# Static Program Analysis

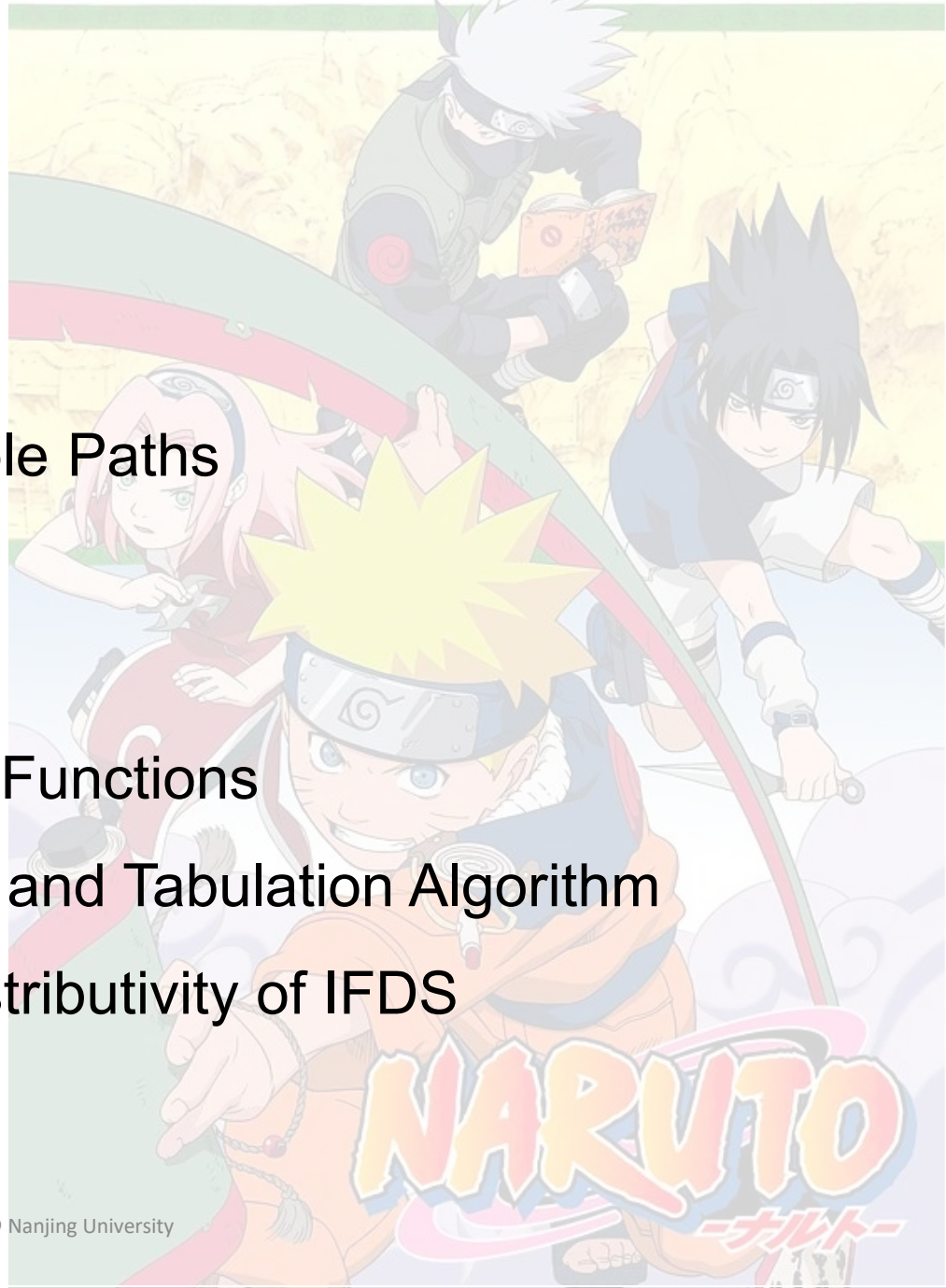## CFL-Reachability and IFDS

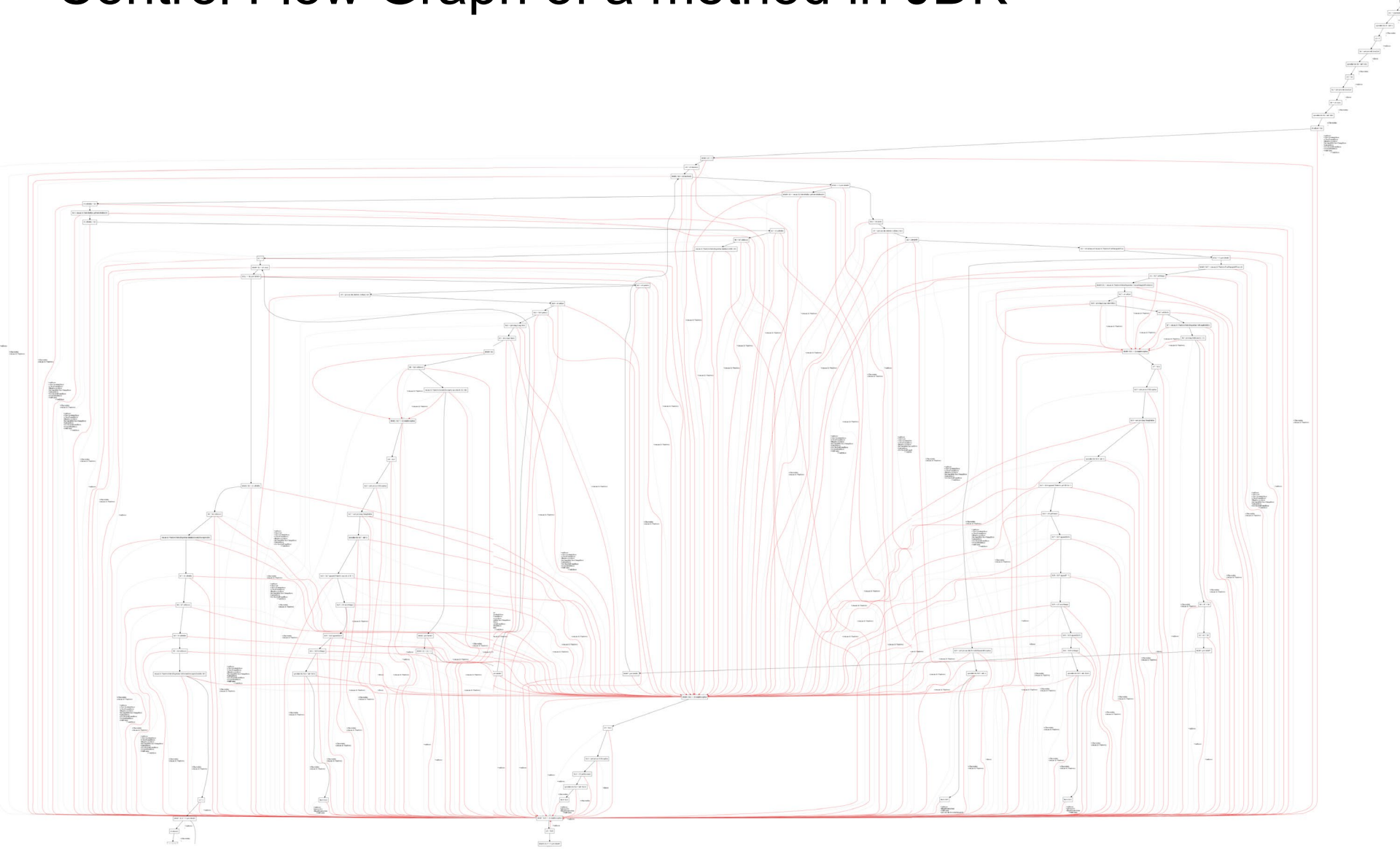Nanjing University

Yue Li

2020

# Contents

1. Feasible and Realizable Paths

2. CFL-Reachability

3. Overview of IFDS

4. Supergraph and Flow Functions

5. Exploded Supergraph and Tabulation Algorithm

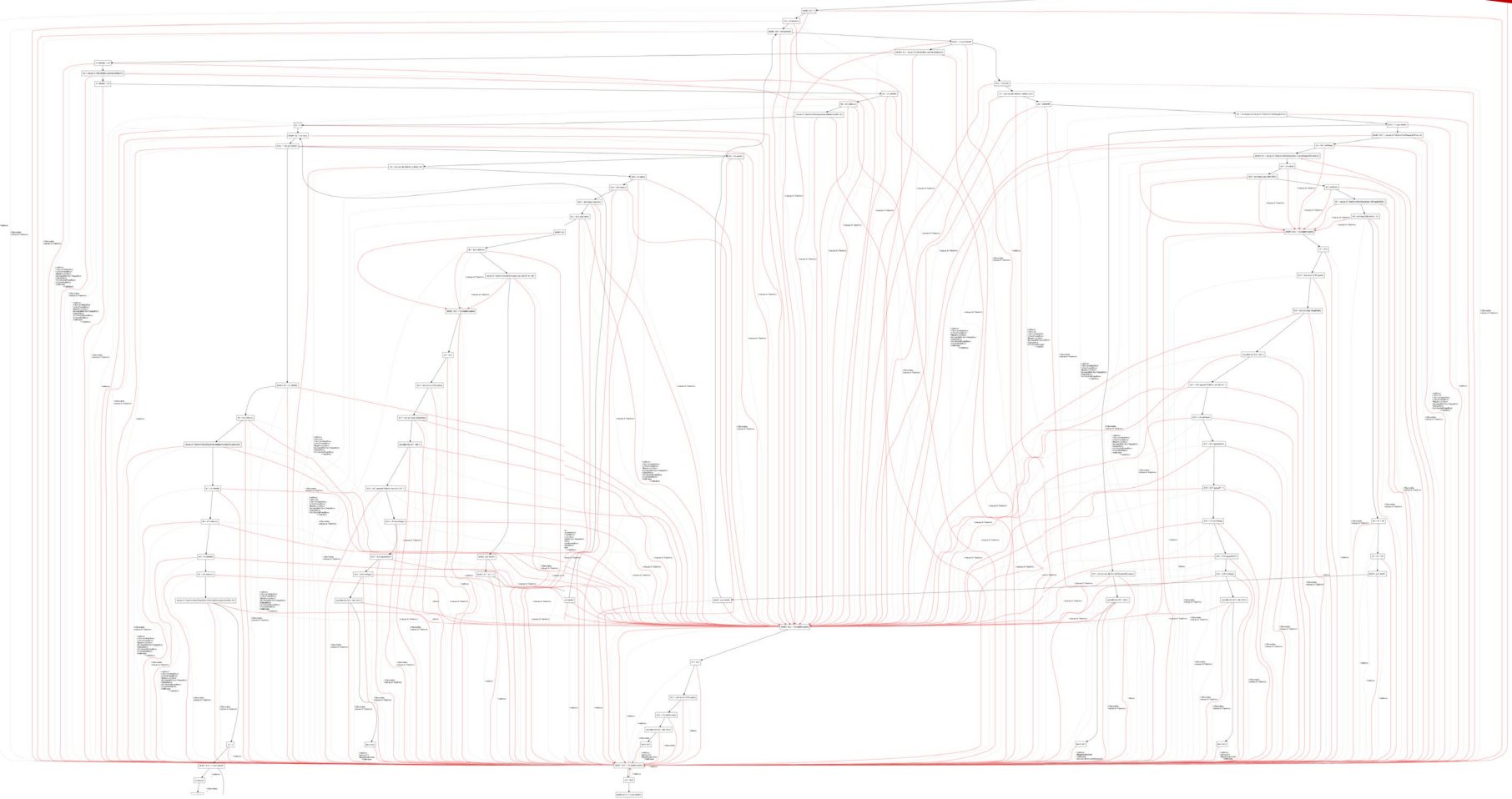6. Understanding the Distributivity of IFDS

# Control Flow Graph of a method in JDK

# Control Flow Graph of a method in JDK
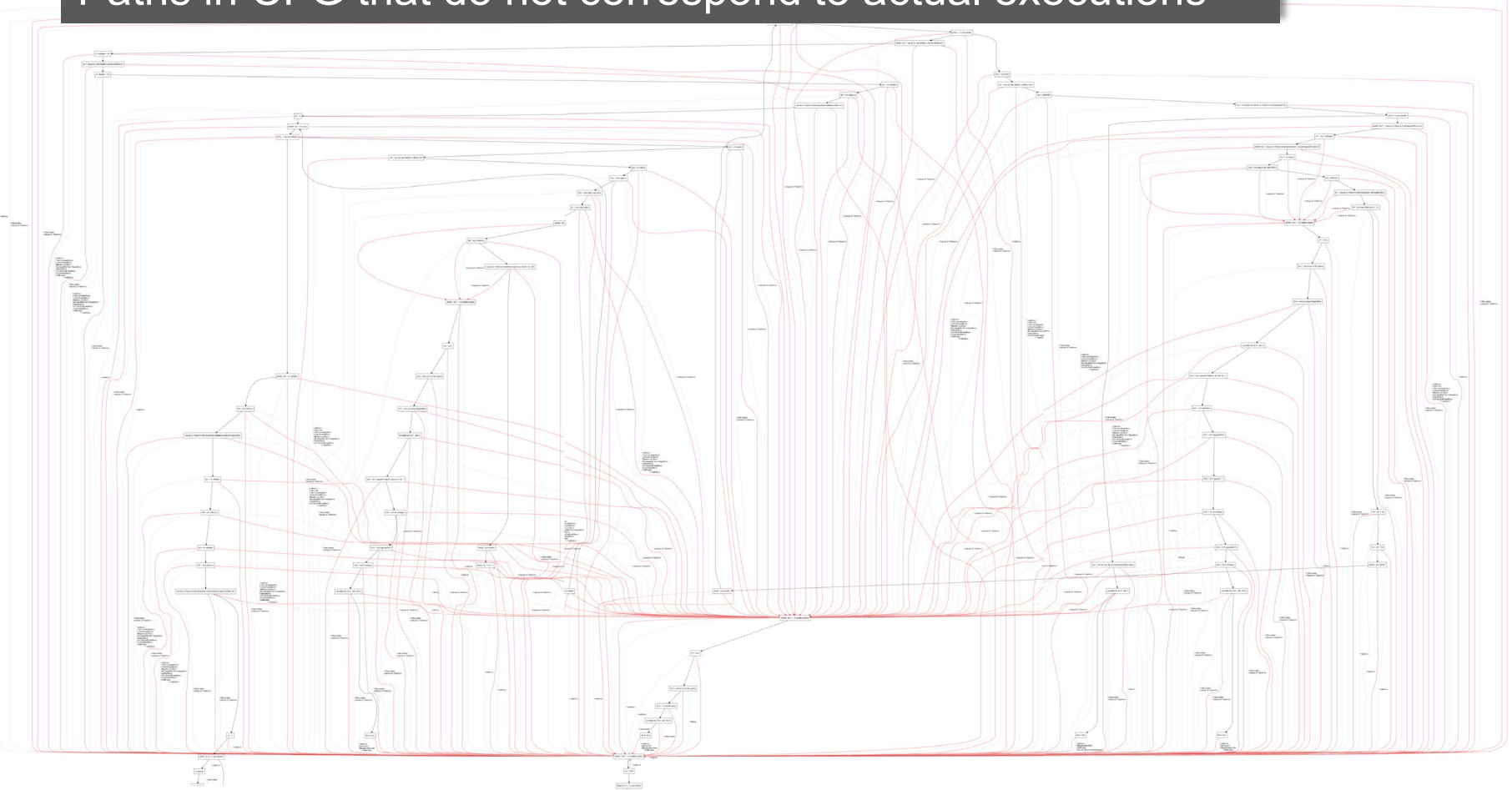


Are all the paths executable?

# Control Flow Graph of a method in JDK

Infeasible Paths:

Paths in CFG that do not correspond to actual executions

Are all the paths executable?

# Control Flow Graph of a method in JDK

**Infeasible Paths:**

Paths in CFG that do not correspond to actual executions
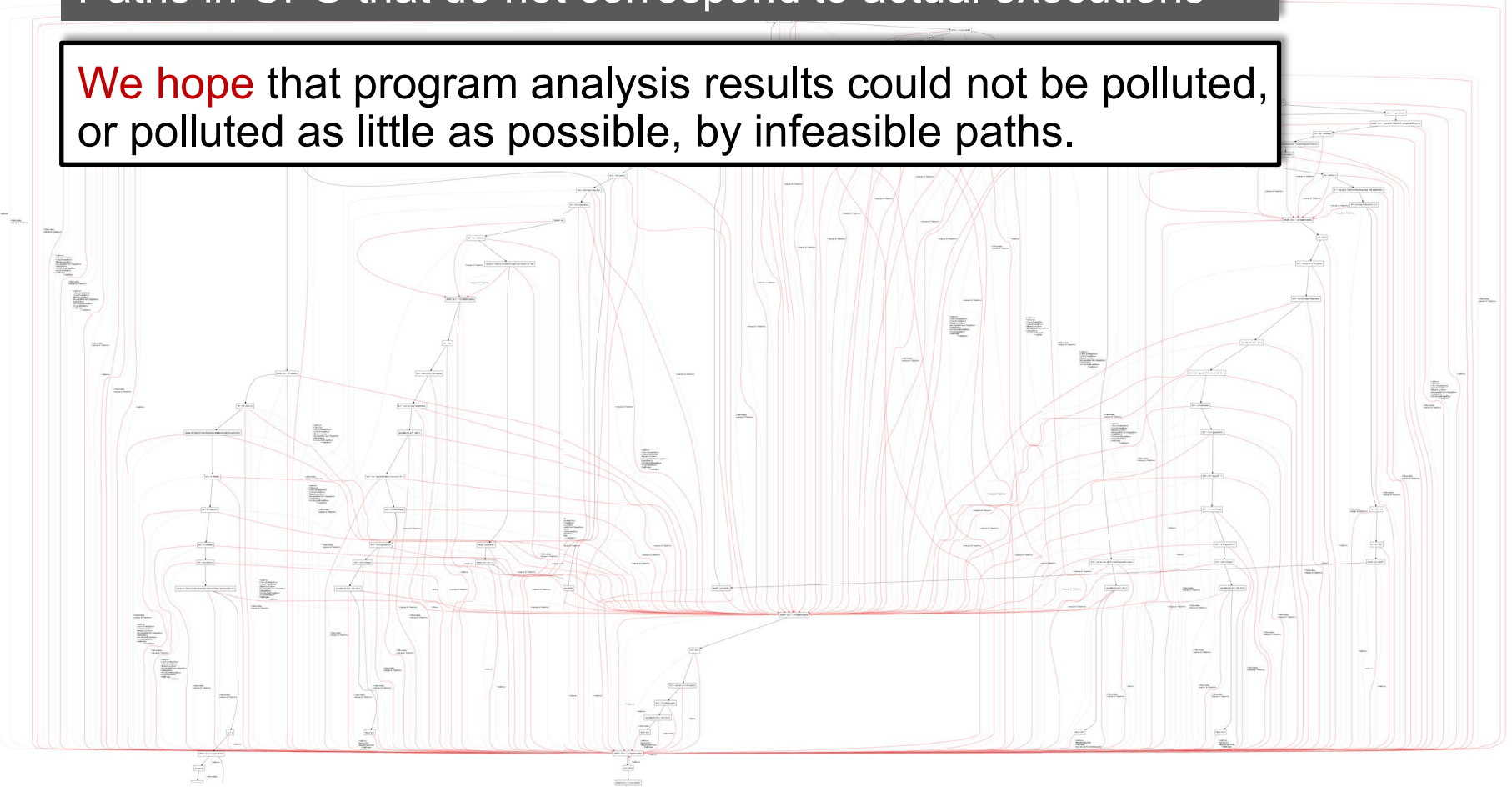
**Are all the paths executable?**

We hope that program analysis results could not be polluted, or polluted as little as possible, by infeasible paths.

# Control Flow Graph of a method in JDK

Are all the paths executable?

**Infeasible Paths:**

Paths in CFG that do not correspond to actual executions

We hope that program analysis results could not be polluted, or polluted as little as possible, by infeasible paths.

But given a path, determine whether it is feasible is, in general, undecidable.

# Control Flow Graph of a method in JDK
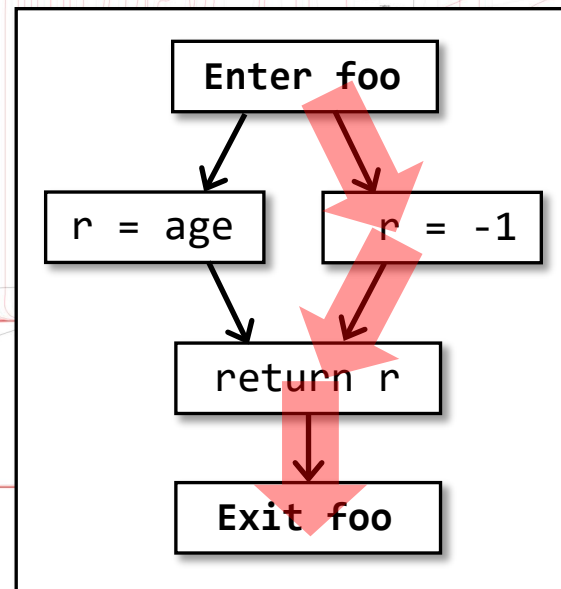
**Infeasible Paths:**
Paths in CFG that do not correspond to actual executions

We hope that program analysis results could not be polluted, or polluted as little as possible, by infeasible paths.

But given a path, determine whether it is feasible is, in general, undecidable.

```
foo(int age) {

    if(age >= 0)

        r = age;
    else
        r = -1;
    return r;
}
```

# Control Flow Graph of a method in JDK

Are all the paths executable?

**Infeasible Paths:**

Paths in CFG that do not correspond to actual executions
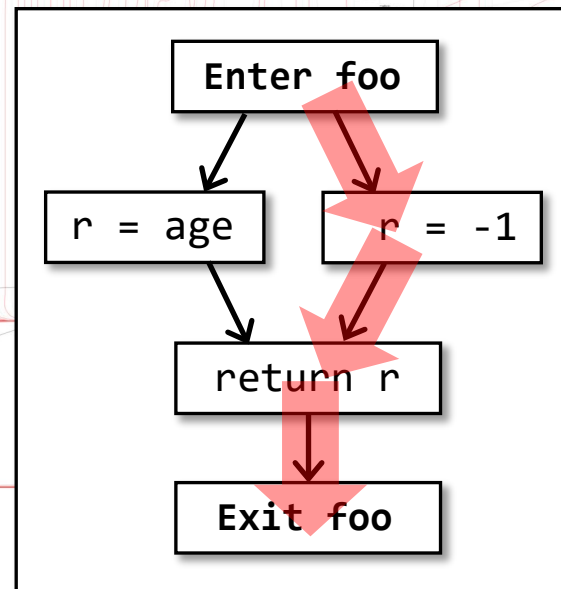
We hope that program analysis results could not be polluted, or polluted as little as possible, by infeasible paths.

No Hope?

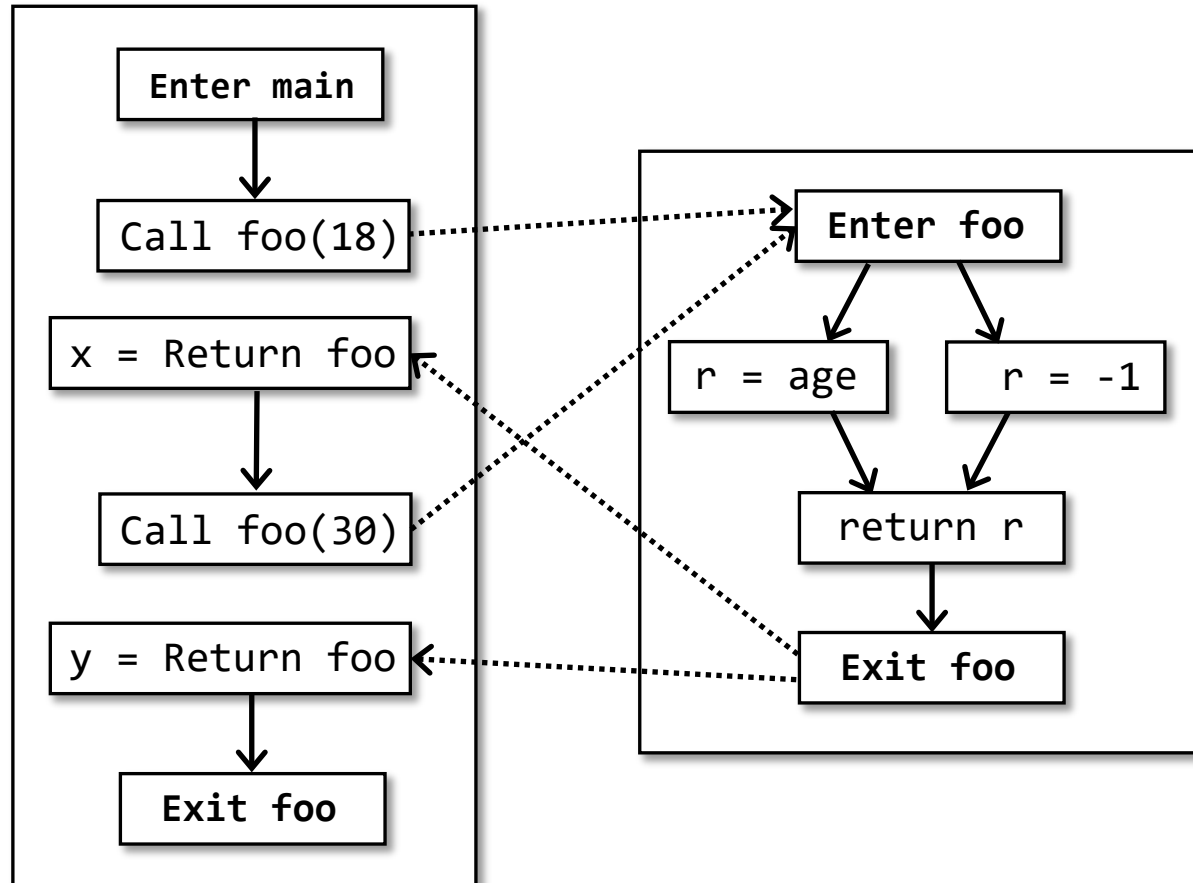But given a path, determine whether it is feasible is, in general, undecidable.

```
foo(int age) {

    if(age >= 0)

        r = age;
    else

        r = -1;

    return r;
}
```

```
Enter foo
r = age        r = -1
return r
Exit foo
```

```
main() {
  x = foo(18);
    ⋮
  y = foo(30);
}
foo(int age) {
  if(age >= 0)
    r = age;
  else
    r = -1;
  return r;
}
```

Enter main → Call foo(18) → x = Return foo → Call foo(30) → y = Return foo → Exit foo

Enter foo → r = age / r = -1 → return r → Exit foo

```
main() {
    x = foo(18);
       ⋮
    y = foo(30);
}
foo(int age) {
    if(age >= 0)
        r = age;
    else
        r = -1;
    return r;
}
```

x=18,30,-1

y=18,30,-1

Enter main

Call foo(18)

x = Return foo

Call foo(30)

y = Return foo

Exit foo

Enter foo

r = age

r = -1

return r

Exit foo

```
main() {
  x = foo(18);
     ⋮
  y = foo(30);
}
foo(int age) {
  if(age >= 0)
    r = age;
  else
    r = -1;
  return r;
}
```

Enter main

Call foo(18)

x = Return foo

Call foo(30)

y = Return foo

Exit foo

x=18,30,-1

y=18,30,-1

Inevitable

Enter foo

r = age        r = -1

return r

Exit foo

```
main() {
  x = foo(18);
    ⋮
  y = foo(30);
}
foo(int age) {
  if(age >= 0)
    r = age;
  else
    r = -1;
  return r;
}
```
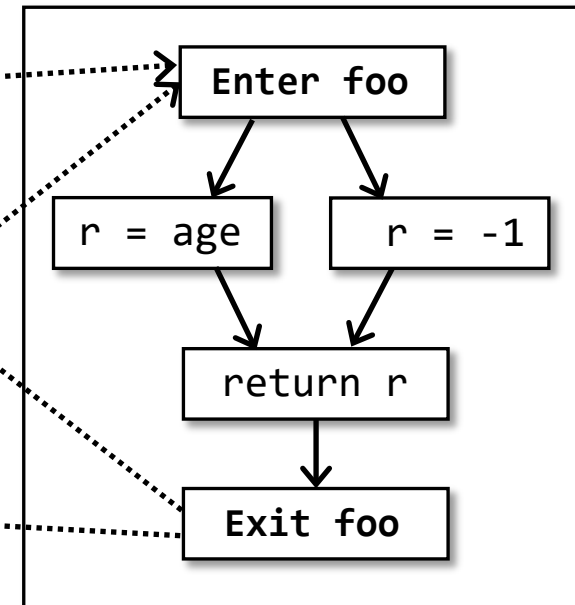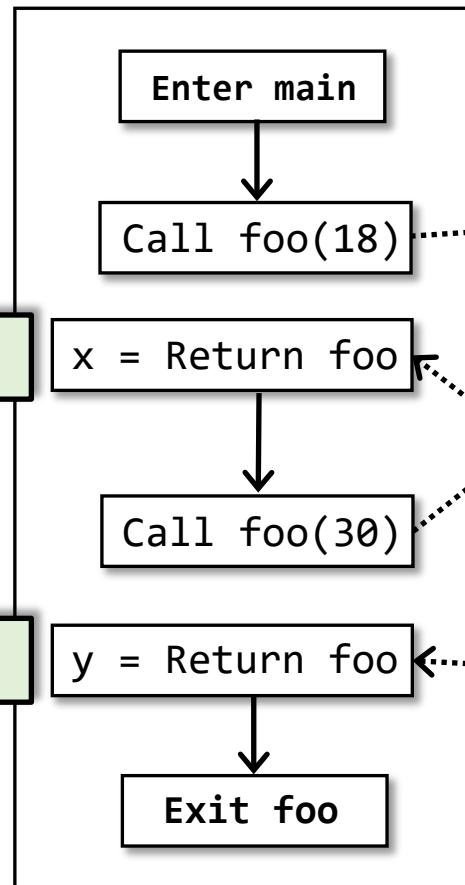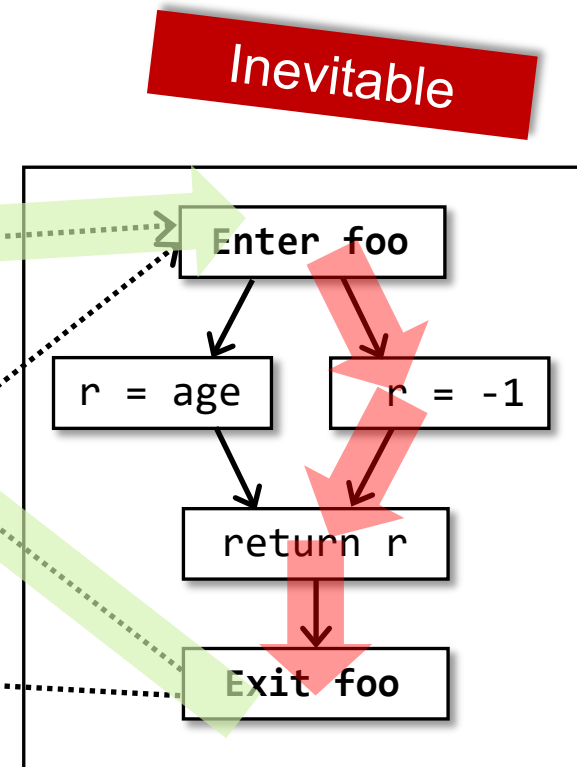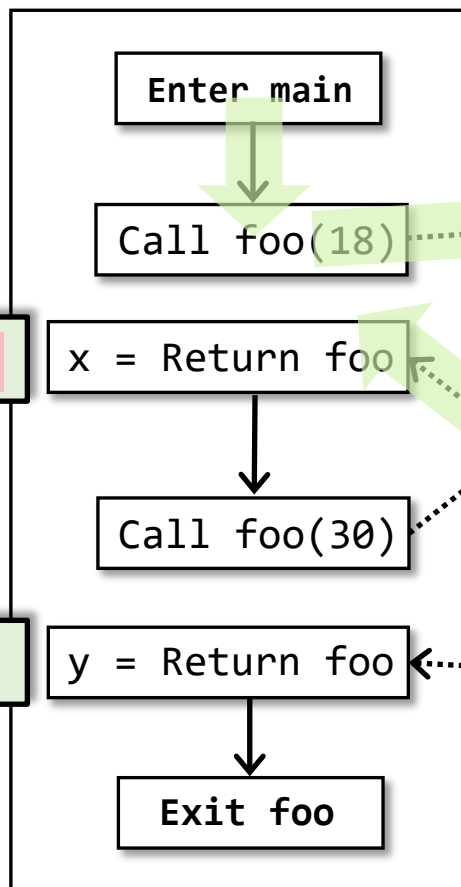
x=18,30,-1

y=18,30,-1

Enter main

Call foo(18)

x = Return foo

Call foo(30)

y = Return foo

Exit foo

Enter foo

r = age

r = -1

return r

Exit foo

```
main() {
  x = foo(18);
    ⋮
  y = foo(30);
}
foo(int age) {
  if(age >= 0)
    r = age;
  else
    r = -1;
  return r;
}
```

Enter main

Call foo(18)

x=18,30,-1

x = Return foo

Call foo(30)

y=18,30,-1

y = Return foo

Exit foo

Enter foo

r = age     r = -1

return r

Exit foo

```
main() {
  x = foo(18);
    ⋮
  y = foo(30);
}
foo(int age) {
  if(age >= 0)
    r = age;
  else
    r = -1;
  return r;
}
```

