

Program Analysis and Verification

0368-4479

Noam Rinetzky

Lecture 11: Shape Analysis + Interprocedural Analysis

Slides credit: Roman Manevich, Mooly Sagiv, Eran Yahav

3-Value logic based shape analysis

Sequential Stack

```
void push (int v) {  
    Node *x = malloc(sizeof(Node));  
    x->d = v;  
    x->n = Top;  
    Top = x;  
}
```

```
int pop() {  
    if (Top == NULL) return EMPTY;  
    Node *s = Top->n;  
    int r = Top->d;  
    Top = s;  
    return r;  
}
```

Want to Verify

No Null Dereference

Underlying list remains acyclic after each operation

Shape Analysis via 3-valued Logic

1) Abstraction

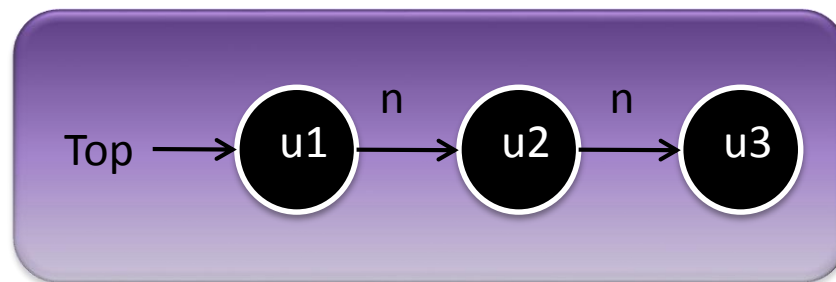
- 3-valued logical structure
- canonical abstraction

2) Transformers

- via logical formulae
- soundness by construction
 - embedding theorem, [SRW02]

Concrete State

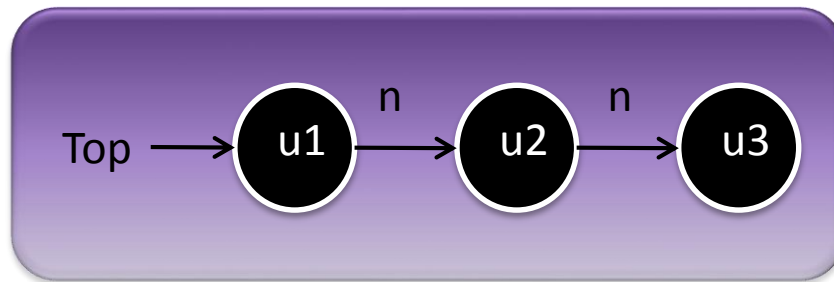
- represent a concrete state as a two-valued logical structure
 - Individuals = heap allocated objects
 - Unary predicates = object properties
 - Binary predicates = relations
- parametric vocabulary



(storeless, no heap addresses)

Concrete State

- $S = \langle U, \iota \rangle$ over a vocabulary P
- U – universe
- ι - interpretation, mapping each predicate from p to its truth value in S



- $U = \{ u1, u2, u3 \}$
- $P = \{ Top, n \}$

$\iota(n)(u1, u2) = 1, \iota(n)(u1, u3) = 0, \iota(n)(u2, u1) = 0, \dots$ ■

$\iota(Top)(u1) = 1, \iota(Top)(u2) = 0, \iota(Top)(u3) = 0$ ■

Formulae for Observing Properties

```
void push (int v) {
```

```
  Node *x =  
    malloc(sizeof(Node));
```

$\exists w: x(w)$

```
x->d = v;
```

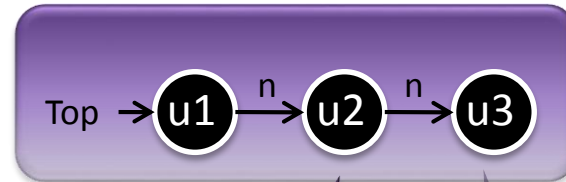
$\exists w: x(w)$

```
x->n = Top;
```

```
  Top = x;
```

```
}  $\neg \exists v1, v2: n(v1, v2) \wedge n^*(v2, v1)$ 
```

```
 $\neg \exists v1, v2: n(v1, v2) \wedge Top(v2)$ 
```



