

Precise Interprocedural Dataflow Analysis via Graph Reachability

Thomas Reps

Susan Horwitz

Mooly Sagiv

POPL 1995

IFDS

Thomas Reps

Susan Horwitz
POPL 1995

Mooly Sagiv

The IFDS "framework"

- a *model* for dataflow problems
- a *uniform solution* to these problems
- a polynomial-time *algorithm*

Reaching Definitions

What statements are affected by a given definition?

*there is a path
from $x := e'$ to e*

*x is not re-defined
along the path*

$x := e'$ reaches e

Available Expressions

What expressions can a statement re-use?

*there is a path
from $x := \dots e' \dots;$
to e*

*no variable in e'
is re-defined along
the path*

e' is available at e

Live Variable Analysis

What variables are referenced at/after a statement?

*there is a path
from e to e'*

v is referenced at e'

v is live at e

Possibly-Uninitialized Variables

Which variables may be null at a given statement?

$x := e'$;

does not appear
on any path to **e**

x may be null at e

$x := e'$;

is the most recent binding

$\exists y \in e'$.

y may be null at e'

x may be null at e

Program Slicing

"The algorithm described in this paper yields an **improved interprocedural-slicing algorithm ...**

6x as fast as the Horwitz-Reps-Binkley algorithm."

Speeding up Slicing

Reps, Horwitz, Sagiv, Rosay; FSE '94

PLDI 1988



Interprocedural Slicing Using Dependence Graphs
Susan Horwitz, Thomas Reps, David Binkley

The Program Summary Graph and Flow-Sensitive Interprocedural Data-Flow Analysis
David Callahan

Interprocedural Side-Effect Analysis in Linear Time
Keith D. Cooper, Ken Kennedy

Possibly-Uninitialized Variables

Which variables may be null at a given statement?

$x := e'$;

does not appear
on any path to **e**

x may be null at e

$x := e'$;

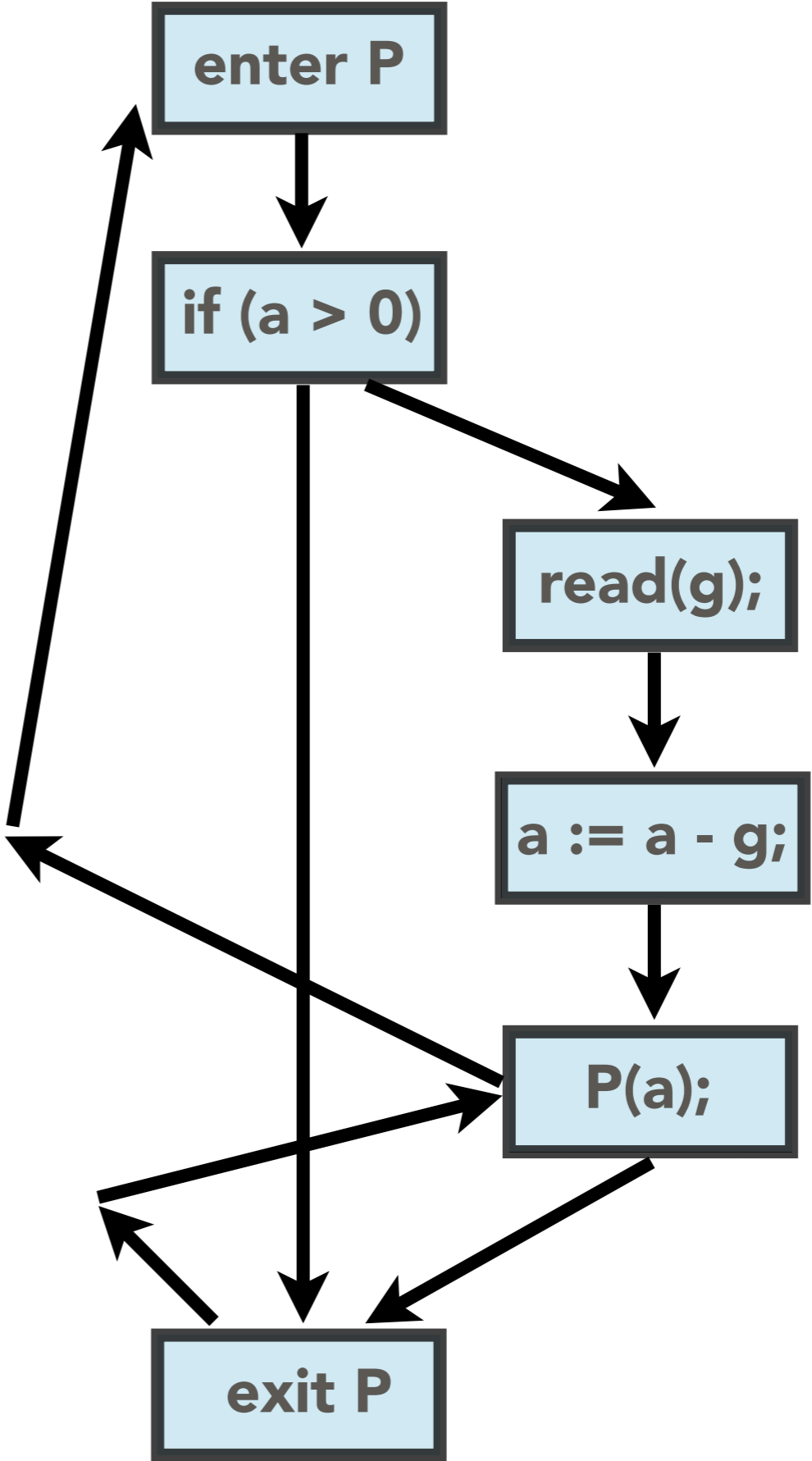
is the most recent binding

$\exists y \in e'$.

y may be null at e'

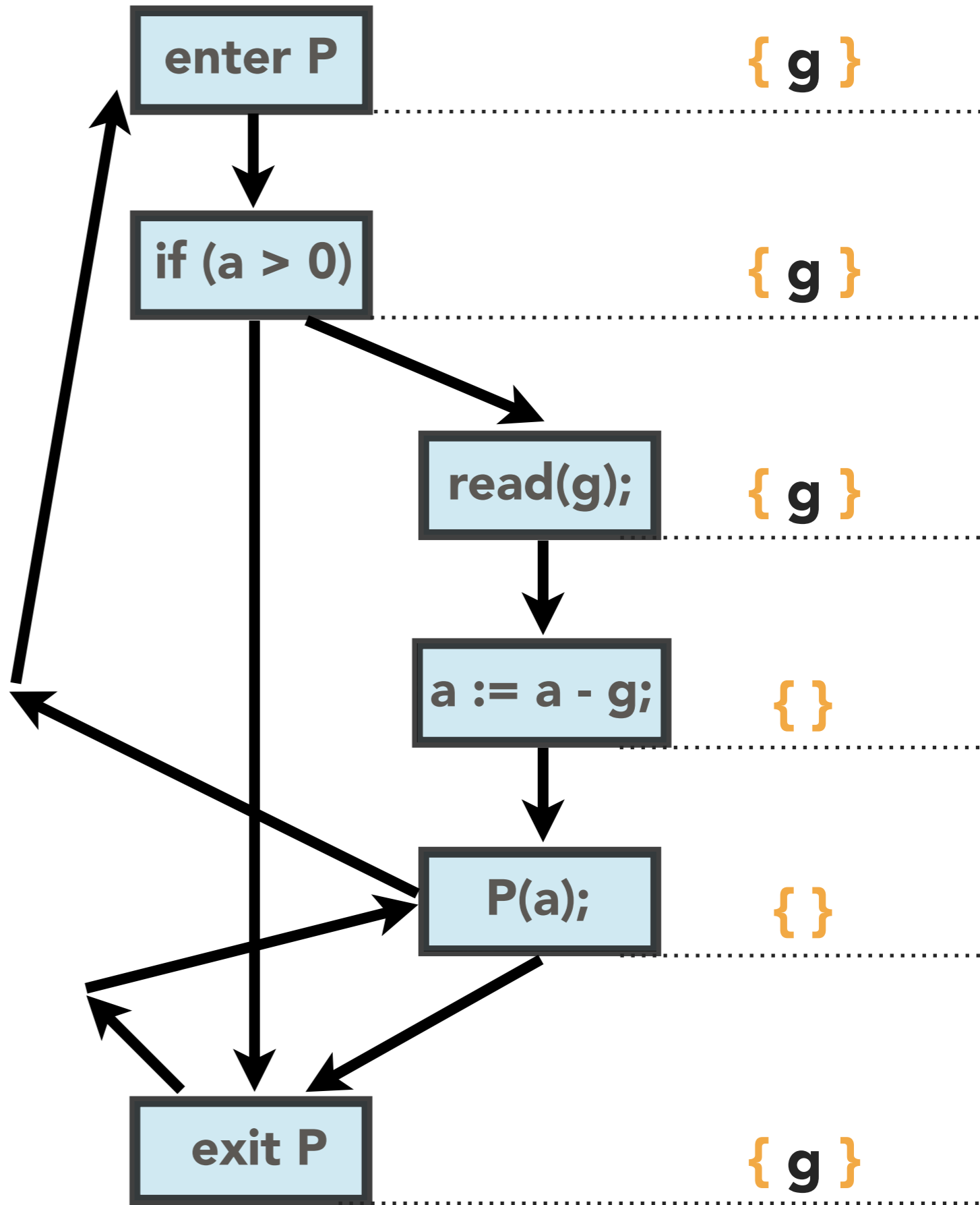
x may be null at e

```
int g;  
  
void P(int a) {  
    if (a > 0) {  
        read(g);  
        a := a - g;  
        P(a);  
    }  
}
```



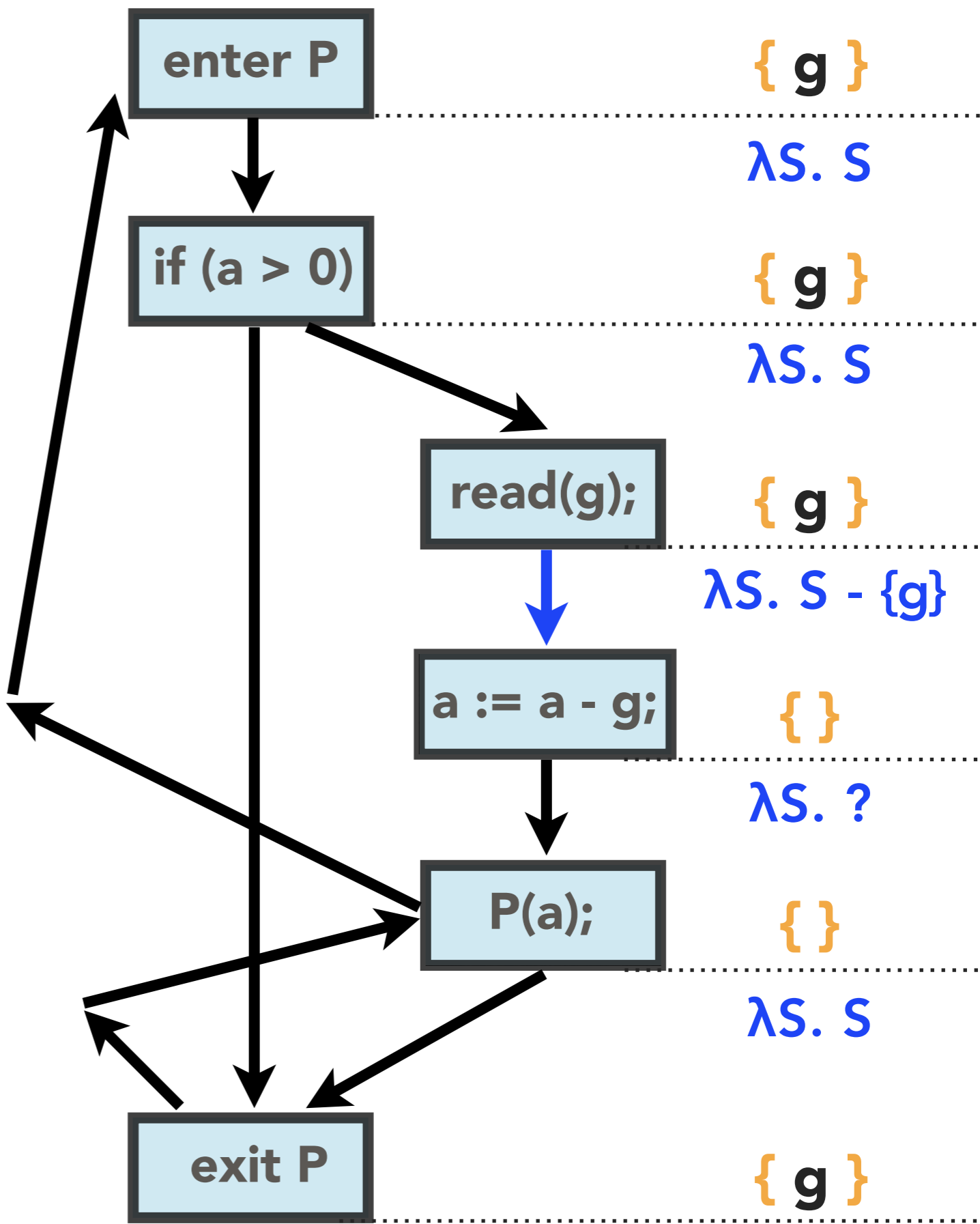
```
int g;

void P(int a) {
  if (a > 0) {
    read(g);
    a := a - g;
    P(a);
  }
}
```



```
int g;

void P(int a) {
  if (a > 0) {
    read(g);
    a := a - g;
    P(a);
  }
}
```