Enter main

Call foo(18)

x = Return foo

Call foo(30)

y = Return foo

Exit foo

Enter foo

r = age     r = -1

return r

Exit foo

x=18,30,-1

y=18,30,-1

```
main() {
  x = foo(18);
    ⋮
  y = foo(30);
}
foo(int age) {
  if(age >= 0)
    r = age;
  else
    r = -1;
  return r;
}
```
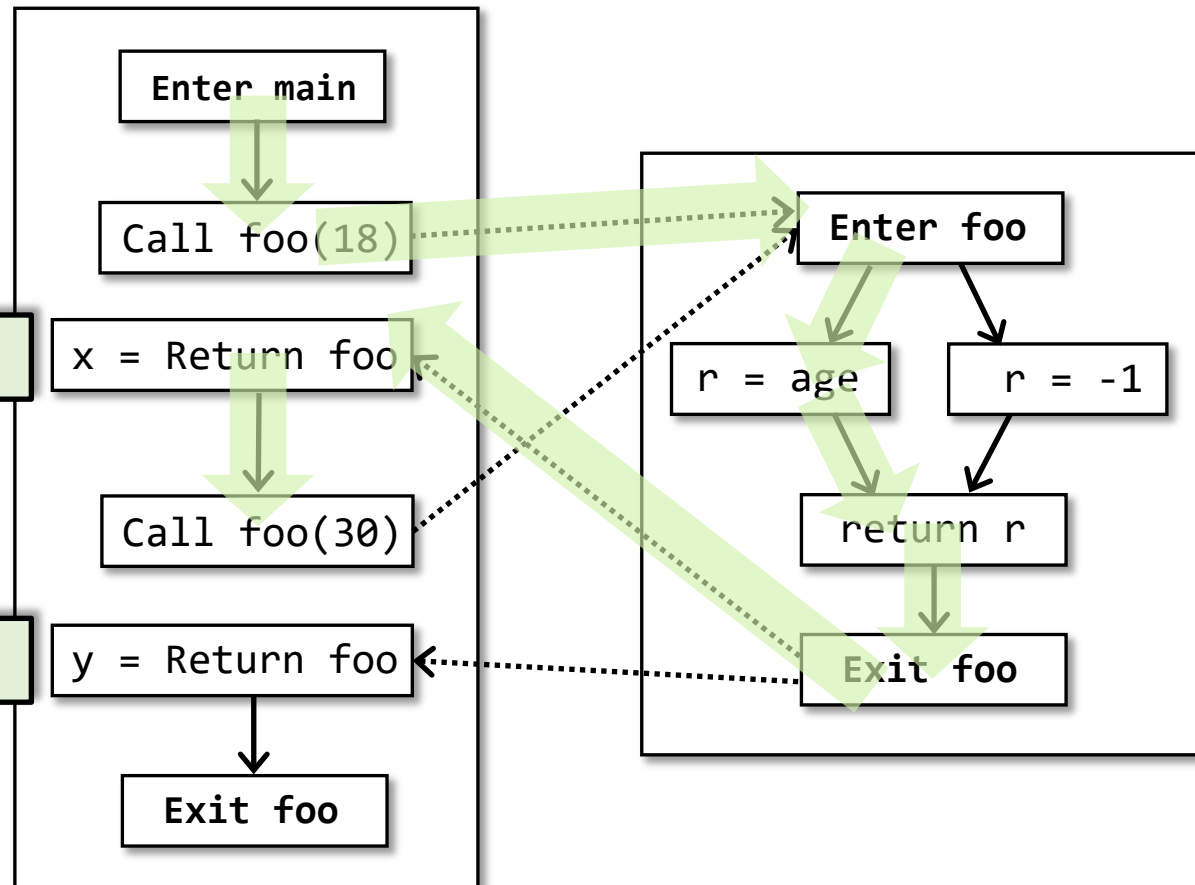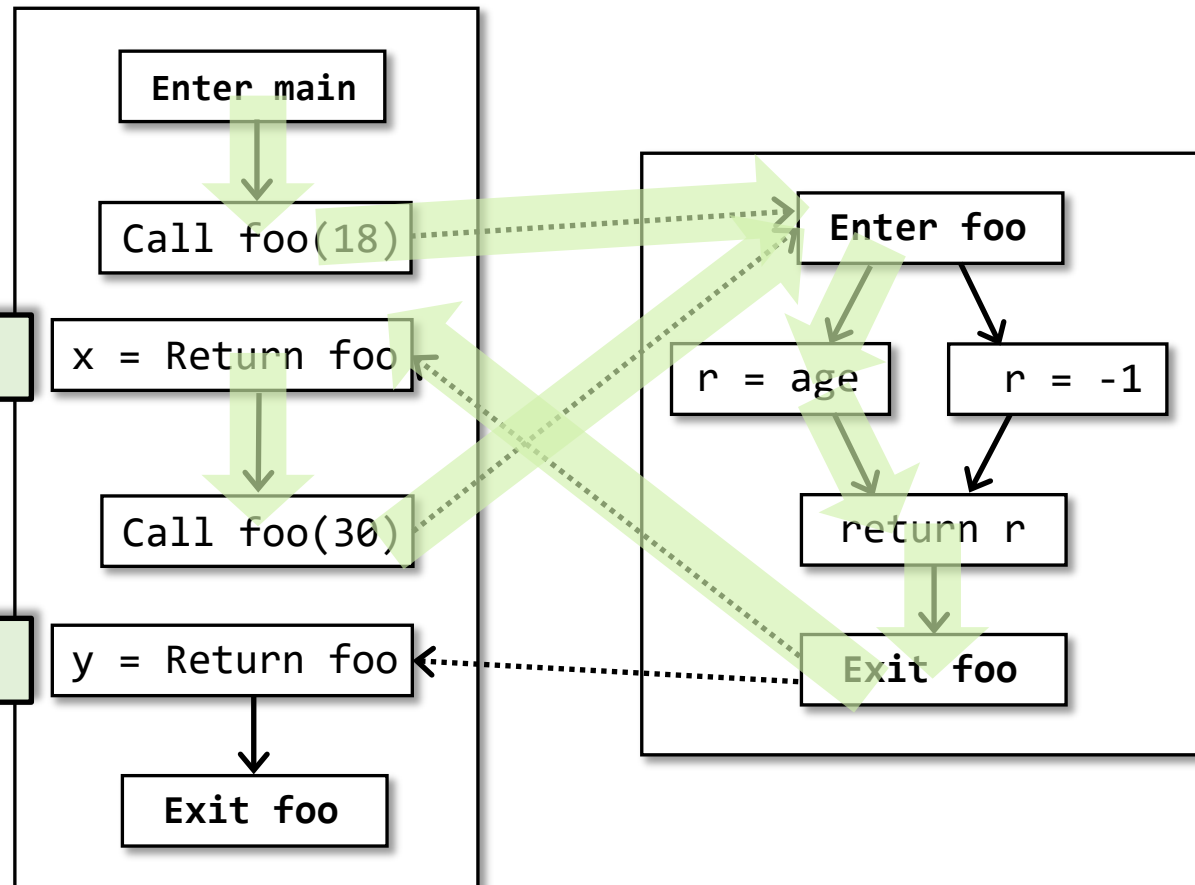
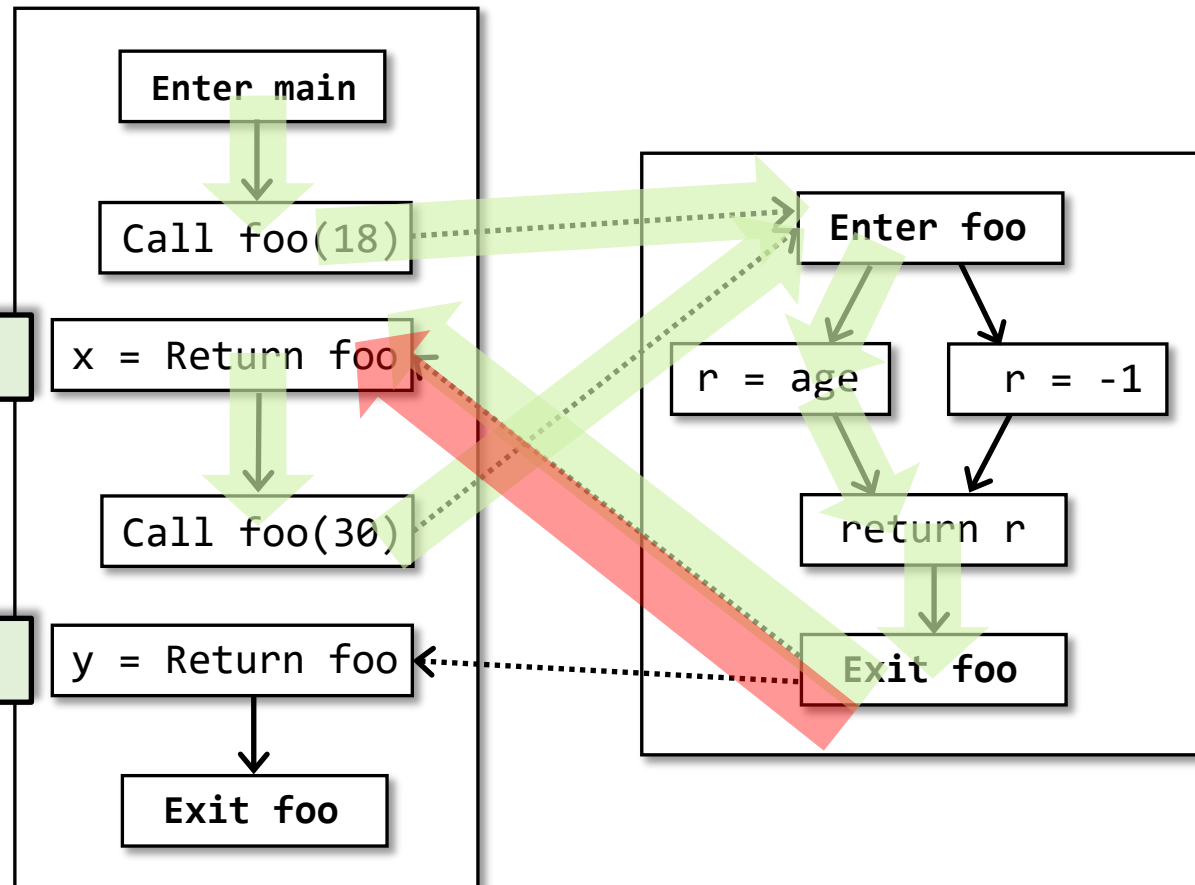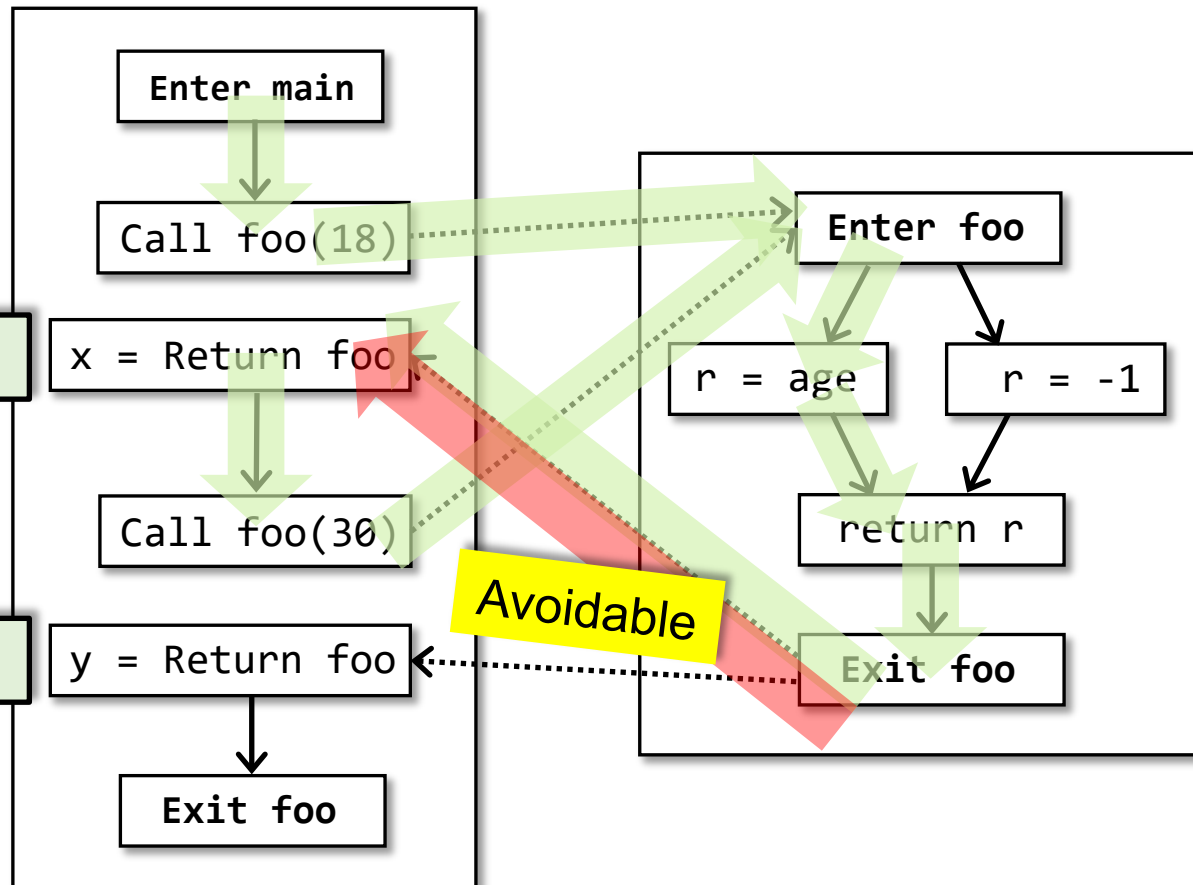x=18,30,-1

y=18,30,-1

Enter main

Call foo(18)

x = Return foo

Call foo(30)

y = Return foo

Exit foo

Enter foo

r = age

r = -1

return r

Exit foo

Avoidable

# Realizable Paths

**Realizable Paths:**

The paths in which "returns" are matched with corresponding "calls"

# Realizable Paths

Realizable Paths:
The paths in which "returns" are matched with corresponding "calls"

- Realizable paths may not be executable, but unrealizable paths must not be executable.

- Our goal is to recognize realizable paths so that we could avoid polluting analysis results along unrealizable paths.
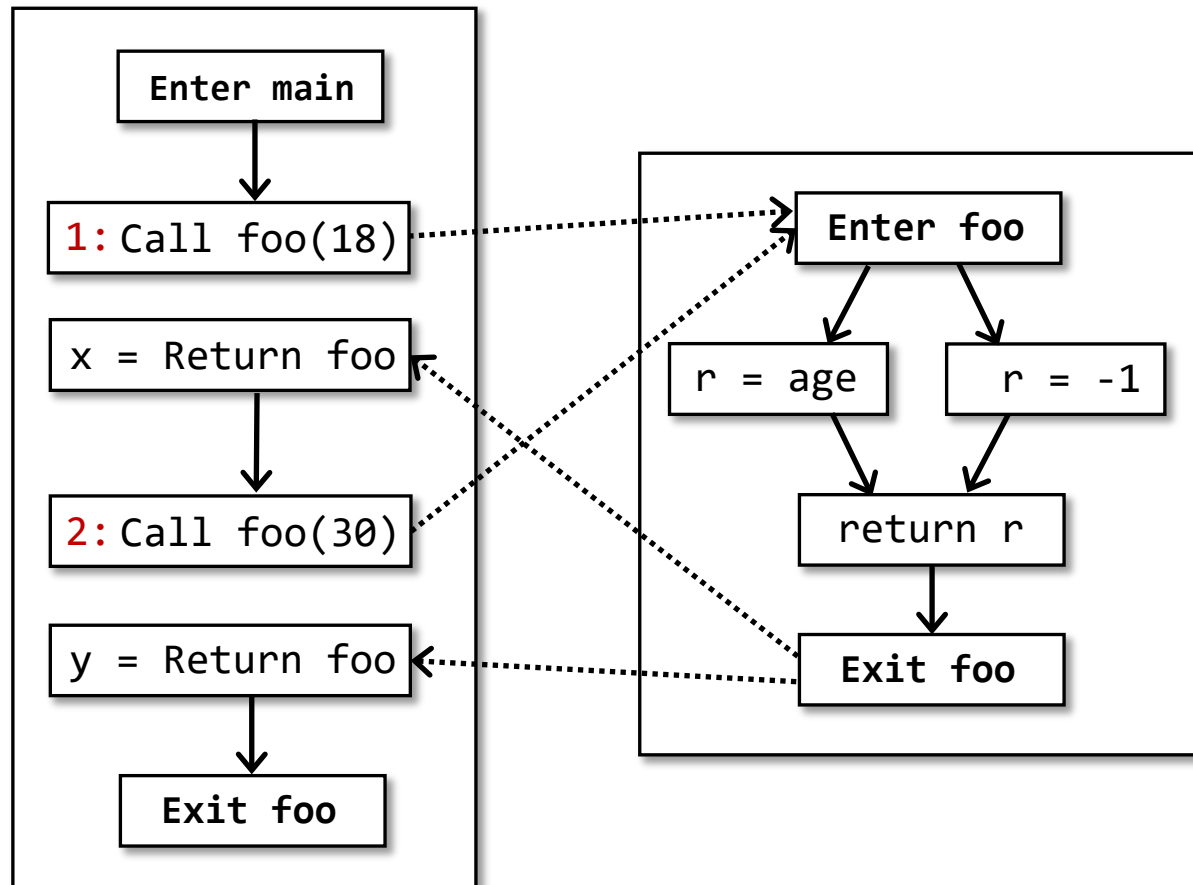
# Realizable Paths

> Realizable Paths:
> The paths in which "returns" are matched with corresponding "calls"

- Realizable paths may not be executable, but unrealizable paths must not be executable.

- Our goal is to recognize realizable paths so that we could avoid polluting analysis results along unrealizable paths.

**How?**

```
main() {
  x = foo(18);
    ⋮
  y = foo(30);
}
foo(int age) {
  if(age >= 0)
    r = age;
  else
    r = -1;
  return r;
}
```
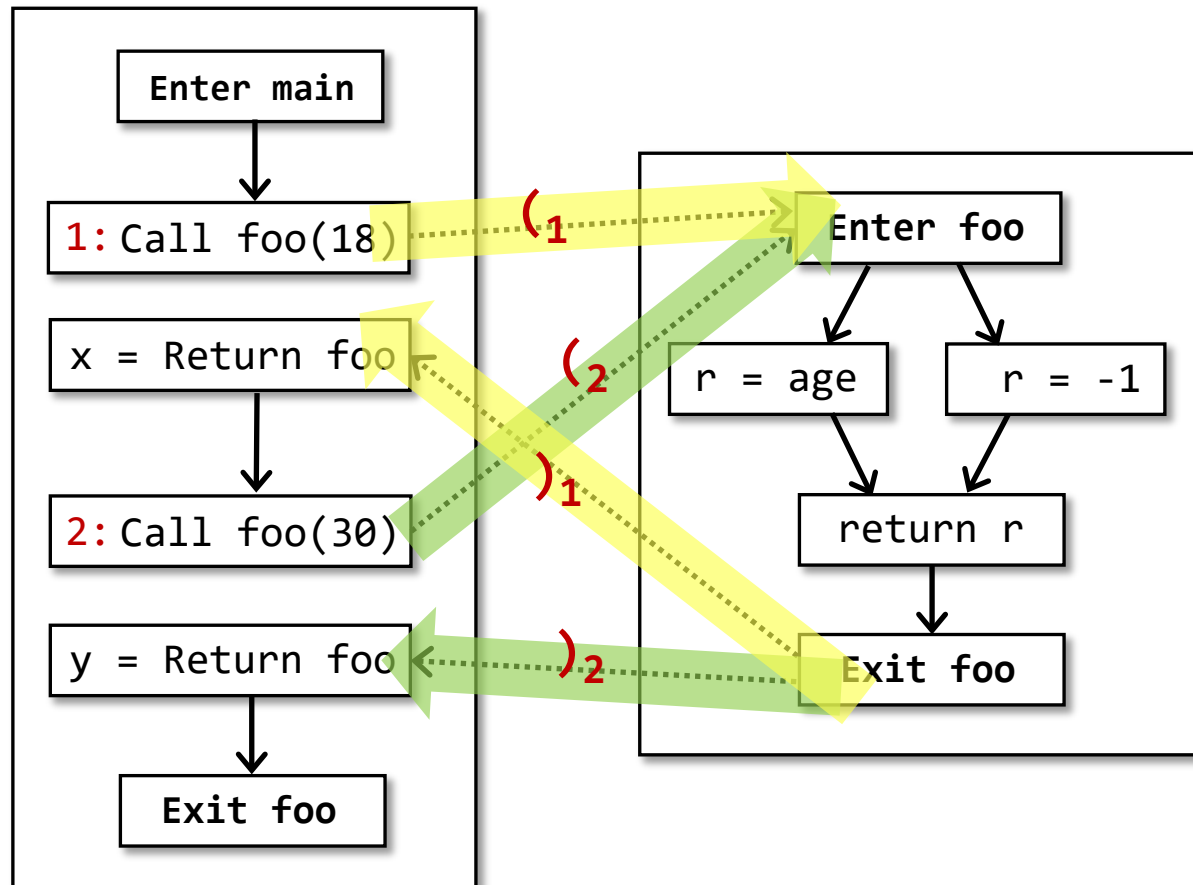
Enter main

1: Call foo(18)

x = Return foo

2: Call foo(30)

y = Return foo

Exit foo

Enter foo

r = age

r = -1

return r

Exit foo

```
main() {
  x = foo(18);
    ⋮
  y = foo(30);
}
foo(int age) {

  if(age >= 0)
    r = age;
  else
    r = -1;
  return r;
}
```
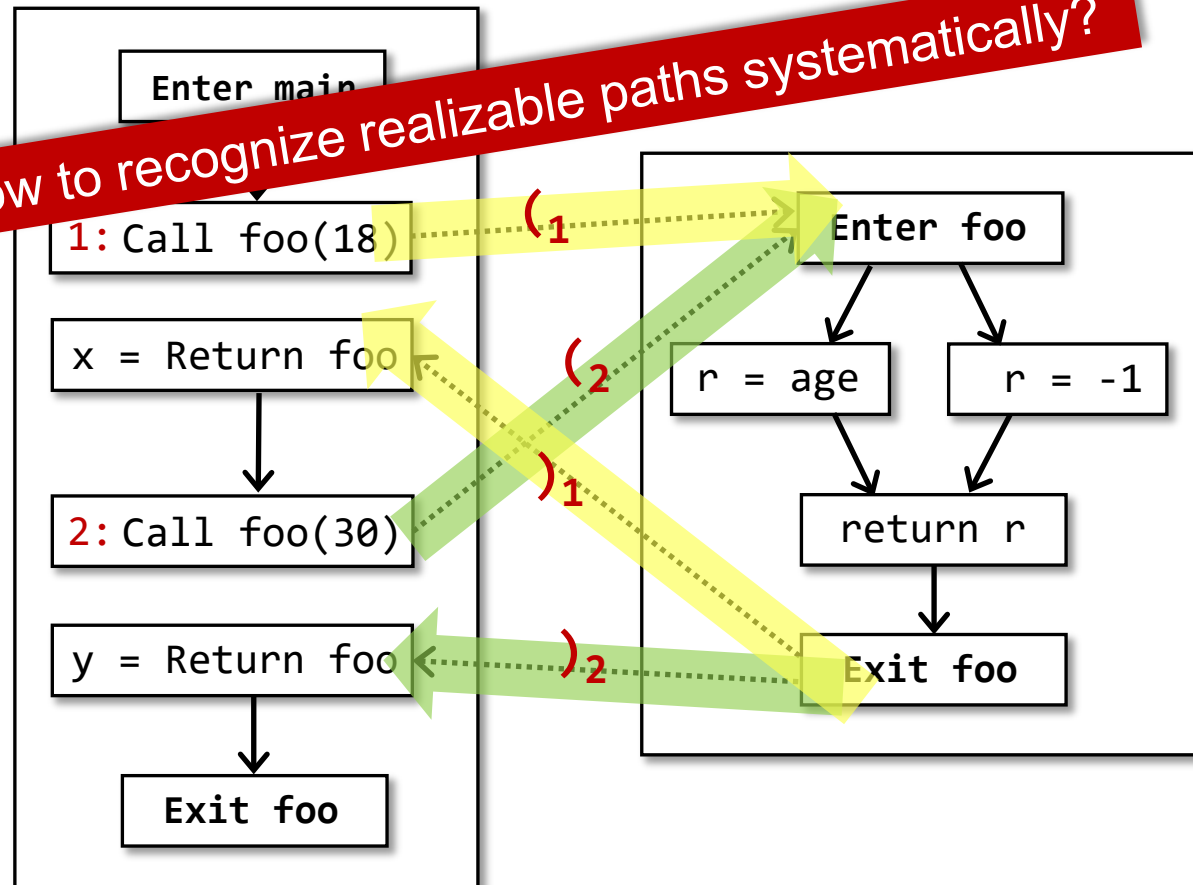
Enter main

1: Call foo(18)

x = Return foo

2: Call foo(30)

y = Return foo

Exit foo

Enter foo

r = age

r = -1

return r

Exit foo

(₁  (₂  )₁  )₂

How to recognize realizable paths systematically?

```
main() {
  x = foo(18);
    ⋮
  y = foo(30);
}

foo(int age) {

  if(age >= 0)

    r = age;
  else
    r = -1;
  return r;
}
```

Enter main

1: Call foo(18)

x = Return foo

2: Call foo(30)

y = Return foo

Exit foo

Enter foo

r = age     r = -1

return r

Exit foo

(₁  (₂  )₁  )₂

# CFL-Reachability

**CFL-Reachability**

A path is considered to connect two nodes A and B, or B is reachable from A, only if the concatenation of the labels on the edges of the path is a word in a specified context-free language.

# CFL-Reachability

> ## CFL-Reachability
>
> A path is considered to connect two nodes A and B, or B is reachable from A, only if the concatenation of the labels on the edges of the path is a word in a specified context-free language.

- A valid sentence in language L must follow L's grammar.
- A context-free language is a language generated by a context-free grammar (CFG).

# CFL-Reachability

> **CFL-Reachability**
>
> A path is considered to connect two nodes A and B, or B is reachable from A, only if the concatenation of the labels on the edges of the path is a word in a specified context-free language.

- A valid sentence in language L must follow L's grammar.
- A context-free language is a language generated by a context-free grammar (CFG).

CFG is a formal grammar in which every production is of the form:
$$S \rightarrow \alpha$$
where S is a single nonterminal and $\alpha$ could be a string of terminals and/or nonterminals, or empty.
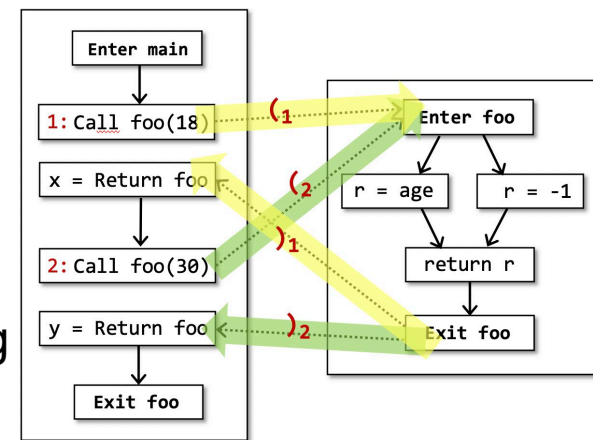
- S $\rightarrow$ aSb
- S $\rightarrow \varepsilon$

Context-free means S could be replaced by aSb/$\varepsilon$ anywhere, regardless of where S occurs.

# CFL-Reachability


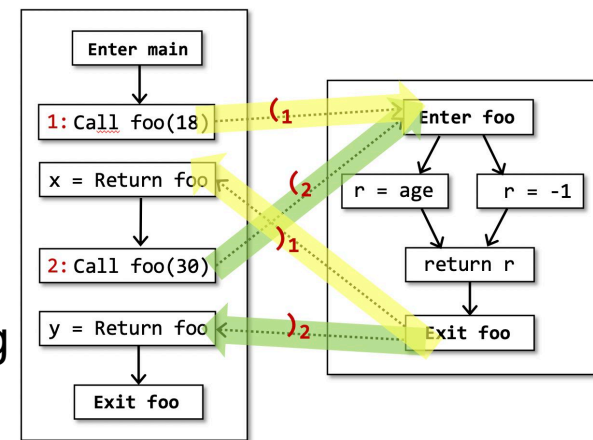
| Partially Balanced-Parenthesis Problem via CFL |
| --- |

- Every right parenthesis "$)_i$" is balanced by a preceding left parenthesis "$(_i$", but the converse needs not hold

- For each call site $i$, label its call edge "$(_i$" and return edge "$)_i$"

- Label all other edges with symbol "e"

# CFL-Reachability



Partially Balanced-Parenthesis Problem via CFL

- Every right parenthesis "$)_i$" is balanced by a preceding left parenthesis "$(_i$", but the converse needs not hold

- For each call site $i$, label its call edge "$(_i$" and return edge "$)_i$"
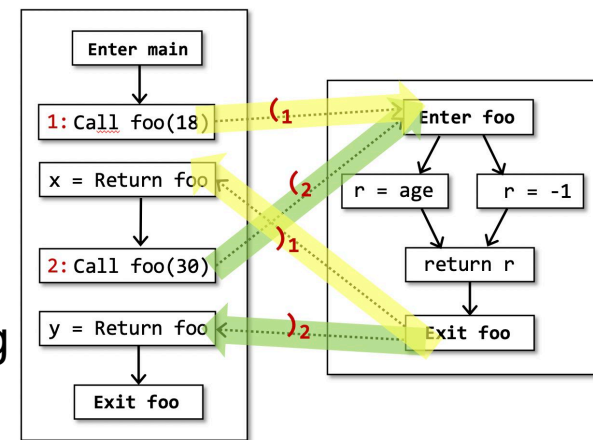
- Label all other edges with symbol "e"

A path is a realizable path iff the path' word is in the language *L(realizable)*

# CFL-Reachability



| Partially Balanced-Parenthesis Problem via CFL |
| --- |

- Every right parenthesis ")$_i$" is balanced by a preceding left parenthesis "($_i$", but the converse needs not hold

- For each call site $i$, label its call edge "($_i$" and return edge ")$_i$"
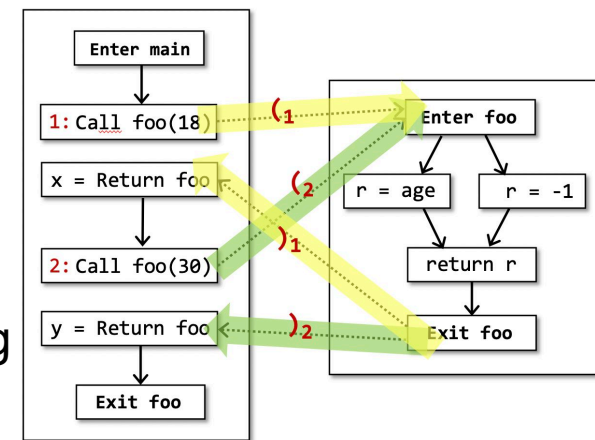
- Label all other edges with symbol "e"

A path is a realizable path iff the path' word is in the language *L(realizable)*

e.g., ($_1$ ($_2$ e )$_2$ )$_1$ ($_3$

# CFL-Reachability



| Partially Balanced-Parenthesis Problem via CFL |
| --- |

- Every right parenthesis "$)_i$" is balanced by a preceding left parenthesis "$(_i$", but the converse needs not hold

- For each call site $i$, label its call edge "$(_i$" and return edge "$)_i$"

- Label all other edges with symbol "e"

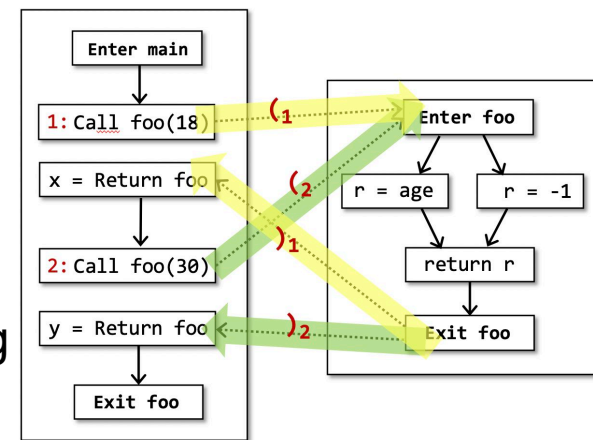A path is a realizable path iff the path' word is in the language *L(realizable)*

*realizable* → *matched  realizable*

e.g., $(_1 (_2 e )_2 )_1 (_3$

# CFL-Reachability



Partially Balanced-Parenthesis Problem via CFL

- Every right parenthesis "$)_i$" is balanced by a preceding left parenthesis "$(_i$", but the converse needs not hold

- For each call site $i$, label its call edge "$(_i$" and return edge "$)_i$"

- Label all other edges with symbol "e"

A path is a realizable path iff the path' word is in the language *L(realizable)*

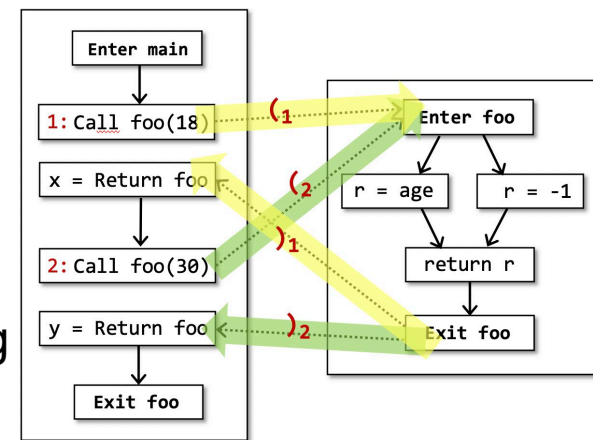*realizable* → *matched realizable*

         → $(_i$

e.g., $(_1 (_2 e )_2 )_1 (_3$

# CFL-Reachability



---
**Partially Balanced-Parenthesis Problem via CFL**
---

- Every right parenthesis "$)_i$" is balanced by a preceding left parenthesis "$(_i$", but the converse needs not hold

- For each call site $i$, label its call edge "$(_i$" and return edge "$)_i$"

- Label all other edges with symbol "e"

A path is a realizable path iff the path' word is in the language *L(realizable)*

$$\textit{realizable} \rightarrow \textit{matched } \textit{realizable}$$
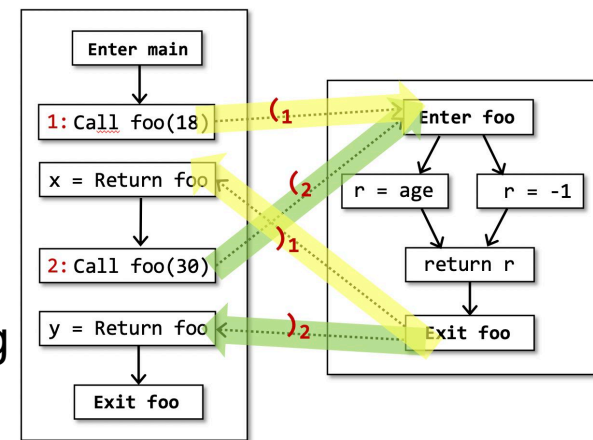$$\rightarrow (_i$$

e.g., $(_1 (_2 \ e \ )_2 )_1 (_3$

e.g., $(_1 (_2 \ e \ )_2 )_1 (_3 (_4$

# CFL-Reachability



| Partially Balanced-Parenthesis Problem via CFL |
|---|

- Every right parenthesis "$)_i$" is balanced by a preceding left parenthesis "$(_i$", but the converse needs not hold

- For each call site $i$, label its call edge "$(_i$" and return edge "$)_i$"

- Label all other edges with symbol "e"

A path is a realizable path iff the path' word is in the language *L(realizable)*

*realizable* → *matched realizable*
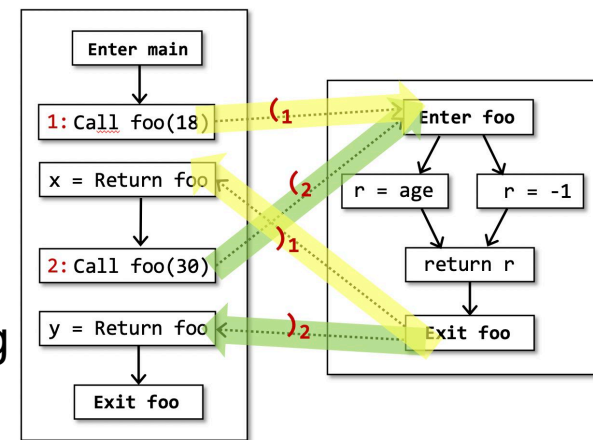         → $(_i$ *realizable*
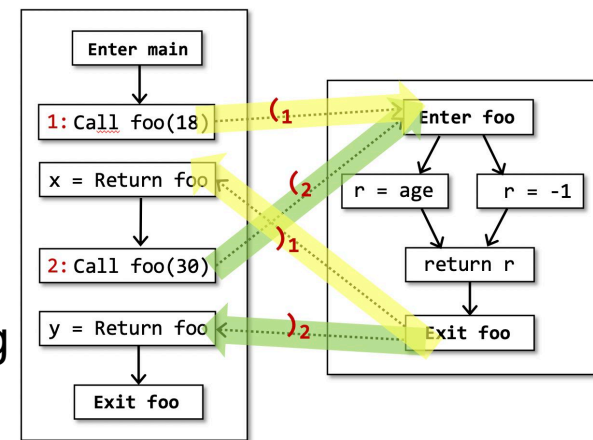
e.g., $(_1 (_2 e )_2 )_1 (_3$

e.g., $(_1 (_2 e )_2 )_1 (_3 (_4$

# CFL-Reachability



Partially Balanced-Parenthesis Problem via CFL

- Every right parenthesis "$)_i$" is balanced by a preceding left parenthesis "$(_i$", but the converse needs not hold

- For each call site $i$, label its call edge "$(_i$" and return edge "$)_i$"

- Label all other edges with symbol "e"

A path is a realizable path iff the path' word is in the language *L(realizable)*

$$realizable \rightarrow matched \ \ realizable$$
$$\rightarrow (_i \ \ realizable$$
$$\rightarrow \varepsilon$$

e.g., $(_1 \, (_2 \, e \, )_2 \, )_1 \, (_3$

e.g., $(_1 \, (_2 \, e \, )_2 \, )_1 \, (_3 \, (_4$

# CFL-Reachability



Partially Balanced-Parenthesis Problem via CFL

- Every right parenthesis ")$_i$" is balanced by a preceding left parenthesis "($_i$", but the converse needs not hold

- For each call site $i$, label its call edge "($_i$" and return edge ")$_i$"

- Label all other edges with symbol "e"

A path is a realizable path iff the path' word is in the language *L(realizable)*

$$realizable \rightarrow matched\ realizable$$
$$\rightarrow (_i\ realizable$$
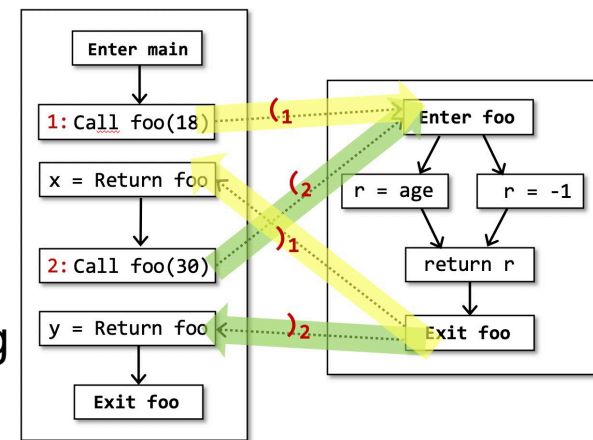$$\rightarrow \varepsilon$$

$$matched \rightarrow (_i\ matched\ )_i$$

e.g., ($_1$ ($_2$ e )$_2$ )$_1$ ($_3$

e.g., ($_1$ ($_2$ e )$_2$ )$_1$ ($_3$ ($_4$

# CFL-Reachability



Partially Balanced-Parenthesis Problem via CFL

- Every right parenthesis "$)_i$" is balanced by a preceding left parenthesis "$(_i$", but the converse needs not hold

