 \rightarrow n_2 \in E^\sharp$:

propagate($n_0 \rightarrow n_2$)

else if b_1 is *Call* node **then begin**

foreach call edge $n_1 \rightarrow n_2 \in E^\sharp$:

propagate($n_2 \rightarrow n_2$)

foreach non-call edge $n_1 \rightarrow n_2 \in E^\sharp \cup \text{SUMMARYINST}$:

propagate($n_0 \rightarrow n_2$)

end else if b_1 is *Exit* node **then begin**

foreach caller/return node pair b_i^c, b_i^r that calls b_0 and vars v_0, v_1 **do**

$n_s = \langle b_i^c, v_0 \rangle$; $n_r = \langle b_i^r, v_1 \rangle$

if $\{n_s \rightarrow n_0, n_0 \rightarrow n_1, n_1 \rightarrow n_r\} \subseteq E^\sharp$ **and not** $n_s \rightarrow n_r \in \text{SUMMARYINST}$ **then**

SUMMARYINST := **SUMMARYINST** $\cup \{n_s \rightarrow n_r\}$

foreach $n_z \rightarrow n_s \in \text{PATHEDGE}$:

propagate(n_z, n_r)

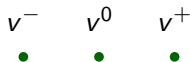
end done end done end

Summary: IFDS Algorithm

- ▶ Computes yes-or-no analysis on all variables
 - ▶ Original notion of 'variables' is slightly broader)
- ▶ Represents facts-of-interest as nodes $\langle b, v \rangle$:
 - ▶ b is node (basic block) in CFG
 - ▶ v is variable that we are interested in
- ▶ Uses
 - ▶ '*Exploded Supergraph*' $G^\#$
 - ▶ All CFGs in program in one graph
 - ▶ Plus interprocedural call edges
 - ▶ *Representation relations*
 - ▶ *Graph reachability*
 - ▶ *A worklist*
- ▶ Distinguishes between *Call* nodes, *Exit* nodes, others
- ▶ **Demand-driven**: only analyses what it needs
- ▶ **Whole-program analysis**
- ▶ **Computes Least Fixpoint on distributive frameworks**

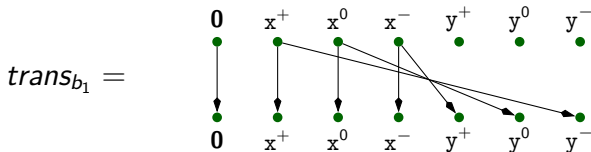
BONUS SLIDES

Beyond True and False



- ▶ What if abstract domain is not boolean?
 - ▶ e.g., $\{\top, A^+, A^-, A^0, \perp\}$
- ▶ Multiple boolean properties per variable
 - ▶ easy for powerset lattice $\mathcal{P}(\{+, -, 0\})$
- ▶ *Limitation*: Transfer functions only depend on one variable
- ▶ Some problems not representable, others must adapt lattice

Consider $b_1 = \boxed{y := 0 - x}$:



This is how the algorithm was originally proposed

Extending IFDS?

- ▶ Not all analyses map well to IFDS
- ▶ Core ideas are appealing:
 - ▶ Automatically compute procedure summaries
 - ▶ Exploit graph reachability + worklist for *dependency tracking*

It is possible to extend this to other classes of problems

Linear Reaching Values

Statement	in_b	out_b
$x := 42$	M	M with $[x \mapsto 42]$
$x := y + 1$	$M = \{[y \mapsto c], \dots\}$	M with $[x \mapsto c + 1]$
$x := y * 7$	$M = \{[y \mapsto c], \dots\}$	M with $[x \mapsto c * 7]$
$x := y + z$	M	M with $[x \mapsto \top]$

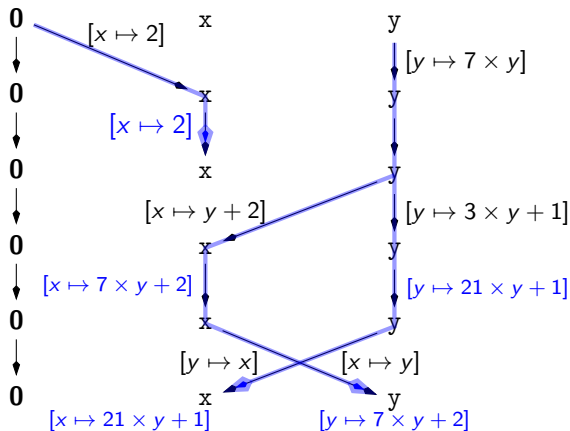
- ▶ “ M with $[x \mapsto e]$ ” means “Remove from M any $[x \mapsto \dots]$ if it exists, and then add $[x \mapsto e]$ ”.
- ▶ The above sketches a *distributive* reaching values analysis
 - ▶ Each annotation of form $v_1 \mapsto c_1 \times v_2 + c_2$
 - ▶ No support for adding / multiplying / ... multiple variables
- ▶ Encode in IFDS?