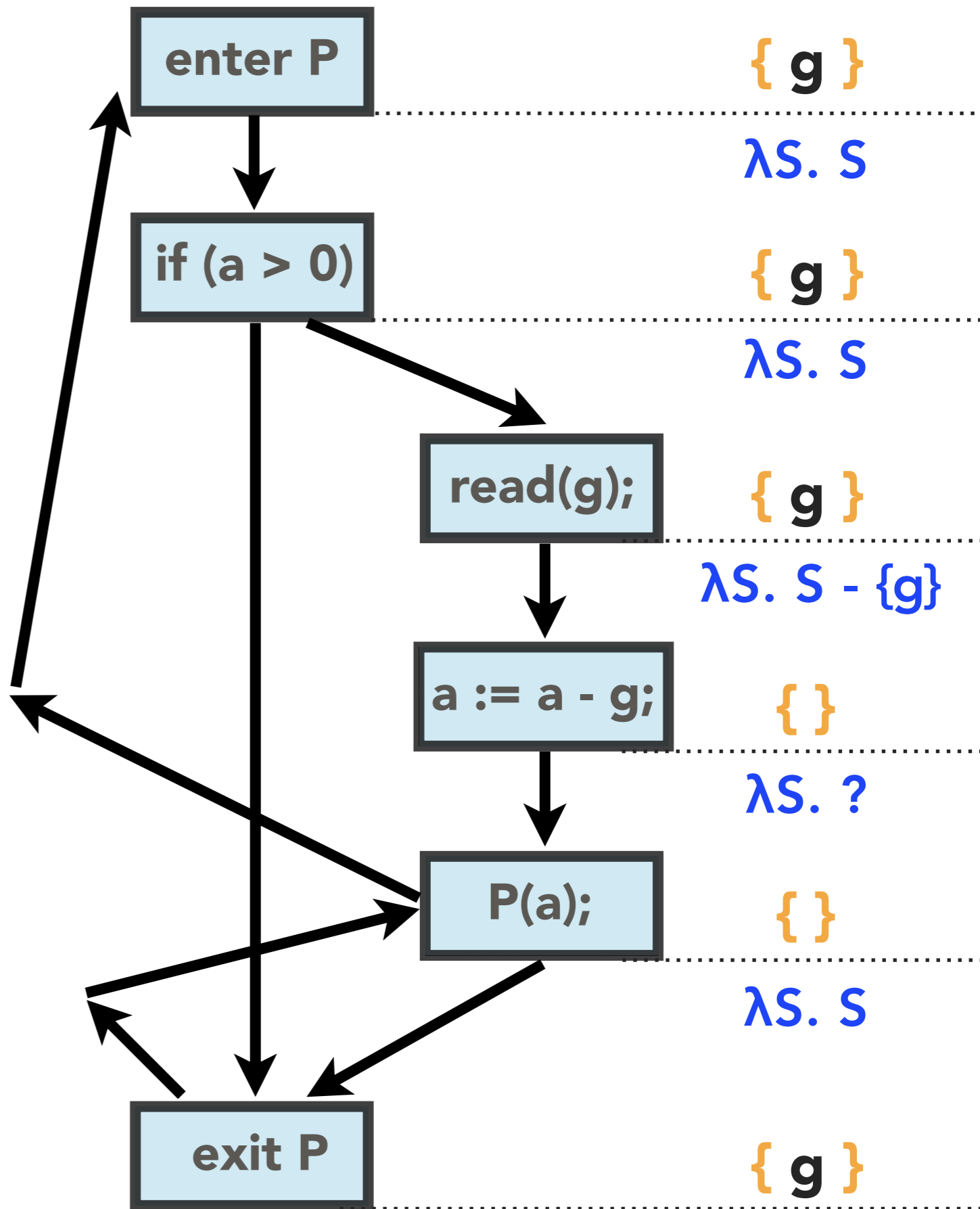
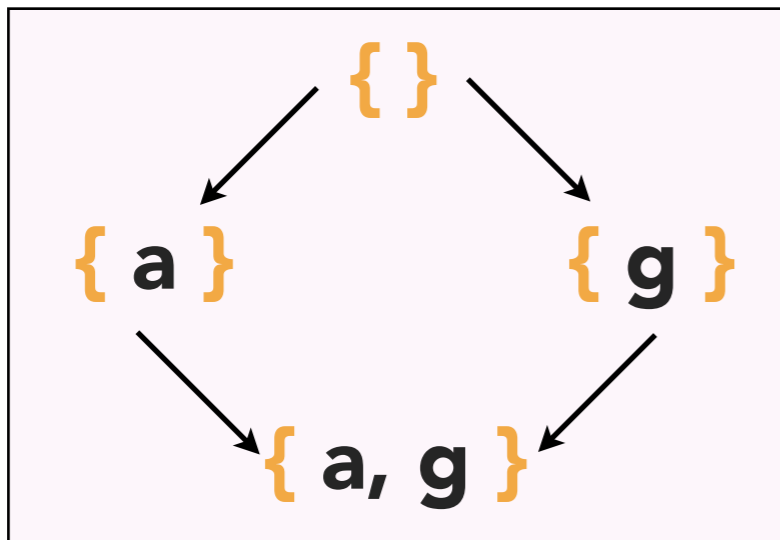
Kildall, POPL 1973

"Meet over all paths"

$$\bigcap \text{fn} (\dots (\text{f1} (\{ \mathbf{g} \})) \dots)$$

$\text{f1} \dots \text{fn} \in \text{AllPaths}$

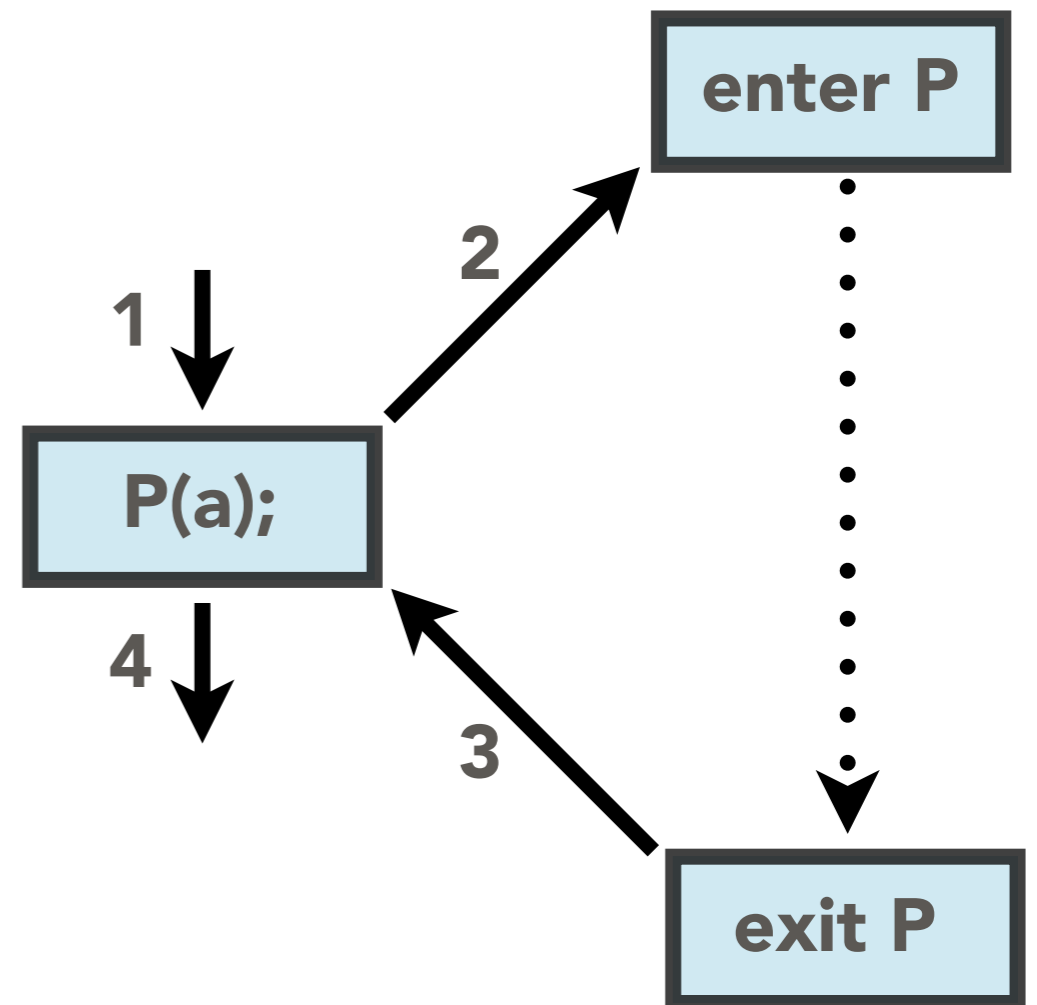
Infinite Set!



IFDS, POPL 1995

Meet over all **valid** paths

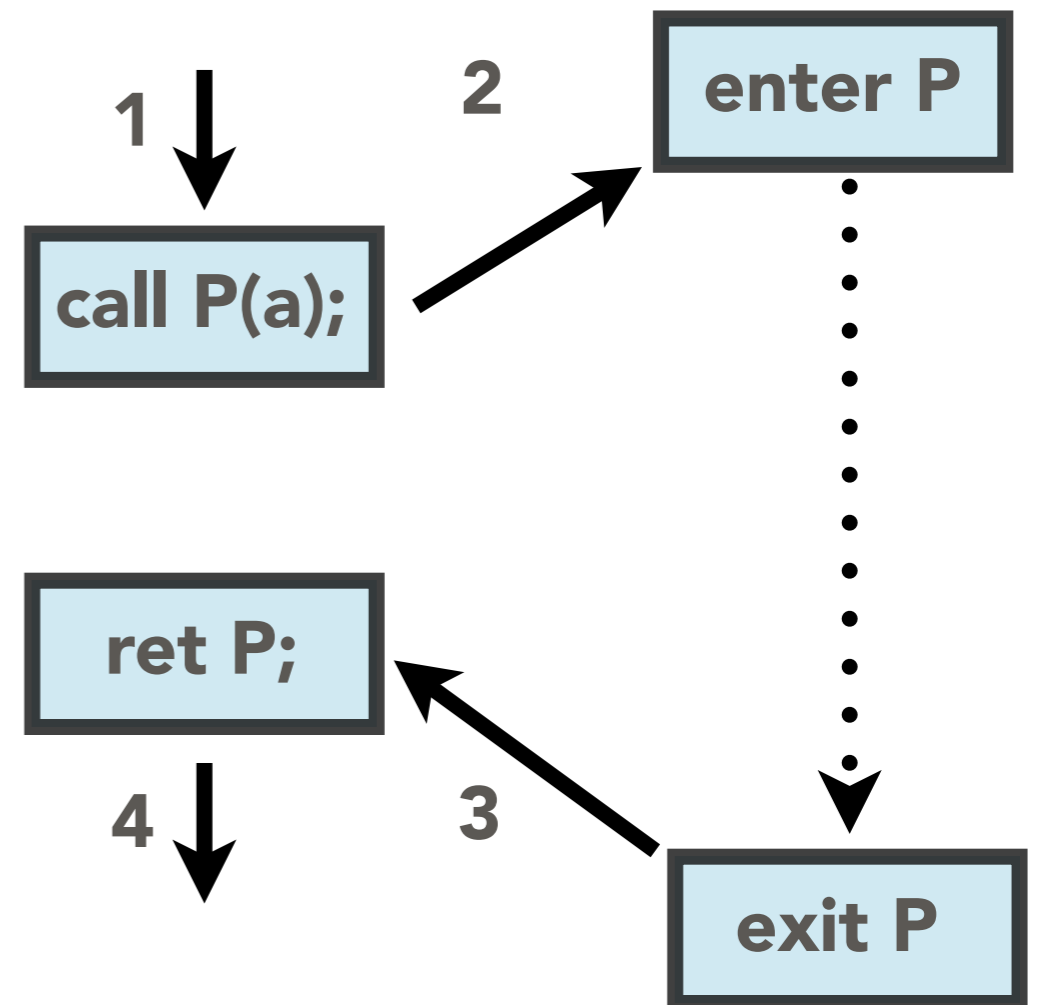
- Calls & Returns must match



IFDS, POPL 1995

Meet over all **valid** paths

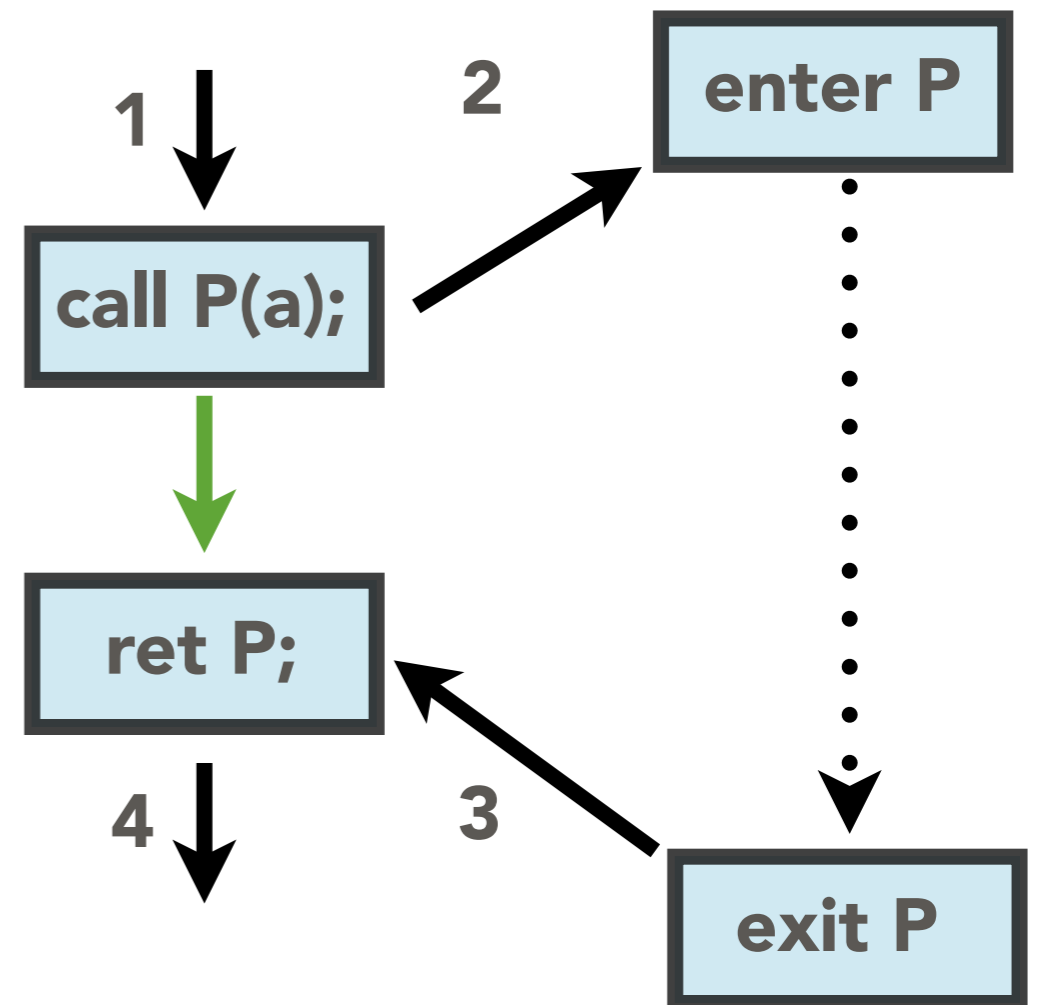
- Calls & Returns must match
- Enforced by **call** & **ret** nodes

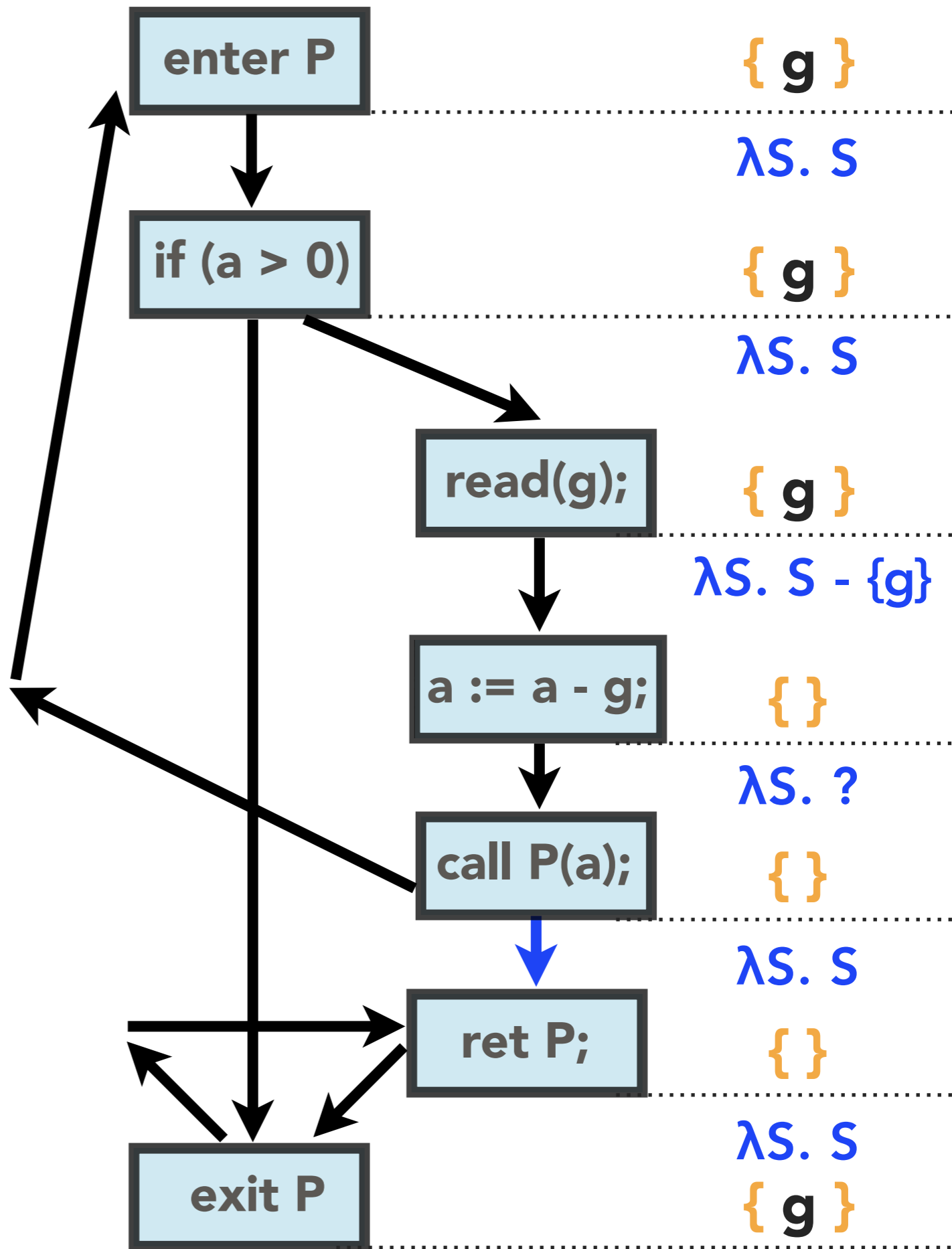


IFDS, POPL 1995

Meet over all **valid** paths

- Calls & Returns must match
- Enforced by **call** & **ret** nodes
- Track local variables with a **call-to-return edge**