- For each call site $i$, label its call edge "$(_i$" and return edge "$)_i$"

- Label all other edges with symbol "e"

A path is a realizable path iff the path' word is in the language *L(realizable)*

*realizable* → *matched realizable*

→ $(_i$ *realizable*

→ $\varepsilon$

*matched* → $(_i$ *matched* $)_i$

→ e

→ $\varepsilon$

e.g., $(_1 (_2 \text{ e } )_2 )_1 (_3$

e.g., $(_1 (_2 \text{ e } )_2 )_1 (_3 (_4$

# CFL-Reachability



Partially Balanced-Parenthesis Problem via CFL

- Every right parenthesis ")$_i$" is balanced by a preceding left parenthesis "($_i$", but the converse needs not hold

- For each call site $i$, label its call edge "($_i$" and return edge ")$_i$"

- Label all other edges with symbol "e"

A path is a realizable path iff the path' word is in the language *L(realizable)*

$$realizable \rightarrow matched\ realizable$$
$$\rightarrow (_i\ realizable$$
$$\rightarrow \varepsilon$$

$$matched \rightarrow (_i\ matched\ )_i$$
$$\rightarrow e$$
$$\rightarrow \varepsilon$$

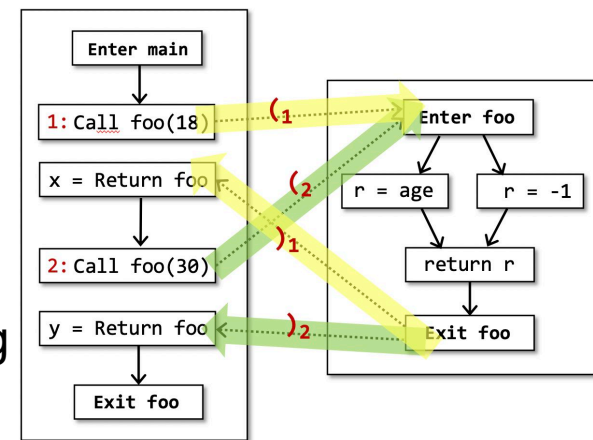e.g., ($_1$ ($_2$ e )$_2$ )$_1$ ($_3$

e.g., ($_1$ ($_2$ e )$_2$ )$_1$ ($_3$ ($_4$

e.g., ($_1$ ($_2$ e e e )$_2$ )$_1$ ($_3$ ($_4$

# CFL-Reachability



| Partially Balanced-Parenthesis Problem via CFL |
|---|

- Every right parenthesis ")$_i$" is balanced by a preceding left parenthesis "($_i$", but the converse needs not hold

- For each call site $i$, label its call edge "($_i$" and return edge ")$_i$"

- Label all other edges with symbol "e"

A path is a realizable path iff the path' word is in the language *L(realizable)*

$$realizable \rightarrow matched\ realizable$$
$$\rightarrow (_i\ realizable$$
$$\rightarrow \varepsilon$$

$$matched \rightarrow (_i\ matched\ )_i$$
$$\rightarrow e$$
$$\rightarrow \varepsilon$$
$$\rightarrow matched\ matched$$
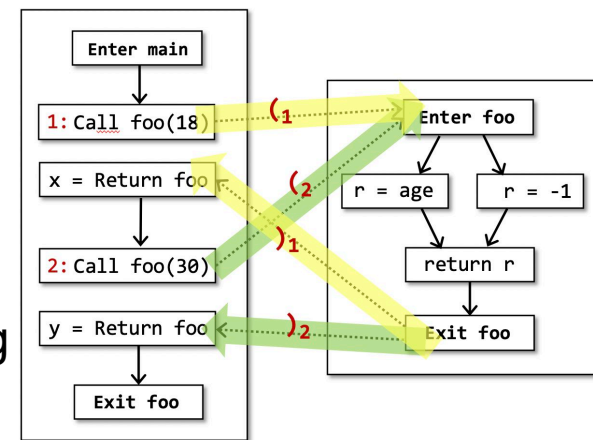
e.g., ($_1$ ($_2$ e )$_2$ )$_1$ ($_3$

e.g., ($_1$ ($_2$ e )$_2$ )$_1$ ($_3$ ($_4$

e.g., ($_1$ ($_2$ e e e )$_2$ )$_1$ ($_3$ ($_4$

# CFL-Reachability



| Partially Balanced-Parenthesis Problem via CFL |
| --- |

- Every right parenthesis ")$_i$" is balanced by a preceding left parenthesis "($_i$", but the converse needs not hold

- For each call site $i$, label its call edge "($_i$" and return edge ")$_i$"

- Label all other edges with symbol "e"

A path is a realizable path iff the path' word is in the language *L(realizable)*

*realizable* → *matched  realizable*

  → ($_i$  *realizable*

  → ε

*matched* → ($_i$ *matched* )$_i$

  → e

  → ε

  → *matched matched*

e.g., ($_1$ ($_2$ e )$_2$ )$_1$ ($_3$

e.g., ($_1$ ($_2$ e )$_2$ )$_1$ ($_3$ ($_4$

e.g., ($_1$ ($_2$ e e e )$_2$ )$_1$ ($_3$ ($_4$

e.g., e e ($_1$ ($_2$ e e e )$_2$ )$_1$ ($_3$ ($_4$ e

$L(realizable):$

$realizable \rightarrow matched\ realizable$

$\rightarrow (_i\ realizable$

$\rightarrow \varepsilon$

$matched \rightarrow (_i\ matched\ )_i$

$\rightarrow e$

$\rightarrow \varepsilon$

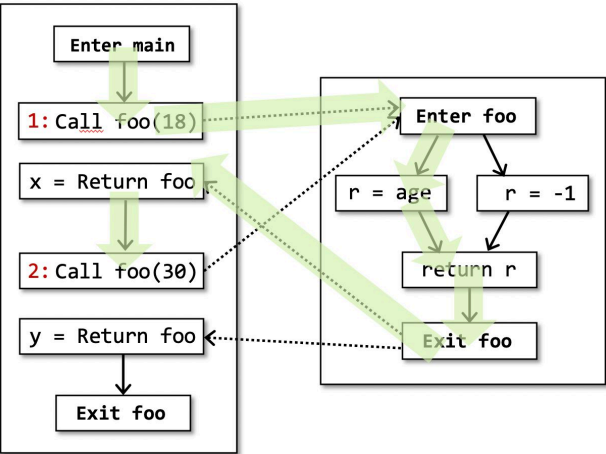$\rightarrow matched\ matched$

$L(realizable):$

$realizable \rightarrow matched\ realizable$
$\rightarrow (_i\ realizable$
$\rightarrow \varepsilon$

$matched \rightarrow (_i\ matched\ )_i$
$\rightarrow e$
$\rightarrow \varepsilon$
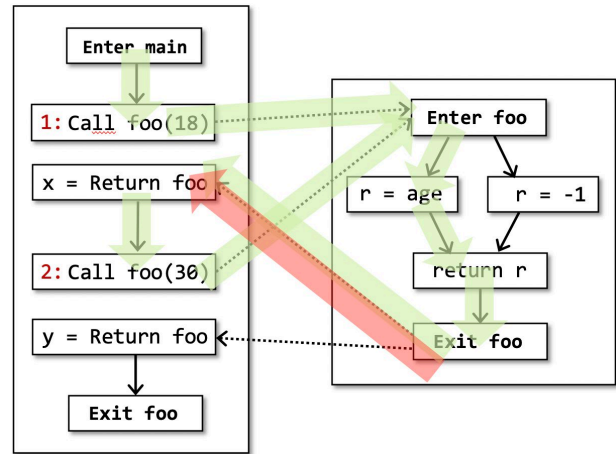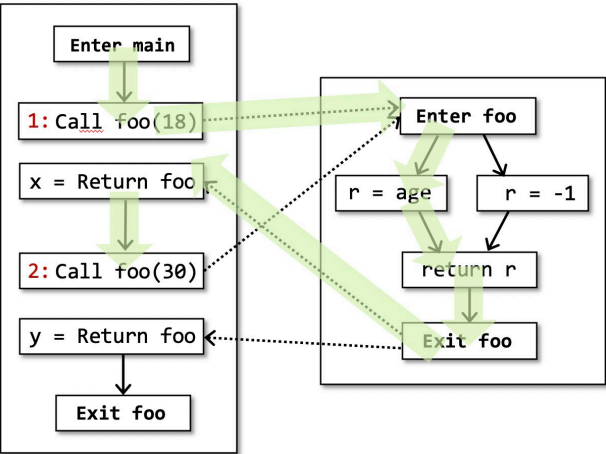$\rightarrow matched\ matched$

$e(_1eee)_1e \in L(realizable)$

$$L(realizable):$$

$$realizable \rightarrow matched\ realizable$$
$$\rightarrow (_i\ realizable$$
$$\rightarrow \varepsilon$$

$$matched \rightarrow (_i\ matched\ )_i$$
$$\rightarrow e$$
$$\rightarrow \varepsilon$$
$$\rightarrow matched\ matched$$

$$e(_1 eee)_1 e \in L(realizable)$$

$$e(_1 eee)_1 e(_2 eee)_1$$

$$\notin L(realizable)$$

# IFDS

A Program Analysis Framework via Graph Reachability

Yue Li @ Nanjing University

# IFDS

"Precise Interprocedural Dataflow Analysis via Graph Reachability"

*Thomas Reps, Susan Horwitz, and Mooly Sagiv, POPL'95*

IFDS (Interprocedural, Finite, Distributive, Subset Problem)

# IFDS

"Precise Interprocedural Dataflow Analysis via Graph Reachability"

*Thomas Reps, Susan Horwitz, and Mooly Sagiv, POPL'95*

IFDS (Interprocedural, Finite, Distributive, Subset Problem)

IFDS is for interprocedural data flow analysis
with distributive flow functions over finite domains.

# IFDS

"Precise Interprocedural Dataflow Analysis via Graph Reachability"

*Thomas Reps, Susan Horwitz, and Mooly Sagiv, POPL'95*

IFDS (Interprocedural, Finite, Distributive, Subset Problem)

IFDS is for interprocedural data flow analysis
with distributive flow functions over finite domains.

Provide meet-over-all-realizable-paths (MRP) solution.

# Meet-Over-All-Realizable-Paths (MRP)

Path function for path $p$, denoted as $pf_p$, is a composition of flow functions for all edges (sometimes nodes) on $p$.

Recall

$$pf_p = f_n \circ \dots \circ f_2 \circ f_1$$

# Meet-Over-All-Realizable-Paths (MRP)

Path function for path $p$, denoted as $pf_p$, is a composition of flow functions for all edges (sometimes nodes) on $p$.

**Recall**

$$pf_p = f_n \circ \dots \circ f_2 \circ f_1$$

$$\text{MOP}_n = \bigsqcup_{p \in \text{Paths}(start, n)} pf_p(\bot)$$

For each node n, MOP$_n$ provides a "meet-over-all-paths" solution where $\text{Paths}(start, n)$ denotes the set of paths in CFG from the start node to n.

# Meet-Over-All-Realizable-Paths (MRP)

Path function for path $p$, denoted as $pf_p$, is a composition of flow functions for all edges (sometimes nodes) on $p$.

$$pf_p = f_n \circ \ldots \circ f_2 \circ f_1$$

$$\text{MOP}_n = \bigsqcup_{p \in \text{Paths}(start, n)} pf_p (\bot)$$

For each node n, $\text{MOP}_n$ provides a "meet-over-all-paths" solution where $\text{Paths}(start, n)$ denotes the set of paths in CFG from the start node to n.

$$\text{MRP}_n = \bigsqcup_{p \in \text{RPaths}(start, n)} pf_p (\bot)$$

For each node n, $\text{MRP}_n$ provides a "meet-over-all-realizable-paths" solution where $\text{RPaths}(start, n)$ denotes the set of realizable paths (the path's word is in the language $L(realizable)$) from the start node to n.

# Meet-Over-All-Realizable-Paths (MRP)

Path function for path $p$, denoted as $pf_p$, is a composition of flow functions for all edges (sometimes nodes) on $p$.

Recall

$$pf_p = f_n \circ \ldots \circ f_2 \circ f_1$$

$$MOP_n = \bigsqcup_{p \in \text{Paths}(start, n)} pf_p(\bot)$$

For each node n, $MOP_n$ provides a "meet-over-all-paths" solution where $\text{Paths}(start, n)$ denotes the set of paths in CFG from the start node to n.

$$MRP_n = \bigsqcup_{p \in \text{RPaths}(start, n)} pf_p(\bot)$$

For each node n, $MRP_n$ provides a "meet-over-all-realizable-paths" solution where $\text{RPaths}(start, n)$ denotes the set of realizable paths (the path's word is in the language $L(realizable)$) from the start node to n.
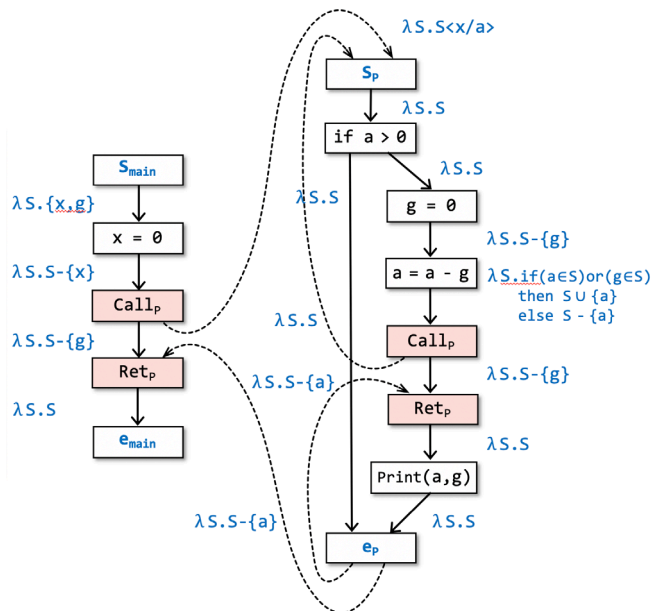
$$MRP_n \sqsubseteq MOP_n$$

# Overview of IFDS

Given a program P, and a dataflow-analysis problem Q

# Overview of IFDS

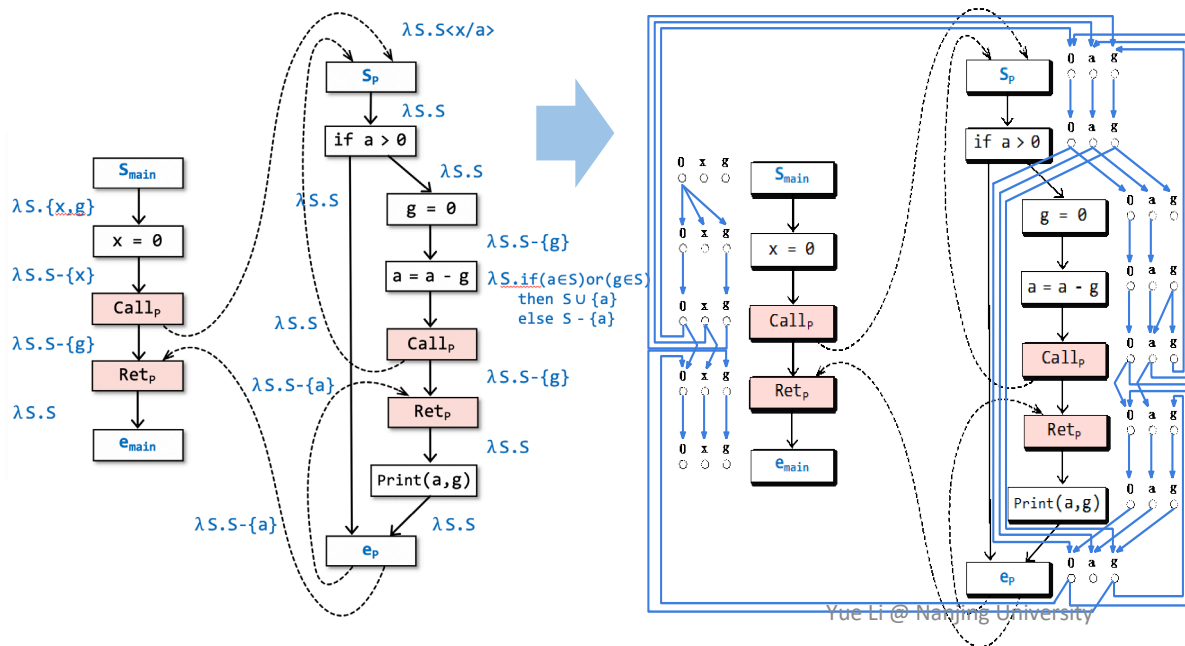Given a program P, and a dataflow-analysis problem Q

- Build a supergraph G* for P and
  define flow functions for edges in G* based on Q

# Overview of IFDS
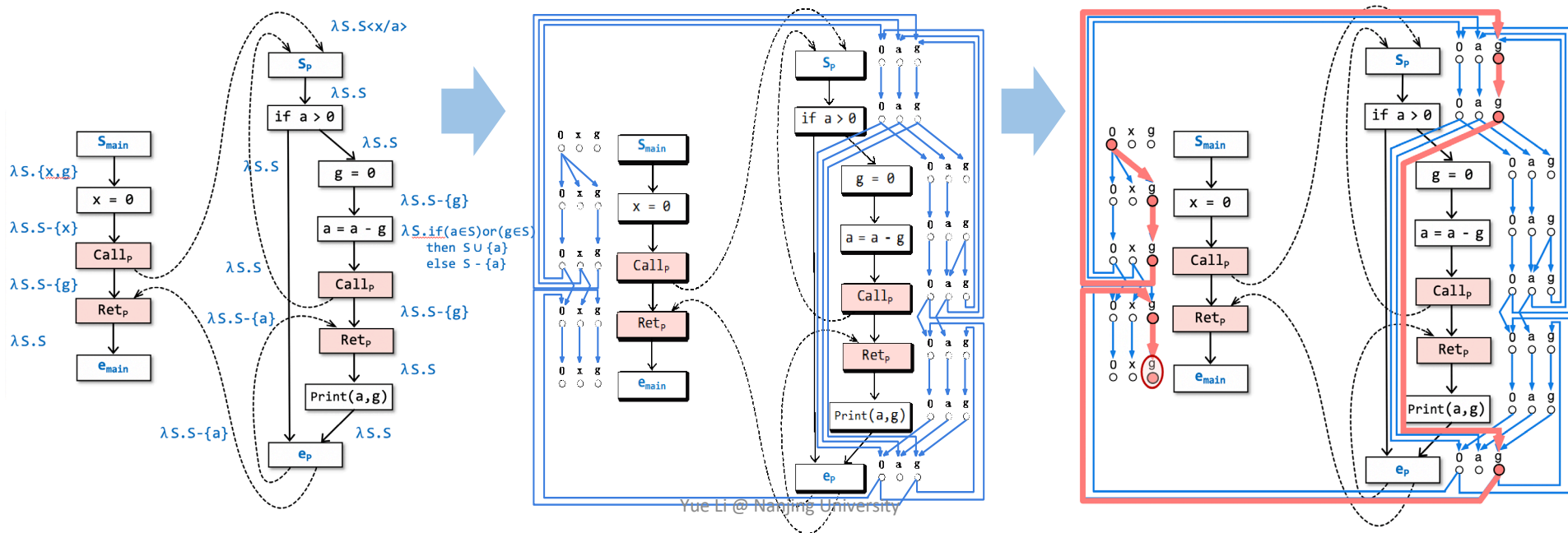
Given a program P, and a dataflow-analysis problem Q

- Build a supergraph G* for P and
  define flow functions for edges in G* based on Q

- Build exploded supergraph G# for P by transforming
  flow functions to representation relations (graphs)

# Overview of IFDS
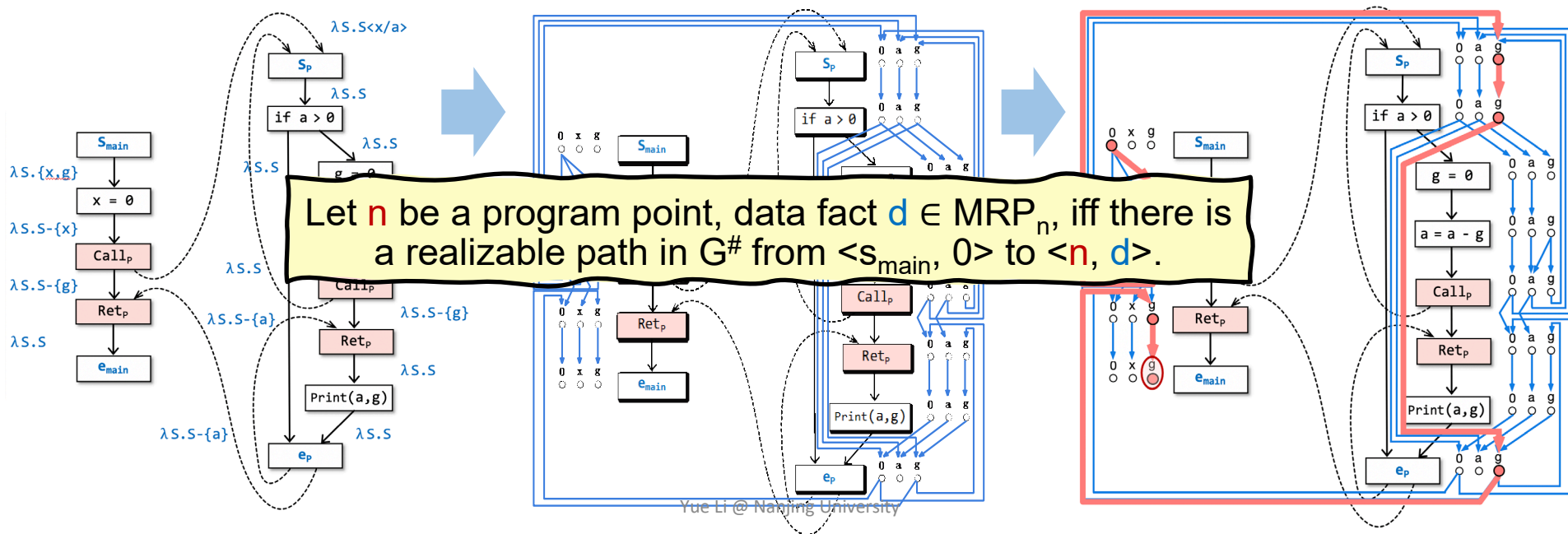
Given a program P, and a dataflow-analysis problem Q

- Build a supergraph G* for P and
  define flow functions for edges in G* based on Q

- Build exploded supergraph G# for P by transforming
  flow functions to representation relations (graphs)

- Q can be solved as graph reachability problems (find out MRP solutions)
  via applying Tabulation algorithm on G#

# Overview of IFDS
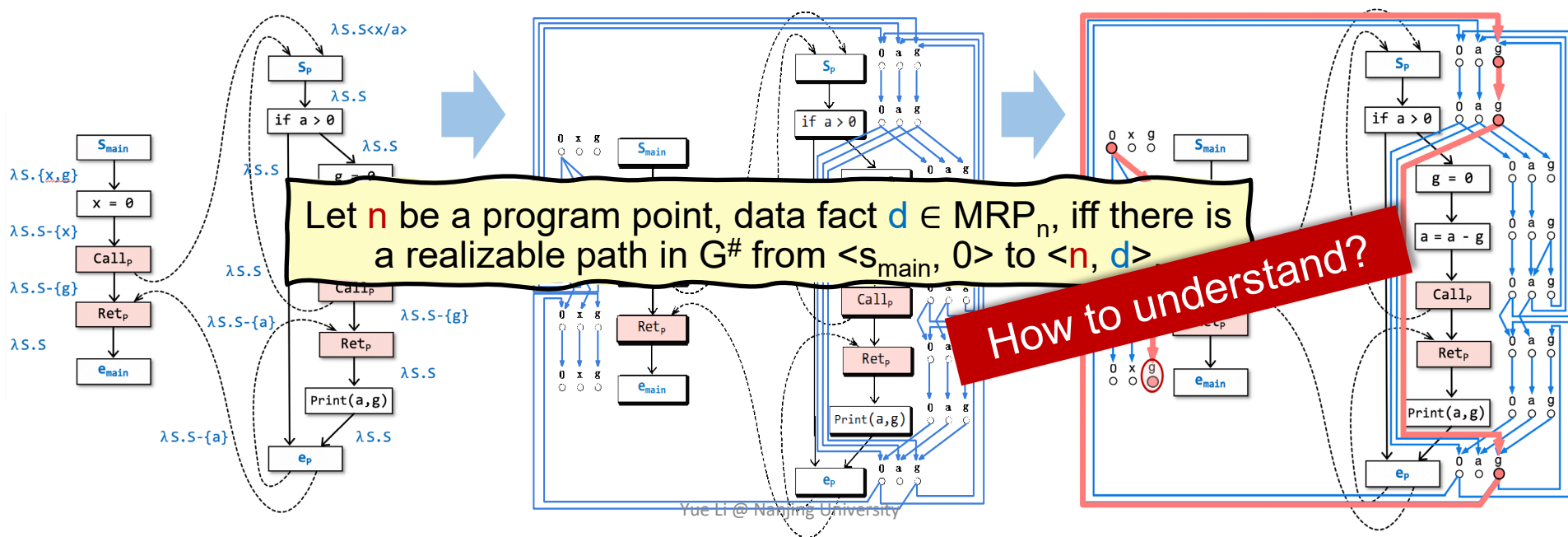
Given a program P, and a dataflow-analysis problem Q

- Build a supergraph G* for P and
  define flow functions for edges in G* based on Q

- Build exploded supergraph G# for P by transforming
  flow functions to representation relations (graphs)

- Q can be solved as graph reachability problems (find out MRP solutions)
  via applying Tabulation algorithm on G#



Let n be a program point, data fact $d \in MRP_n$, iff there is
a realizable path in G# from $<s_{main}, 0>$ to $<n, d>$.

# Overview of IFDS

Given a program P, and a dataflow-analysis problem Q
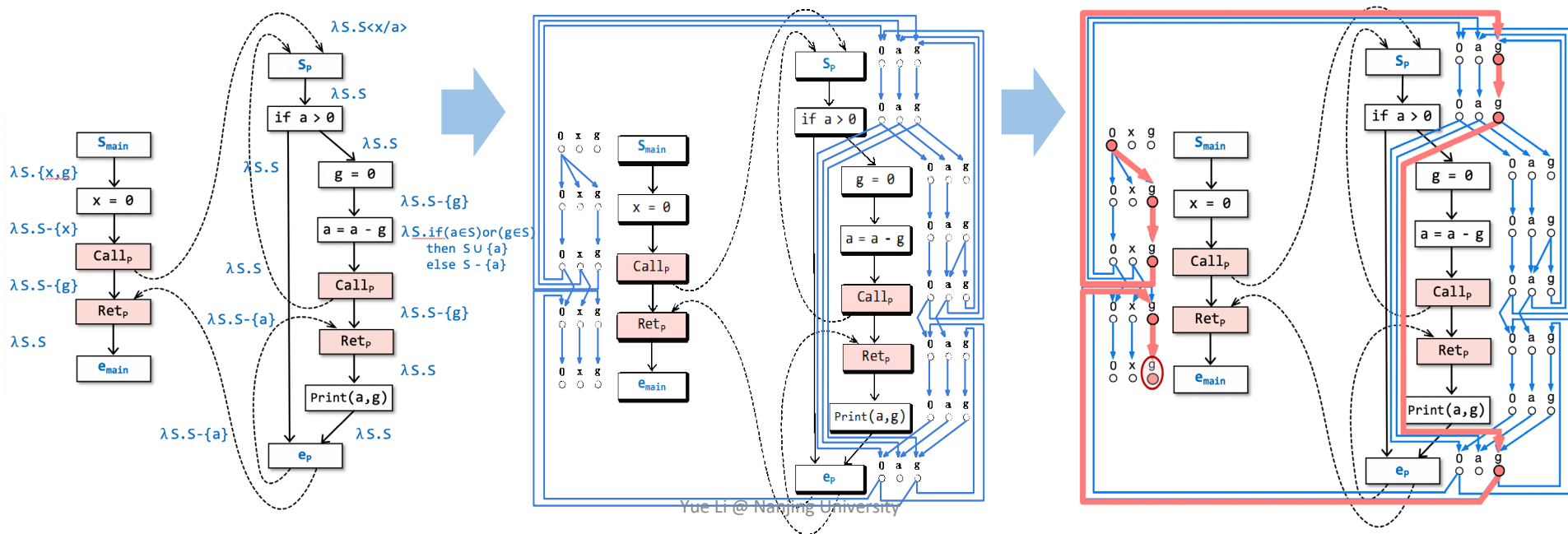
- Build a supergraph G* for P and
  define flow functions for edges in G* based on Q

- Build exploded supergraph G# for P by transforming
  flow functions to representation relations (graphs)

- Q can be solved as graph reachability problems (find out MRP solutions)
  via applying Tabulation algorithm on G#



Let n be a program point, data fact $d \in MRP_n$, iff there is
a realizable path in G# from $\langle s_{main}, 0 \rangle$ to $\langle n, d \rangle$.

How to understand?

# Overview of IFDS

Given a program P, and a dataflow-analysis problem Q
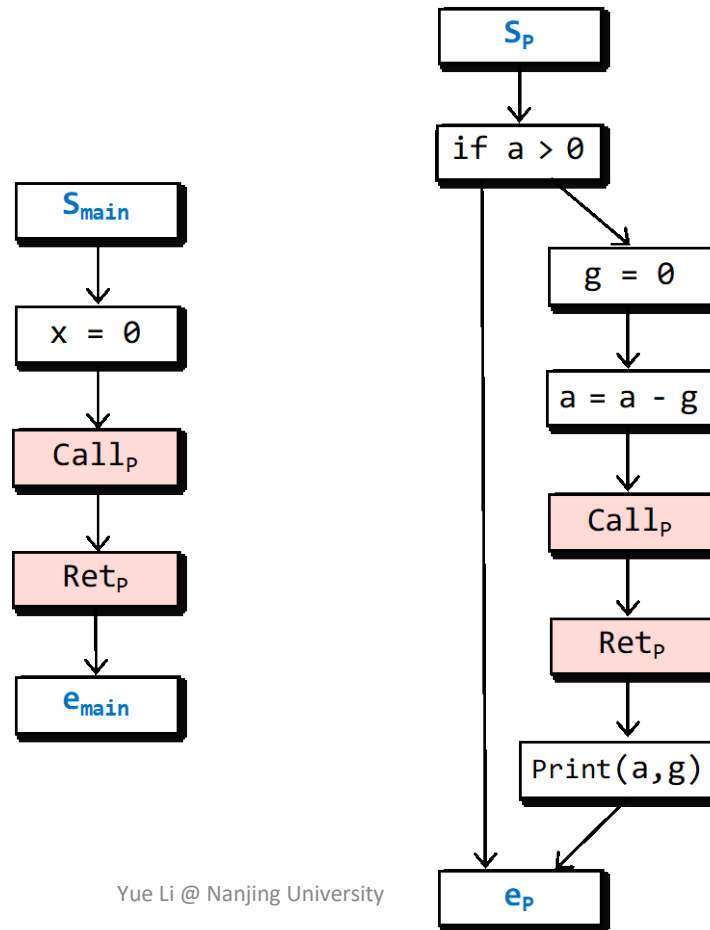
- Build a supergraph G* for P and
  define flow functions for edges in G* based on Q

- Build exploded supergraph G# for P by transforming
  flow functions to representation relations (graphs)

- Q can be solved as graph reachability problems (find out MRP solutions)
  via applying Tabulation algorithm on G#

# Supergraph

In IFDS, a program is represented by G* = (N*, E*) called a supergraph.

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a – g;
    P(a);
    Print(a,g);
  }
}
```
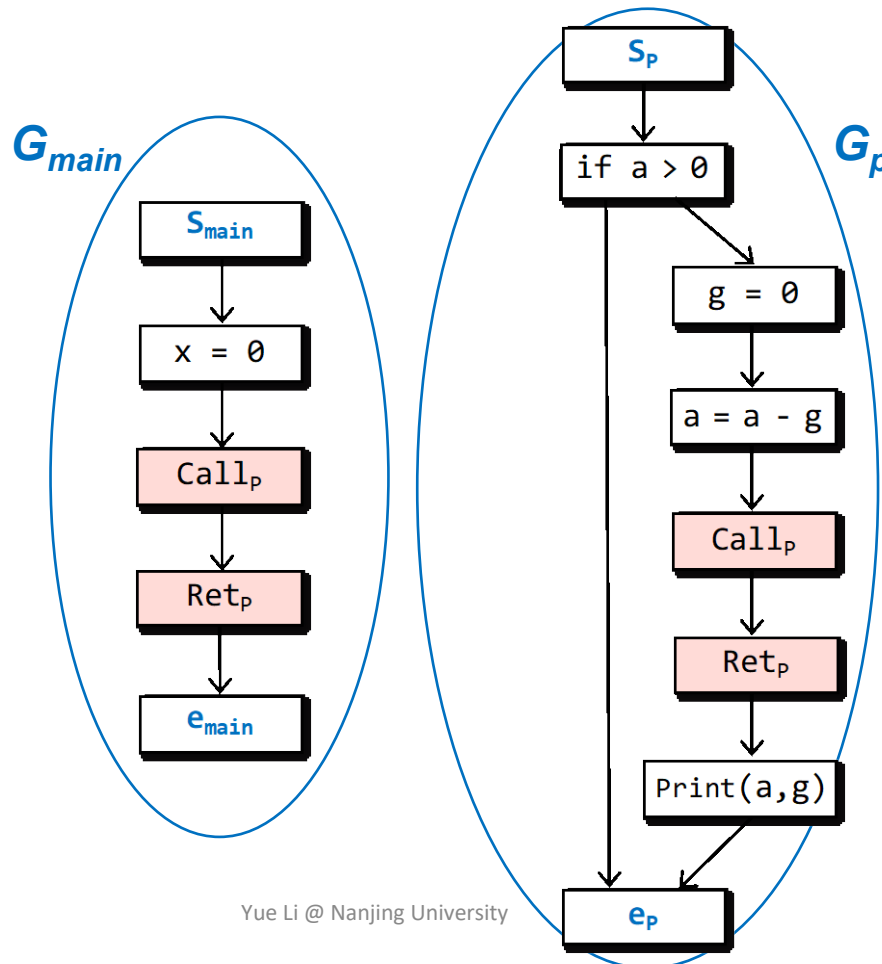
# Supergraph

In IFDS, a program is represented by G* = (N*, E*) called a supergraph.

- *G* consists of a collection of flow graphs $G_1$, $G_2$,... (one for each procedure)*

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a – g;
    P(a);
    Print(a,g);
  }
}
```

# Supergraph

In IFDS, a program is represented by G* = (N*, E*) called a supergraph.