Labelling Graph Edges

```
x := 2;  
y := y * 7;
```

```
x := y + 2;  
y := 3 * y + 1;
```

```
{ t := x;  
  x := y;  
  y := t; }
```



- ▶ Extending IFDS to support information processing
- ▶ Carrying over key techniques:
 - ▶ *Track dependencies*
 - ▶ *Generate procedure summaries on the fly*

Representation

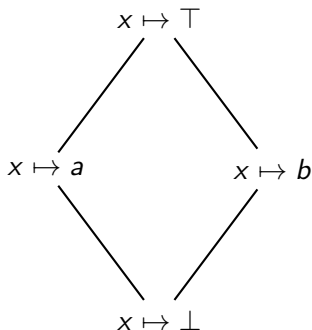
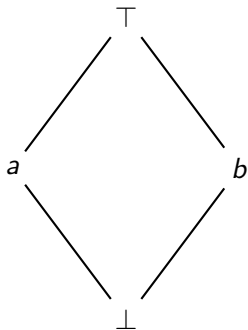
$$\left\{ \begin{array}{l} [x \mapsto c_{x,2} \times x + d_{x,2}] \\ [y \mapsto c_{y,2} \times y + d_{y,2}] \end{array} \right\} \circ \left\{ \begin{array}{l} [x \mapsto c_{x,1} \times v + d_{x,1}] \\ [y \mapsto c_{y,1} \times w + d_{y,1}] \end{array} \right\} \\ = \\ \left\{ \begin{array}{l} [x \mapsto (c_{x,1} \times c_{x,2}) \times v + (d_{x,2} + c_{x,2} \times d_{x,1})] \\ [y \mapsto (c_{y,1} \times c_{y,2}) \times w + (d_{y,2} + c_{y,2} \times d_{y,1})] \end{array} \right\}$$

- ▶ c_i, d_i : constants
- ▶ v, w : program variables
- ▶ (Maps of) linear functions are closed under composition
- ▶ Must support \sqcup to merge, map to \perp on mismatch

$$\left\{ \begin{array}{l} [x \mapsto c_{x,1} \times v_1 + d_{x,1}] \\ [y \mapsto c_{y,1} \times v_3 + d_{y,1}] \end{array} \right\} \sqcup \left\{ \begin{array}{l} [x \mapsto c_{x,1} \times v_1 + d_{x,1}] \\ [y \mapsto c_{y,2} \times v_2 + d_{y,2}] \end{array} \right\} \\ = \\ \left\{ \begin{array}{l} [x \mapsto c_{x,1} \times x + d_{x,1}] \\ [y \mapsto \perp] \end{array} \right\}$$

Micro-Functions and Lattices

- ▶ Extend lattices to such 'Micro-Functions':



Micro-Functions, Efficient Representation

- ▶ Micro-Functions must support:

Encoding		$O(1)$ space
Computation	$f(x)$	$O(1)$ time
Equality testing	$f = f'$	$O(1)$ time
Composition	$f \circ f'$	$O(1)$ time
Meet	$f \sqcup f'$	$O(1)$ time

- ▶ Micro-functions are **efficiently representable** if they satisfy space / time constraints
 - ▶ Required for the algorithm's time bounds
- ▶ Other examples:
 - ▶ IFDS problems
 - ▶ Value bounds analysis

The IDE Algorithm (1/1)

- ▶ Interprocedural **D**istributive **E**nvironments algorithm
- ▶ Extends IFDS to 'labelled' edges as described above
- ▶ Assumes distributive framework over micro-functions
- ▶ Algorithmic changes:
 - ▶ First phase analogous to IFDS
 - ▶ Second phase applies computed functions to read out results
- ▶ Maintain/update mapping from path edges to micro-functions f :

$$\text{PATHEDGE} = \{ \langle b_0, v_0 \rangle \xrightarrow{f_0} \langle b_1, v_1 \rangle, \dots \}$$

- ▶ 'Missing edges' equivalent to $x \mapsto \perp$
- ▶ Initialise:

$$\text{PATHEDGE} = \{ \langle b_0, v_0 \rangle \xrightarrow{v_1 \mapsto \perp} \langle b_1, v_1 \rangle, \dots \}$$

- ▶ Always exactly one f per $\{ \langle b_0, v_0 \rangle \xrightarrow{f} \langle b_1, v_1 \rangle \} \in \text{PATHEDGE}$

The IDE Algorithm (2/2)

Procedure propagate($n_1 \rightarrow n_2$): -- IFDS version

begin

if $n_1 \rightarrow n_2 \in \text{PATHEDGE}$ **then**

return

$\text{PATHEDGE} := \text{PATHEDGE} \cup \{n_1 \rightarrow n_2\}$

$\text{WORKLIST} := \text{WORKLIST} \cup \{n_1 \rightarrow n_2\}$

end

⇓

Procedure propagate_{IDE}($n_1 \xrightarrow{f} n_2$): -- IDE version

begin

let $n_1 \xrightarrow{f'} n_2 \in \text{PATHEDGE}$

$f_{\text{upd}} := f \sqcup f'$

if $f_{\text{upd}} = f'$ **then**

return

$\text{PATHEDGE} := (\text{PATHEDGE} \setminus \{n_1 \xrightarrow{f'} n_2\}) \cup \{n_1 \xrightarrow{f_{\text{upd}}} n_2\}$

$\text{WORKLIST} := \text{WORKLIST} \cup \{n_1 \rightarrow n_2\}$

end

Summary

- ▶ IDE strictly generalises IFDS
- ▶ Utilises **Micro-Functions** to ensure efficient summaries:
 - ▶ Intra-procedural summaries via `PATHEDGE`
 - ▶ Inter-procedural procedure summaries via `SUMMARYINST`
- ▶ Runtime is $O(LED^3)$ if micro-functions are **efficiently representable**
 - ▶ L : Lattice height
 - ▶ IFDS: 1
 - ▶ IDE: length of longest descending chain
 - ▶ E : Number of control-flow edges
 - ▶ D : Number of variables
- ▶ IFDS supported by many popular dataflow frameworks