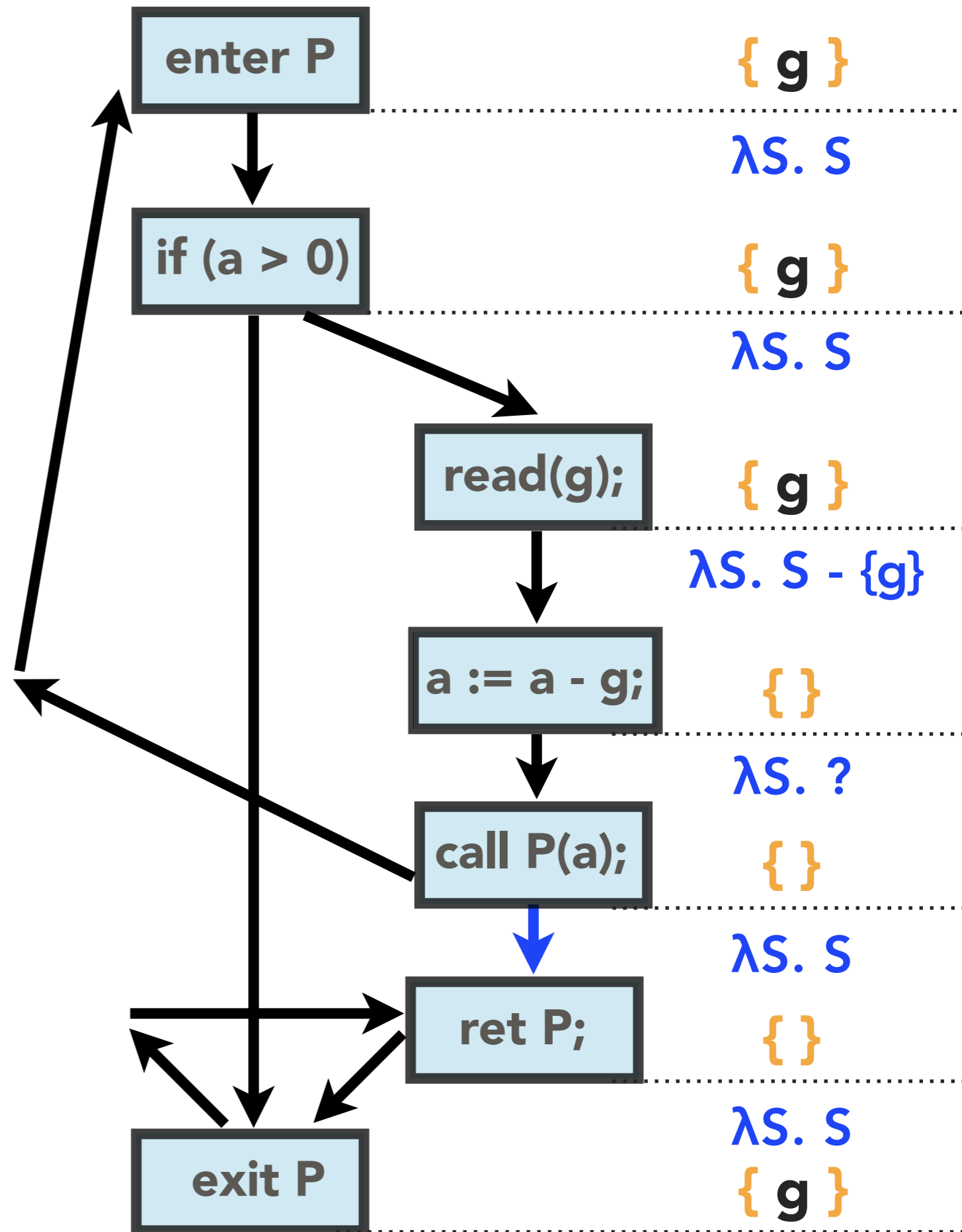




The "supergraph"

```
int g;  
  
void main(void) {  
  int x;  
  read(x);  
  P(x);  
}
```

Adds similar CFGs for other procedures



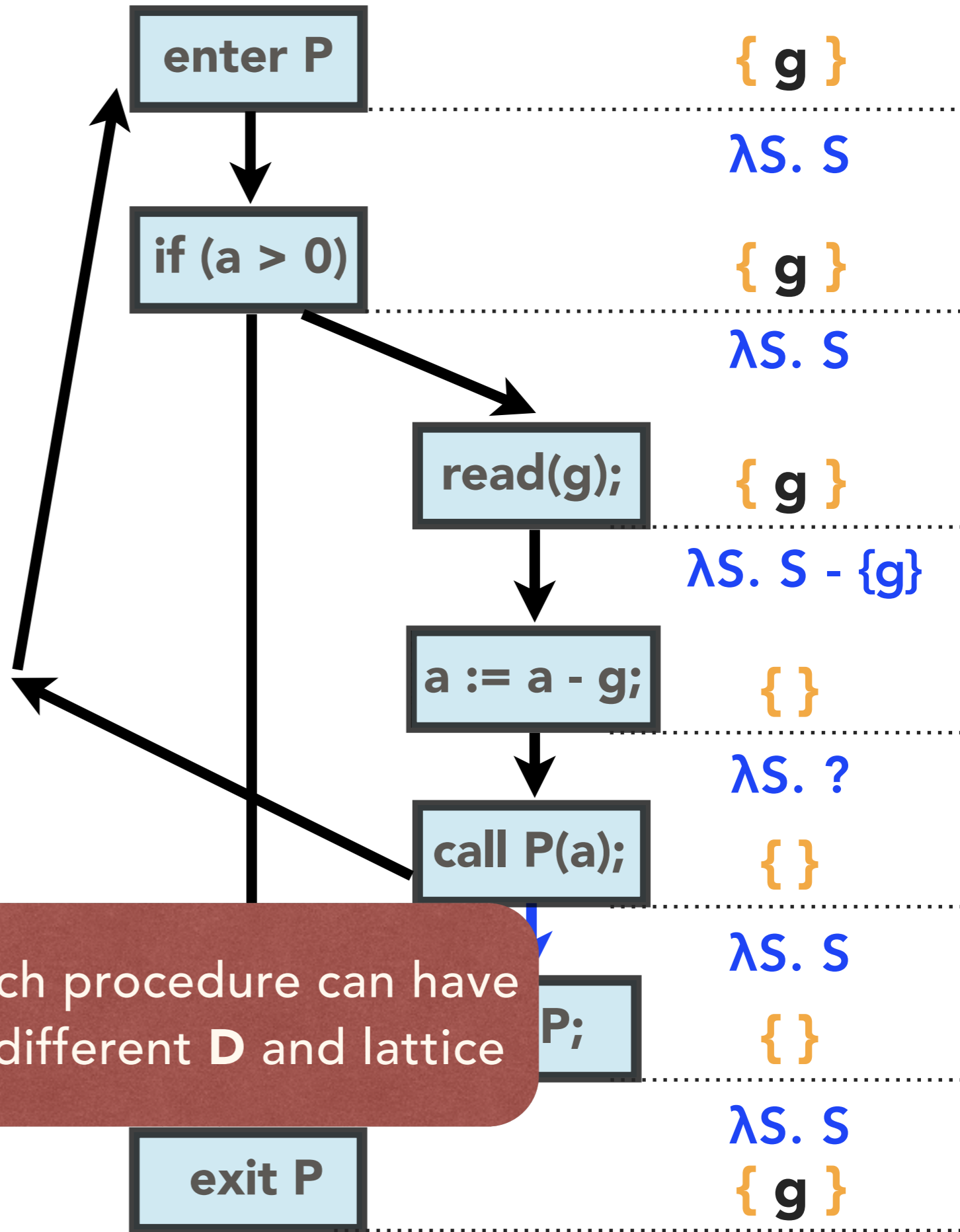
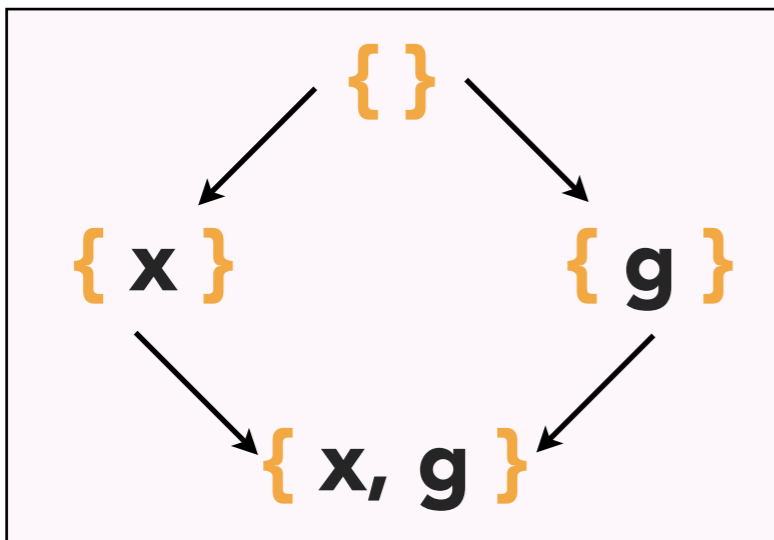
The "supergraph"

```

int g;

void main(void) {
  int x;
  read(x);
  P(x);
}

```



Each procedure can have a different **D** and lattice

An IFDS problem instance

- G = a supergraph
- D = a finite set (determines a lattice)
- F = a set of **distributive** functions over the lattice
- M = a map from edges in G to functions in F
- \sqcap = meet operator on the lattice

Interprocedural

Finite

Distributive

Subset

A few "IFDS" problems

D

- Reaching definitions All Variables
- Available Expressions All Expressions
- Live Variable Analysis All Variables
- Possibly-Uninitialized Variables All Variables
- Type Analysis Variables \times Types

In general, apply "path function" to a subset of D

$$\prod_{f_1 \dots f_n \in \text{All Valid Paths}} \text{fn} (\dots (\text{f1} (\{ \mathbf{g} \})) \dots)$$

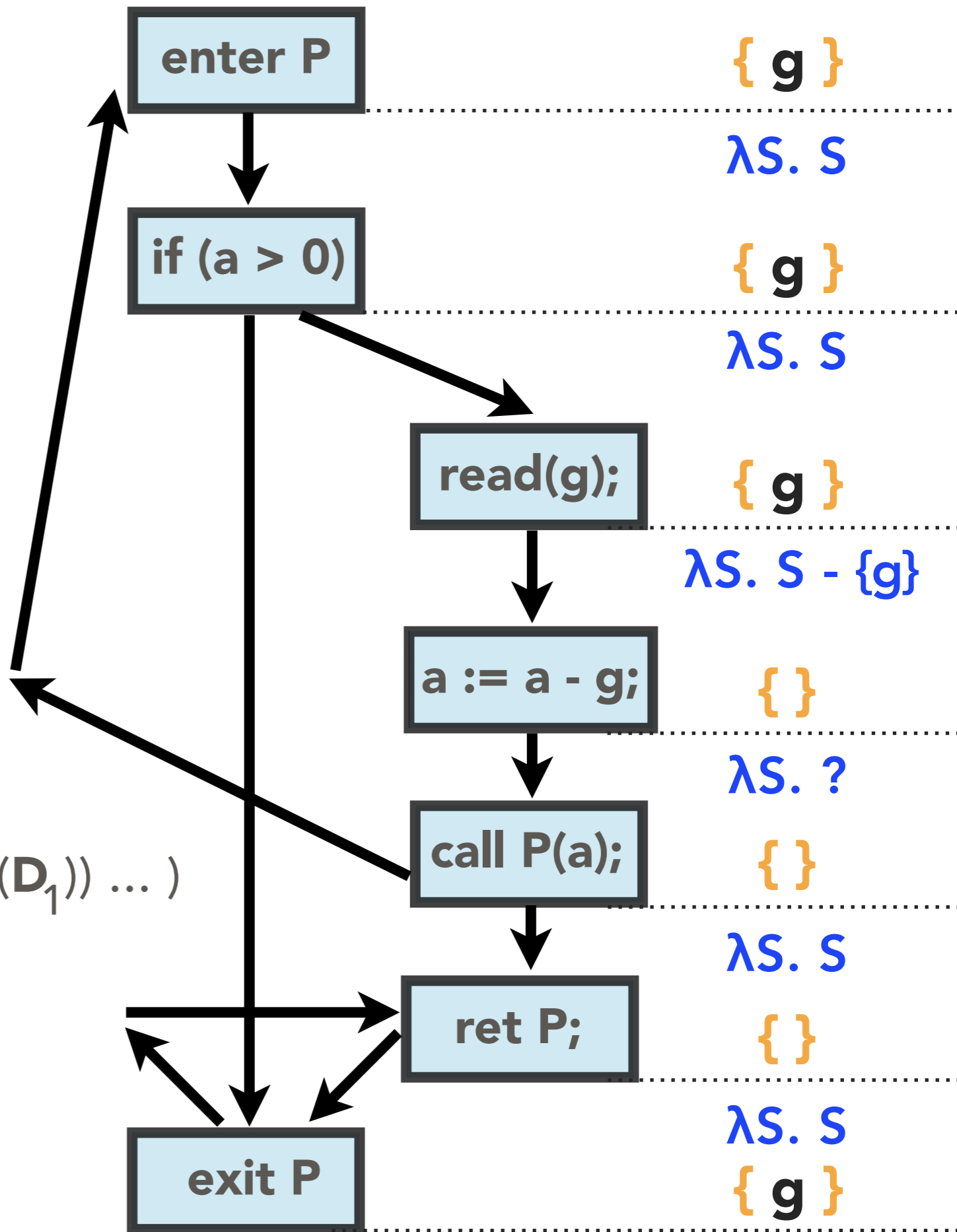
f1 ... fn ∈ All Valid Paths

... since each fi is distributive

$$f_i (D) = \{ f_i (D_1), \dots, f_i (D_k) \}$$

... therefore

$$\text{fn} (\dots (\text{f1} (D)) \dots) = \{ \text{fn} (\dots (\text{f1} (D_1)) \dots), \dots, \text{fn} (\dots (\text{f1} (D_k)) \dots) \}$$