- *G\* consists of a collection of flow graphs $G_1$, $G_2$, … (one for each procedure)*
- *Each flowgraph $G_p$ has a unique start node $s_p$, and a unique exit node $e_p$*

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a - g;
    P(a);
    Print(a,g);
  }
}
```
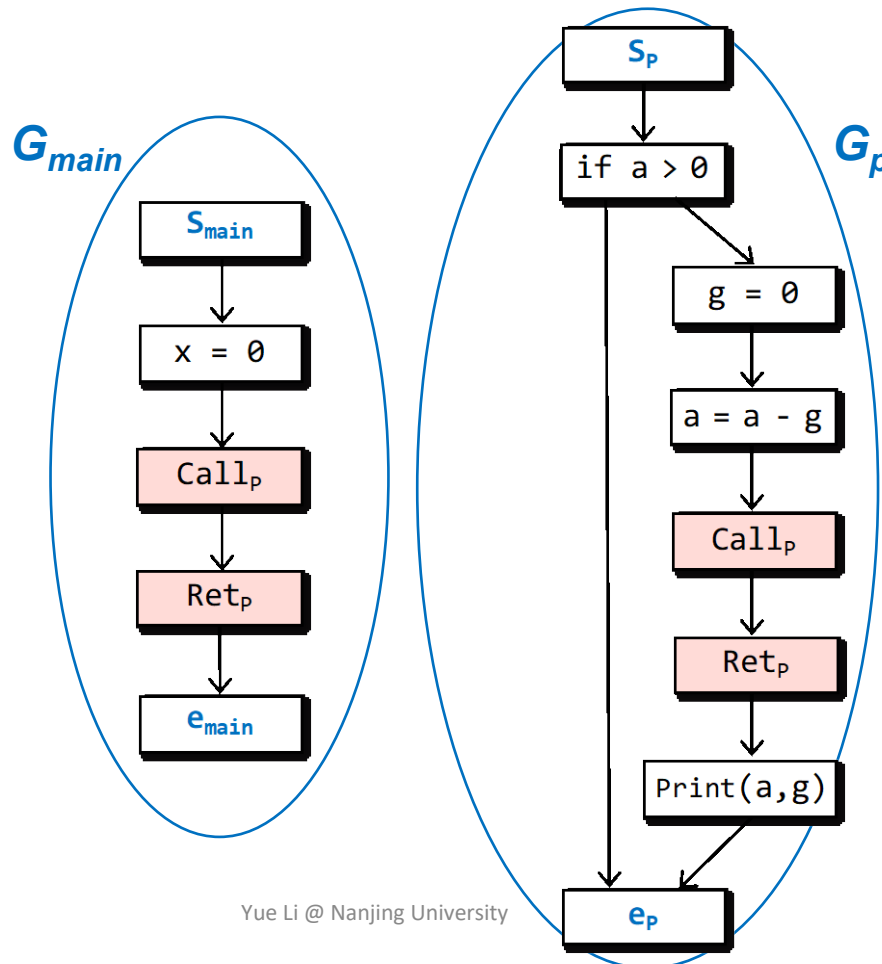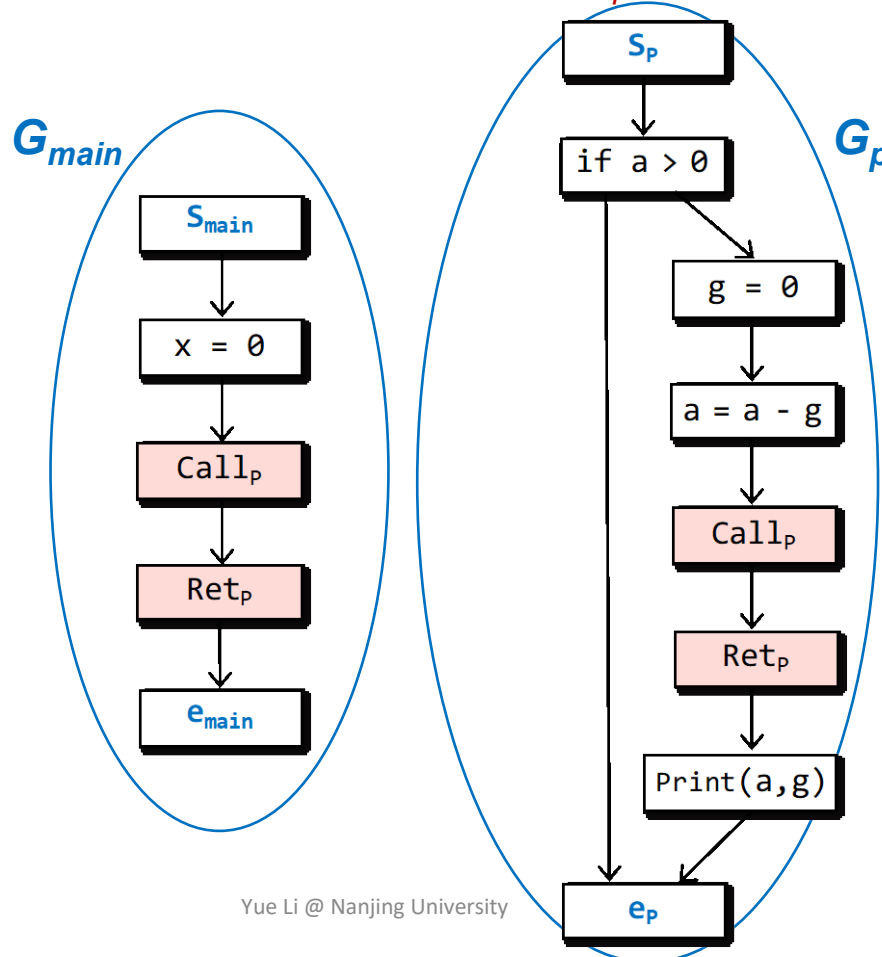
# Supergraph

In IFDS, a program is represented by G* = (N*, E*) called a supergraph.

- *G* consists of a collection of flow graphs $G_1$, $G_2$,… (one for each procedure)*
- *Each flowgraph $G_p$ has a unique start node $s_p$, and a unique exit node $e_p$*
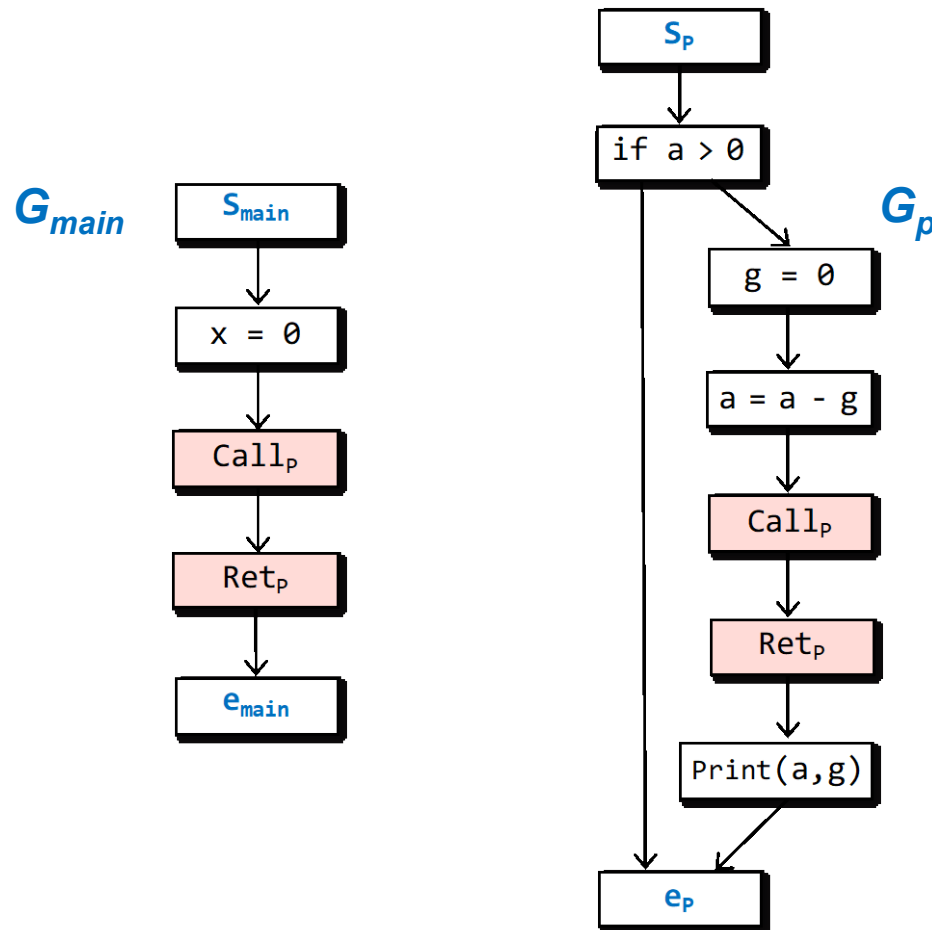- *A procedure call is represented by a call node $Call_p$, and a return-site node $Ret_p$*

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a – g;
    P(a);
    Print(a,g);
  }
}
```

# Supergraph

*G\* has three edges for each procedure call:*

```
int g;
main(){
    int x;
    x = 0;
    P(x);
}
P(int a){
    if(a > 0){
        g = 0;
        a = a – g;
        P(a);
        Print(a,g);
    }
}
```

**G<sub>main</sub>**

- S<sub>main</sub>
- x = 0
- Call<sub>P</sub>
- Ret<sub>P</sub>
- e<sub>main</sub>

**G<sub>P</sub>**

- S<sub>P</sub>
- if a > 0
- g = 0
- a = a - g
- Call<sub>P</sub>
- Ret<sub>P</sub>
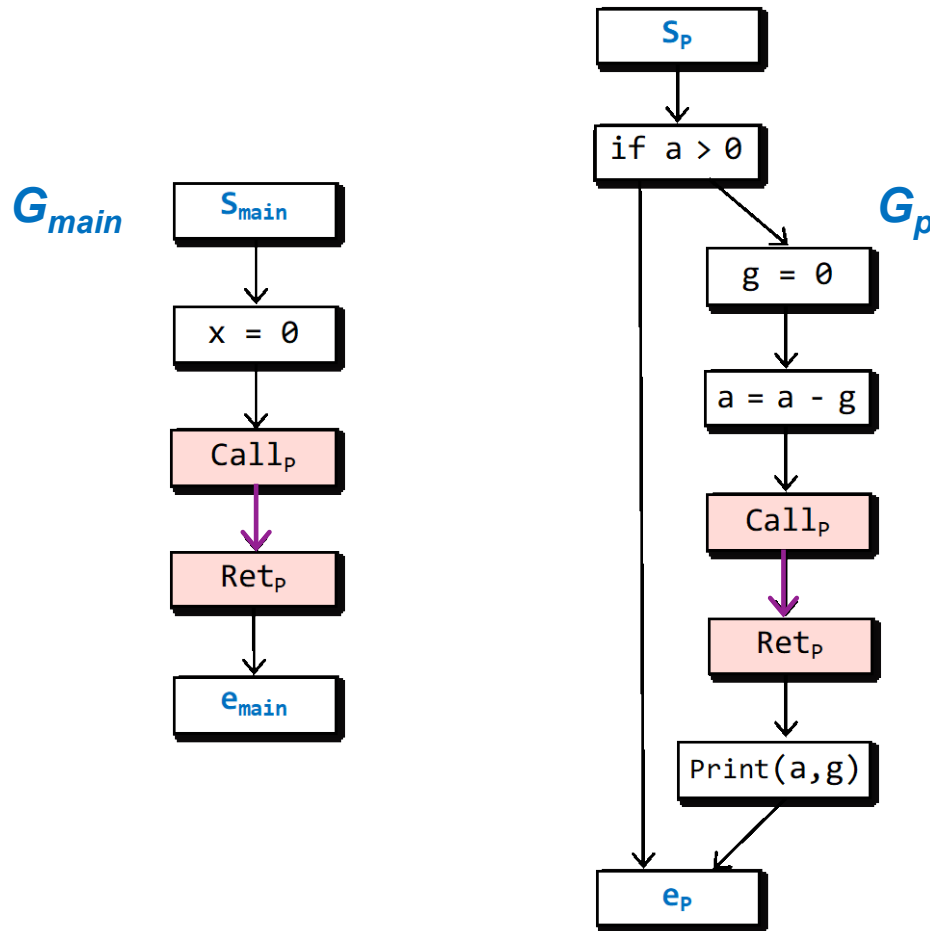- Print(a,g)
- e<sub>P</sub>

# Supergraph

*G\* has three edges for each procedure call:*

- *An intraprocedural* *call-to-return-site edge* *from* *Call$_p$* *to* *Ret$_p$*

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a – g;
    P(a);
    Print(a,g);
  }
}
```

$G_{main}$

$S_{main}$ → x = 0 → Call$_P$ → Ret$_P$ → $e_{main}$

$G_p$

S$_P$ → if a > 0 → g = 0 → a = a - g → Call$_P$ → Ret$_P$ → Print(a,g) → e$_P$

# Supergraph

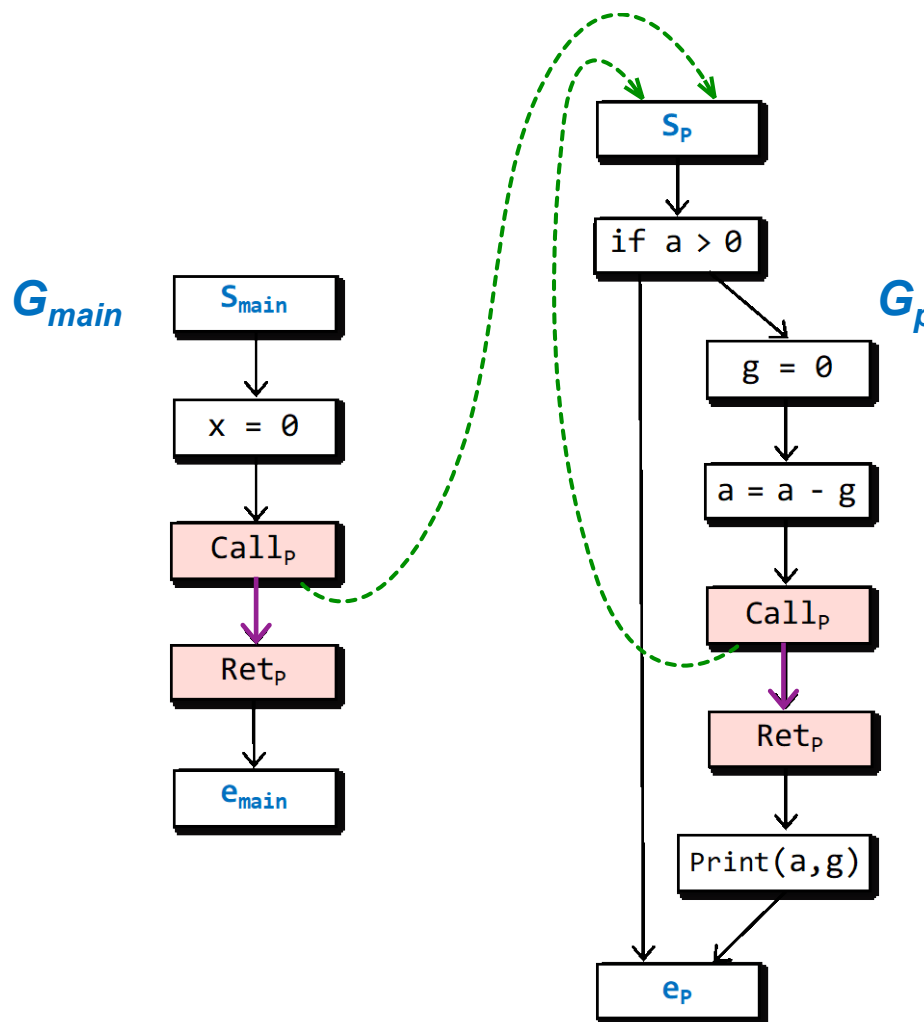*G\* has three edges for each procedure call:*

- *An intraprocedural call-to-return-site edge from $Call_p$ to $Ret_p$*
- *An interprocedural call-to-start edge from $Call_p$ to $s_p$ of the called procedure*

```
int g;
main(){
    int x;
    x = 0;
    P(x);
}
P(int a){
    if(a > 0){
        g = 0;
        a = a – g;
        P(a);
        Print(a,g);
    }
}
```
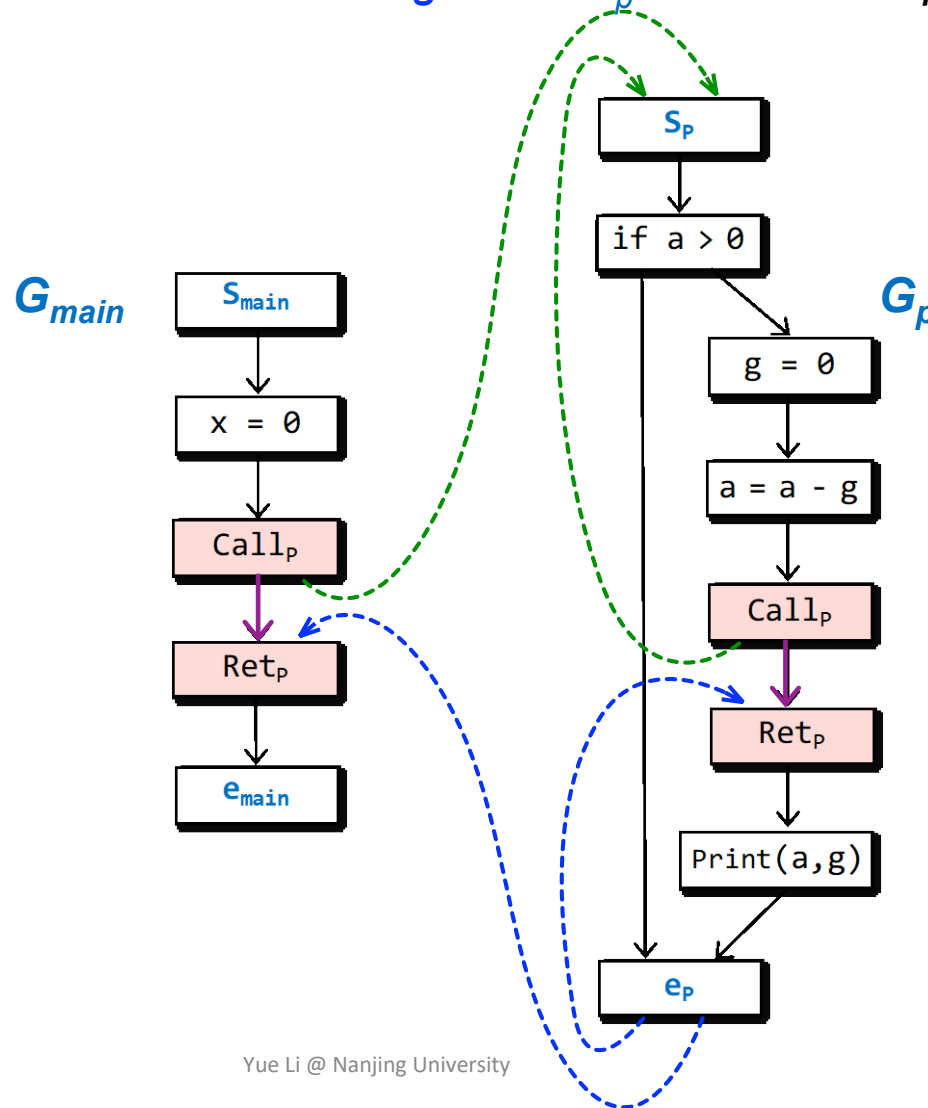
**$G_{main}$**

| $S_{main}$ |
| x = 0 |
| $Call_P$ |
| $Ret_P$ |
| $e_{main}$ |

**$G_p$**

| $S_P$ |
| if a > 0 |
| g = 0 |
| a = a - g |
| $Call_P$ |
| $Ret_P$ |
| Print(a,g) |
| $e_P$ |

# Supergraph

*G\* has three edges for each procedure call:*

- *An intraprocedural call-to-return-site edge from $Call_p$ to $Ret_p$*
- *An interprocedural call-to-start edge from $Call_p$ to $s_p$ of the called procedure*
- *An interprocedural exit-to-return-site edge from $e_p$ of the called procedure to $Ret_p$*

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a - g;
    P(a);
    Print(a,g);
  }
}
```



$G_{main}$

$G_p$

# Design Flow Functions

"Possibly-uninitialized variables": for each node n ∈ N*, determine the set of variables that may be uninitialized before execution reaches n.
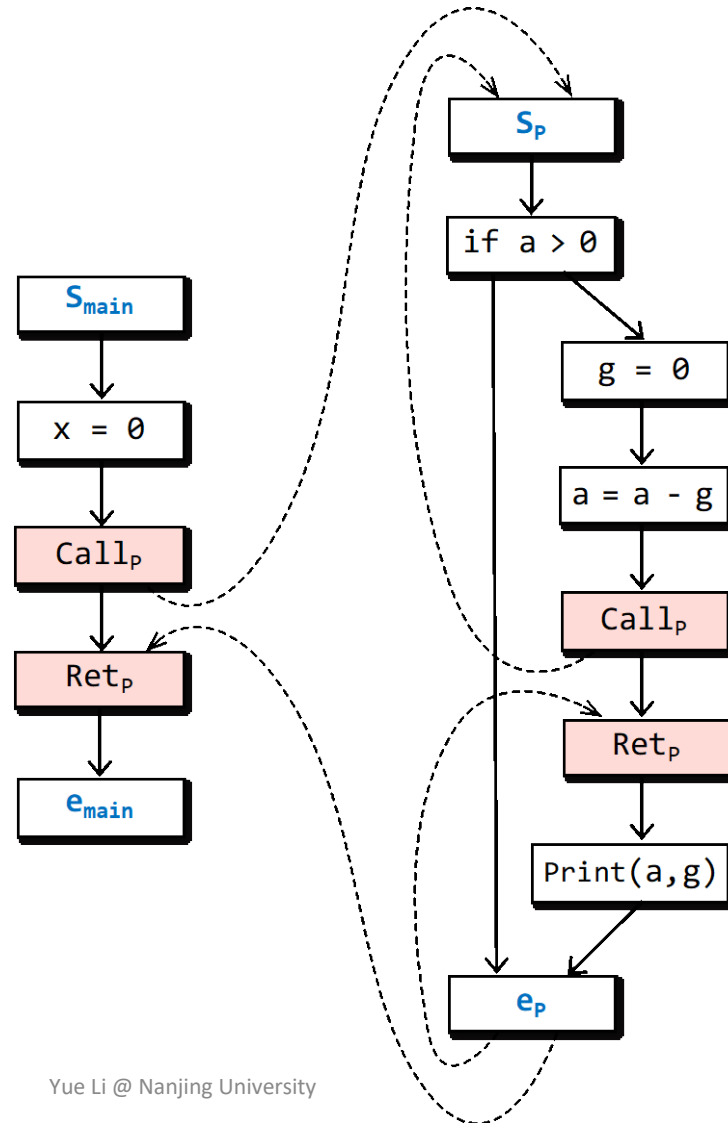


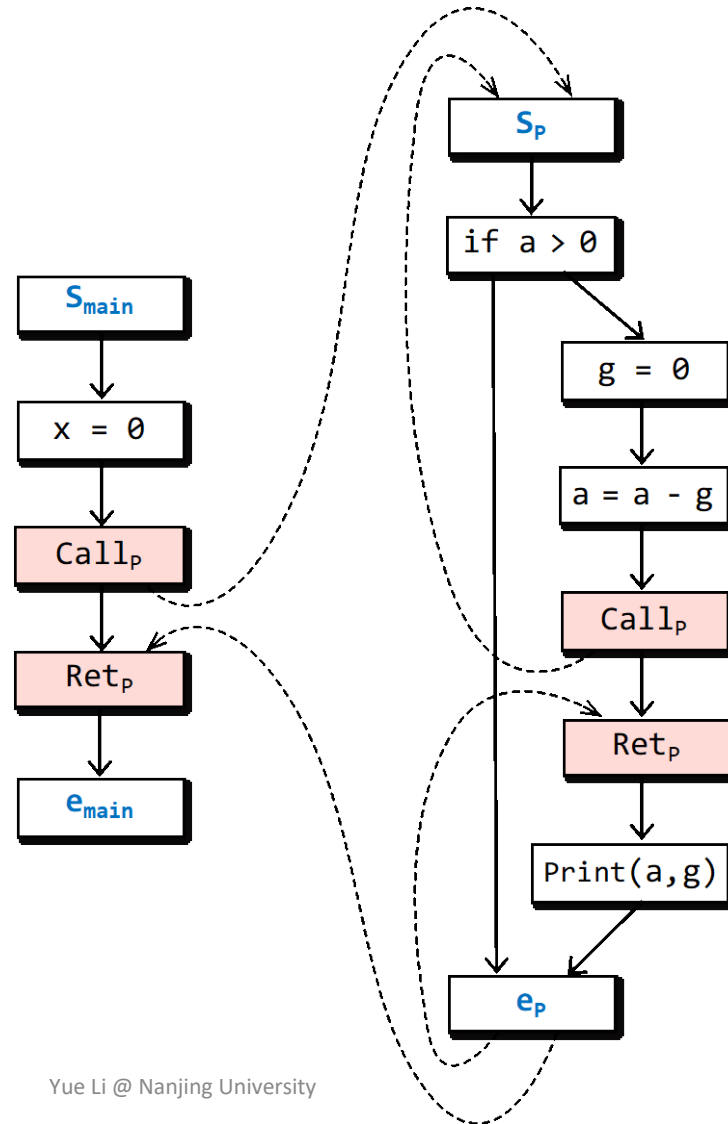```
int g;
main(){
   int x;
   x = 0;
   P(x);
}
P(int a){
   if(a > 0){
      g = 0;
      a = a – g;
      P(a);
      Print(a,g);
   }
}
```

# Design Flow Functions

"Possibly-uninitialized variables": for each node n ∈ N*, determine the set of variables that may be uninitialized before execution reaches n.

$$\lambda\, e_{param} \cdot e_{body}$$

# Design Flow Functions

"Possibly-uninitialized variables": for each node n ∈ N*, determine the set of variables that may be uninitialized before execution reaches n.

$$\lambda\, e_{param} \cdot e_{body}$$

*e.g.,* $\lambda\, x \cdot x+1$

# Design Flow Functions

"Possibly-uninitialized variables": for each node n ∈ N*, determine the set of variables that may be uninitialized before execution reaches n.
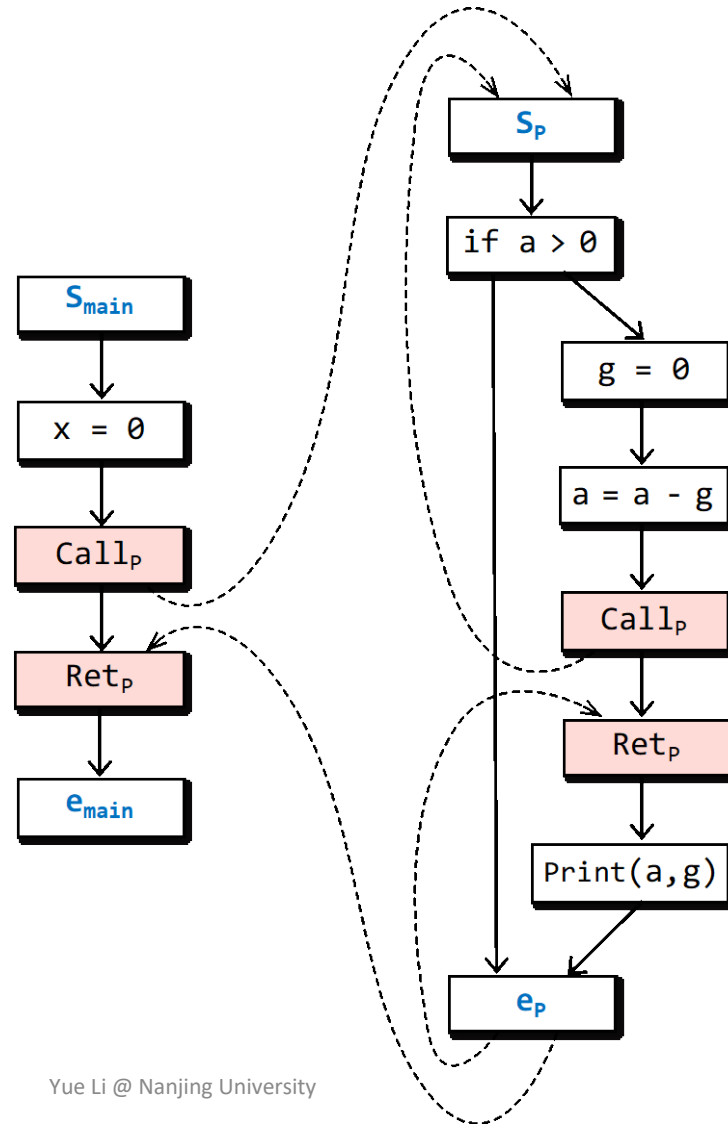
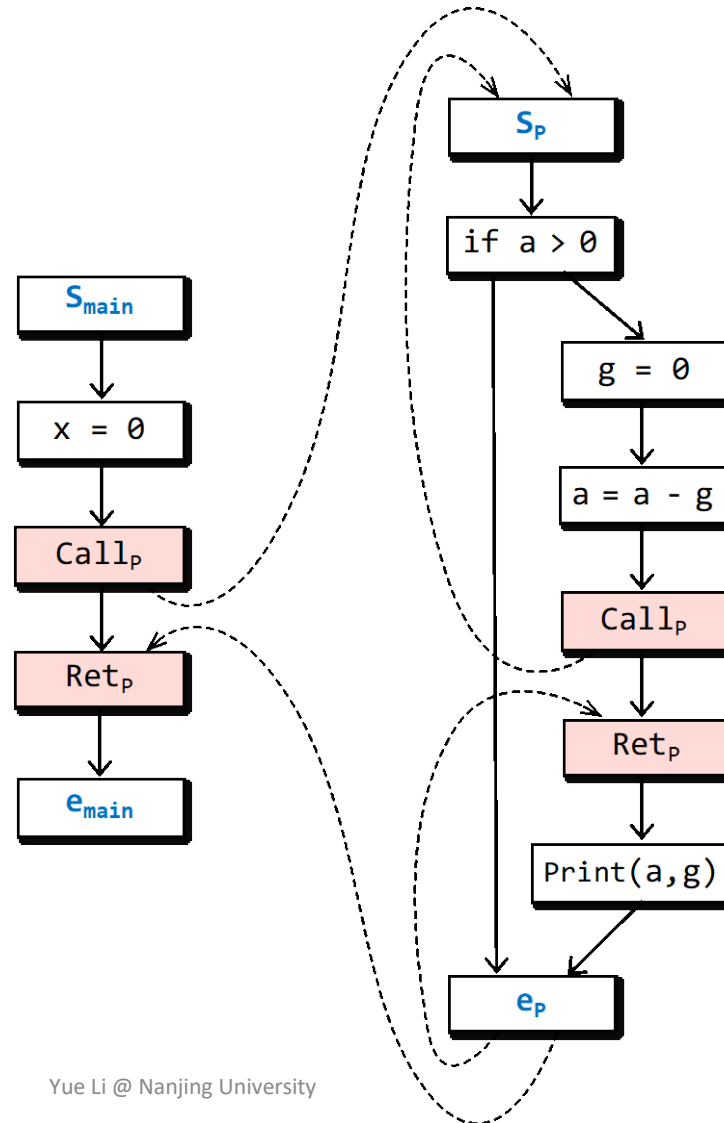$$\lambda\, e_{param} \cdot e_{body}$$

*e.g.,* $\lambda\, \text{x.x+1}$

$(\lambda\, \text{x.x+1})\text{3}$

# Design Flow Functions

"Possibly-uninitialized variables": for each node n ∈ N*, determine the set of variables that may be uninitialized before execution reaches n.
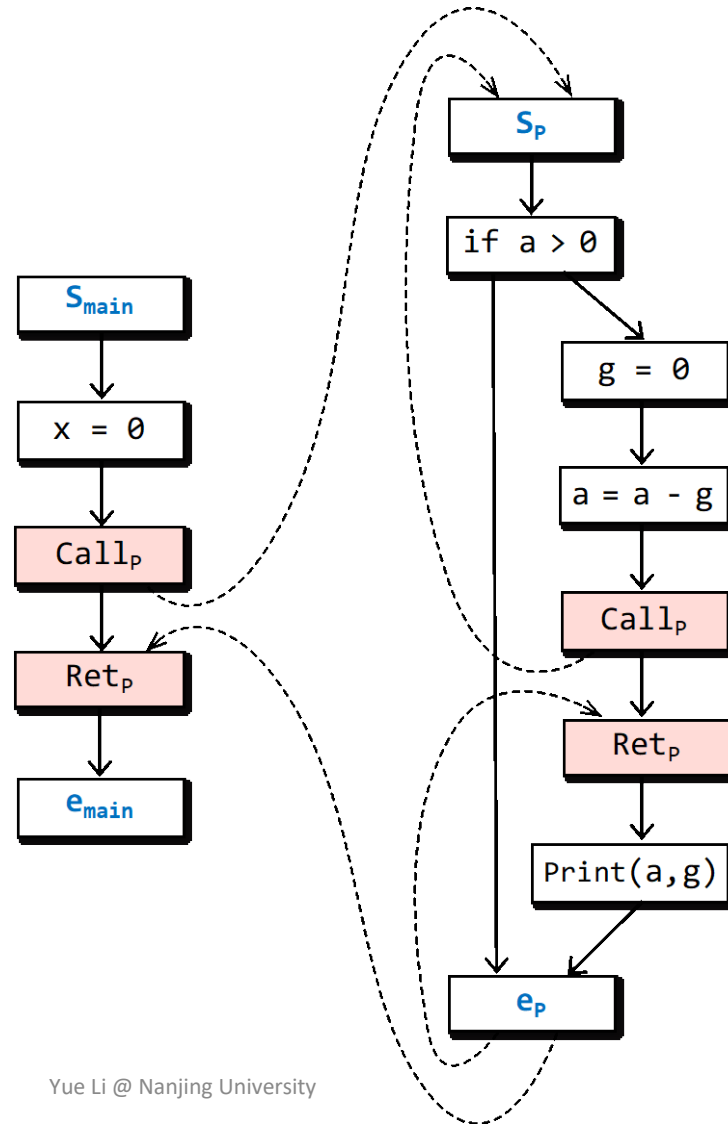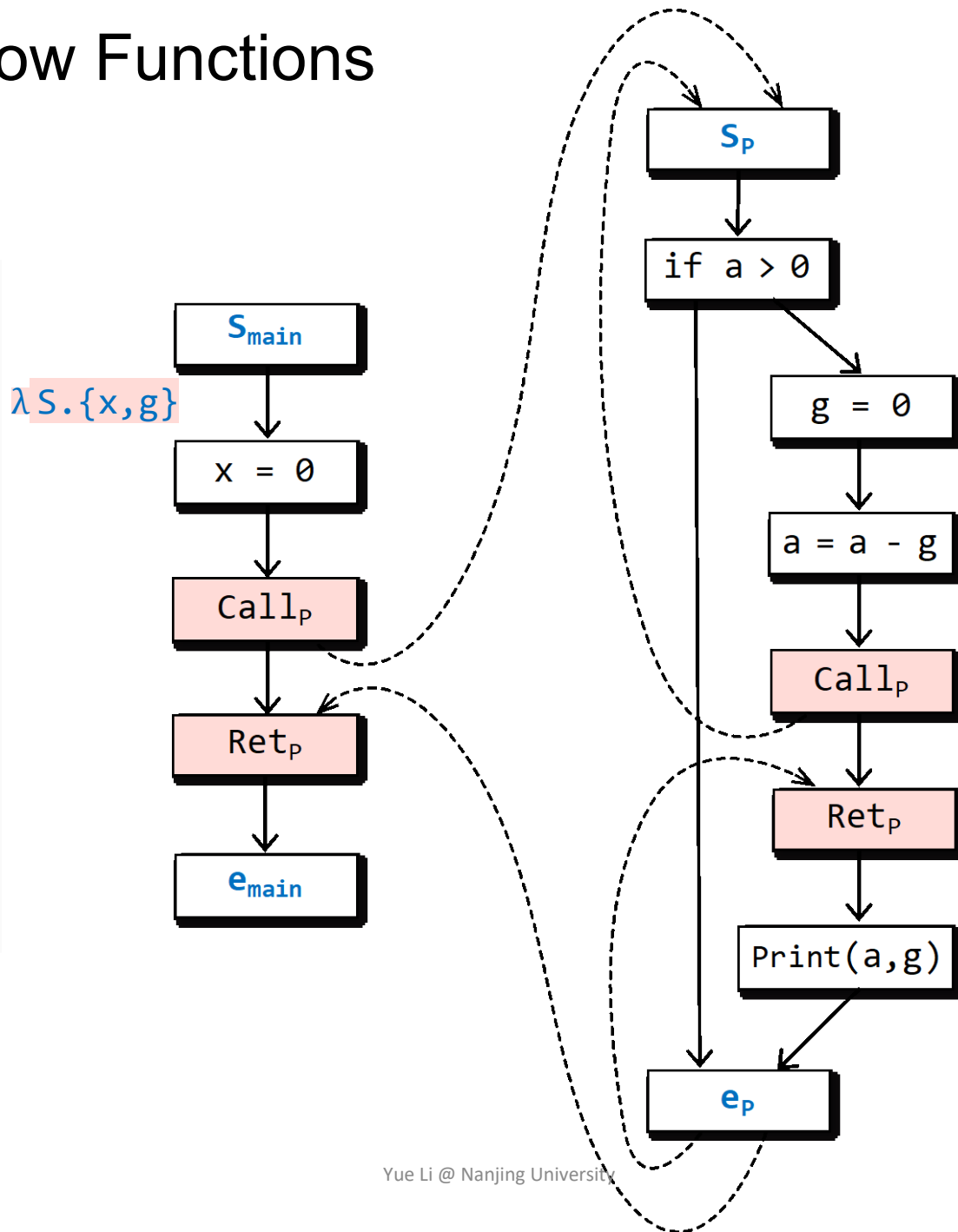
$$\lambda\, e_{param} \cdot e_{body}$$

$e.g.,$  $\lambda$ x.x+1

$(\lambda$ x.x+1$)$3
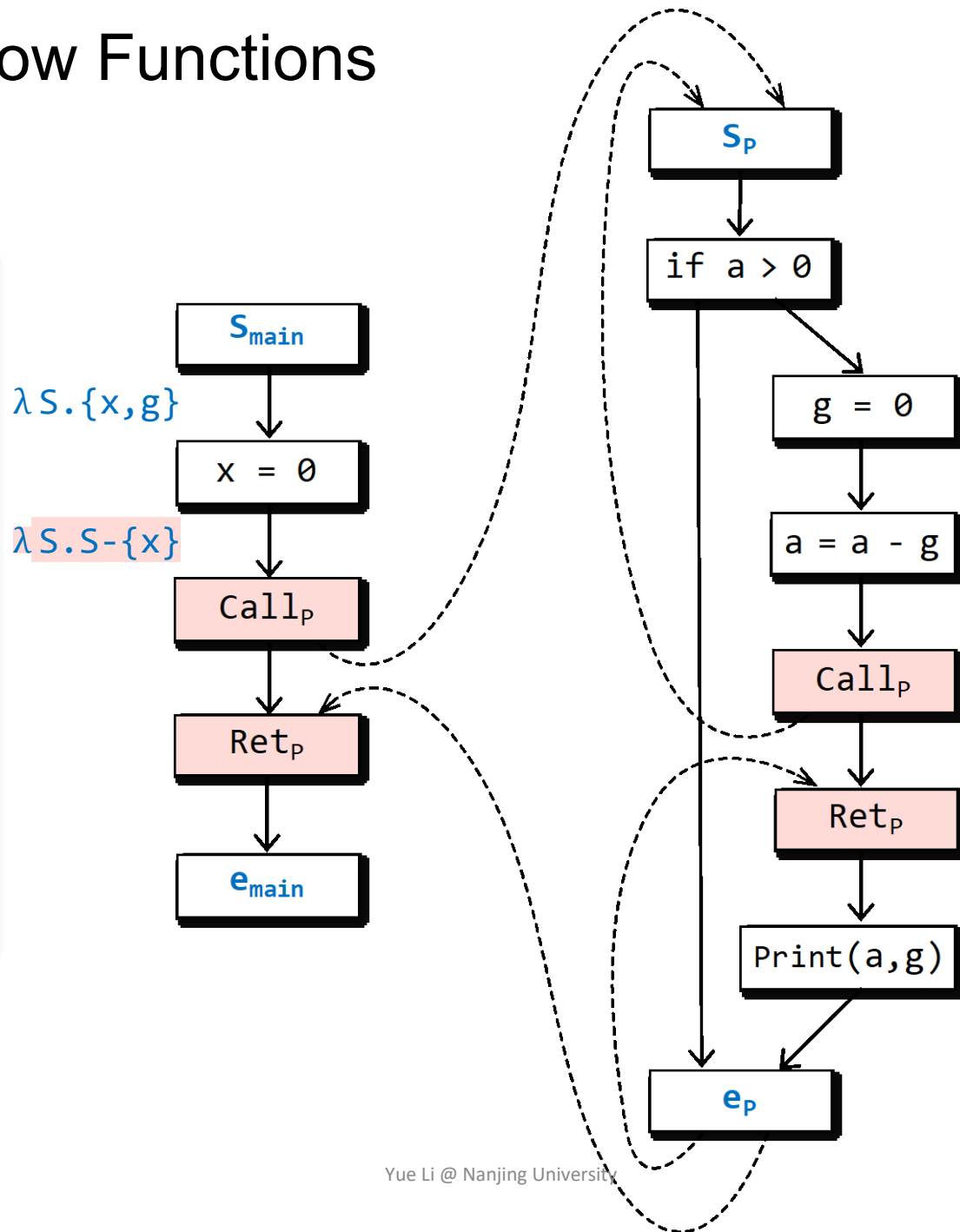
$\Rightarrow$ 3+1

$\Rightarrow$ 4

# Design Flow Functions



```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a – g;
    P(a);
    Print(a,g);
  }
}
```
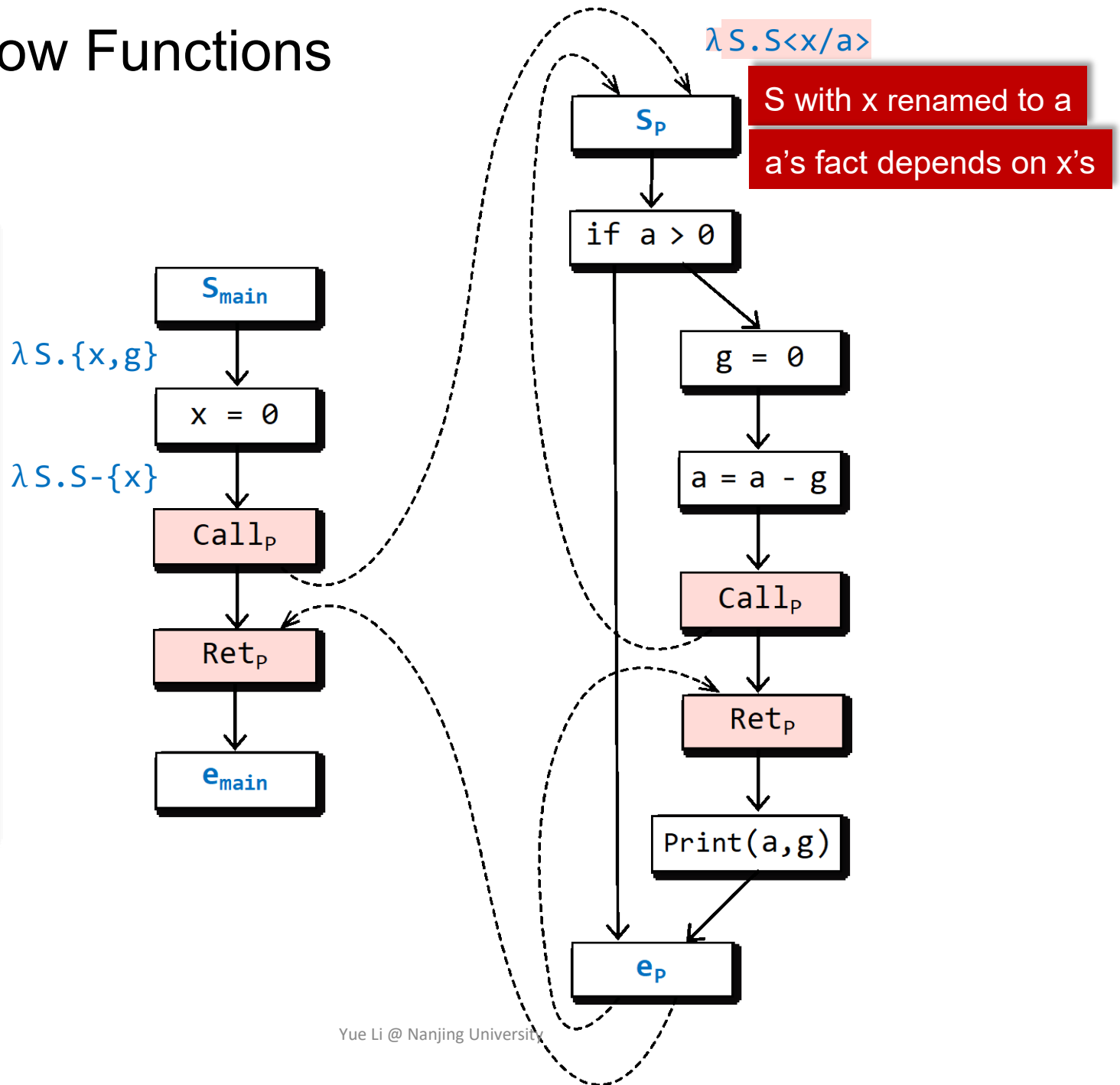
$\lambda S.\{x,g\}$

$S_{main}$

x = 0

$Call_P$

$Ret_P$

$e_{main}$

$S_P$

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

# Design Flow Functions



```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a - g;
    P(a);
    Print(a,g);
  }
}
```

$\lambda S.\{x,g\}$

$\lambda S.S-\{x\}$

S_main

x = 0

Call_P

Ret_P

e_main

S_P

if a > 0

g = 0

a = a - g

Call_P

Ret_P

Print(a,g)

e_P

# Design Flow Functions

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a – g;
    P(a);
    Print(a,g);
  }
}
```

$\lambda S.\{x,g\}$

$\lambda S.S-\{x\}$

$\lambda S.S\langle x/a\rangle$

S with x renamed to a

a's fact depends on x's



$S_{main}$

x = 0

$Call_P$

$Ret_P$

$e_{main}$

$S_P$

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

# Design Flow Functions

```
int g;
main(){
    int x;
    x = 0;
    P(x);
}
P(int a){
    if(a > 0){
        g = 0;
        a = a – g;
        P(a);
        Print(a,g);
    }
}
```
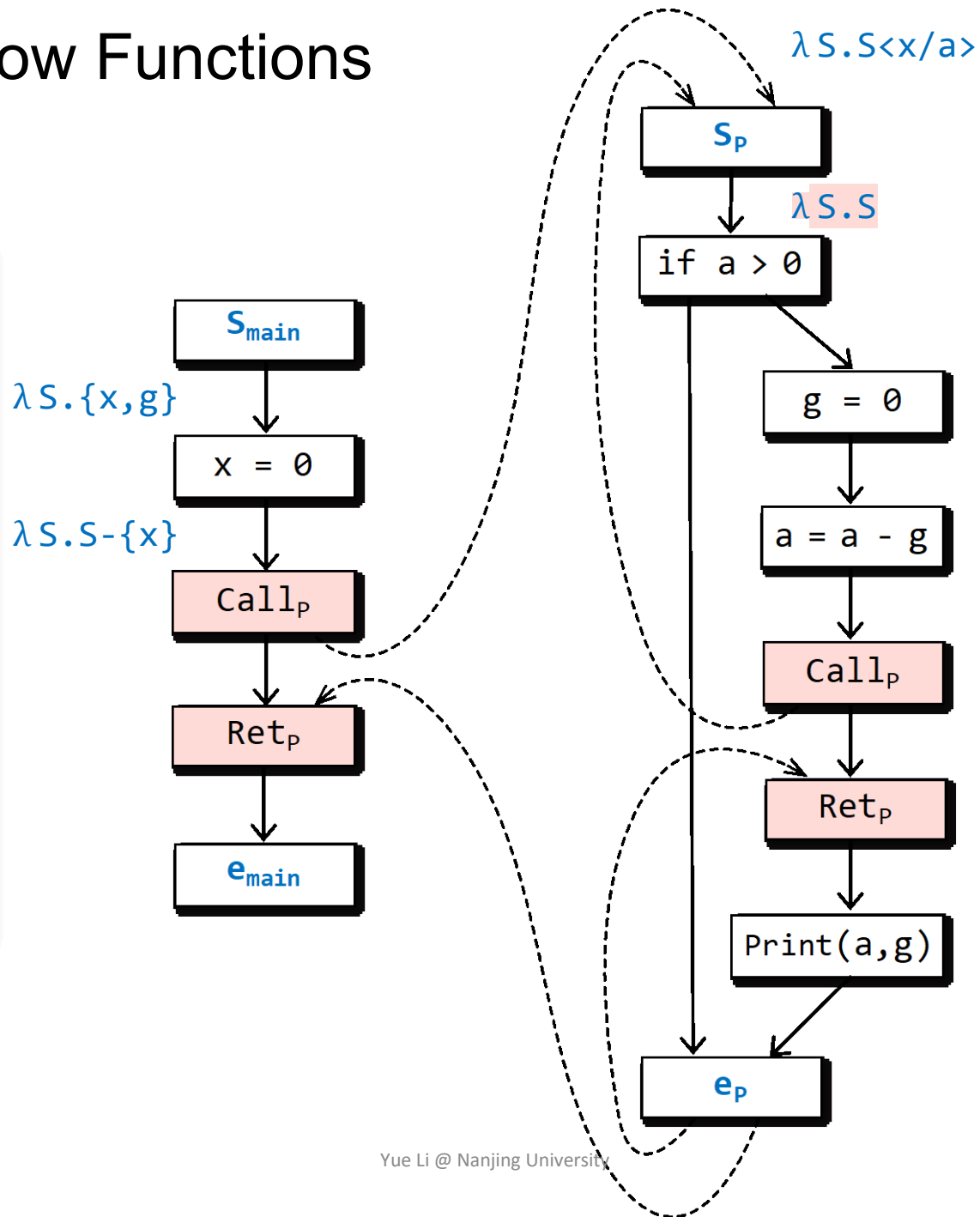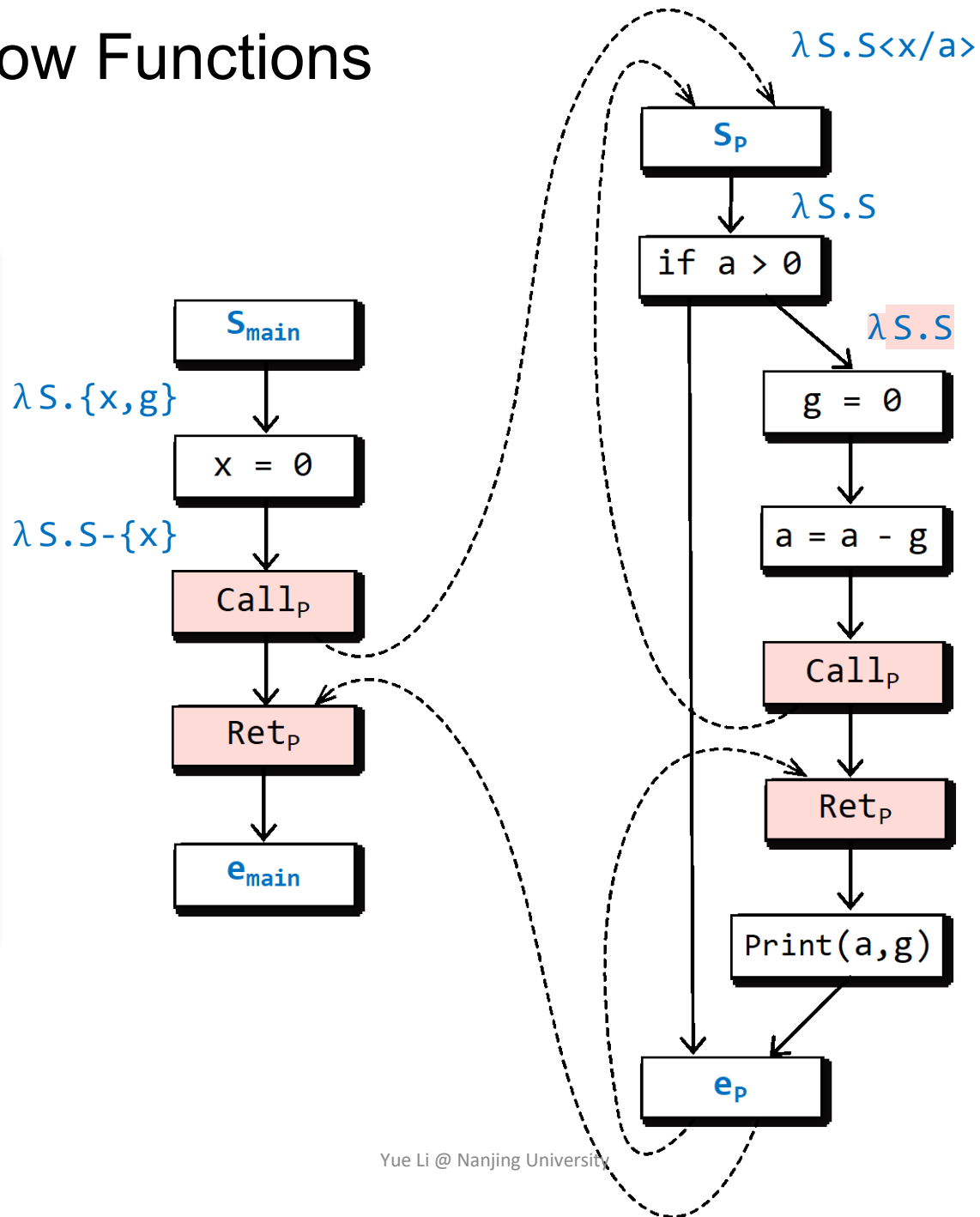
$\lambda S.S\langle x/a\rangle$

$S_P$

$\lambda S.S$

`if a > 0`

$S_{main}$

$\lambda S.\{x,g\}$

`x = 0`

`g = 0`

$\lambda S.S-\{x\}$

`a = a - g`

$Call_P$

$Call_P$

$Ret_P$

$Ret_P$

$e_{main}$

`Print(a,g)`

$e_P$

# Design Flow Functions

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a – g;
    P(a);
    Print(a,g);
  }
}
```
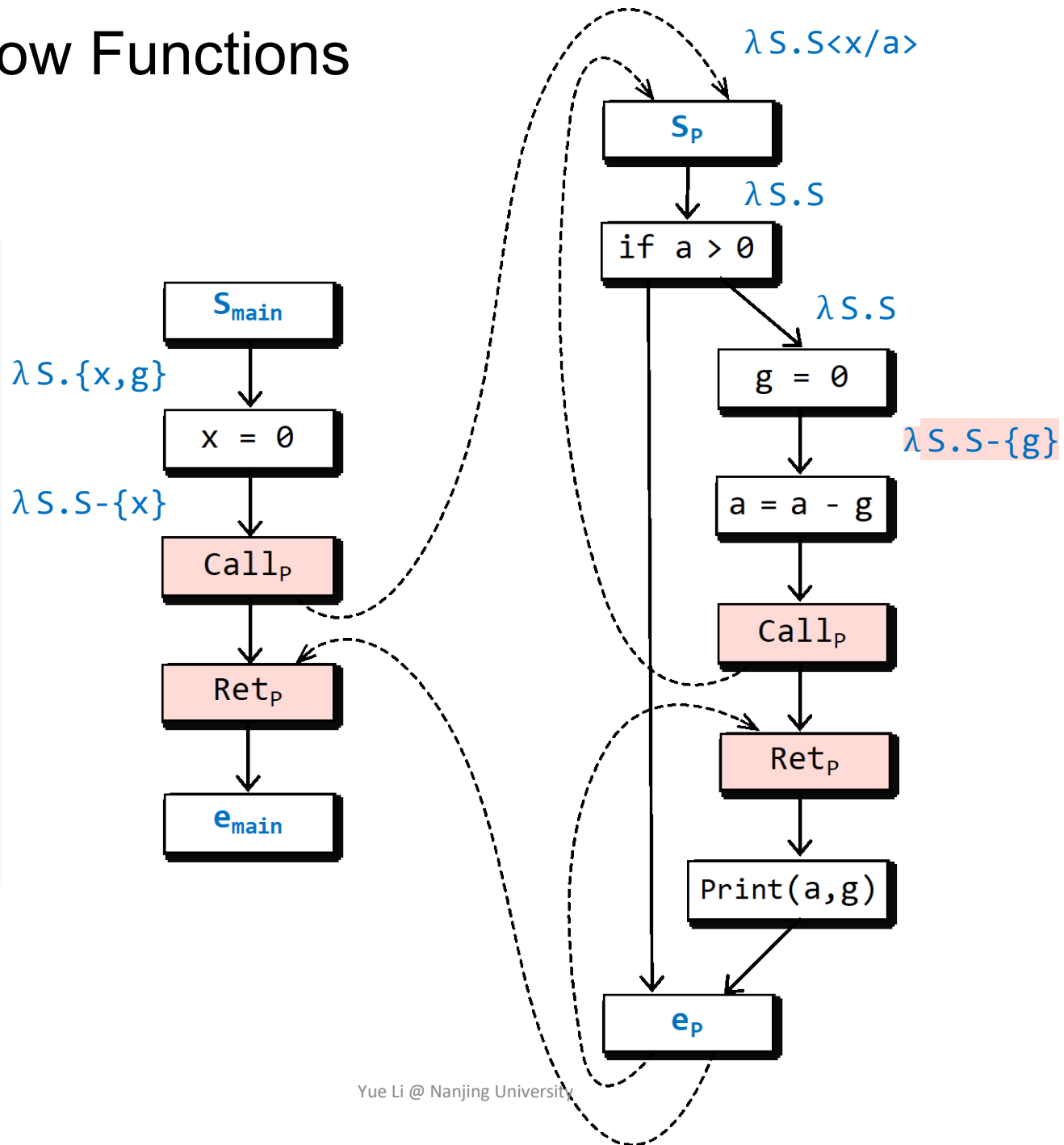
$\lambda S.S\langle x/a\rangle$

$S_P$

$\lambda S.S$

if a > 0

$\lambda S.S$

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

$S_{main}$

$\lambda S.\{x,g\}$

x = 0

$\lambda S.S-\{x\}$

$Call_P$

$Ret_P$

$e_{main}$

# Design Flow Functions

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a – g;
    P(a);
    Print(a,g);
  }
}
```
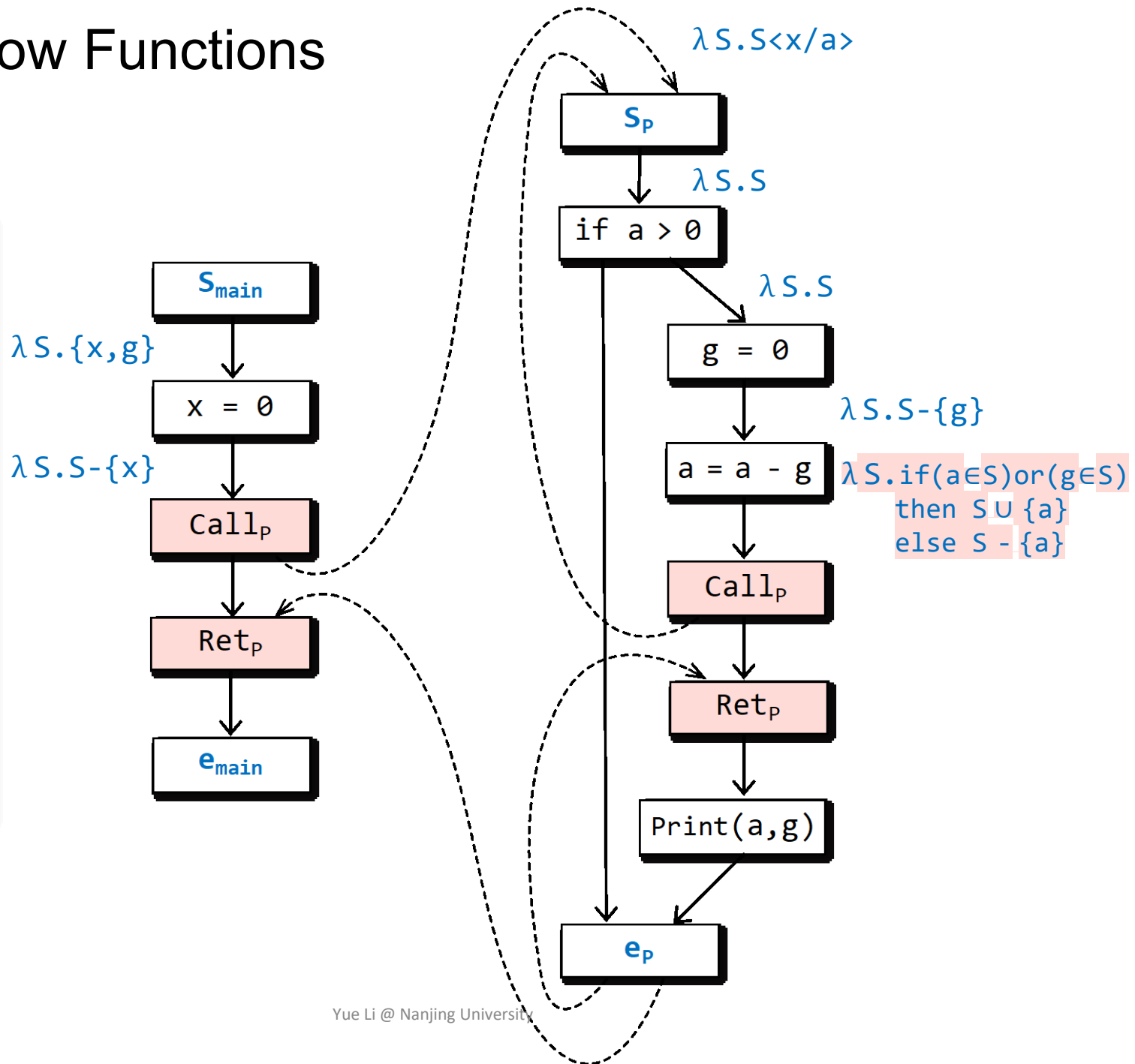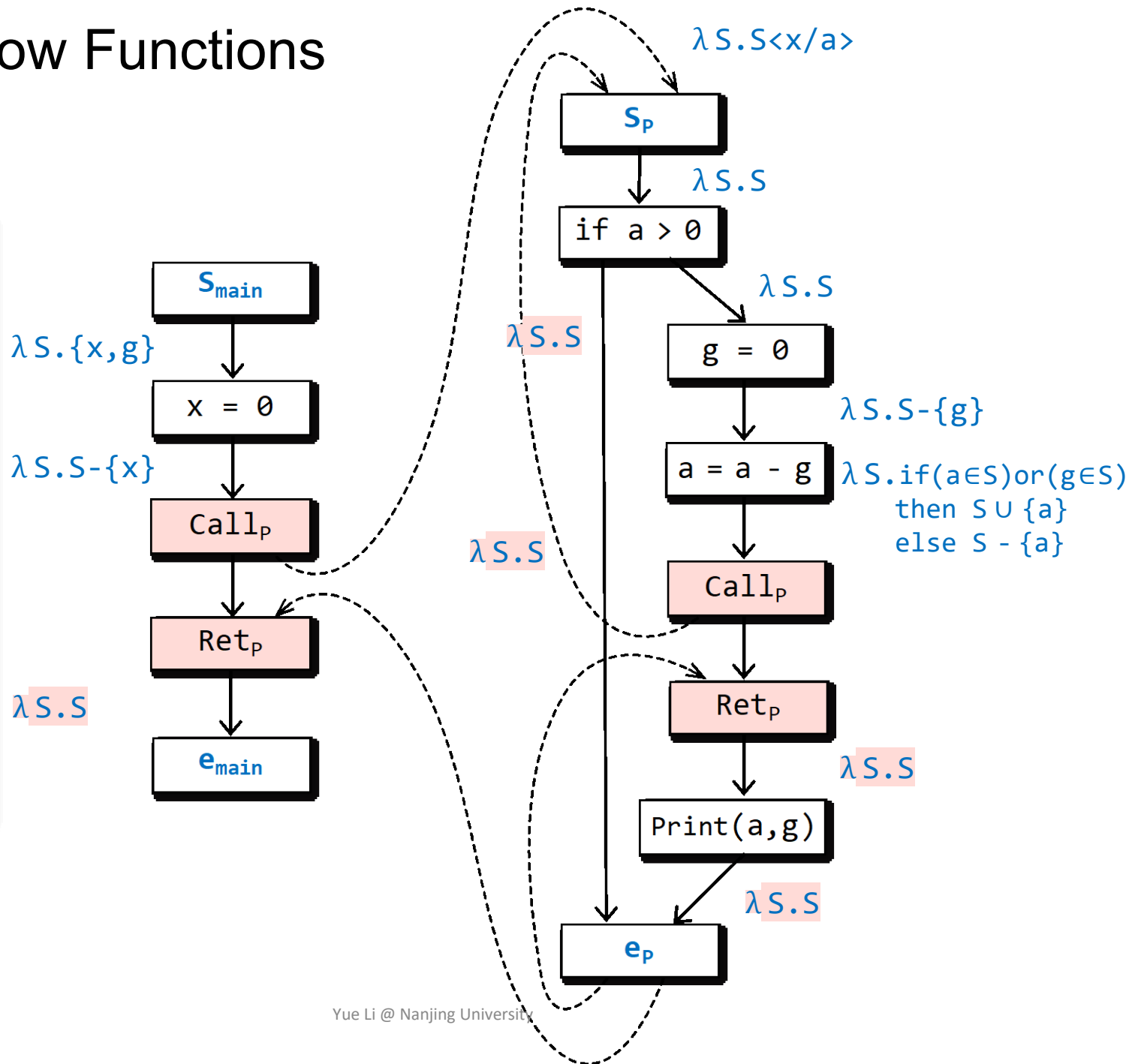
$\lambda S.\{x,g\}$

$\lambda S.S\text{-}\{x\}$

$\lambda S.S\text{<}x/a\text{>}$

$\lambda S.S$

$\lambda S.S$

$\lambda S.S\text{-}\{g\}$

$S_{main}$

x = 0

$Call_P$

$Ret_P$

$e_{main}$

$S_P$

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

# Design Flow Functions

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a – g;
    P(a);
    Print(a,g);
  }
}
```
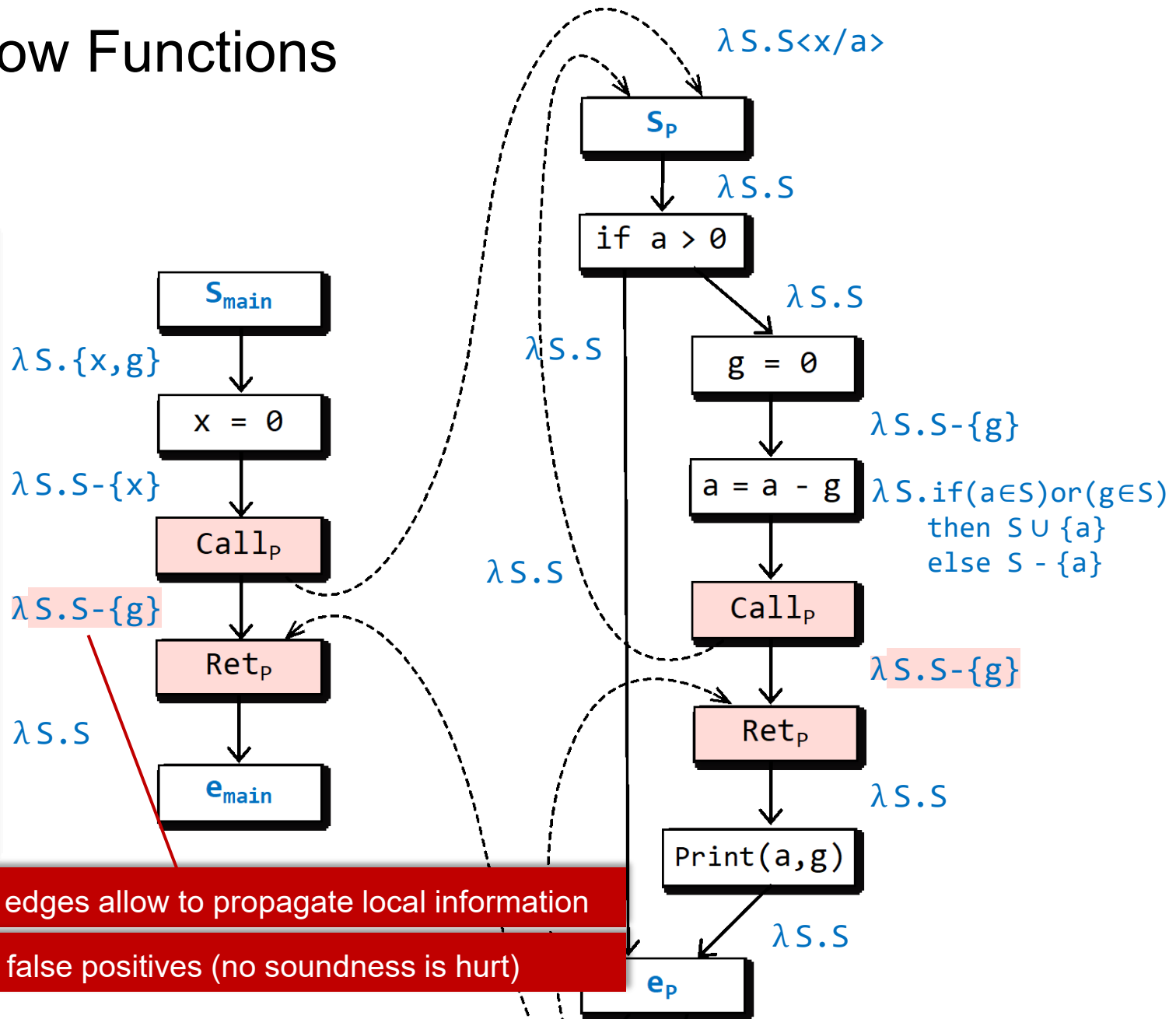
$\lambda S.S<x/a>$

$S_P$

$\lambda S.S$

`if a > 0`

$\lambda S.S$

`g = 0`

$\lambda S.S-\{g\}$

`a = a - g`

$\lambda S.if(a\in S)or(g\in S)$
  $then\ S\cup\{a\}$
  $else\ S-\{a\}$

$Call_P$

$Ret_P$

`Print(a,g)`

$e_P$

$S_{main}$

$\lambda S.\{x,g\}$

`x = 0`

$\lambda S.S-\{x\}$

$Call_P$

$Ret_P$

$e_{main}$

# Design Flow Functions

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a - g;
    P(a);
    Print(a,g);
  }
}
```
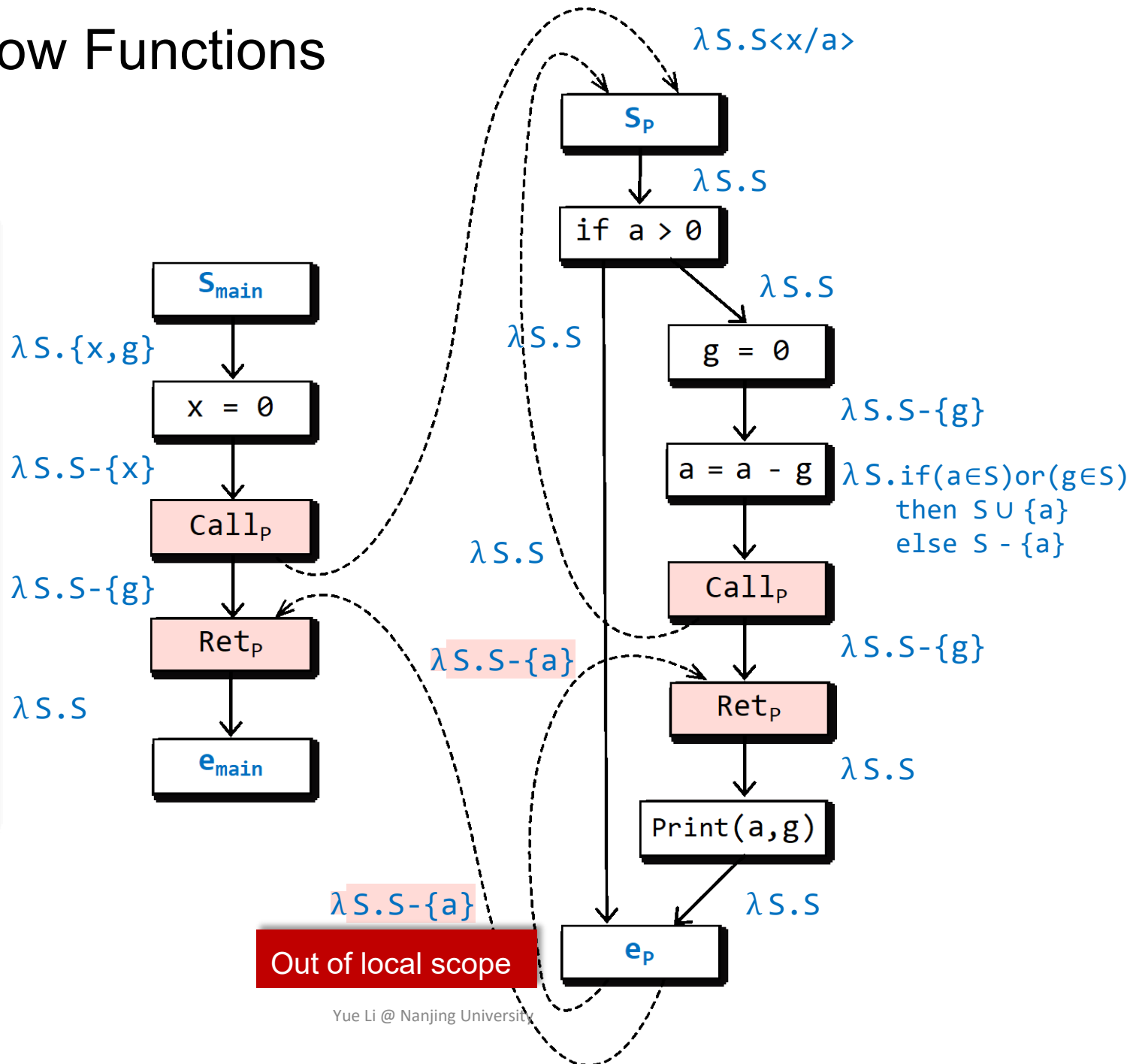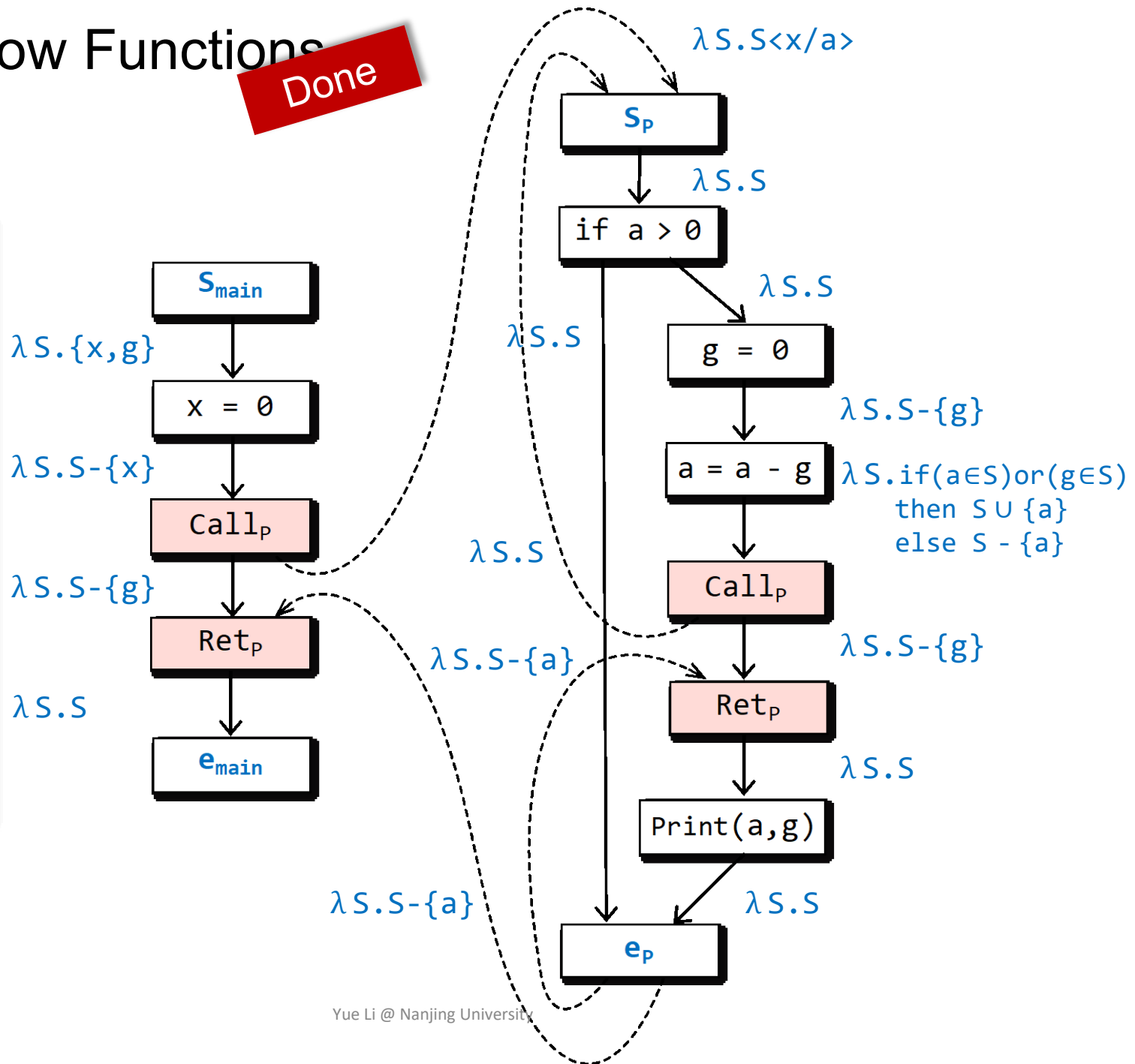
$\lambda S.S<x/a>$