$$\bigcap \text{fn} (\dots (\text{f1} (\mathbf{D})) \dots)$$

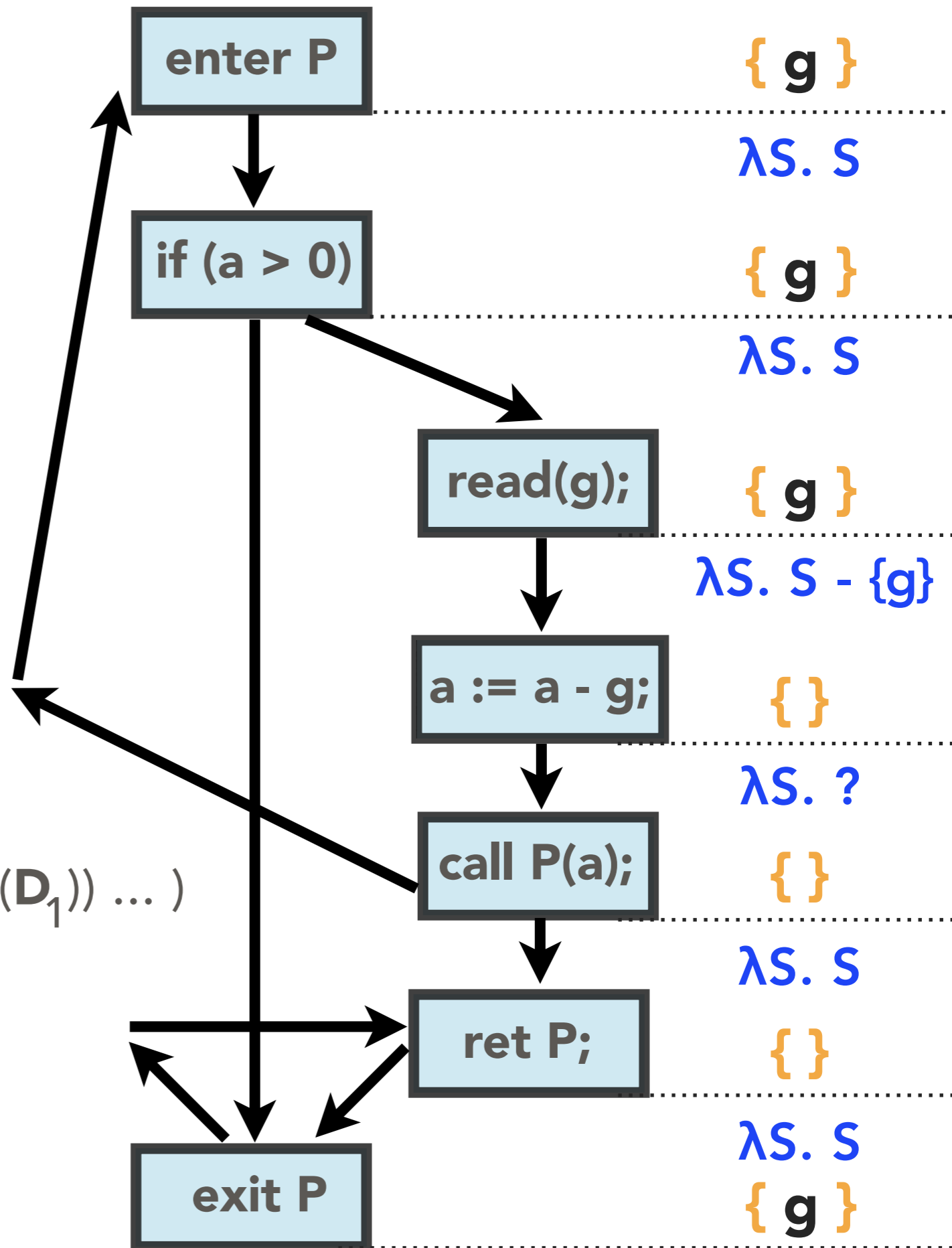
$\text{f1} \dots \text{fn} \in \text{All Valid Paths}$

... since each fi is distributive

$$\text{fi} (\mathbf{D}) = \{ \text{fi} (\mathbf{D}_1) , \dots , \text{fi} (\mathbf{D}_k) \}$$

... therefore

$$\text{fn} (\dots (\text{f1} (\mathbf{D})) \dots) = \{ \text{fn} (\dots (\text{f1} (\mathbf{D}_1)) \dots) , \dots , \text{fn} (\dots (\text{f1} (\mathbf{D}_k)) \dots) \}$$



$$\bigcap \text{fn} (\dots (\text{f1} (\mathbf{D})) \dots)$$

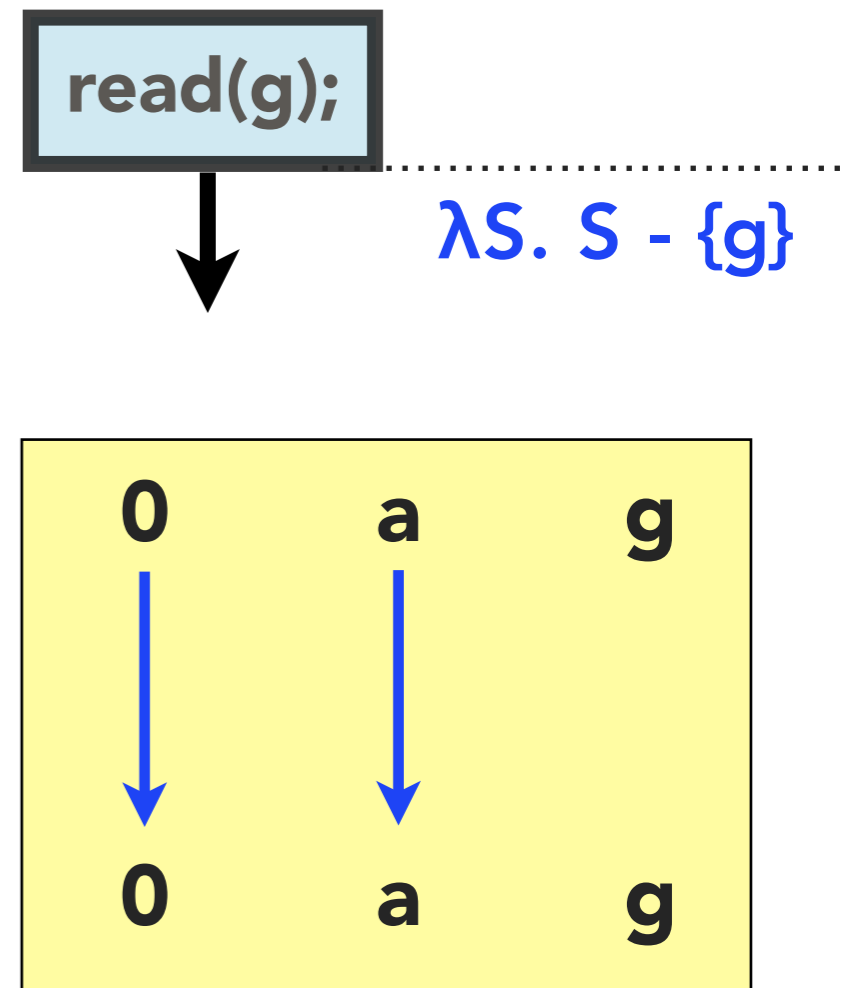
f1 ... **fn** ∈ All Valid Paths

... since each **fi** is distributive

$$\text{fi} (\mathbf{D}) = \{ \text{fi} (\mathbf{D}_1) , \dots , \text{fi} (\mathbf{D}_k) \}$$

... therefore

$$\text{fn} (\dots (\text{f1} (\mathbf{D})) \dots) = \{ \text{fn} (\dots (\text{f1} (\mathbf{D}_1)) \dots) , \dots , \text{fn} (\dots (\text{f1} (\mathbf{D}_k)) \dots) \}$$



$$\bigcap \text{fn} (\dots (\text{f1} (\mathbf{D})) \dots)$$

f1 ... **fn** ∈ All Valid Paths

... since each **fi** is distributive

$$\text{fi} (\mathbf{D}) = \{ \text{fi} (\mathbf{D}_1) , \dots , \text{fi} (\mathbf{D}_k) \}$$

... therefore

$$\text{fn} (\dots (\text{f1} (\mathbf{D})) \dots) = \{ \text{fn} (\dots (\text{f1} (\mathbf{D}_1)) \dots) , \dots , \text{fn} (\dots (\text{f1} (\mathbf{D}_k)) \dots) \}$$

a := a - g;



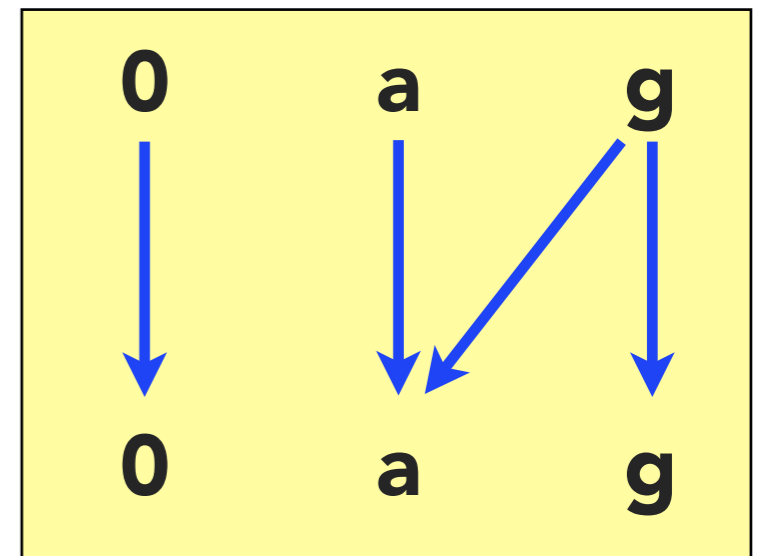
$\lambda S.$

if $a \in S$

or $g \in S$

then $S \cup \{a\}$

else $S - \{a\}$



$$\bigcap \text{fn} (\dots (\mathbf{f1} (\mathbf{D})) \dots)$$

$\mathbf{f1} \dots \mathbf{fn} \in$ All Valid Paths

... since each \mathbf{fi} is distributive

$$\mathbf{fi} (\mathbf{D}) = \{ \mathbf{fi} (\mathbf{D}_1) , \dots , \mathbf{fi} (\mathbf{D}_k) \}$$

... therefore

$$\text{fn} (\dots (\mathbf{f1} (\mathbf{D})) \dots) = \{ \text{fn} (\dots (\mathbf{f1} (\mathbf{D}_1)) \dots) , \dots_k \}$$

For any \mathbf{f} ,

$$\mathbf{0} \longrightarrow \mathbf{0}$$

$$\mathbf{0} \longrightarrow \mathbf{y} \quad \text{if } \mathbf{f}(\mathbf{0}) = \mathbf{y}$$

$$\mathbf{x} \longrightarrow \mathbf{y} \quad \text{if } \mathbf{f}(\mathbf{x}) = \mathbf{y}$$

and $\mathbf{f}(\mathbf{0}) \neq \mathbf{y}$

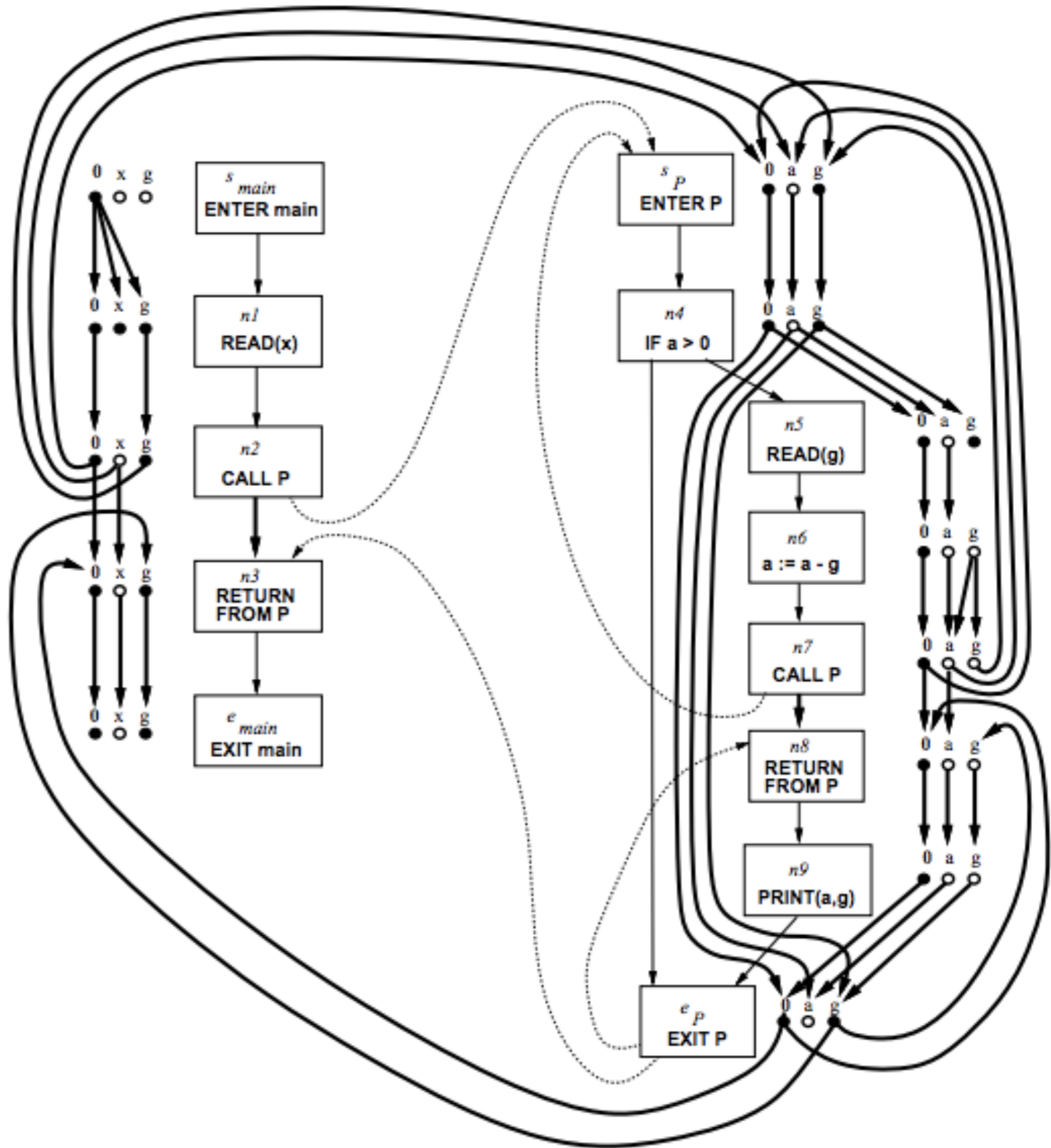
Exploded supergraph

```

int g;

void main(void) {
  int x;
  read(x);
  P(x);
}

void P(int a) {
  if (a > 0) {
    read(g);
    a := a - g;
    P(a);
  }
}
  