$S_P$

$\lambda S.S$

if a > 0

$\lambda S.S$

$\lambda S.S$

g = 0

$\lambda S.S-\{g\}$

a = a - g

$\lambda S.\text{if}(a \in S)\text{or}(g \in S)$
$\quad\text{then } S \cup \{a\}$
$\quad\text{else } S - \{a\}$

$Call_P$

$Ret_P$

$\lambda S.S$

Print(a,g)

$\lambda S.S$

$e_P$

$\lambda S.S$

$S_{main}$

$\lambda S.\{x,g\}$

x = 0

$\lambda S.S-\{x\}$

$Call_P$

$Ret_P$

$\lambda S.S$

$e_{main}$

$\lambda S.S$

# Design Flow Functions

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a - g;
    P(a);
    Print(a,g);
  }
}
```
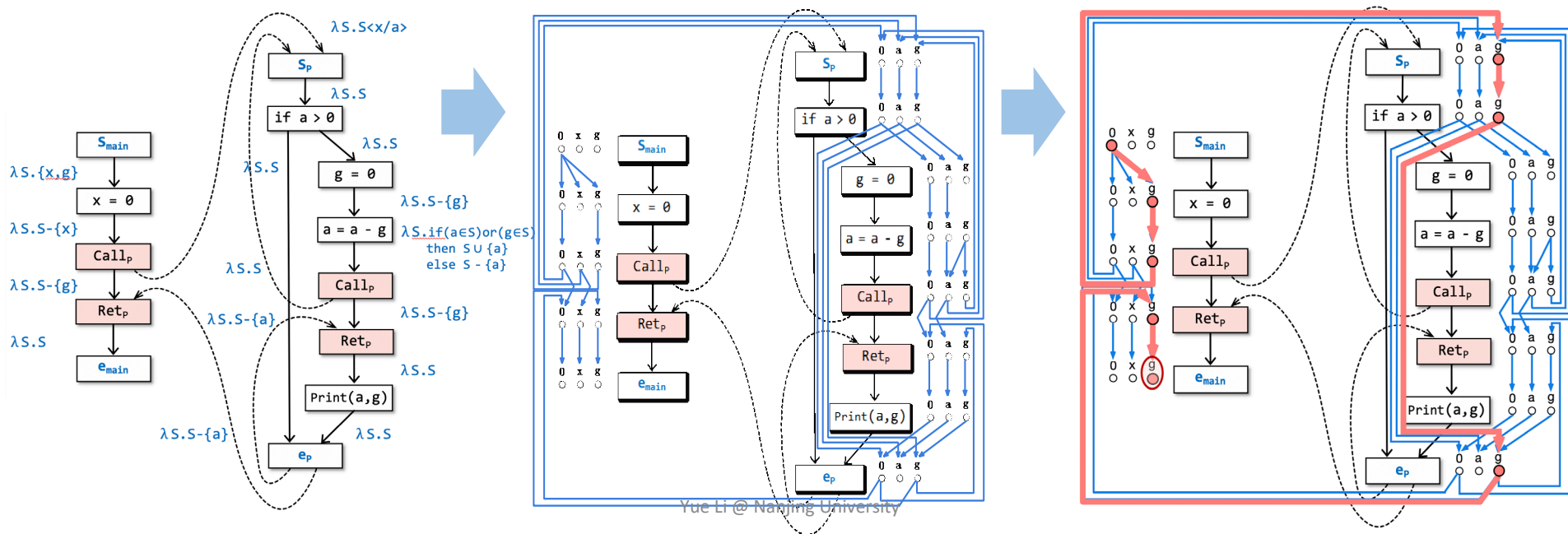
$\lambda S.S\langle x/a\rangle$

**S**$_P$

$\lambda S.S$

if a > 0

$\lambda S.S$

g = 0

$\lambda S.S\text{-}\{g\}$

a = a - g

$\lambda S.\text{if}(a\in S)\text{or}(g\in S)$
    then $S\cup\{a\}$
    else $S - \{a\}$

**Call**$_P$

$\lambda S.S\text{-}\{g\}$

**Ret**$_P$

$\lambda S.S$

Print(a,g)

$\lambda S.S$

**e**$_P$

**S**$_{main}$

$\lambda S.\{x,g\}$

x = 0

$\lambda S.S\text{-}\{x\}$

**Call**$_P$

$\lambda S.S\text{-}\{g\}$

**Ret**$_P$

$\lambda S.S$

**e**$_{main}$

$\lambda S.S$

$\lambda S.S$

$\lambda S.S$

"call-to-return-site" edges allow to propagate local information

S-{g} helps reduce false positives (no soundness is hurt)

# Design Flow Functions

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a – g;
    P(a);
    Print(a,g);
  }
}
```

$\lambda S.S\langle x/a\rangle$

$S_P$

$\lambda S.S$

`if a > 0`

$\lambda S.S$

`g = 0`

$\lambda S.S.S-\{g\}$

`a = a - g`

$\lambda S.\text{if}(a\in S)\text{or}(g\in S)$
$\quad\text{then } S\cup\{a\}$
$\quad\text{else } S-\{a\}$

$Call_P$

$\lambda S.S-\{g\}$

$Ret_P$

$\lambda S.S$

`Print(a,g)`

$\lambda S.S$

$e_P$

$S_{main}$

$\lambda S.\{x,g\}$

`x = 0`

$\lambda S.S-\{x\}$

$Call_P$

$\lambda S.S-\{g\}$

$Ret_P$

$\lambda S.S$

$e_{main}$

$\lambda S.S$

$\lambda S.S$

$\lambda S.S-\{a\}$

$\lambda S.S-\{a\}$

Out of local scope

# Design Flow Functions

**Done**

```
int g;
main(){
  int x;
  x = 0;
  P(x);
}
P(int a){
  if(a > 0){
    g = 0;
    a = a – g;
    P(a);
    Print(a,g);
  }
}
```

$\lambda S.S\langle x/a\rangle$

**S$_P$**

$\lambda S.S$

if a > 0

$\lambda S.S$

**S$_{main}$**

$\lambda S.\{x,g\}$

x = 0

$\lambda S.S-\{x\}$

$\lambda S.S$

g = 0

$\lambda S.S-\{g\}$

Call$_P$

a = a - g

$\lambda S.if(a\in S)or(g\in S)$
then $S\cup\{a\}$
else $S - \{a\}$

$\lambda S.S-\{g\}$

Ret$_P$

$\lambda S.S$

Call$_P$

$\lambda S.S-\{a\}$

$\lambda S.S-\{g\}$

$\lambda S.S$

**e$_{main}$**

Ret$_P$

$\lambda S.S$

Print(a,g)

$\lambda S.S$

$\lambda S.S-\{a\}$

**e$_P$**

$\lambda S.S$

# Overview of IFDS

Given a program P, and a dataflow-analysis problem Q

- Build a supergraph G* for P and
  define flow functions for edges in G* based on Q

- Build exploded supergraph G# for P by transforming
  flow functions to representation relations (graphs)

- Q can be solved as graph reachability problems (find out MRP solutions)
  via applying Tabulation algorithm on G#

# Build Exploded Supergraph

- Build exploded supergraph $G^\#$ for a program by transforming flow functions to representation relations (graphs)

- Each flow function can be represented as a graph with $2(D+1)$ nodes (at most $(D+1)^2$ edges), where D is a finite set of dataflow facts

# Build Exploded Supergraph

- Build exploded supergraph $G^{\#}$ for a program by transforming flow functions to representation relations (graphs)

- Each flow function can be represented as a graph with $2(D+1)$ nodes (at most $(D+1)^2$ edges), where D is a finite set of dataflow facts

The representation relation of flow function f, $R_f \subseteq (D \cup 0) \times (D \cup 0)$ is a binary relation (or graph) defined as follows:

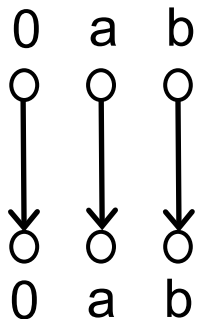$\quad R_f = \{ (0,0) \}$                     Edge: $0 \rightarrow 0$

$\quad\quad \cup \{ (0,y) \mid y \in f(\emptyset) \}$          Edge: $0 \rightarrow d_1$

$\quad\quad \cup \{ (x,y) \mid y \notin f(\emptyset) \text{ and } y \in f(\{x\}) \}$ Edge: $d_1 \rightarrow d_2$

# Build Exploded Supergraph

- Build exploded supergraph $G^\#$ for a program by transforming flow functions to representation relations (graphs)

- Each flow function can be represented as a graph with $2(D+1)$ nodes (at most $(D+1)^2$ edges), where D is a finite set of dataflow facts

The representation relation of flow function f, $R_f \subseteq (D \cup 0) \times (D \cup 0)$ is a binary relation (or graph) defined as follows:

$R_f = \{ (0,0) \}$        Edge: $0 \rightarrow 0$

$\cup \{ (0,y) \mid y \in f(\emptyset) \}$        Edge: $0 \rightarrow d_1$

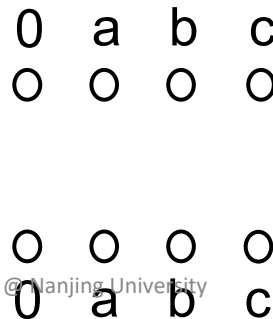$\cup \{ (x,y) \mid y \notin f(\emptyset) \text{ and } y \in f(\{x\}) \}$ Edge: $d_1 \rightarrow d_2$
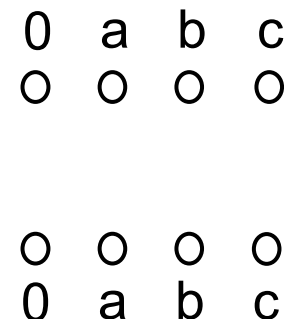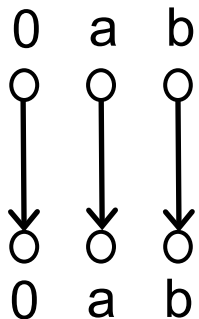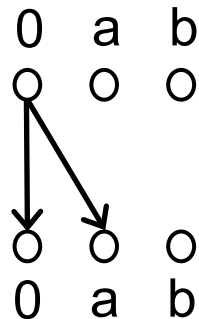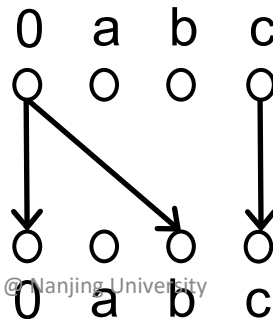
$\lambda S.S$      $\lambda S.\{a\}$      $\lambda S.(S-\{a\})\cup\{b\}$      $\lambda S.\text{if } a \in S$ then $S \cup \{b\}$ else $S - \{b\}$

| 0 | a | b |
|---|---|---|
| O | O | O |

| 0 | a | b |
|---|---|---|
| O | O | O |

| 0 | a | b | c |
|---|---|---|---|
| O | O | O | O |

| 0 | a | b | c |
|---|---|---|---|
| O | O | O | O |

| O | O | O |
|---|---|---|
| 0 | a | b |

| O | O | O |
|---|---|---|
| 0 | a | b |

| O | O | O | O |
|---|---|---|---|
| 0 | a | b | c |

| O | O | O | O |
|---|---|---|---|
| 0 | a | b | c |

# Build Exploded Supergraph

- Build exploded supergraph $G^\#$ for a program by transforming flow functions to <span style="color:red">representation relations</span> (graphs)

- Each flow function can be represented as a graph with $2(D+1)$ nodes (at most $(D+1)^2$ edges), where D is a finite set of dataflow facts

The <span style="color:red">representation relation</span> of <span style="color:blue">flow function f</span>, $R_f \subseteq (D \cup 0) \times (D \cup 0)$ is a binary relation (or graph) defined as follows:

$R_f = \{ (0,0) \}$                    Edge: $0 \rightarrow 0$

$\cup \{ (0,y) \mid y \in f(\emptyset) \}$            Edge: $0 \rightarrow d_1$

$\cup \{ (x,y) \mid y \notin f(\emptyset) \text{ and } y \in f(\{x\}) \}$ Edge: $d_1 \rightarrow d_2$

$\lambda S.S$                 $\lambda S.\{a\}$              $\lambda S.(S-\{a\})\cup\{b\}$              $\lambda S.\text{if } a \in S$
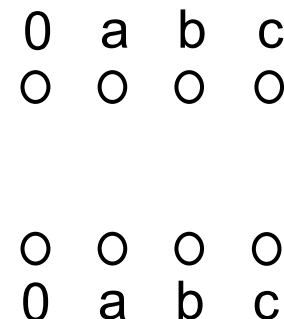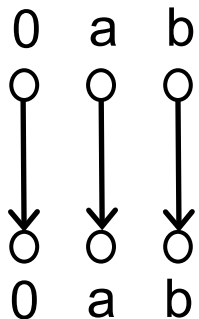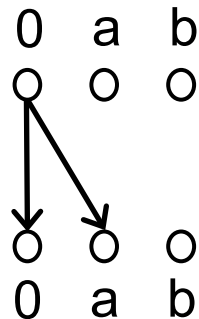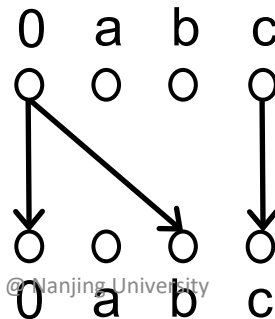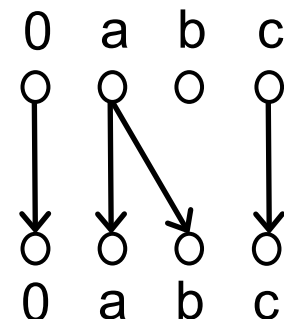                                                                                                      $\text{then } S \cup \{b\}$
                                                                                                      $\text{else } S - \{b\}$

# Build Exploded Supergraph

- Build exploded supergraph $G^\#$ for a program by transforming flow functions to representation relations (graphs)

- Each flow function can be represented as a graph with 2(D+1) nodes (at most $(D+1)^2$ edges), where D is a finite set of dataflow facts

The representation relation of flow function f, $R_f \subseteq (D \cup 0) \times (D \cup 0)$ is a binary relation (or graph) defined as follows:

$R_f$ = { (0,0) }        Edge: $0 \rightarrow 0$

     ∪ { (0,y) | y ∈ f(∅) }       Edge: $0 \rightarrow d_1$

     ∪ { (x,y) | y ∉ f(∅) and y ∈ f({x}) } Edge: $d_1 \rightarrow d_2$



λ S.S

λ S.{a}

λ S.(S-{a})∪{b}

λ S.if a ∈ S
    then S ∪{b}
    else S - {b}

# Build Exploded Supergraph

- Build exploded supergraph $G^\#$ for a program by transforming flow functions to representation relations (graphs)

- Each flow function can be represented as a graph with 2(D+1) nodes (at most $(D+1)^2$ edges), where D is a finite set of dataflow facts

---

The representation relation of flow function f, $R_f \subseteq (D \cup 0) \times (D \cup 0)$ is a binary relation (or graph) defined as follows:

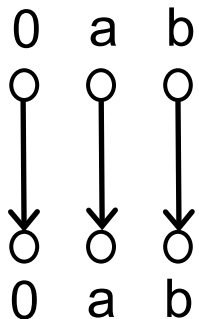$R_f$ = { (0,0) }                               Edge: $0 \rightarrow 0$

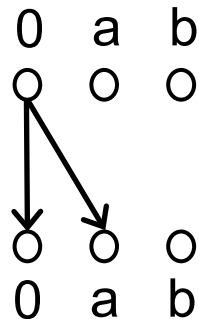$\cup$ { (0,y) | y $\in$ f($\emptyset$) }                    Edge: $0 \rightarrow d_1$

$\cup$ { (x,y) | y $\notin$ f($\emptyset$) and y $\in$ f({x}) } Edge: $d_1 \rightarrow d_2$
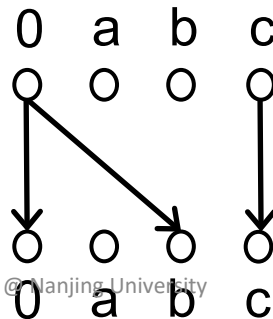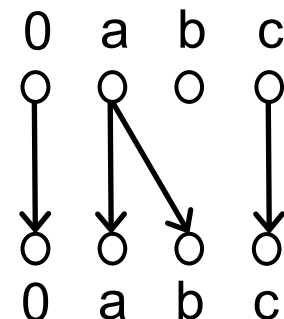
---

$\lambda S.S$   $\lambda S.\{a\}$   $\lambda S.(S-\{a\})\cup\{b\}$   $\lambda S.\text{if } a \in S$
$\text{then } S \cup \{b\}$
$\text{else } S - \{b\}$

# Build Exploded Supergraph

- Build exploded supergraph $G^{\#}$ for a program by transforming flow functions to <span style="color:red">representation relations</span> (graphs)

- Each flow function can be represented as a graph with $2(D+1)$ nodes (at most $(D+1)^2$ edges), where D is a finite set of dataflow facts

---

The <span style="color:red">representation relation</span> of <span style="color:blue">flow function f</span>, $R_f \subseteq (D \cup 0) \times (D \cup 0)$ is a binary relation (or graph) defined as follows:

$R_f$ = { (0,0) }             Edge: $0 \rightarrow 0$
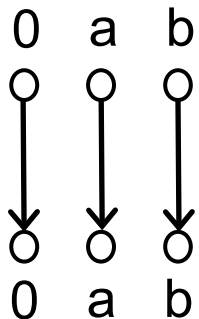
$\cup$ { (0,y) | y $\in$ f(∅) }        Edge: $0 \rightarrow d_1$
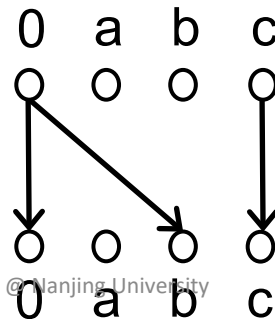
$\cup$ { (x,y) | y $\notin$ f(∅) and y $\in$ f({x}) } Edge: $d_1 \rightarrow d_2$

---

$\lambda S.S$

$\lambda S.\{a\}$

$\lambda S.(S-\{a\}) \cup \{b\}$

$\lambda S.\text{if } a \in S$
$\quad \text{then } S \cup \{b\}$
$\quad \text{else } S - \{b\}$

# Build Exploded Supergraph

The representation relation of flow function f, $R_f \subseteq (D \cup 0) \times (D \cup 0)$ is a binary relation (or graph) defined as follows:

$$R_f = \{ (0,0) \} \qquad\qquad \text{Edge: } 0 \rightarrow 0$$

$$\cup \{ (0,y) \mid y \in f(\emptyset) \} \qquad \text{Edge: } 0 \rightarrow d_1$$

$$\cup \{ (x,y) \mid y \notin f(\emptyset) \text{ and } y \in f(\{x\}) \} \text{ Edge: } d_1 \rightarrow d_2$$



$\lambda\, S.S$

$\lambda\, S.\{a\}$

$\lambda\, S.(S-\{a\})\cup\{b\}$

$\lambda\, S.\text{if } a \in S$
    $\text{then } S \cup \{b\}$
    $\text{else } S - \{b\}$

# Build Exploded Supergraph

The representation relation of flow function f, $R_f \subseteq (D \cup 0) \times (D \cup 0)$ is a binary relation (or graph) defined as follows:

$R_f$ = { (0,0) }                  Edge: $0 \rightarrow 0$

∪ { (0,y) | y ∈ f(∅) }              Edge: $0 \rightarrow d_1$

∪ { (x,y) | y ∉ f(∅) ... }          Edge: $d_1 \rightarrow d_2$

**Why we need 0 → 0 edges?**

$\lambda S.S$

$\lambda S.\{a\}$

$\lambda S.(S-\{a\}) \cup \{b\}$

$\lambda S.\text{if } a \in S$
$\quad\text{then } S \cup \{b\}$
$\quad\text{else } S - \{b\}$

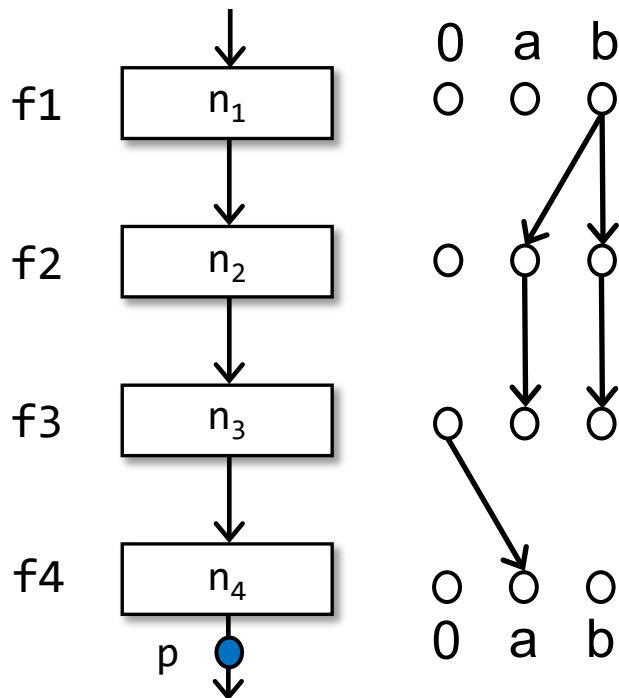# Why We Need Edge 0 → 0?

In traditional data flow analysis, to see whether data fact a holds at program point p, we check if a is in OUT[$n_4$] after the analysis finishes

$$\text{OUT}[n_4] = f4 \circ f3 \circ f2 \circ f1(\text{IN}[n1])$$

f1 | $n_1$

f2 | $n_2$

f3 | $n_3$

f4 | $n_4$

p

# Why We Need Edge 0 → 0?

In traditional data flow analysis, to see whether data fact a holds at program point p, we check if a is in OUT[$n_4$] after the analysis finishes

$$OUT[n_4] = f4 \circ f3 \circ f2 \circ f1(IN[n1])$$

Data facts are propagated via the composition of flow functions. In this case, the "reachability" is directly retrieved from the final result in OUT[$n_4$] .
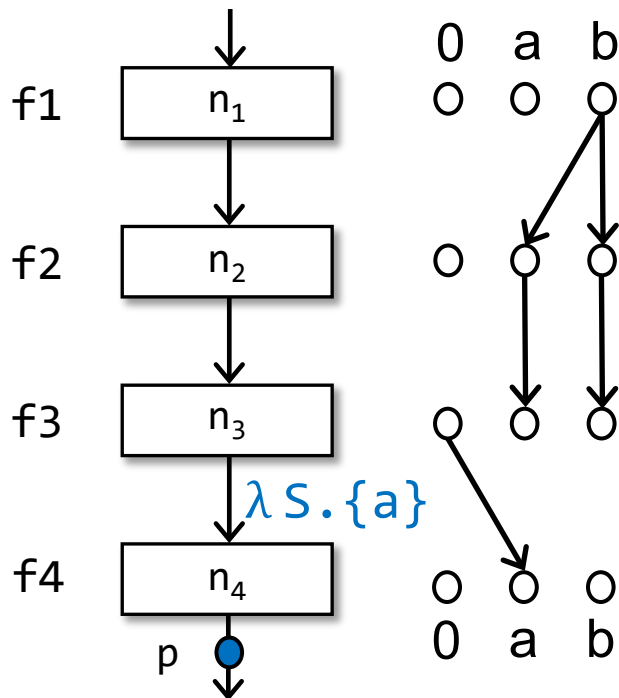
f1    $n_1$

f2    $n_2$

f3    $n_3$

f4    $n_4$

p

# Why We Need Edge 0 → 0?

In traditional data flow analysis, to see whether data fact a holds at program point p, we check if a is in OUT[$n_4$] after the analysis finishes

$$\texttt{OUT[}n_4\texttt{]} = \texttt{f4} \circ \texttt{f3} \circ \texttt{f2} \circ \texttt{f1(IN[n1])}$$

Data facts are propagated via the composition of flow functions. In this case, the "reachability" is directly retrieved from the final result in OUT[$n_4$].



For the same case, in IFDS, whether data fact a holds at p depends on if there is a path from <$s_{main}$, 0> to <$n_4$,a>, and the "reachability" is retrieved by connecting the edges (finding out a path) on the "pasted" representation relations

# Why We Need Edge 0 → 0?

In traditional data flow analysis, to see whether data fact a holds at program point p, we check if a is in OUT[$n_4$] after the analysis finishes

$$\texttt{OUT[}n_4\texttt{]} = \texttt{f4} \circ \texttt{f3} \circ \texttt{f2} \circ \texttt{f1(IN[n1])}$$

Data facts are propagated via the composition of flow functions. In this case, the "reachability" is directly retrieved from the final result in OUT[$n_4$] .
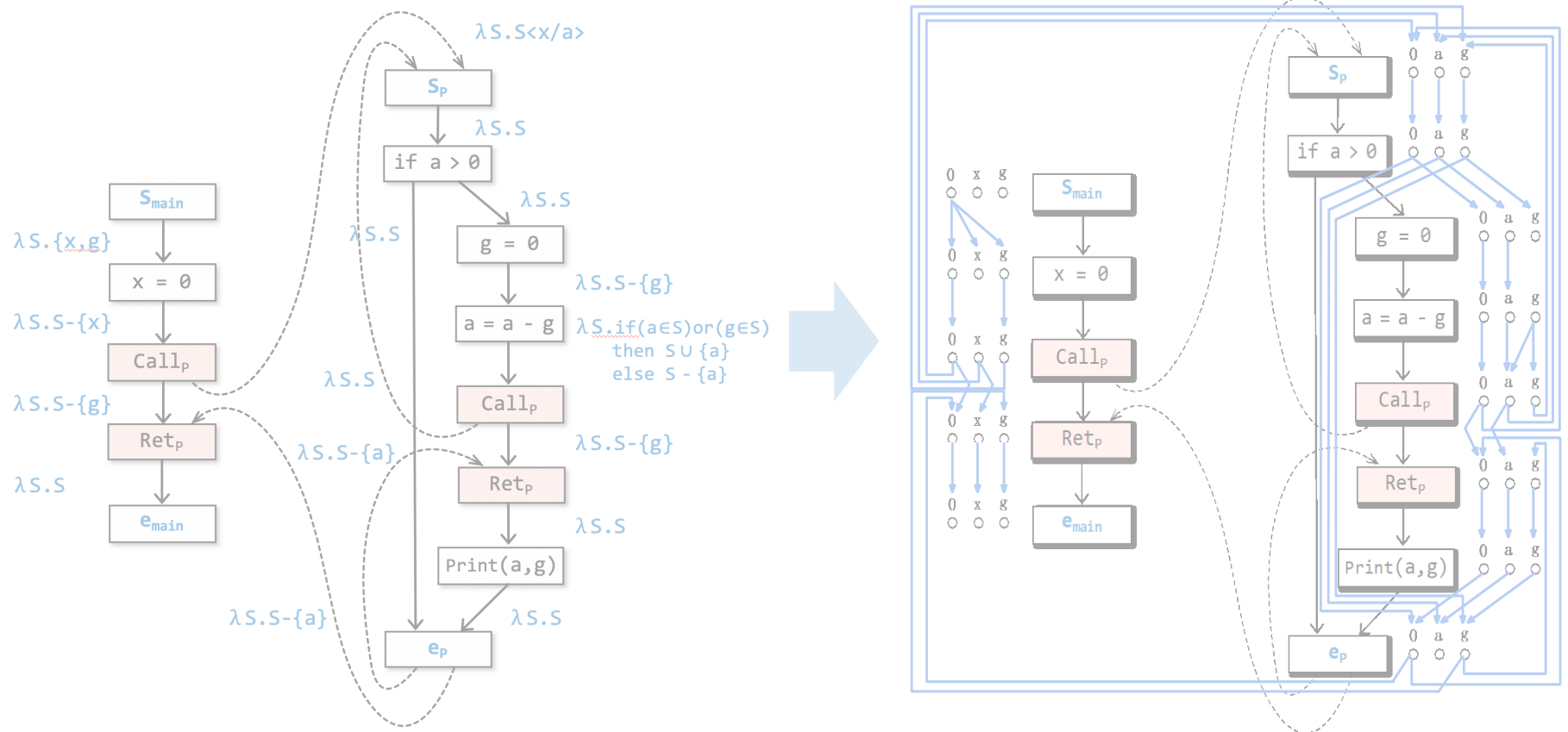


For the same case, in IFDS, whether data fact a holds at p depends on if there is a path from <$s_{main}$, 0> to <$n_4$,a>, and the "reachability" is retrieved by connecting the edges (finding out a path) on the "pasted" representation relations

$\lambda$S.{a} says a holds at p regardless of input S; however, *without edge 0→0,*
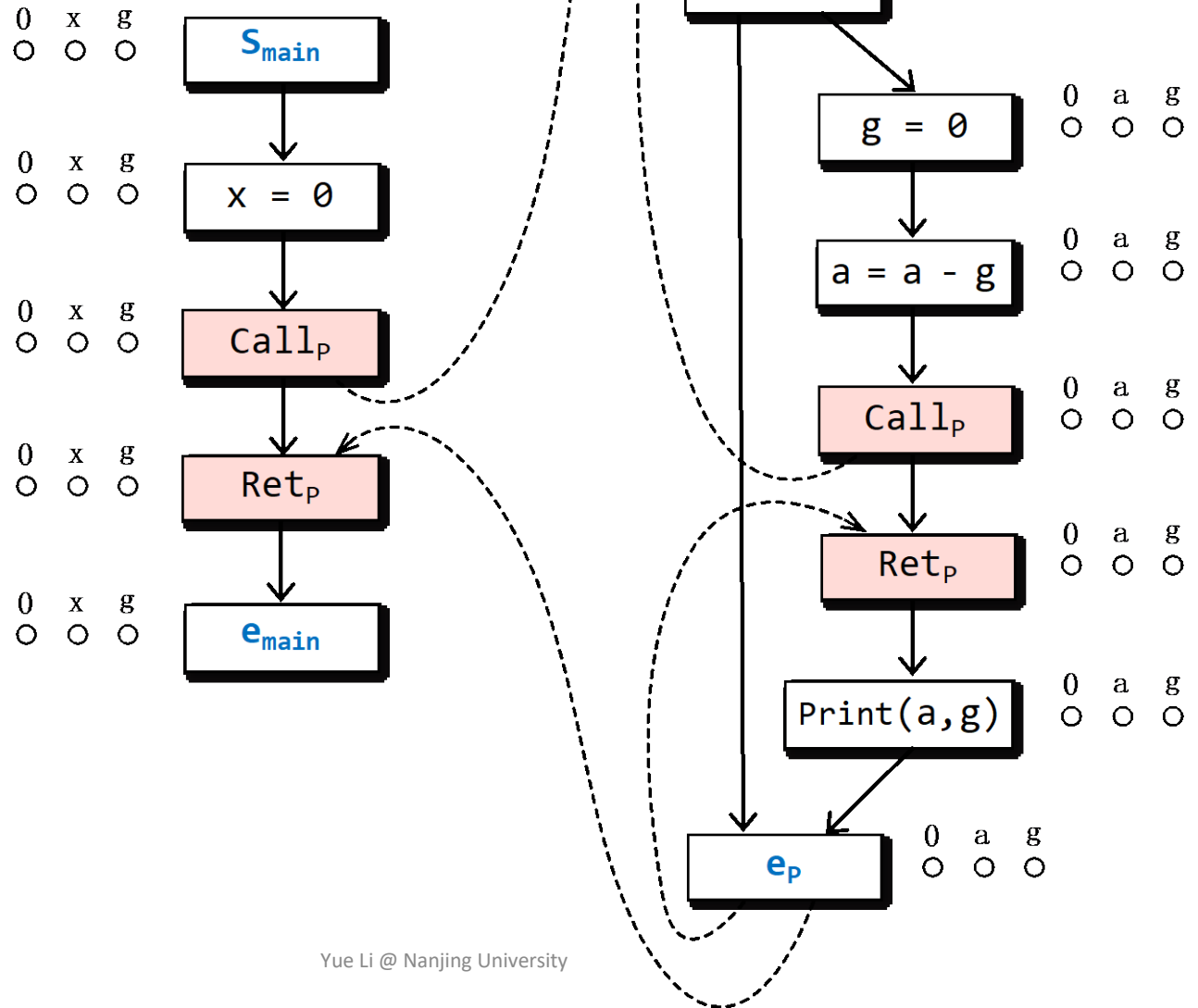
> *the representation relation for each edge cannot be connected or "pasted" together, like flow functions cannot be composed together in traditional data flow analysis.*

Thus IFDS cannot produce correct solutions via such disconnected representation relations.
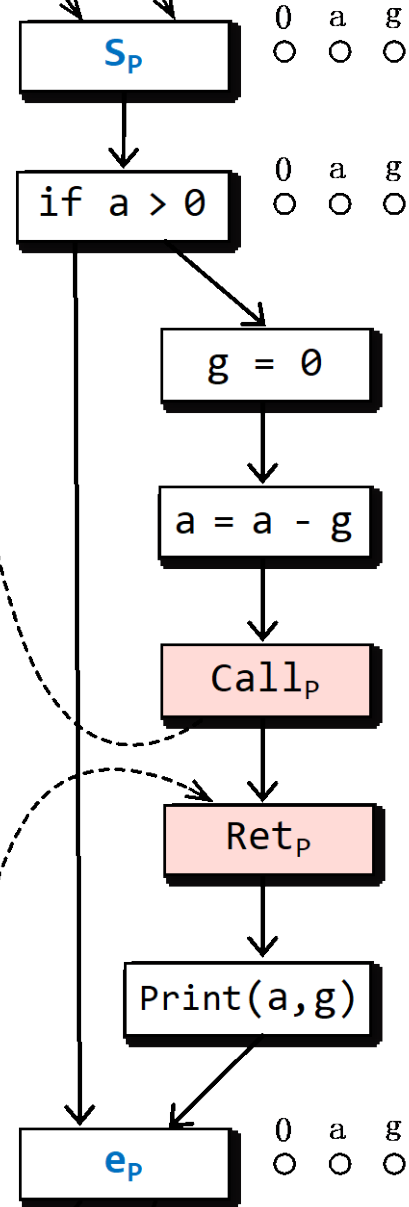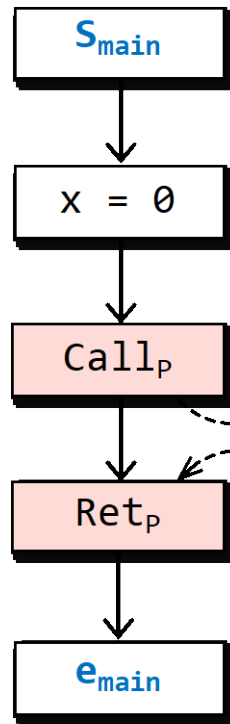
# So We Need the "Glue Edge" $0 \rightarrow 0$!

In traditional data flow analysis, to see whether data fact a holds at program point p, we check if a is in OUT[$n_4$] after the analysis finishes

$$\text{OUT}[n_4] = \text{f4} \circ \text{f3} \circ \text{f2} \circ \text{f1(IN[n1])}$$

Data facts are propagated via the composition of flow functions. In this case, the "reachability" is directly retrieved from the final result in OUT[$n_4$] .

For the same case, in IFDS, whether data fact a holds at p depends on if there is a path from <$s_{main}$, 0> to <$n_4$,a>, and the "reachability" is retrieved by connecting the edges (finding out a path) on the "pasted" representation relations

$\lambda S.\{a\}$ says a holds at p regardless of input S; however, *without edge 0→0,*

*the representation relation for each edge cannot be connected or "pasted" together, like flow functions cannot be composed together in traditional data flow analysis.*

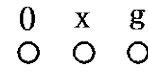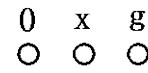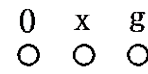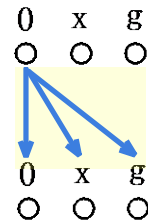Thus IFDS cannot produce correct solutions via such disconnected representation relations.

# Now, let's build an exploded supergraph

$\lambda S.\{x,g\}$

0 x g

$S_{main}$

x = 0

$Call_P$

$Ret_P$

$e_{main}$

0 a g

$S_P$

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

λS.{x,g}

λS.S-{x}

λS.S-{x}

λ S.S<x/a>

0   x   g
○   ○   ○

S_main

0   x   g
○   ○   ○

x = 0

0   x   g
○   ○   ○

Call_P

0   x   g
○   ○   ○

Ret_P

0   x   g
○   ○   ○

e_main

0   a   g
○   ○   ○

S_P

0   a   g
○   ○   ○

if a > 0

0   a   g
○   ○   ○

g = 0

0   a   g
○   ○   ○

a = a - g

0   a   g
○   ○   ○

Call_P

0   a   g
○   ○   ○

Ret_P

0   a   g
○   ○   ○

Print(a,g)

0   a   g
○   ○   ○

e_P

λS.S<x/a>

0 x g
○ ○ ○

Smain

0 x g
○ ○ ○

x = 0

0 x g
○ ○ ○

CallP

0 x g
○ ○ ○

RetP

0 x g
○ ○ ○

emain

0 a g
○ ○ ○

SP

0 a g
○ ○ ○

if a > 0

0 a g
○ ○ ○

g = 0

0 a g
○ ○ ○

a = a - g

0 a g
○ ○ ○

CallP

0 a g
○ ○ ○

RetP

0 a g
○ ○ ○

Print(a,g)

0 a g
○ ○ ○

eP

λS.S

0  a  g
○  ○  ○

**S_P**

if a > 0    0  a  g
            ○  ○  ○

g = 0    0  a  g
         ○  ○  ○

a = a - g    0  a  g
             ○  ○  ○

**Call_P**    0  a  g
              ○  ○  ○

**Ret_P**    0  a  g
             ○  ○  ○

Print(a,g)    0  a  g
              ○  ○  ○

**e_P**    0  a  g
           ○  ○  ○

0  x  g
○  ○  ○

**S_main**

x = 0

0  x  g
○  ○  ○

**Call_P**

0  x  g
○  ○  ○

0  x  g
○  ○  ○

**Ret_P**

0  x  g
○  ○  ○

**e_main**

$\lambda S.S$

$S_{main}$

$x = 0$

$Call_P$

$Ret_P$

$e_{main}$

$S_P$

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

$\lambda S.S$

$\lambda S.S$

$\lambda S.S\text{-}\{g\}$

$S_{main}$

x = 0

$Call_P$

$Ret_P$

$e_{main}$

$S_P$

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

$\lambda S.S\text{-}\{g\}$

S_P

if a > 0

g = 0

a = a - g

Call_P

Ret_P

Print(a,g)

e_P

S_main

x = 0

Call_P

Ret_P

e_main

$$\lambda\, S.\, \text{if}\,(a \in S)\,\text{or}\,(g \in S)$$
$$\text{then}\ S \cup \{a\}$$
$$\text{else}\ S - \{a\}$$

0 x g

0 x g

$S_{main}$

x = 0

$Call_P$

$Ret_P$

$e_{main}$

0 x g

0 x g

0 x g

$S_P$

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

0 a g

0 a g

0 a g

0 a g

0 a g

0 a g

0 a g

0 a g

$\lambda S.if(a \in S) or(g \in S)$
$then \ S \cup \{a\}$
$else \ S - \{a\}$

$\lambda S.S$

$S_P$

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

$S_{main}$

x = 0

$Call_P$

$Ret_P$

$e_{main}$

$\lambda S.S$

$\lambda S.S$

$S_P$

if a > 0

$S_{main}$

x = 0

$Call_P$

$Ret_P$

$e_{main}$

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$\lambda S.S$

$e_P$

$\lambda S.S-\{a\}$

$\lambda S.S-\{a\}$

S_main

x = 0

Call_P

Ret_P

e_main

S_P

if a > 0

g = 0

a = a - g

Call_P

Ret_P

Print(a,g)

e_P

λS.S

$S_{main}$

$x = 0$

$Call_P$

$Ret_P$

$e_{main}$

$S_P$

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

$\lambda S.S$

$\lambda S.S\text{-}\{a\}$

$\lambda\,S.S\text{-}\{a\}$

0 x g

0 a g

$S_P$

0 a g

$S_{main}$

if a > 0

x = 0

g = 0

0 a g

$Call_P$

a = a - g

$\lambda S.S-\{g\}$

0 x g

0 a g

$Ret_P$

$Call_P$

$e_{main}$

0 x g

$Ret_P$

0 a g

Print(a,g)

0 a g

$e_P$

0 a g

Smain
x = 0
CallP
RetP
λS.S
emain

SP
if a > 0
g = 0
a = a - g
CallP
RetP
Print(a,g)
eP

0 x g
0 a g

$S_{main}$

x = 0

$Call_P$

$Ret_P$

$\lambda S.S$

$e_{main}$

$S_P$

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

0    x    g

$S_{main}$

x = 0

$Call_P$

$Ret_P$

$e_{main}$

Is this g reachable?

$S_P$

0    a    g

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

That g is reachable along **realizable paths** from $\langle S_{main}, 0 \rangle$

0    x    g

$S_{main}$

x = 0

$Call_P$

$Ret_P$

$e_{main}$

$S_P$

if a > 0

g = 0

a = a - g

$Call_P$

$Ret_P$

Print(a,g)

$e_P$

0    a    g

Is this g reachable?

Yue Li @ NANJING University

S_P

if a > 0

S_main

x = 0

g = 0

Call_P

a = a - g

Ret_P

Call_P

e_main

Ret_P

Print(a,g)

e_P

That g is reachable only along **non-realizable paths** from <S_main,0>

Given an exploded supergraph, we apply Tabulation algorithm to identify MRP solutions

Given an exploded supergraph, we apply Tabulation algorithm to identify MRP solutions

$S_P$

0   a   g

0   a   g

if a > 0

0   x   g

$S_{main}$

x = 0

$Call_P$

$Ret_P$

$e_{main}$

g = 0

0   a   g

a = a - g

0   a   g

$Call_P$

0   a   g

$Ret_P$

0   a   g

Print(a,g)

0   a   g

**Blue circles (final results)** denote the nodes that are reachable along realizable paths from $<S_{main},0>$

$e_P$

0   a   g

How?

Given an exploded supergraph, we apply Tabulation algorithm to identify MRP solutions

S_main
x = 0
Call_P
Ret_P
e_main

S_P
if a > 0
g = 0
a = a - g
Call_P
Ret_P
Print(a,g)
e_P

**Blue circles (final results)** denote the nodes that are reachable along realizable paths from $<S_{main}, 0>$

Yue Li @ Nanjing University

# Overview of IFDS

Given a program P, and a dataflow-analysis problem Q

- Build a supergraph G* for P and
  define flow functions for edges in G* based on Q

- Build exploded supergraph G# for P by transforming
  flow functions to representation relations (graphs)

- Q can be solved as graph reachability problems (find out MRP solutions)
  via applying Tabulation algorithm on G#

# Tabulation Algorithm

Given an exploded supergraph G#, Tabulation algorithm determines the MRP solution by finding out all realizable paths starting from <$s_{main}$, 0>



Tabulation Algorithm

# Tabulation Algorithm

Given an exploded supergraph G#, Tabulation algorithm determines the MRP solution by finding out all realizable paths starting from $<s_{main}, 0>$

Let $n$ be a program point, data fact $d \in MRP_n$, iff there is a realizable path in G# from $<s_{main}, 0>$ to $<n, d>$. (then $d$'s white circle turns to blue)



Tabulation Algorithm

# Tabulation Algorithm

**declare** PathEdge, WorkList, SummaryEdge: **global** edge set

**algorithm** Tabulate($G_{IP}^{\#}$)
**begin**
[1]    Let $(N^{\#}, E^{\#}) = G_{IP}^{\#}$
[2]    PathEdge := { $\langle s_{main}, \mathbf{0} \rangle \rightarrow \langle s_{main}, \mathbf{0} \rangle$ }
[3]    WorkList := { $\langle s_{main}, \mathbf{0} \rangle \rightarrow \langle s_{main}, \mathbf{0} \rangle$ }
[4]    SummaryEdge := $\varnothing$
[5]    ForwardTabulateSLRPs()
[6]    **for** each $n \in N^{*}$ **do**
[7]      $X_n := \{ d_2 \in D \mid \exists d_1 \in (D \cup \{ \mathbf{0} \})$ such that $\langle s_{procOf(n)}, d_1 \rangle \rightarrow \langle n, d_2 \rangle \in$ PathEdge $\}$
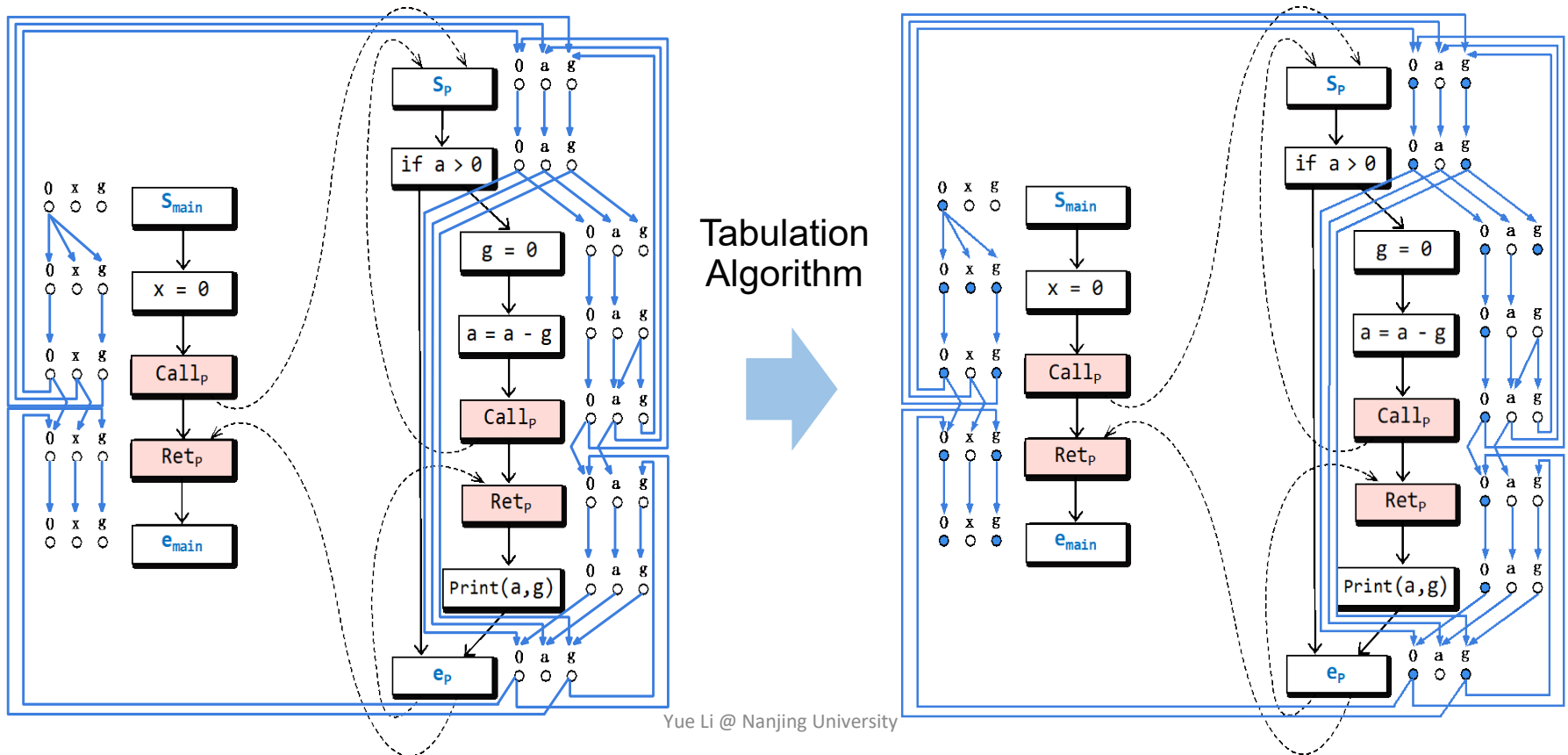[8]    **od**
**end**

**procedure** Propagate($e$)
**begin**
[9]    **if** $e \notin$ PathEdge **then** Insert $e$ into PathEdge; Insert $e$ into WorkList **fi**
**end**

**procedure** ForwardTabulateSLRPs()
**begin**
[10]    **while** WorkList $\neq \varnothing$ **do**
[11]      Select and remove an edge $\langle s_p, d_1 \rangle \rightarrow \langle n, d_2 \rangle$ from WorkList
[12]      **switch** $n$

[13]        **case** $n \in Call_p$ :
[14]          **for** each $d_3$ such that $\langle n, d_2 \rangle \rightarrow \langle s_{calledProc(n)}, d_3 \rangle \in E^{\#}$ **do**
[15]            Propagate($\langle s_{calledProc(n)}, d_3 \rangle \rightarrow \langle s_{calledProc(n)}, d_3 \rangle$)
[16]          **od**
[17]          **for** each $d_3$ such that $\langle n, d_2 \rangle \rightarrow \langle returnSite(n), d_3 \rangle \in (E^{\#} \cup$ SummaryEdge) **do**
[18]            Propagate($\langle s_p, d_1 \rangle \rightarrow \langle returnSite(n), d_3 \rangle$)
[19]          **od**
[20]        **end case**

[21]        **case** $n = e_p$ :
[22]          **for** each $c \in callers(p)$ **do**
[23]            **for** each $d_4, d_5$ such that $\langle c, d_4 \rangle \rightarrow \langle s_p, d_1 \rangle \in E^{\#}$ and $\langle e_p, d_2 \rangle \rightarrow \langle returnSite(c), d_5 \rangle \in E^{\#}$ **do**
[24]              **if** $\langle c, d_4 \rangle \rightarrow \langle returnSite(c), d_5 \rangle \notin$ SummaryEdge **then**
[25]                Insert $\langle c, d_4 \rangle \rightarrow \langle returnSite(c), d_5 \rangle$ into SummaryEdge
[26]                **for** each $d_3$ such that $\langle s_{procOf(c)}, d_3 \rangle \rightarrow \langle c, d_4 \rangle \in$ PathEdge **do**
[27]                  Propagate($\langle s_{procOf(c)}, d_3 \rangle \rightarrow \langle returnSite(c), d_5 \rangle$)
[28]                **od**
[29]              **fi**
[30]            **od**
[31]          **od**
[32]        **end case**

[33]        **case** $n \in (N_p - Call_p - \{ e_p \})$ :
[34]          **for** each $\langle m, d_3 \rangle$ such that $\langle n, d_2 \rangle \rightarrow \langle m, d_3 \rangle \in E^{\#}$ **do**
[35]            Propagate($\langle s_p, d_1 \rangle \rightarrow \langle m, d_3 \rangle$)
[36]          **od**
[37]        **end case**

[38]      **end switch**
[39]    **od**
**end**

$$O(\mathrm{ED}^3)$$

# Tabulation Algorithm

**declare** PathEdge, WorkList, SummaryEdge: **global** edge set

**algorithm** Tabulate($G_{IP}^\#$)
**begin**
[1]    Let $(N^\#, E^\#) = G_{IP}^\#$
[2]    PathEdge := { $\langle s_{main}, \mathbf{0} \rangle \to \langle s_{main}, \mathbf{0} \rangle$ }
[3]    WorkList := { $\langle s_{main}, \mathbf{0} \rangle \to \langle s_{main}, \mathbf{0} \rangle$ }
[4]    SummaryEdge := $\varnothing$
[5]    ForwardTabulateSLRPs()
[6]    **for** each $n \in N^*$ **do**
[7]      $X_n := \{ d_2 \in D \mid \exists d_1 \in (D \cup \{ \mathbf{0} \}) \text{ such that } \langle s_{procOf(n)}, d_1 \rangle \to \langle n, d_2 \rangle \in \text{PathEdge} \}$
[8]    **od**
**end**

**procedure** Propagate($e$)
**begin**
[9]    **if** $e \notin$ PathEdge **then** Insert $e$ into PathEdge; Insert $e$ into WorkList **fi**
**end**

**procedure** ForwardTabulateSLRPs()
**begin**
[10]  **while** WorkList $\neq \varnothing$ **do**
[11]    Select and remove an edge $\langle s_p, d_1 \rangle \to \langle n, d_2 \rangle$ from WorkList
[12]    **switch** $n$
[13]      **case** $n \in Call_p$ :
[14]        **for** each $d_3$ such that $\langle n, d_2 \rangle \to \langle s_{calledProc(n)}, d_3 \rangle \in E^\#$ **do**
[15]          Propagate($\langle s_{calledProc(n)}, d_3 \rangle \to \langle s_{calledProc(n)}, d_3 \rangle$)
[16]        **od**
[17]        **for** each $d_3$ such that $\langle n, d_2 \rangle \to \langle returnSite(n), d_3 \rangle \in (E^\# \cup \text{SummaryEdge})$ **do**
[18]          Propagate($\langle s_p, d_1 \rangle \to \langle returnSite(n), d_3 \rangle$)
[19]        **od**
[20]      **end case**
[21]      **case** $n = e_p$ :
[22]        **for** each $c \in callers(p)$ **do**
[23]          **for** each $d_4, d_5$ such that $\langle c, d_4 \rangle \to \langle s_p, d_1 \rangle \in E^\#$ and $\langle e_p, d_2 \rangle \to \langle returnSite(c), d_5 \rangle \in E^\#$ **do**
[24]            **if** $\langle c, d_4 \rangle \to \langle returnSite(c), d_5 \rangle \notin$ SummaryEdge **then**
[25]              Insert $\langle c, d_4 \rangle \to \langle returnSite(c), d_5 \rangle$ into SummaryEdge
[26]              **for** each $d_3$ such that $\langle s_{procOf(c)}, d_3 \rangle \to \langle c, d_4 \rangle \in$ PathEdge **do**
[27]                Propagate($\langle s_{procOf(c)}, d_3 \rangle \to \langle returnSite(c), d_5 \rangle$)
[28]              **od**
[29]            **fi**
[30]          **od**
[31]        **od**
[32]      **end case**
[33]      **case** $n \in (N_p - Call_p - \{ e_p \})$ :
[34]        **for** each $\langle m, d_3 \rangle$ such that $\langle n, d_2 \rangle \to \langle m, d_3 \rangle \in E^\#$ **do**
[35]          Propagate($\langle s_p, d_1 \rangle \to \langle m, d_3 \rangle$)
[36]        **od**
[37]      **end case**
[38]    **end switch**
[39]  **od**
**end**

$$O(\text{ED}^3)$$

**No time to cover the whole algorithm**

# Tabulation Algorithm

**declare** PathEdge, WorkList, SummaryEdge: **global** edge set

**algorithm** Tabulate($G_{IP}^{\#}$)
**begin**
[1]    Let $(N^{\#}, E^{\#}) = G_{IP}^{\#}$
[2]    PathEdge := { $\langle s_{main}, \mathbf{0} \rangle \rightarrow \langle s_{main}, \mathbf{0} \rangle$ }
[3]    WorkList := { $\langle s_{main}, \mathbf{0} \rangle \rightarrow \langle s_{main}, \mathbf{0} \rangle$ }
[4]    SummaryEdge := $\varnothing$
[5]    ForwardTabulateSLRPs()
[6]    **for** each $n \in N^{*}$ **do**
[7]      $X_n$ := { $d_2 \in D \mid \exists d_1 \in (D \cup \{ \mathbf{0} \})$ such that $\langle s_{procOf(n)}, d_1 \rangle \rightarrow \langle n, d_2 \rangle \in$ PathEdge }
[8]    **od**
**end**

**procedure** Propagate($e$)
**begin**
[9]   **if** $e \notin$ PathEdge **then** Insert $e$ into PathEdge; Insert $e$ into WorkList **fi**
**end**

**procedure** ForwardTabulateSLRPs()
**begin**
[10]   **while** WorkList $\neq \varnothing$ **do**
[11]    Select and remove an edge $\langle s_p, d_1 \rangle \rightarrow \langle n, d_2 \rangle$ from WorkList
[12]    **switch** $n$
[13]      **case** $n \in Call_p$ :
[14]        **for** each $d_3$ such that $\langle n, d_2 \rangle \rightarrow \langle s_{calledProc(n)}, d_3 \rangle \in E^{\#}$ **do**
[15]          Propagate($\langle s_{calledProc(n)}, d_3 \rangle \rightarrow \langle s_{calledProc(n)}, d_3 \rangle$)
[16]        **od**
[17]        **for** each $d_3$ such that $\langle n, d_2 \rangle \rightarrow \langle returnSite(n), d_3 \rangle \in (E^{\#} \cup$ SummaryEdge) **do**
[18]          Propagate($\langle s_p, d_1 \rangle \rightarrow \langle returnSite(n), d_3 \rangle$)
[19]        **od**
[20]      **end case**
[21]      **case** $n = e_p$ :
[22]        **for** each $c \in callers(p)$ **do**
[23]          **for** each $d_4, d_5$ such that $\langle c, d_4 \rangle \rightarrow \langle s_p, d_1 \rangle \in E^{\#}$ and $\langle e_p, d_2 \rangle \rightarrow \langle returnSite(c), d_5 \rangle \in E^{\#}$ **do**
[24]            **if** $\langle c, d_4 \rangle \rightarrow \langle returnSite(c), d_5 \rangle \notin$ SummaryEdge **then**
[25]              Insert $\langle c, d_4 \rangle \rightarrow \langle returnSite(c), d_5 \rangle$ into SummaryEdge
[26]              **for** each $d_3$ such that $\langle s_{procOf(c)}, d_3 \rangle \rightarrow \langle c, d_4 \rangle \in$ PathEdge **do**
[27]                Propagate($\langle s_{procOf(c)}, d_3 \rangle \rightarrow \langle returnSite(c), d_5 \rangle$)
[28]              **od**
[29]            **fi**
[30]          **od**
[31]        **od**
[32]      **end case**
[33]      **case** $n \in (N_p - Call_p - \{ e_p \})$ :
[34]        **for** each $\langle m, d_3 \rangle$ such that $\langle n, d_2 \rangle \rightarrow \langle m, d_3 \rangle \in E^{\#}$ **do**
[35]          Propagate($\langle s_p, d_1 \rangle \rightarrow \langle m, d_3 \rangle$)
[36]        **od**
[37]      **end case**
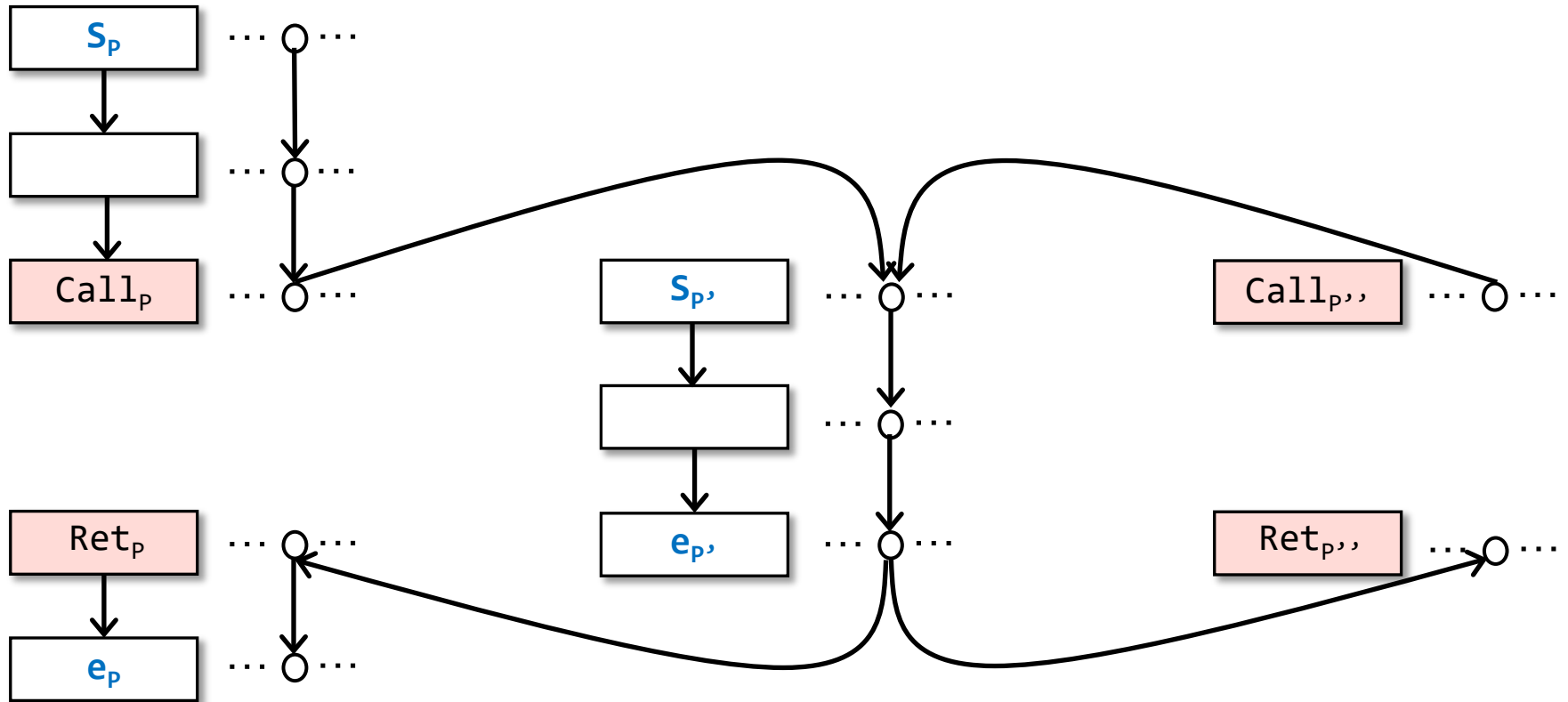[38]    **end switch**
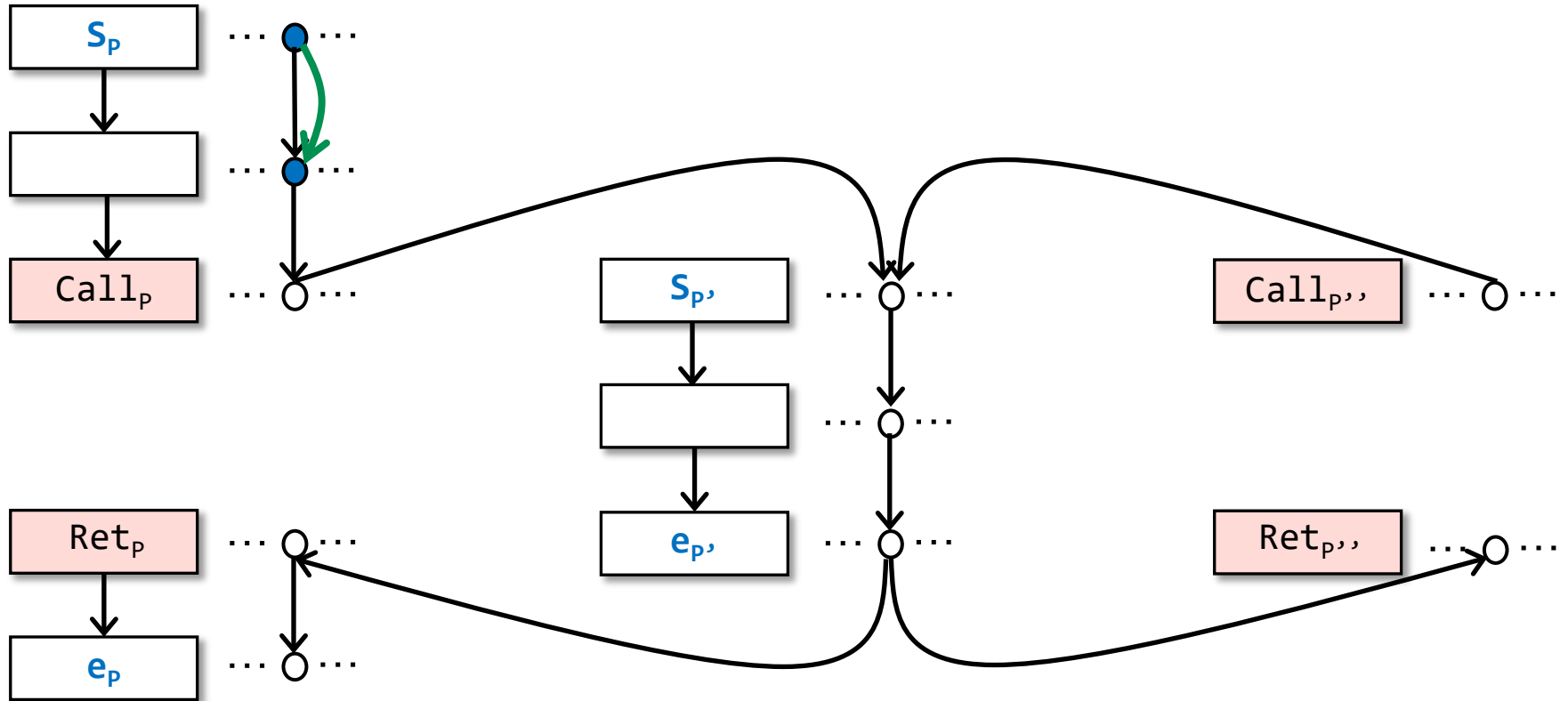[39]   **od**
**end**

$$O(\text{ED}^3)$$

No time to cover the whole algorithm

But we will introduce its core working mechanism by a simple example
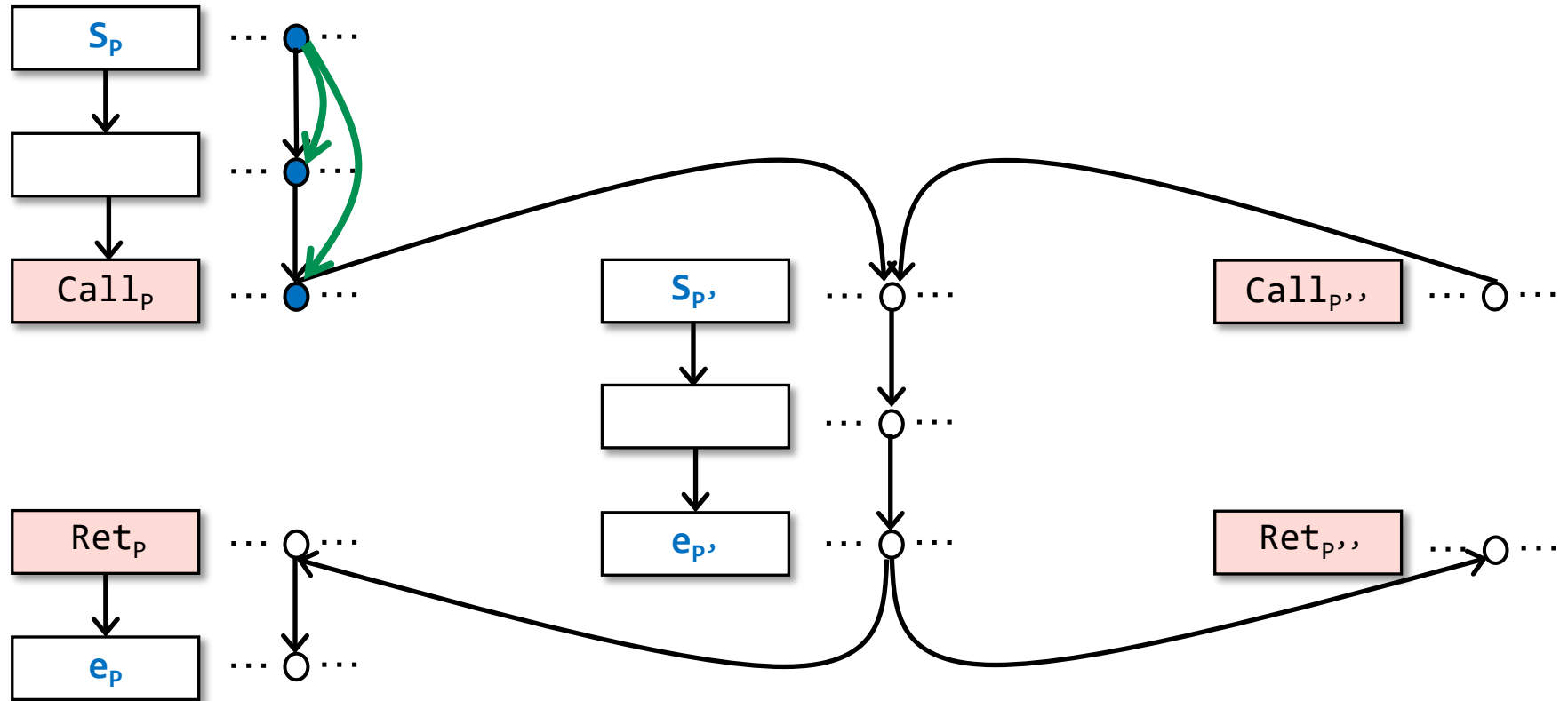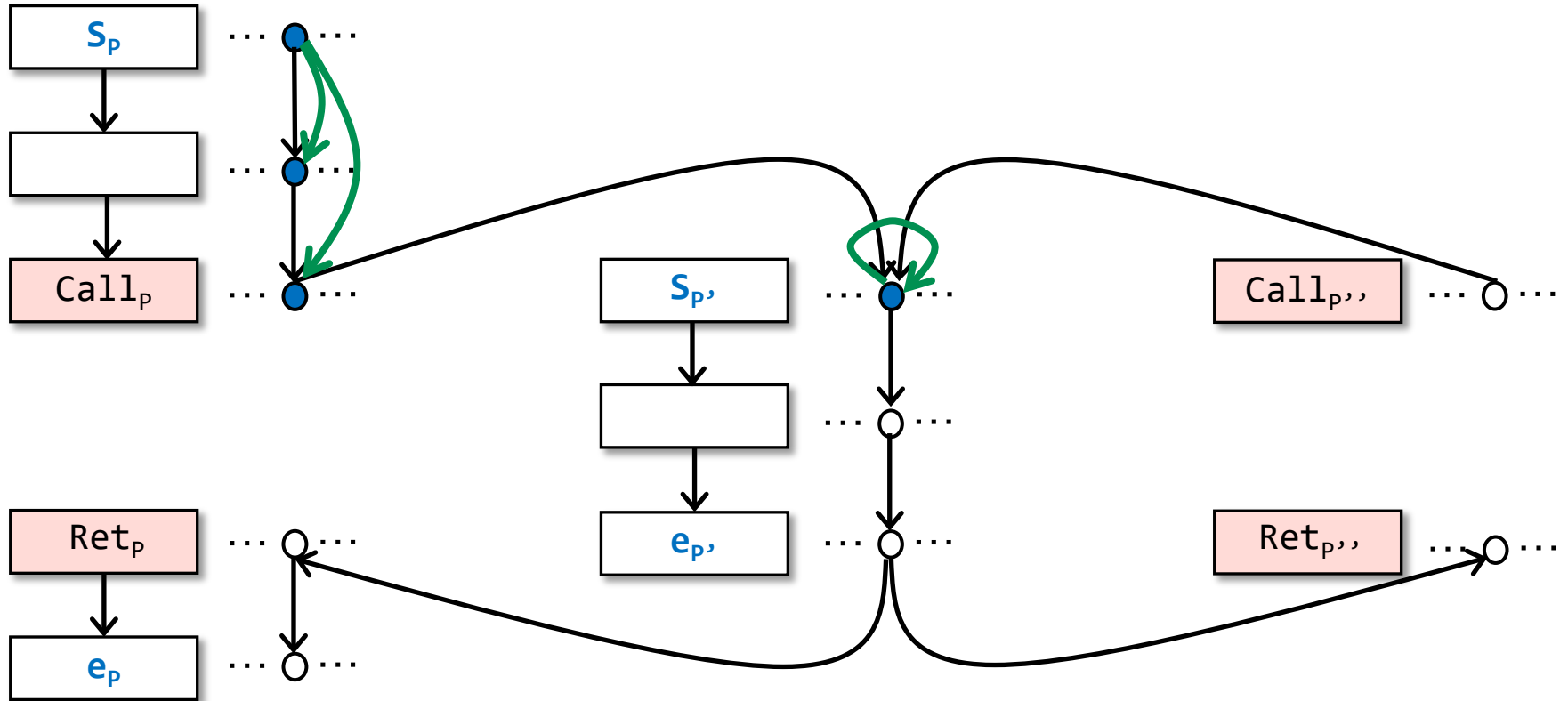
# Core Working Mechanism of Tabulation Algorithm