```



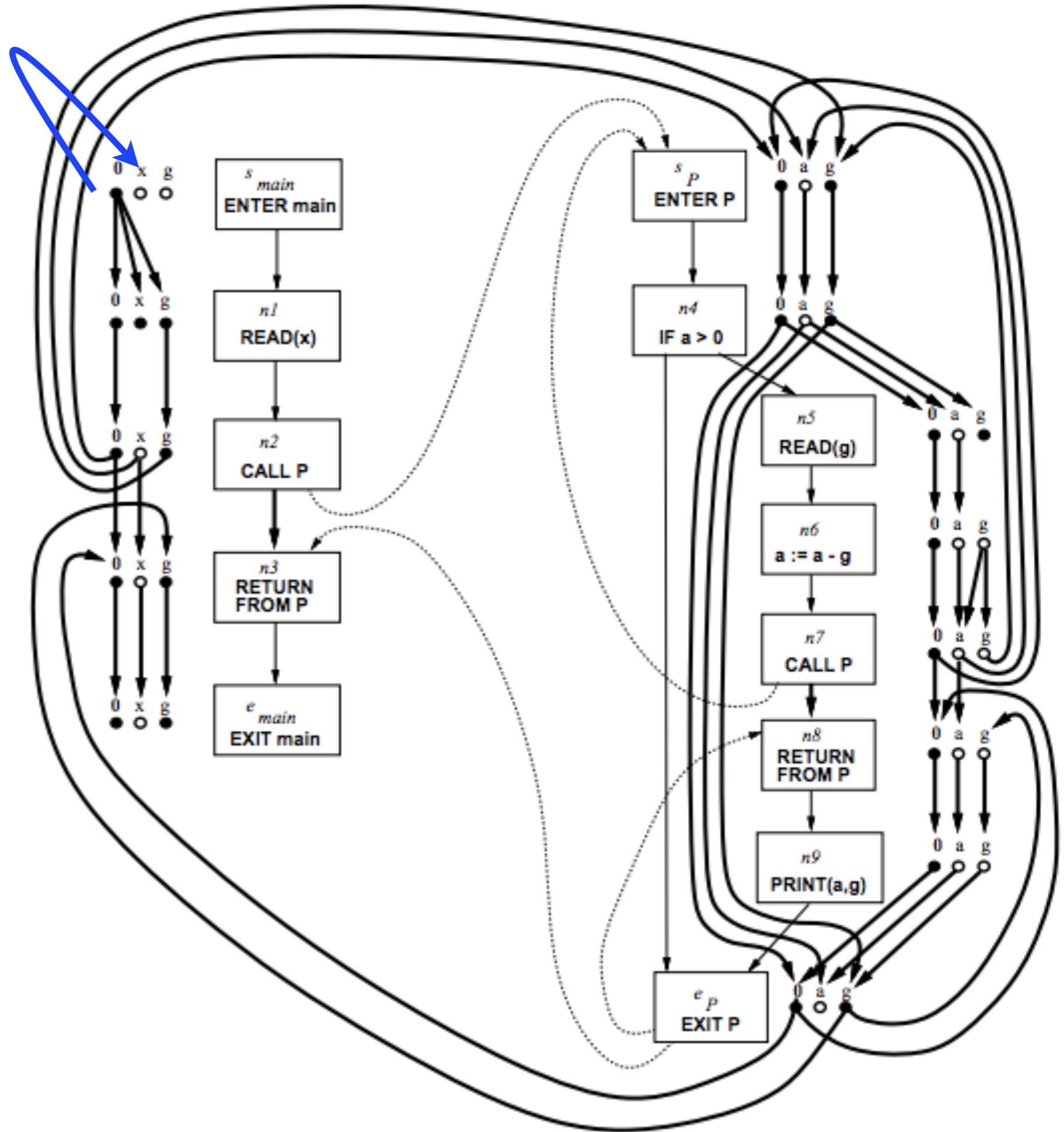
"Tabulation" Algorithm

1. keep a worklist of **Path Edges**
 - (suffixes of valid paths)
2. build set of **Summary Edges**
 - (side effects of a procedure call)
3. result = meet over valid paths

Init
(lines 1-4)

Path Edge

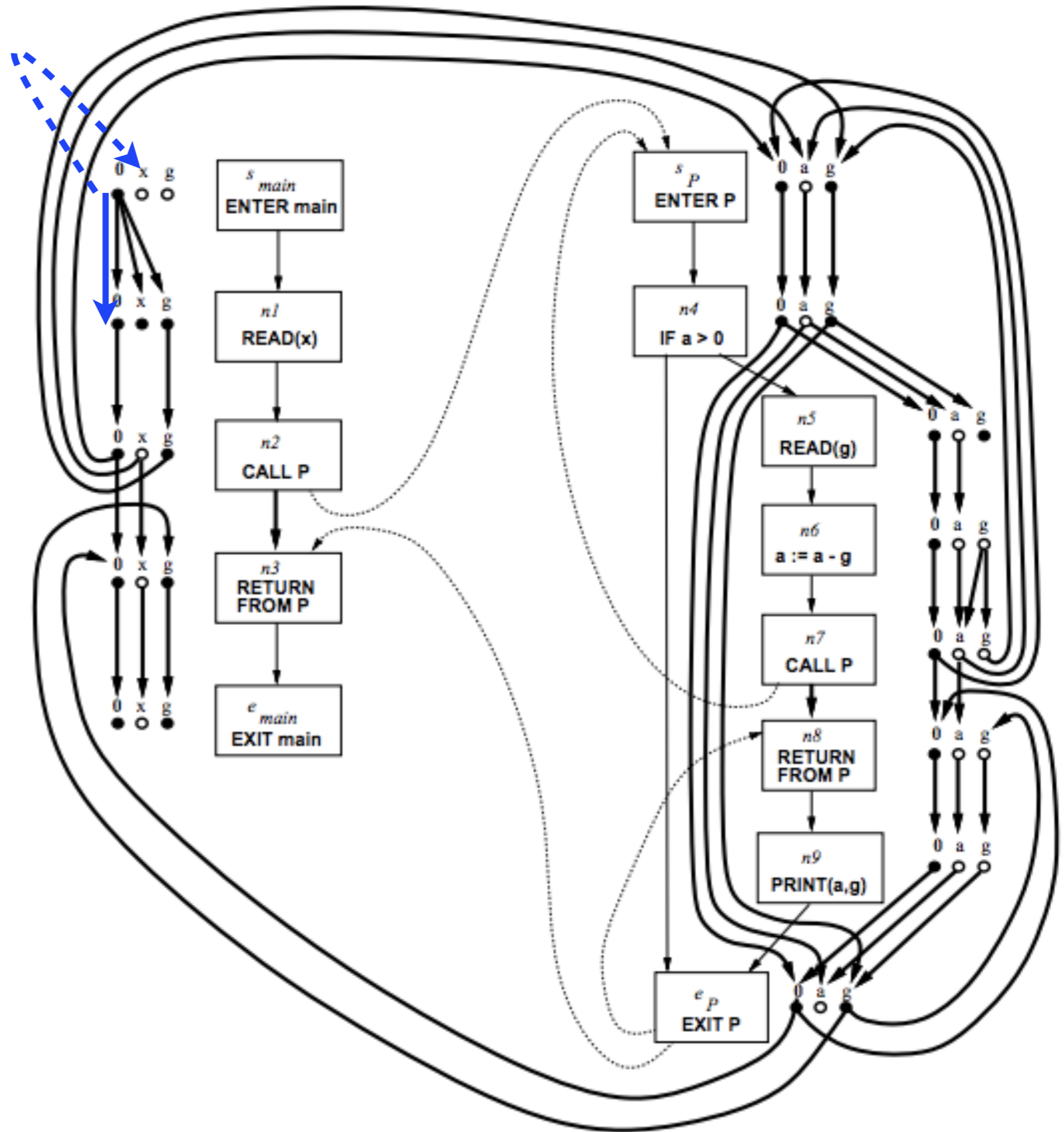
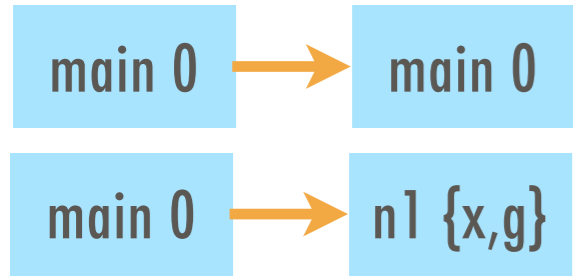
main 0 → main 0



Summary Edge:

Case $n \notin \text{Call}, n \notin \text{Exit}$
 (lines 31-33)

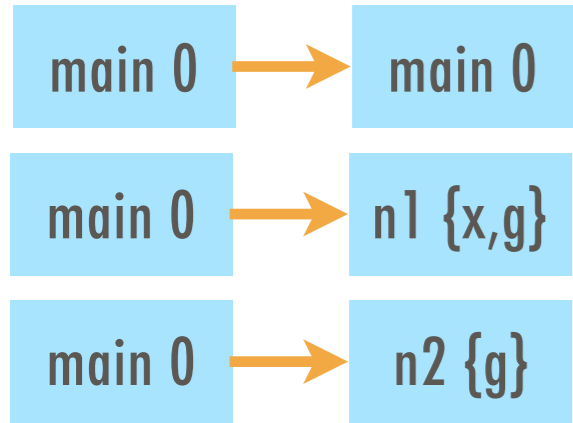
Path Edge



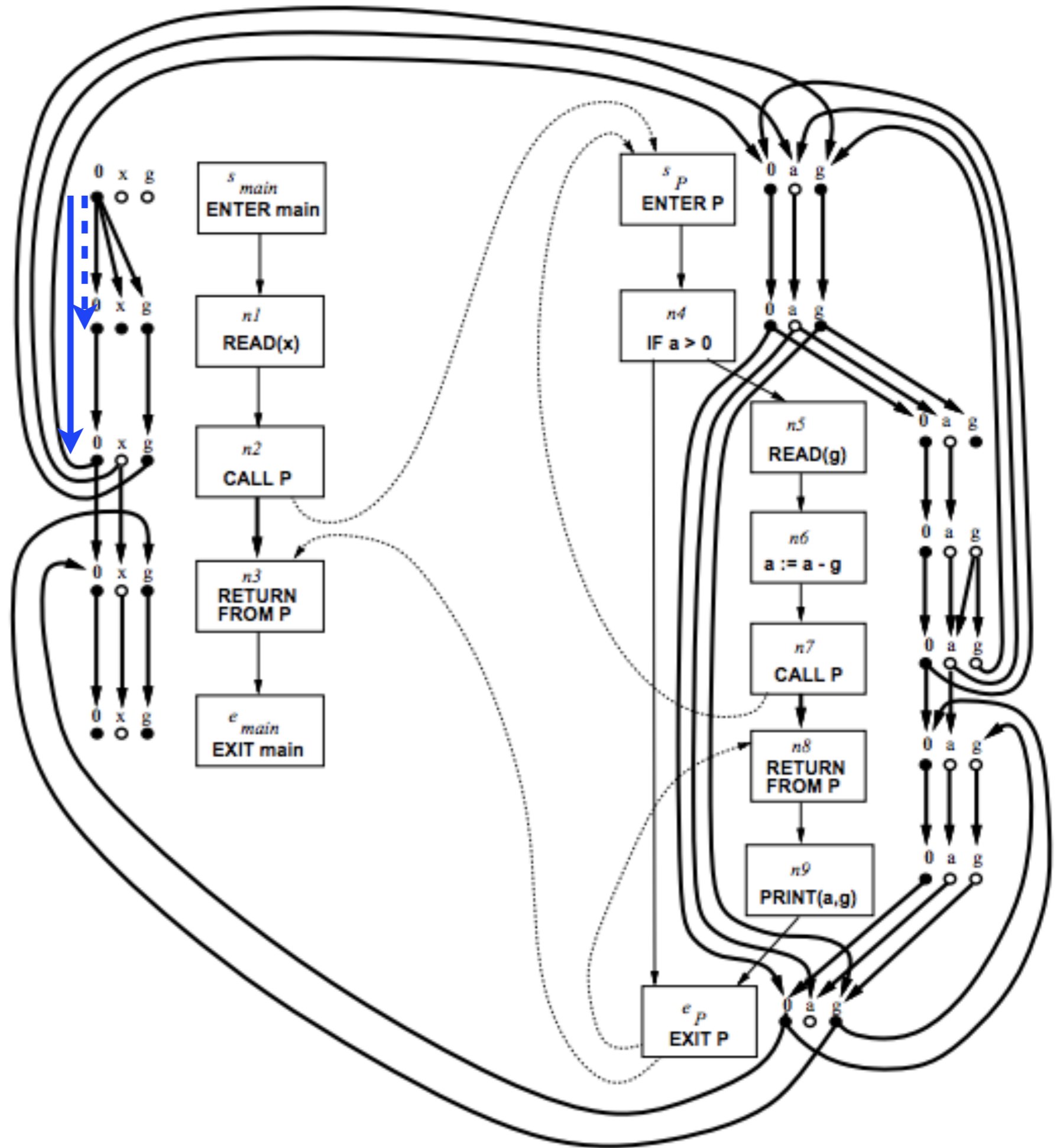
Summary Edge:

Case $n \notin \text{Call}, n \notin \text{Exit}$
 (lines 31-33)

Path Edge

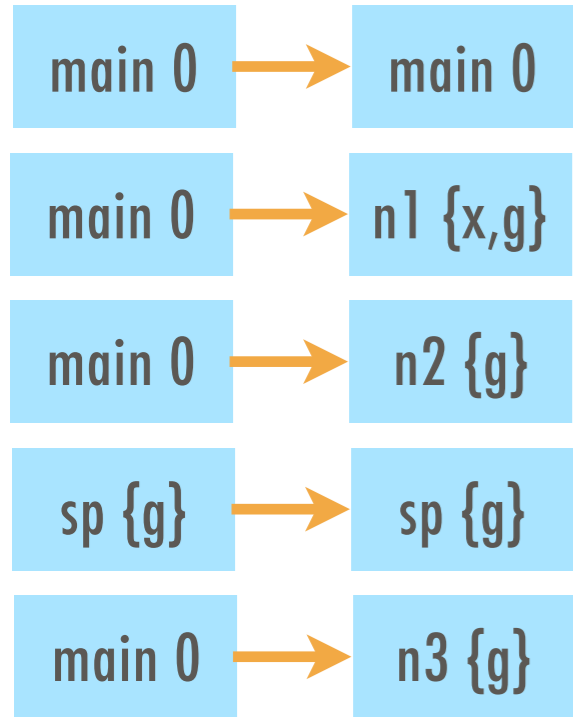


Summary Edge:

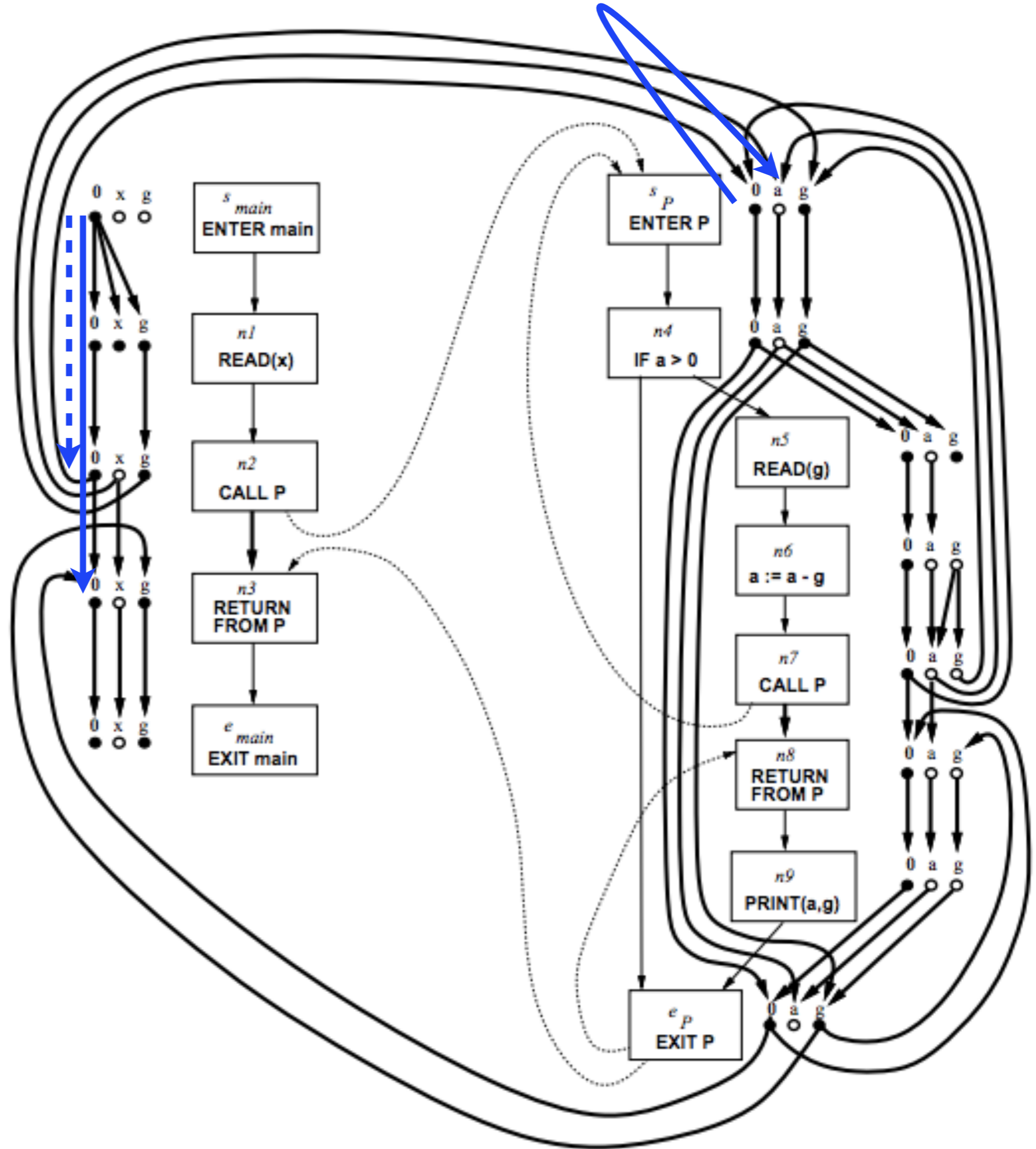


Case $n \in \text{Call}$
 (lines 13-20)

Path Edge

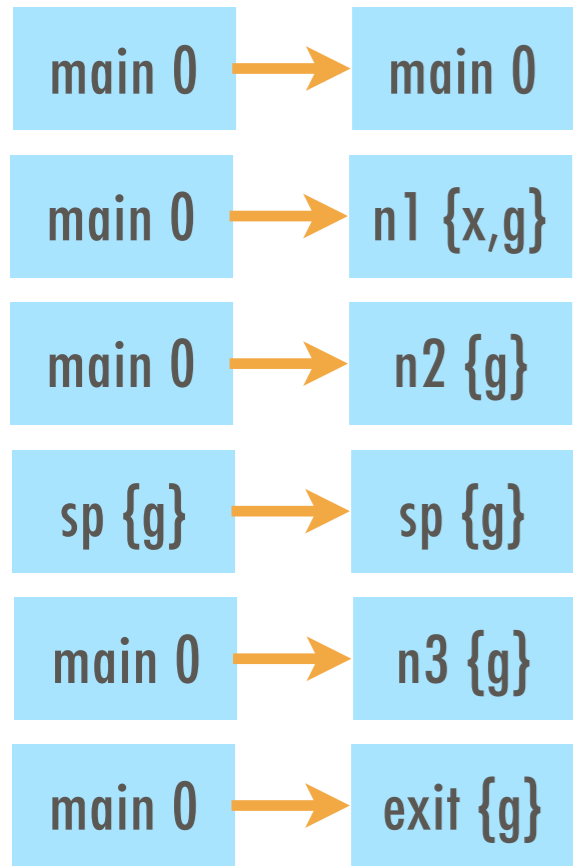


Summary Edge:

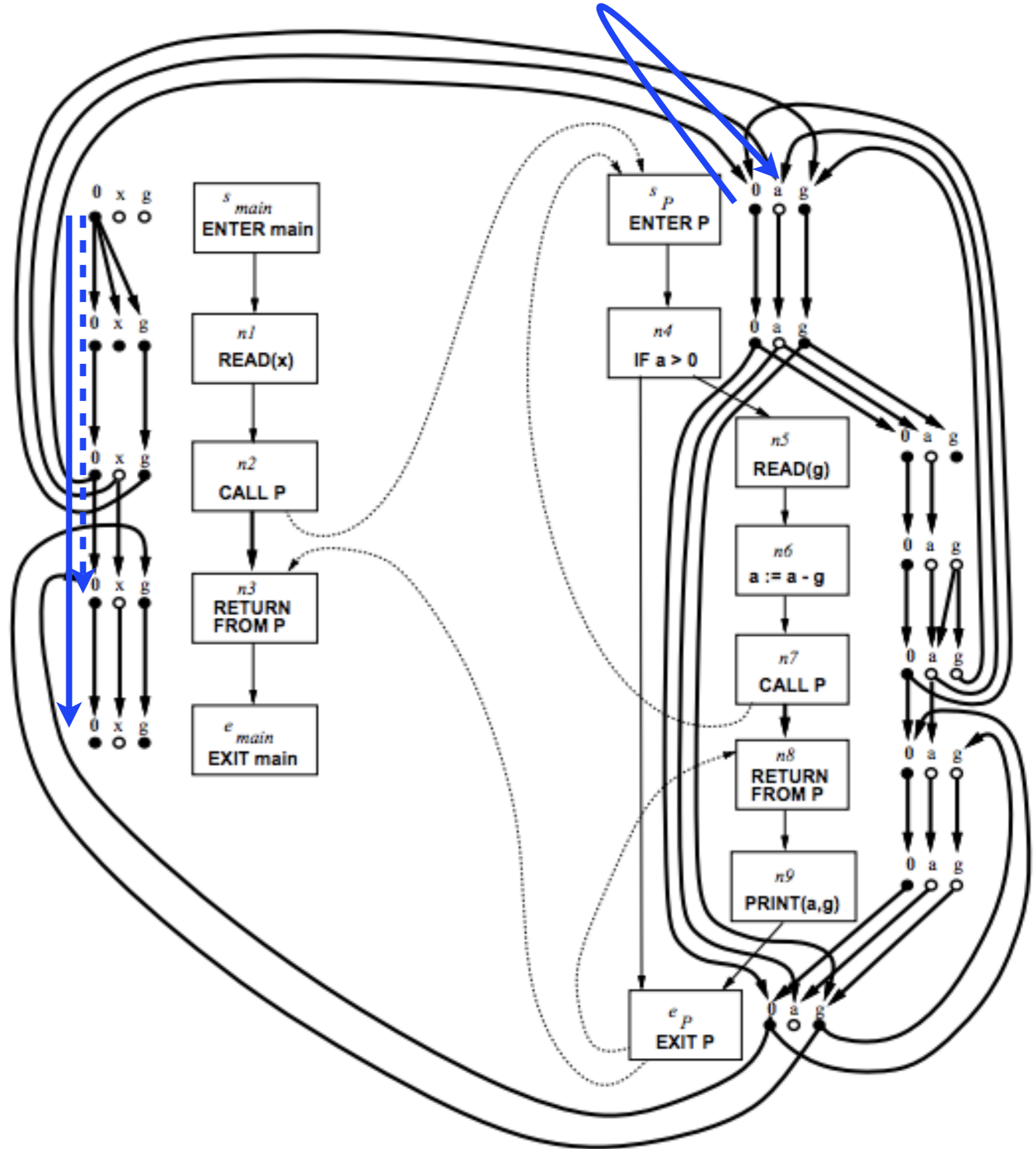


Case $n \notin \text{Call}, n \notin \text{Exit}$
 (lines 31-33)

Path Edge

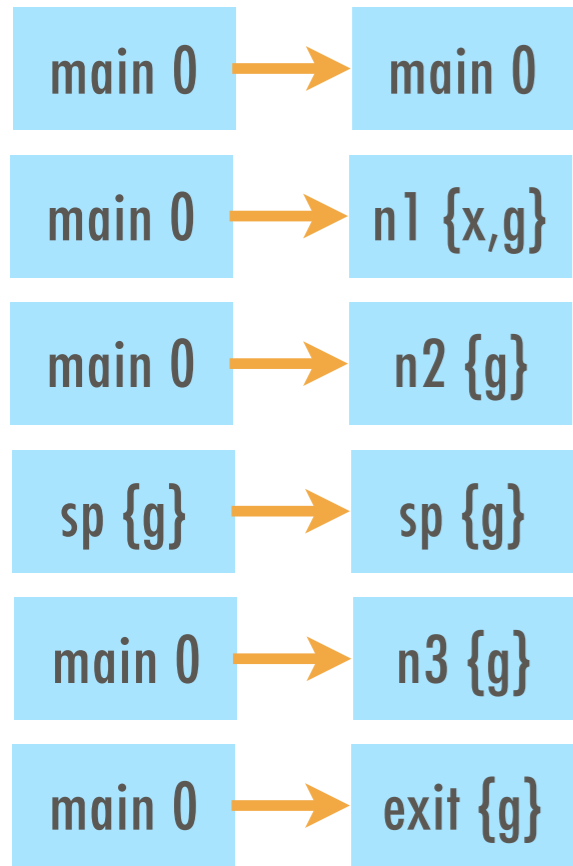


Summary Edge:

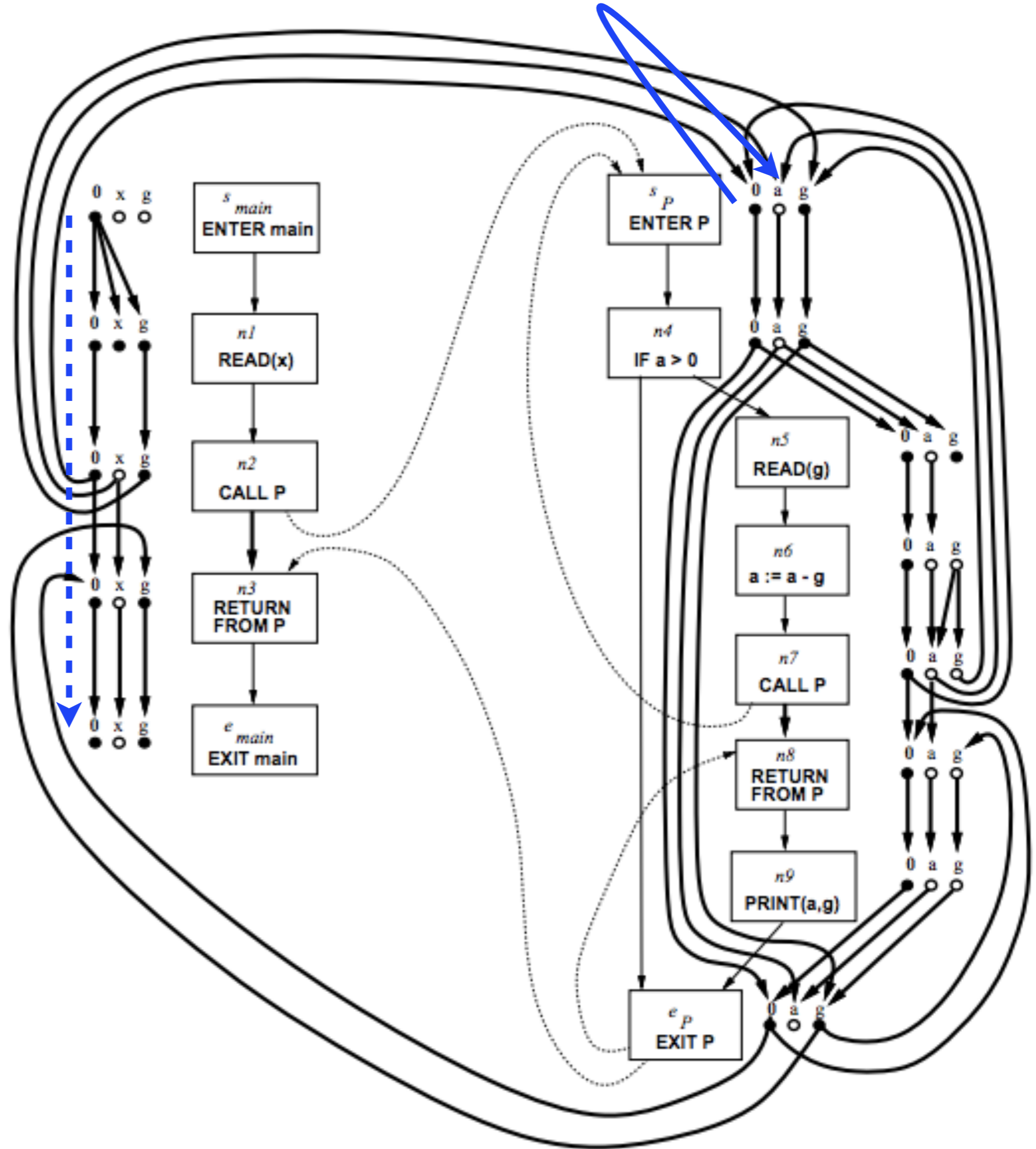


Case $n \notin \text{Call}, n \notin \text{Exit}$
 (lines 31-33)

Path Edge

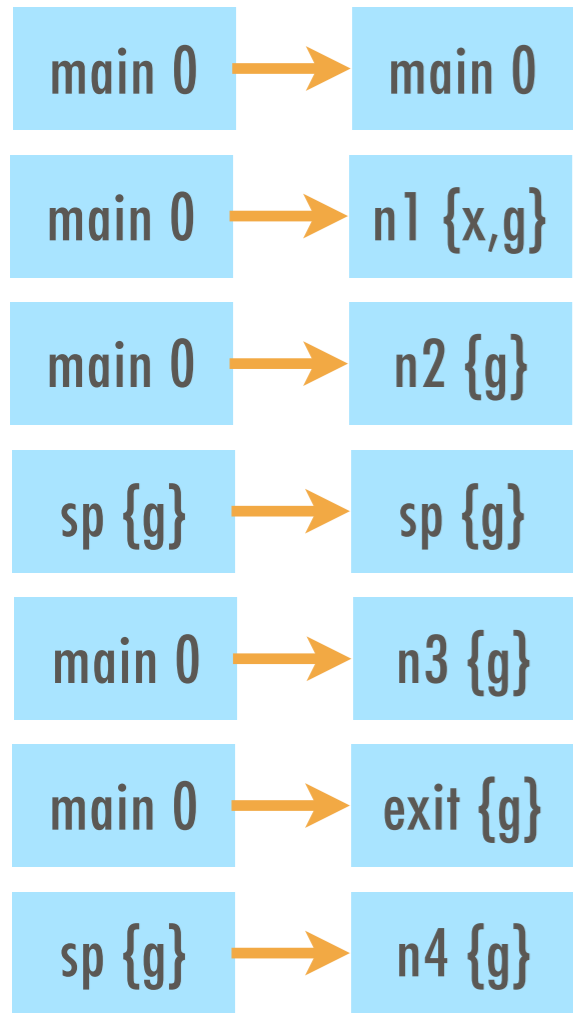


Summary Edge:

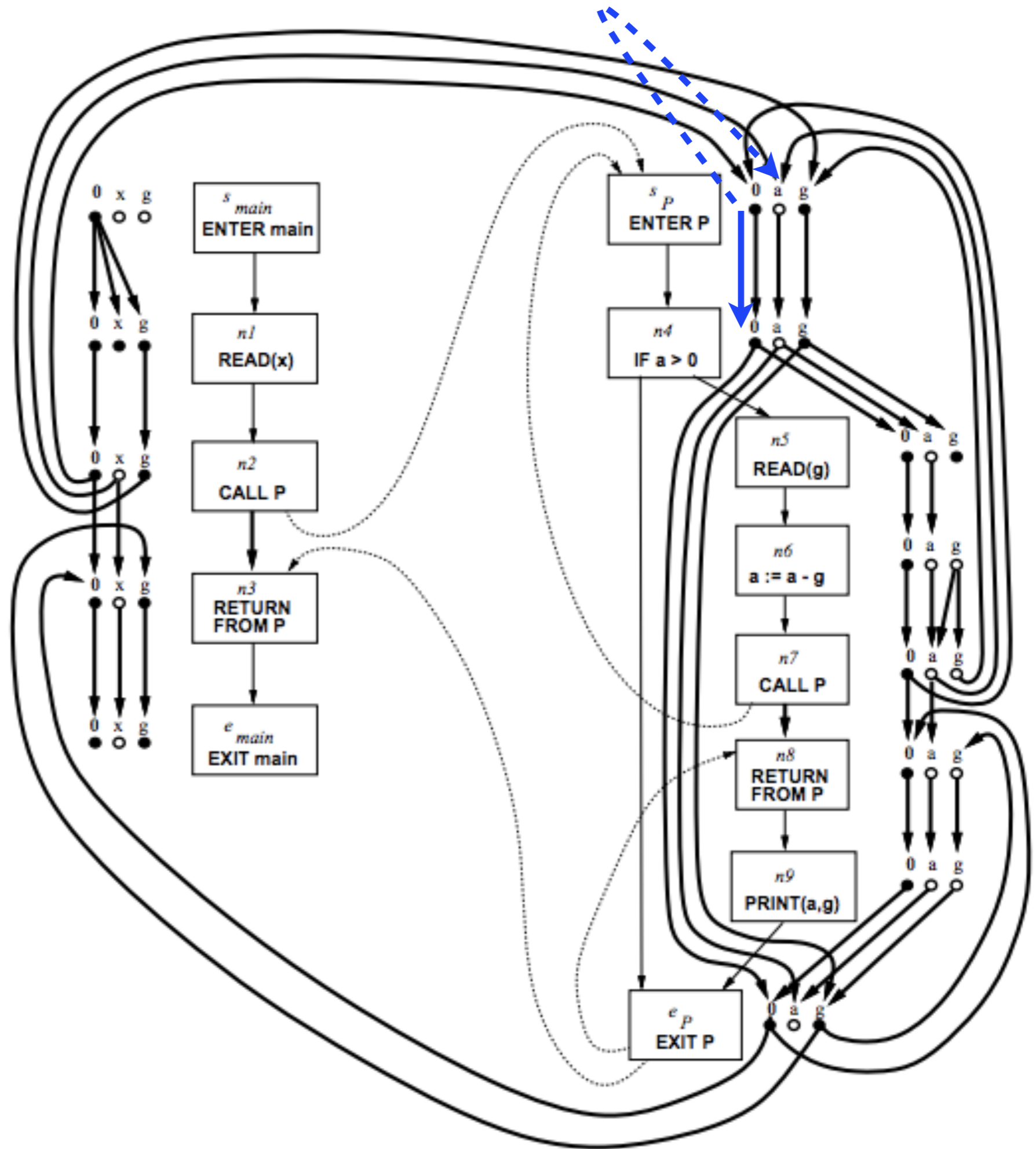


Case $n \notin \text{Call}, n \notin \text{Exit}$
 (lines 31-33)

Path Edge

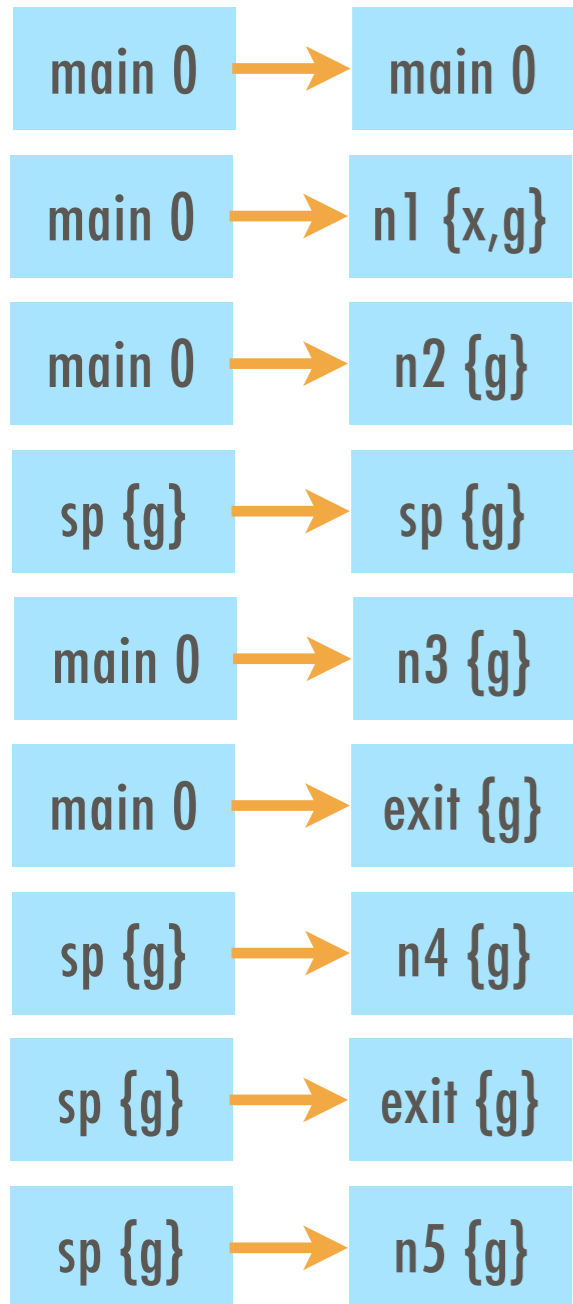


Summary Edge:

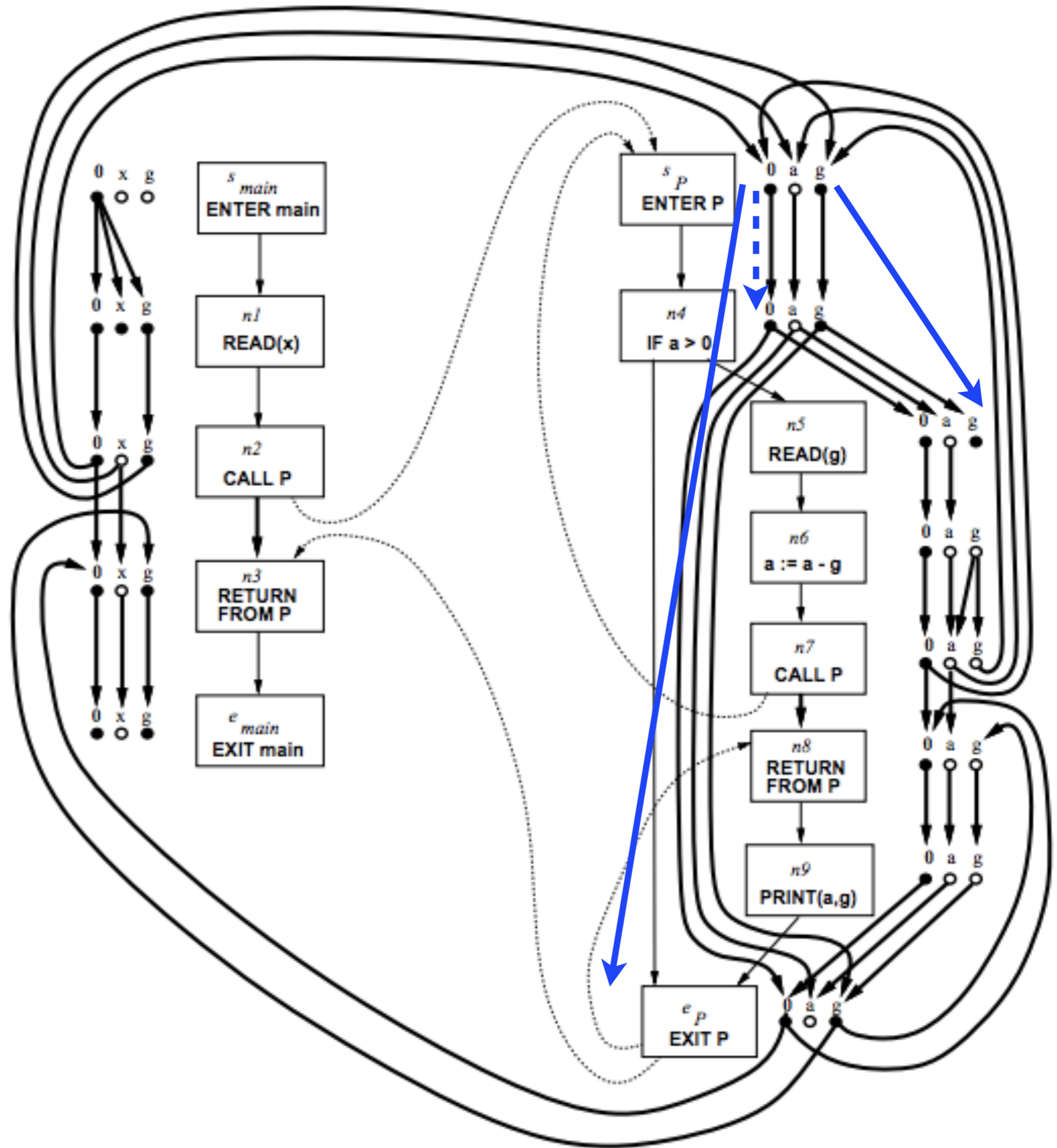


Case $n \notin \text{Call}, n \notin \text{Exit}$
 (lines 31-33)

Path Edge

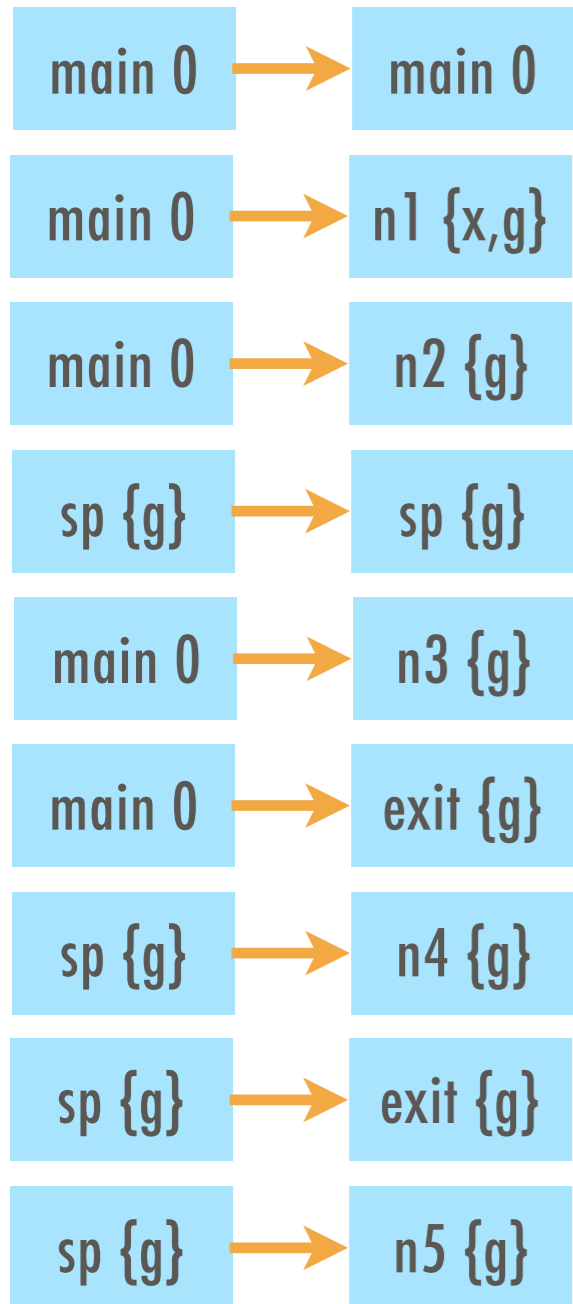


Summary Edge:

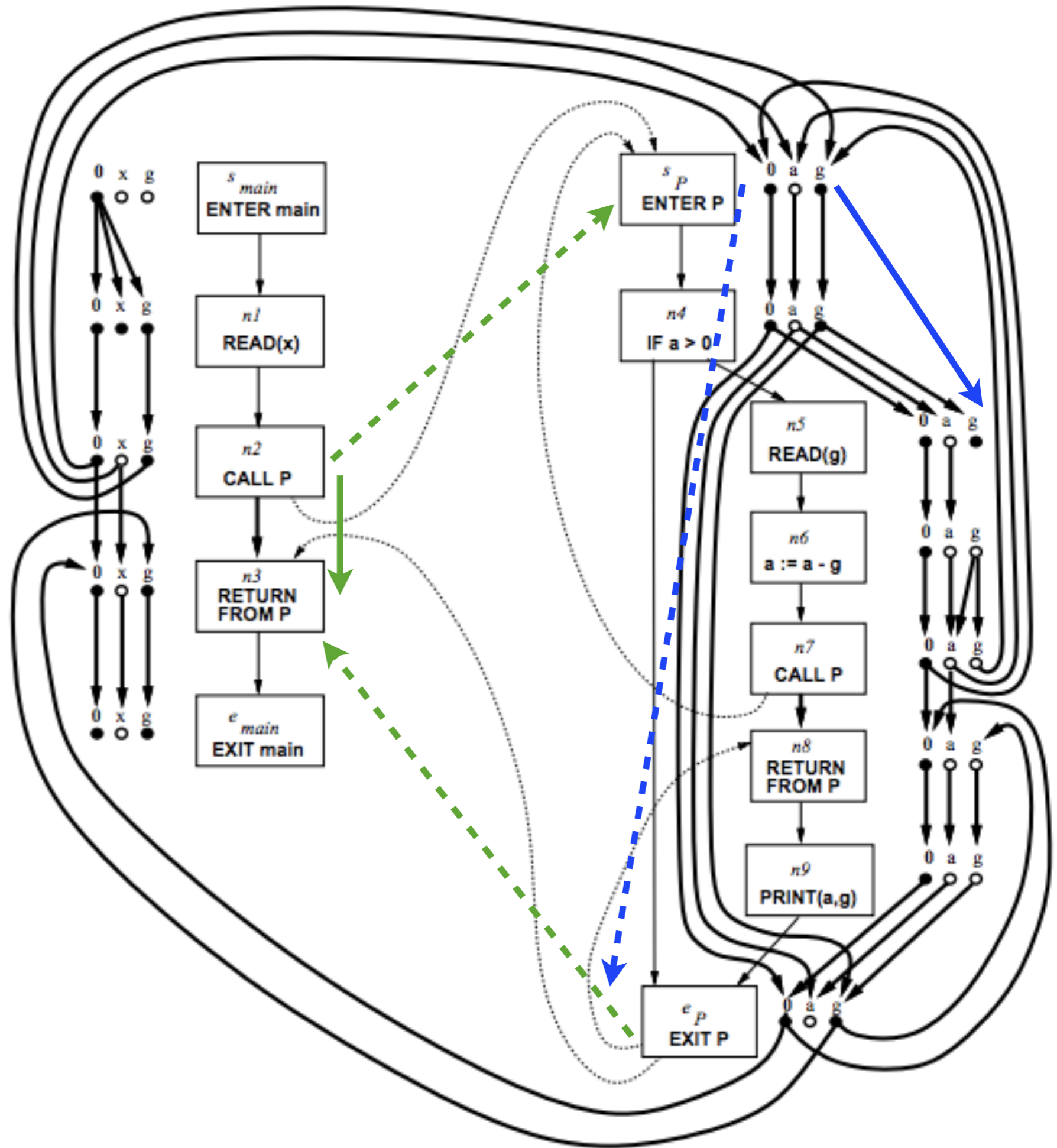
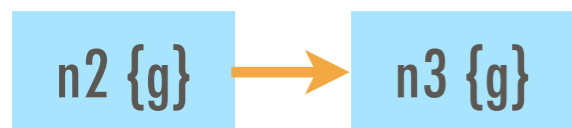