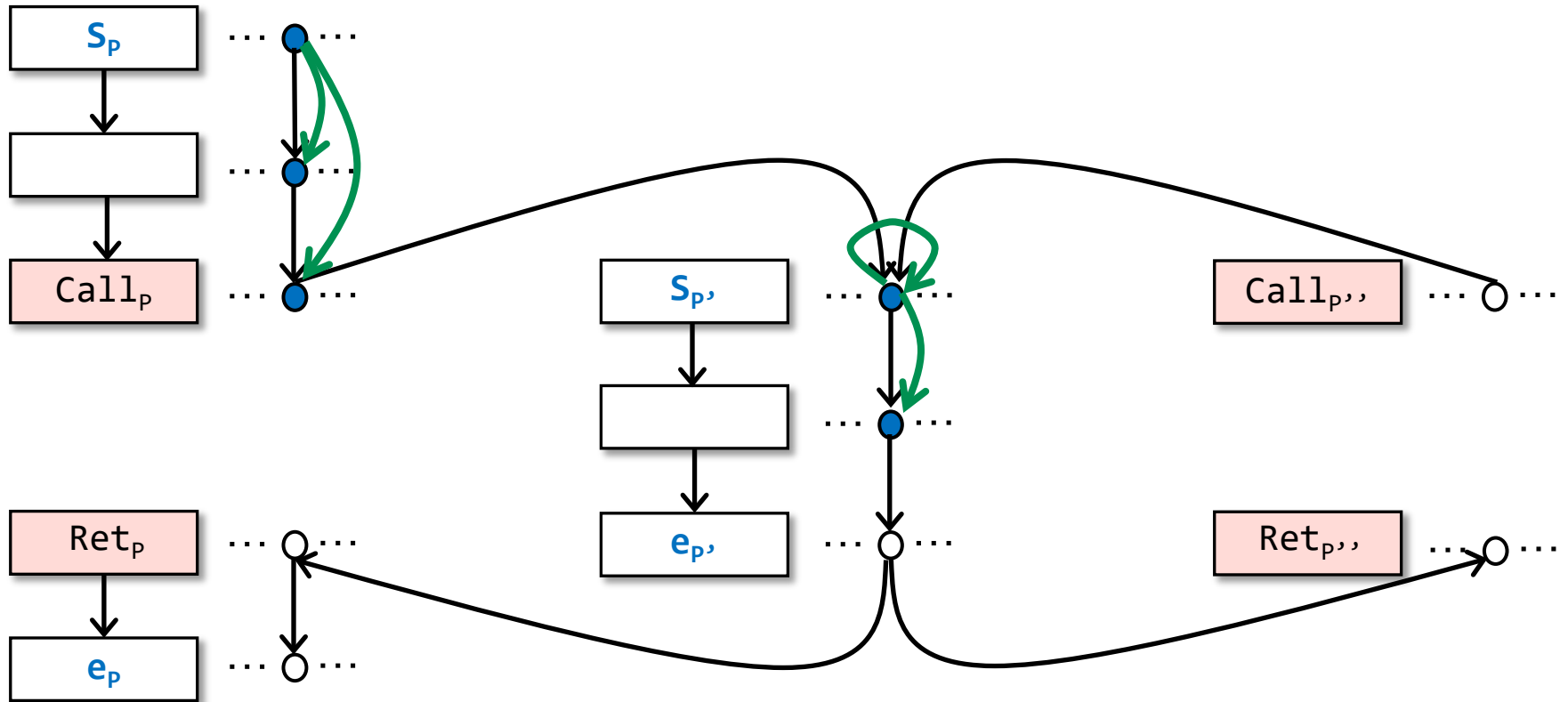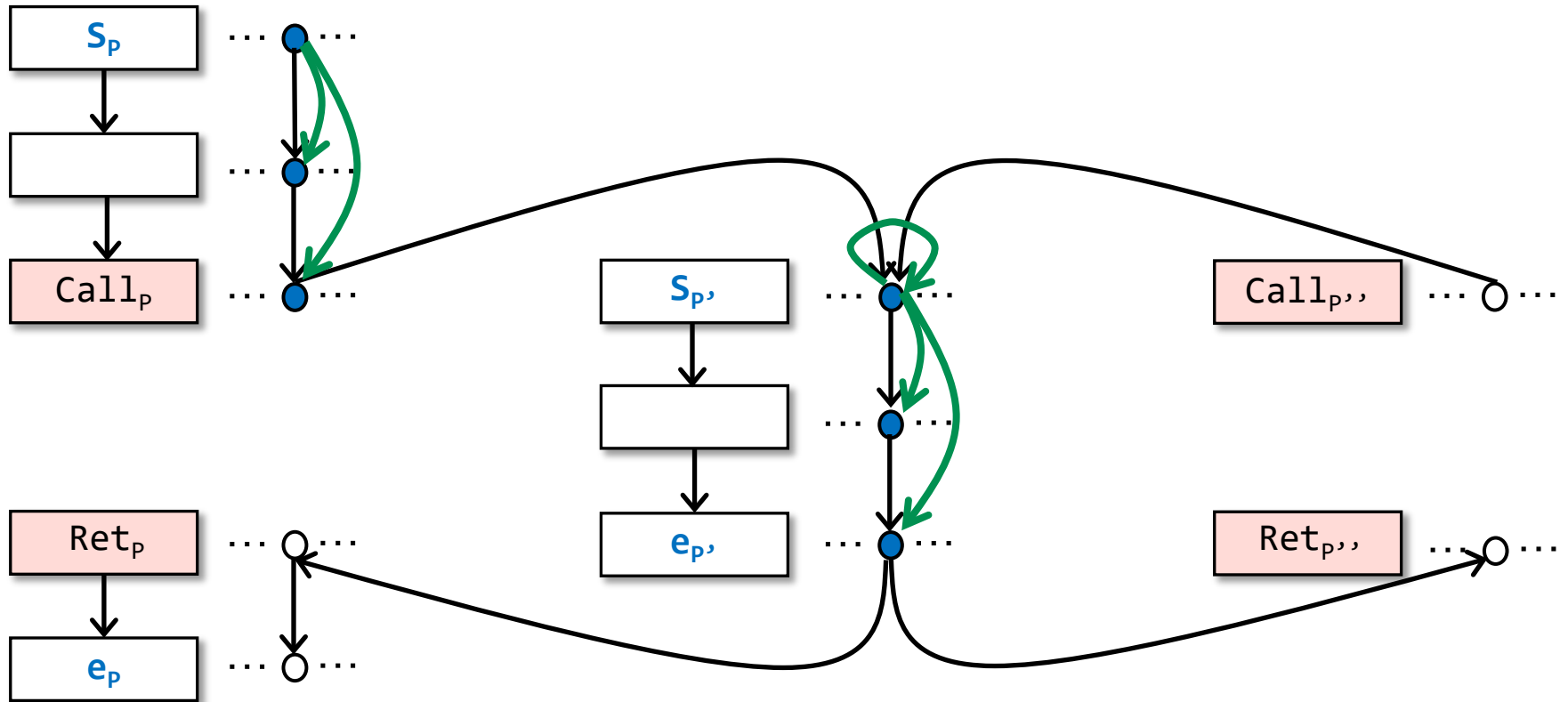# Core Working Mechanism of Tabulation Algorithm

# Core Working Mechanism of Tabulation Algorithm

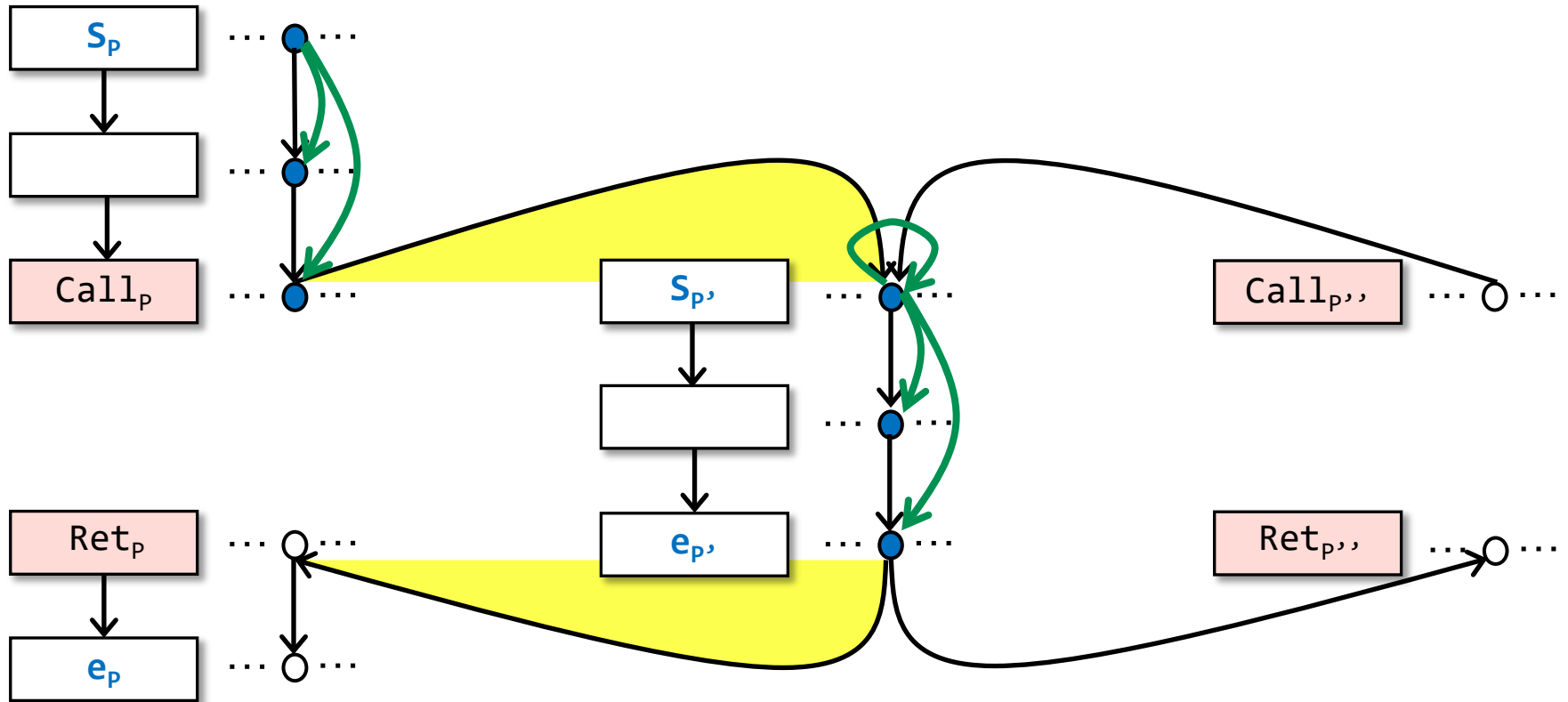# Core Working Mechanism of Tabulation Algorithm

# Core Working Mechanism of Tabulation Algorithm

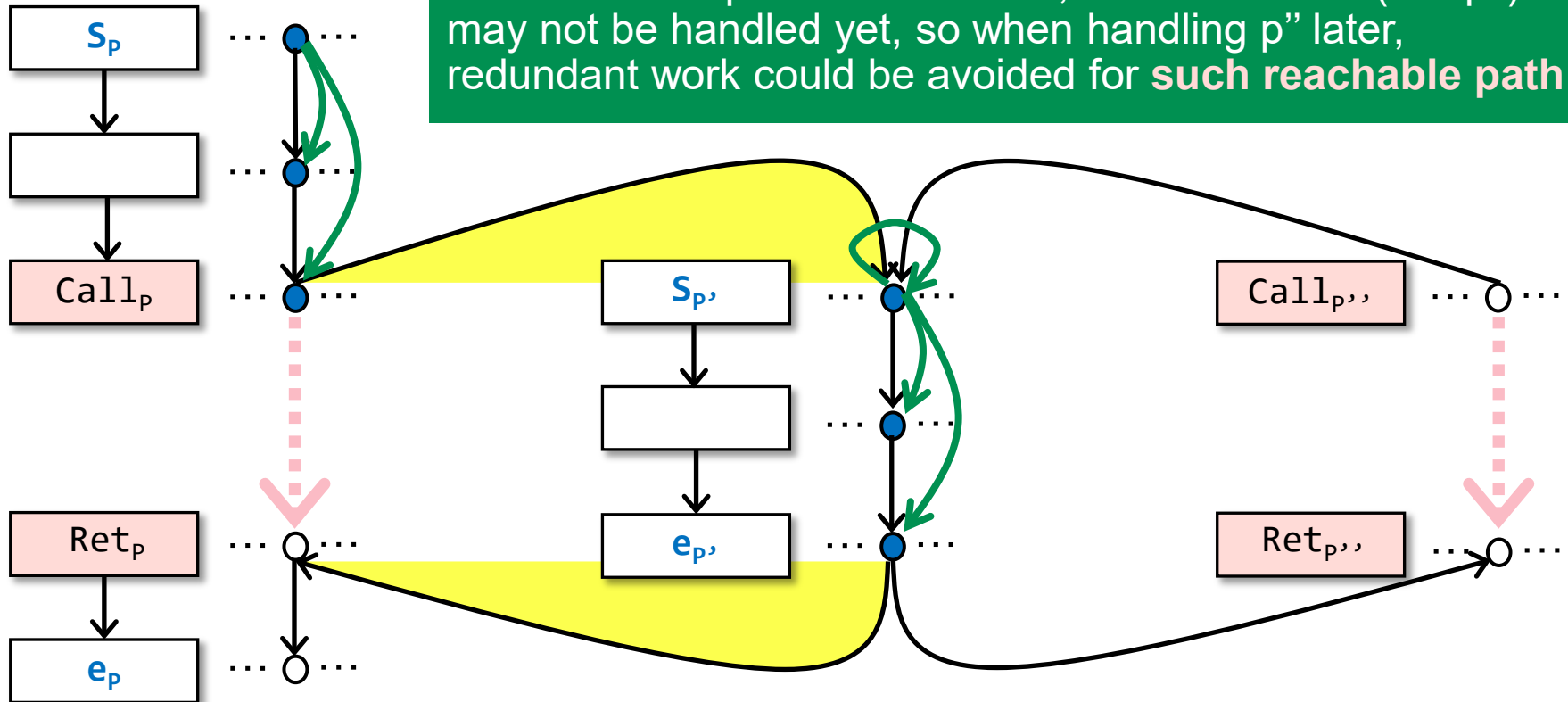# Core Working Mechanism of Tabulation Algorithm

# Core Working Mechanism of Tabulation Algorithm



When handling each exit node ($e_{p'}$), **call-to-return matching** begins: find out the call-sites calling p' ($Call_p$, $Call_{p''}$) and then find out their corresponding return-sites ($Ret_p$, $Ret_{p''}$).

# Core Working Mechanism of Tabulation Algorithm

Actually, here **a summary edge** from <Call,$d_m$> to <Ret,$d_n$> is added to indicate that $d_n$ is **reachable** from $d_m$ through the called method p'. At the moment, some methods (like p'') may not be handled yet, so when handling p'' later, redundant work could be avoided for **such reachable path**.

$S_P$

$Call_P$

$Ret_P$

$e_P$

$S_{P'}$

$e_{P'}$

$Call_{P''}$

$Ret_{P''}$

When handling each exit node ($e_{p'}$), *call-to-return matching* begins: find out the call-sites calling p' ($Call_p$, $Call_{p''}$) and then find out their corresponding return-sites ($Ret_p$, $Ret_{p''}$).
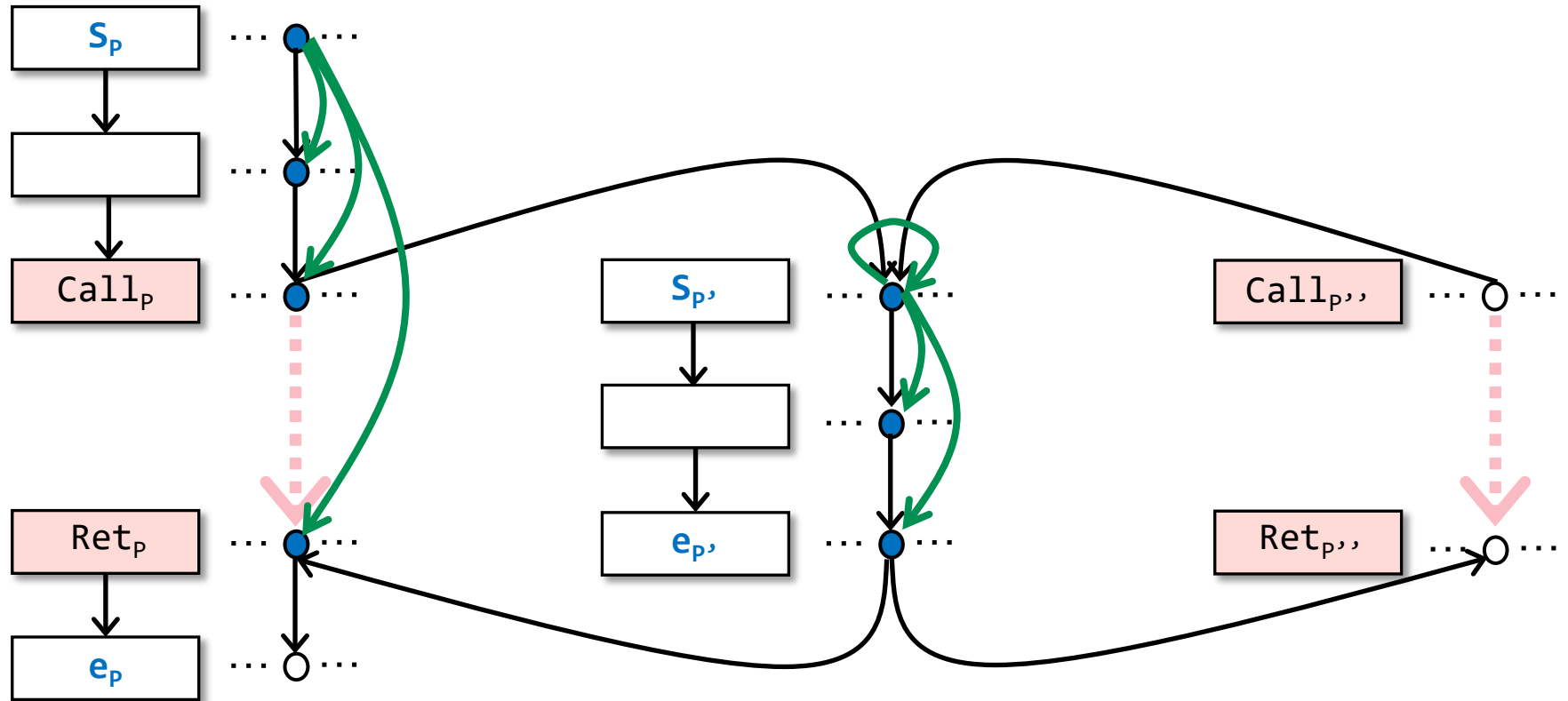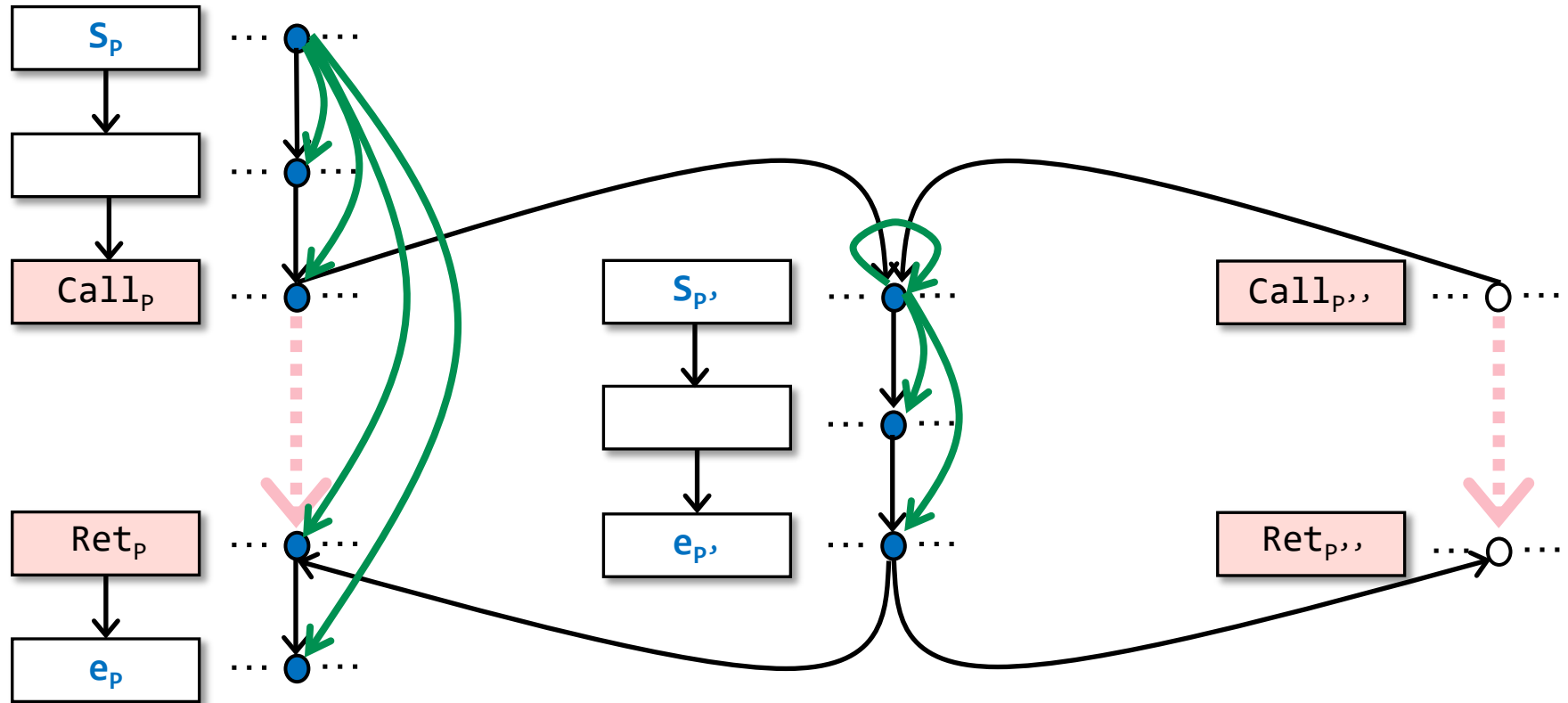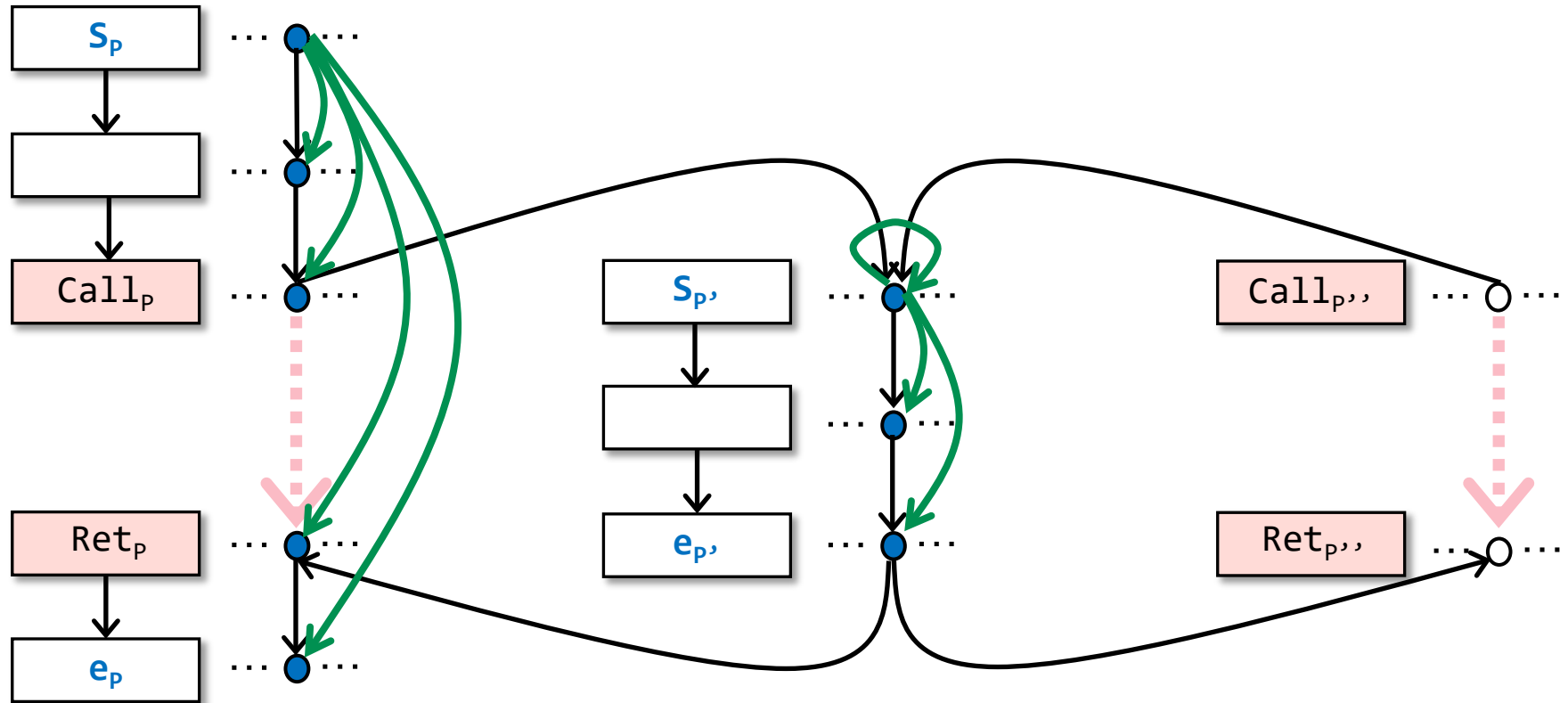
# Core Working Mechanism of Tabulation Algorithm

# Core Working Mechanism of Tabulation Algorithm

# Core Working Mechanism of Tabulation Algorithm



When a data fact (at node n) d's circle is turned to blue, it means that $\langle n, d \rangle$ is reachable from $\langle S_{main}, 0 \rangle$

# Understanding the Distributivity of IFDS

# Understanding the Distributivity of IFDS

- Can we do constant propagation using IFDS?

  Constant propagation has infinite domain, but what if we only deal with finite constant values? Can we still do it using IFDS?

- Can we do pointer analysis using IFDS?

# Understanding the Distributivity of IFDS

- Can we do constant propagation using IFDS? **No**

  Constant propagation has infinite domain, but what if we only deal with finite constant values? Can we still do it using IFDS?

- Can we do pointer analysis using IFDS? **No**

# Understanding the Distributivity of IFDS

- **Distributivity**

  $F(x \wedge y) = F(x) \wedge F(y)$

- **Constant Propagation**

  |         | x | y | z |
  |---------|---|---|---|
  | z = x + y | o | o | o |

  z's value depends on both y's and x's

# Understanding the Distributivity of IFDS

- **Distributivity**

$$F(x \wedge y) = F(x) \wedge F(y)$$

Each flow function in IFDS handles one input data fact per time

| z = x + y | x | y | z |
|-----------|---|---|---|
|           | o | o | o |

z's value depends on both y's and x's

# Understanding the Distributivity of IFDS

- Distributivity

$F(x \wedge y) = F(x) \wedge F(y)$

Each flow function in IFDS handles one input data fact per time

| | x | y | z |
|---|---|---|---|
| z = x + y | o | o | o |

Each representation relation indicates "if x exists, then …", "if y exists then …" But when we need "if both x and y exist", how to draw the representation relation?

# Understanding the Distributivity of IFDS

- Distributivity

Each flow function in IFDS handles one input data fact per time

For constant propagation, we cannot define F if we only know x's (or y's) value

$$z = x + y$$

| x | y | z |
|---|---|---|
| o | o | o |

Each representation relation indicates "if x exists, then …", "if y exists then …" But when we need "if both x and y exist", how to draw the representation relation?

# Understanding the Distributivity of IFDS

- Distributivity

Each flow function in IFDS handles one input data fact per time

For constant propagation, we cannot define F if we only know x's (or y's) value

z = x + y

| x | y | z |
|---|---|---|
| o | o | o |

Each representation relation indicates "if x exists, then …", "if y exists then …" But when we need "if both x and y exist", how to draw the representation relation?

Given a statement S, besides S itself, if we need to consider **multiple** input data facts to create correct outputs, then the analysis is not distributive and should not be expressed in IFDS.

In IFDS, each data fact (circle) and its propagation (edges) could be handled **independently**, and doing so will not affect the correctness of the final results.

# Understanding the Distributivity of IFDS

- Distributivity

For constant propagation, we cannot define F if we only know x's (or y's) value

| | x | y | z |
|---|---|---|---|
| z = x + y | o | o | o |

Each representation relation indicates "if x exists, then …", "if y exists then …" But when we need "if both x and y exist", how to draw the representation relation?

A simple rule to determine whether your analysis could be expressed in IFDS

Given a statement S, besides S itself, if we need to consider **multiple** input data facts to create correct outputs, then the analysis is not distributive and should not be expressed in IFDS.
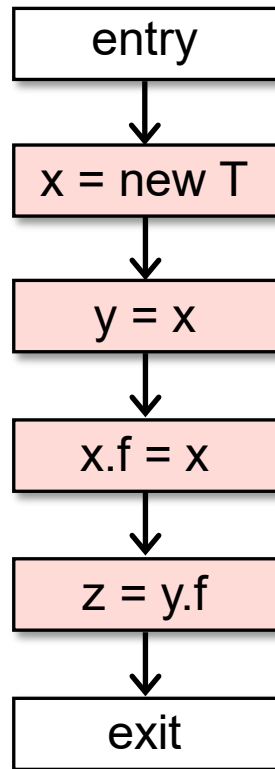
In IFDS, each data fact (circle) and its propagation (edges) could be handled **independently**, and doing so will not affect the correctness of the final results.

# Understanding the Distributivity of IFDS

- Distributivity

Each flow function in IFDS handles one input data fact per time

For constant propagation, we cannot define F if we only know x's (or y's) value

|         | x | y | z |
|---------|---|---|---|
| z = x + y | o | o | o |

Each representation relation indicates "if x exists, then …", "if y exists then …" But when we need "if both x and y exist", how to draw the representation relation?

A simple rule to determine whether your analysis could be expressed in IFDS

Given a statement S, besides S itself, if we need to consider **multiple** input data facts to create correct outputs, then the analysis is not distributive and should not be expressed in IFDS.

In IFDS, each data fact (circle) and its propagation (edges) could be handled **independently**, and doing so will not affect the correctness of the final results.

Regardless of the infinite domain issue, think about whether we could do *linear constant propagation*, e,g., y = 2x + 3, or *copy constant propagation*, e.g., x = 2, y = x, using IFDS-style analysis?

# Understanding the Distributivity of IFDS

- Pointer Analysis



```
entry

x = new T

y = x

x.f = x

z = y.f

exit
```

# Understanding the Distributivity of IFDS

- Pointer Analysis

|  | 0 | x | y | x.f | y.f | z |
|---|---|---|---|---|---|---|
| entry | O | O | O | O | O | O |
| x = new T | O | O | O | O | O | O |
| y = x | O | O | O | O | O | O |
| x.f = x | O | O | O | O | O | O |
| z = y.f | O | O | O | O | O | O |
| exit | O | O | O | O | O | O |

# Understanding the Distributivity of IFDS

- Pointer Analysis

# Understanding the Distributivity of IFDS

- Pointer Analysis

# Understanding the Distributivity of IFDS

- Pointer Analysis

# Understanding the Distributivity of IFDS

- Pointer Analysis

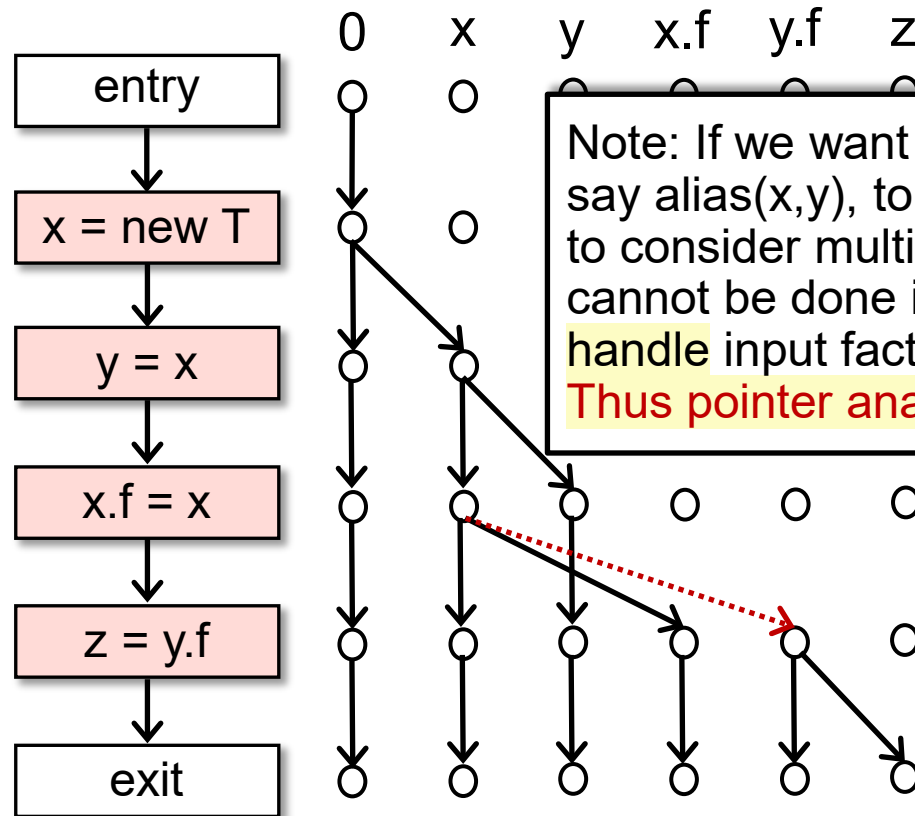# Understanding the Distributivity of IFDS

- Pointer Analysis



z and y.f should have pointed to object [new T]. However, flow function's input data facts lack of the alias information, alias(x,y), alias(x.f,y.f), and we need alias information to produce correct outputs.

# Understanding the Distributivity of IFDS

- Pointer Analysis



Note: If we want to obtain alias information in IFDS, say alias(x,y), to produce correct outputs, we need to consider multiple input data facts, x and y, which cannot be done in *standard* IFDS as flow functions handle input facts independently (one fact per time). Thus pointer analysis is non-distributive.

z and y.f should have pointed to object [new T]. However, flow function's input data facts lack of the alias information, alias(x,y), alias(x.f,y.f), and we need alias information to produce correct outputs.

# Contents

1. Feasible and Realizable Paths

2. CFL-Reachability

3. Overview of IFDS

4. Supergraph and Flow Functions

5. Exploded Supergraph and Tabulation Algorithm

6. Understanding the Distributivity of IFDS

# The X You Need To Understand in This Lecture

- Understand CFL-Reachability

- Understand the basic idea of IFDS

- Understand what problems can be solved by IFDS

注意注意!
划重点了!

感谢陪伴，期待你们的未来

祝同学们在新的一年有新的收获

再见

# 软件分析

南京大学

计算机科学与技术系

程序设计语言与
静态分析研究组

李樾 谭添