Case $n \in \text{Exit}$
 (lines 21-25)

Path Edge

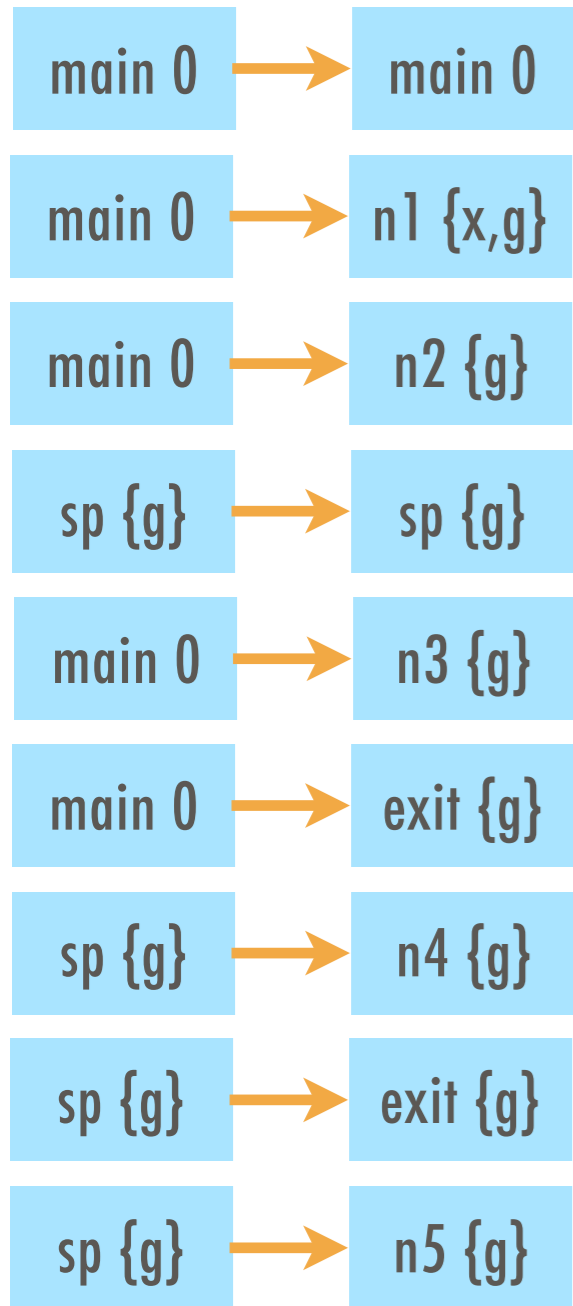


Summary Edge:

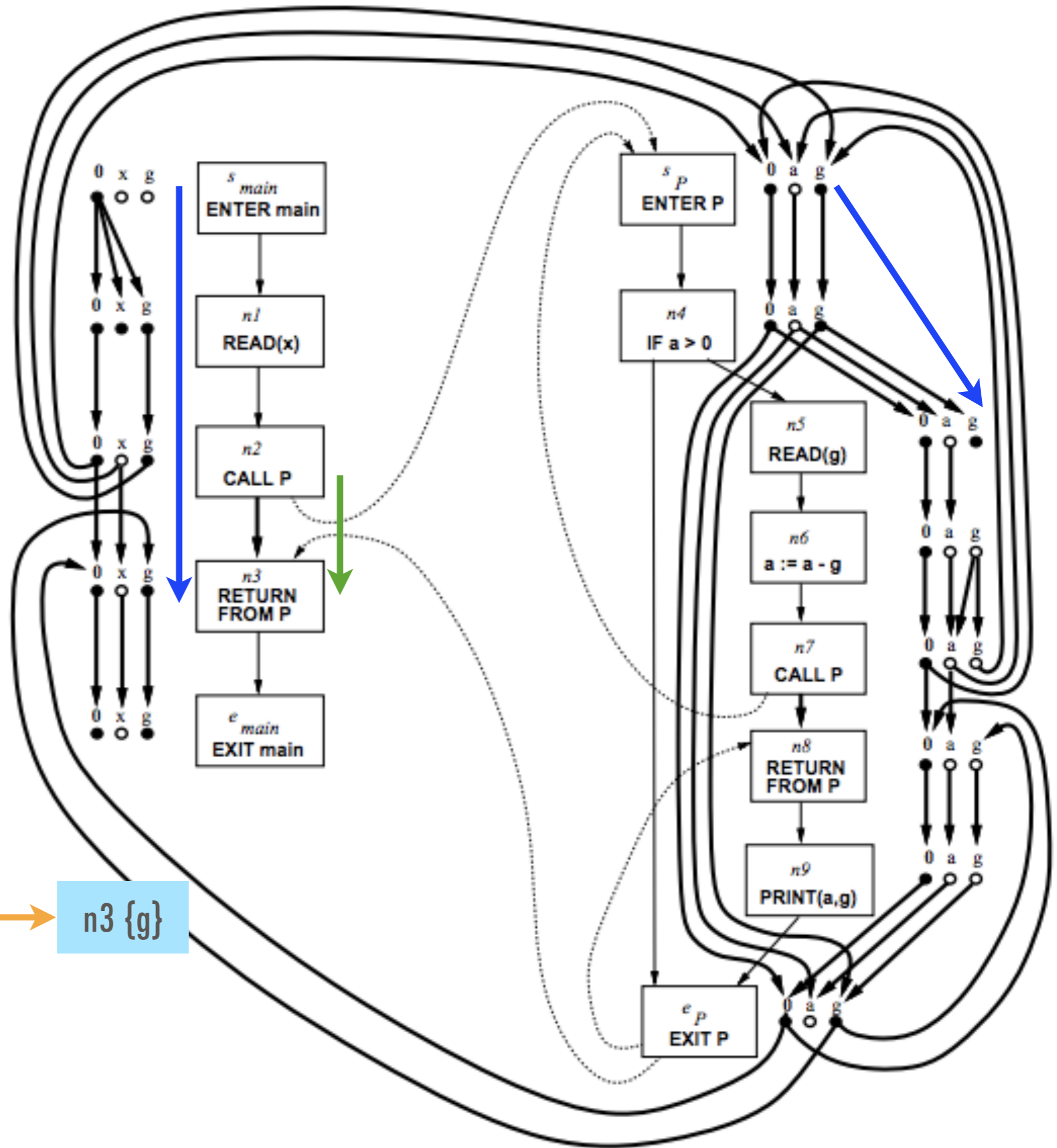
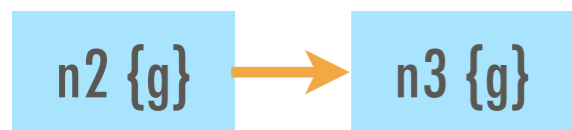


Case $n \in \text{Exit}$
 (lines 25-32)

Path Edge

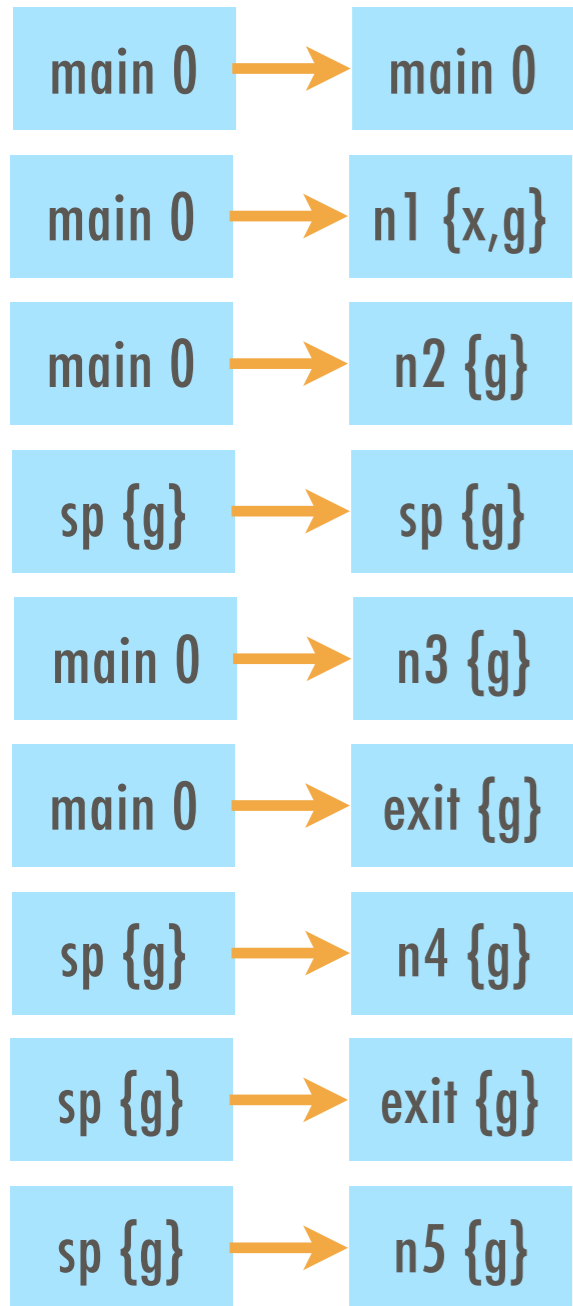


Summary Edge:

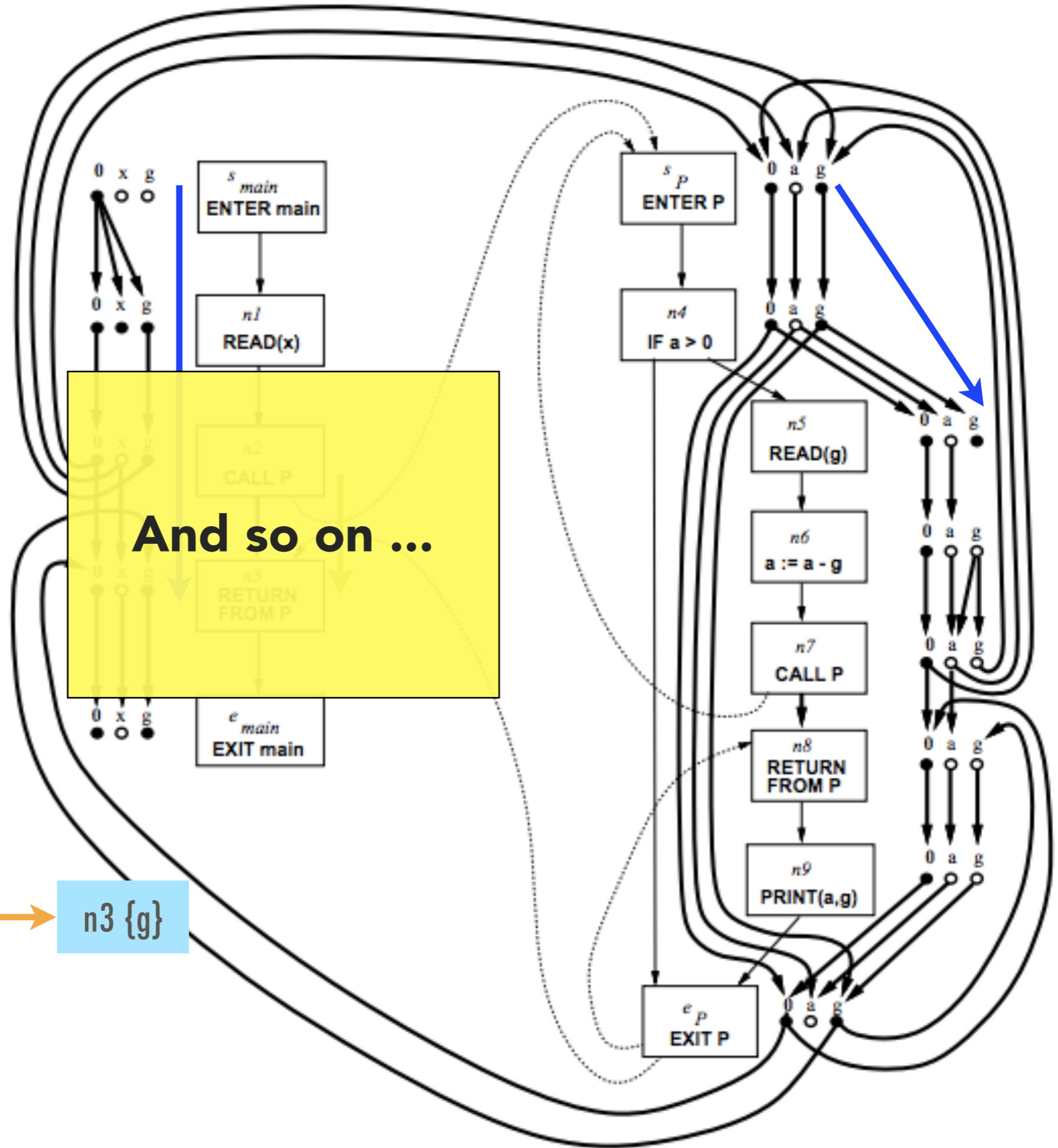
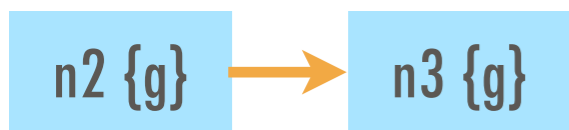


Case $n \in \text{Exit}$
 (lines 25-32)

Path Edge



Summary Edge:



Algorithm II

- "4" ways to find **Path Edges**
 1. **call** edge
 2. **return** edge / **Summary Edge**
 3. normal edge

Running Time

- **E** supergraph edges to explore
- **D** sources to explore from
- **D²** exploded edges for each edge

Class of F	Running Time
Distributive	$O(ED^3)$
<i>h</i> -sparse	$O(\text{Call } D^3 + hED^2)$
Locally Separable	$O(ED)$

Evaluation

exploded supergraph stats

supergraph stats

Program	# lines	# proc.	# calls	# nodes	# edges	D	# n++	# e++
struct-beauty	897	36	214	2188	2860	90	184k	221k
C-parser	1224	48	78	1637	1992	70	104k	112k
ratfor	1345	52	266	2239	2991	87	180k	218k
twig	2388	81	221	3692	4439	142	492k	561k

Program	naive time (s)	naive # null	ifds time (s)	ifds # null
struct-beauty	1.58	583	4.83 (+ 3.25)	543 (- 40)
C-parser	0.54	127	0.7 (+ 0.16)	11 (- 116)
ratfor	1.46	998	3.15 (+ 1.69)	894 (- 104)
twig	5.04	775	5.45 (+ 0.41)	767 (- 8)

Evaluation

Side effects,
live variables,
type analysis, ...

Tabulation
Algorithm

Precise Interprocedural Dataflow Analysis via Graph Reachability

Exploded Supergraph

$F \rightarrow$ bipartite graph

Discussion

- What static analysis problems are / are not IFDS ?
- The uninitialized variables problem is **cubic** in the # global variables, even if these are rarely used.
Can we avoid this overhead?
- Could we allow a (restricted) GOTO?
- Can we add more context-sensitivity?
(Naeem, Lhoták, Rodriguez; CC'10)

Influence

- WALA
- SOOT (Bodden; SOAP '12)
- FLIX (Madsen, Yee, Lhoták; PLDI 2016)
- FlowDroid (Arzt, Rasthofer, Fritz, Bodden, Bartel, Klein, Le Traon, Octeau, McDaniel; PLDI 2014)

END