Top != null
 $\exists w: Top(w)$ **1**

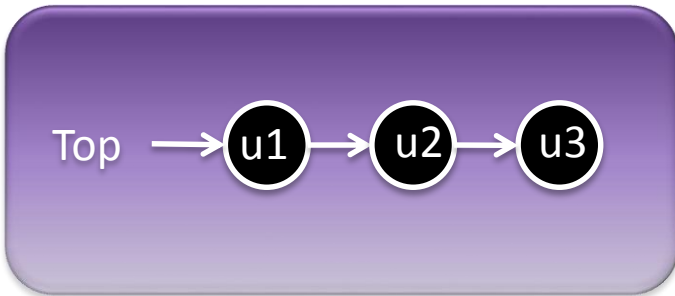
No node precedes Top
 $\neg \exists v1, v2: n(v1, v2) \wedge Top(v2)$ **1**

No Cycles
 $\neg \exists v1, v2: n(v1, v2) \wedge n^*(v2, v1)$ **1**

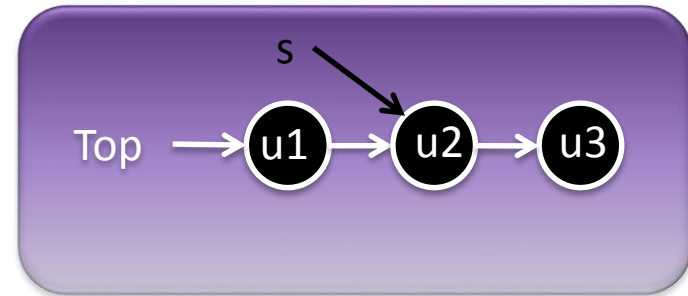
Concrete Interpretation Rules

Statement	Update formula
$x = \text{NULL}$	$x'(v) = 0$
$x = \text{malloc}()$	$x'(v) = \text{IsNew}(v)$
$x = y$	$x'(v) = y(v)$
$x = y \rightarrow \text{next}$	$x'(v) = \exists w: y(w) \wedge n(w, v)$
$x \rightarrow \text{next} = y$	$n'(v, w) = (\neg x(v) \wedge n(v, w)) \vee (x(v) \wedge y(w))$

Example: $s = Top \rightarrow n$



$$s'(v) = \exists v1: Top(v1) \wedge n(v1, v)$$



Top	
u1	1
u2	0
u3	0

n	u1	u2	u3
u1	0	1	0
u2	0	0	1
u3	0	0	0

Top	
u1	1
u2	0
u3	0

n	u1	u2	u3
u1	0	1	0
u2	0	0	1
u3	0	0	0

s	
u1	0
u2	0
u3	0

s	
u1	0
u2	1
u3	0

Collecting Semantics

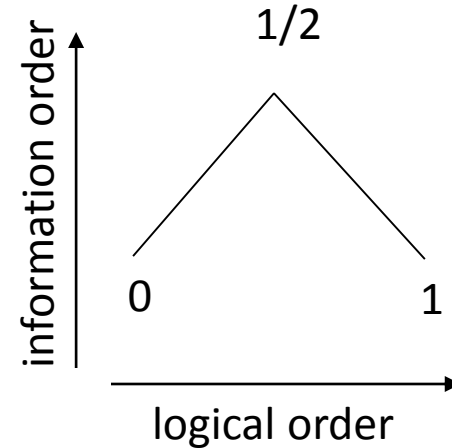
$$\text{CSS}[v] = \begin{cases} \{ \langle \emptyset, \emptyset \rangle \} & \text{if } v = \text{entry} \\ \bigcup \{ \llbracket \text{st}(w) \rrbracket(S) \mid S \in \text{CSS}[w] \} \cup \\ \quad (w,v) \in E(G), \\ \quad w \in \text{Assignments}(G) \\ \bigcup \{ S \mid S \in \text{CSS}[w] \} \cup \\ \quad (w,v) \in E(G), \\ \quad w \in \text{Skip}(G) \\ \bigcup \{ S \mid S \in \text{CSS}[w] \text{ and } S \models \text{cond}(w) \} \cup \\ \quad (w,v) \in \text{True-Branched}(G) \\ \bigcup \{ S \mid S \in \text{CSS}[w] \text{ and } S \models \neg \text{cond}(w) \} \\ \quad (w,v) \in \text{False-Branched}(G) \end{cases} \quad \text{otherwise}$$

Collecting Semantics

- At every program point – a **potentially infinite** set of two-valued logical structures
- Representing (at least) all possible heaps that can arise at the program point
- Next step:
find a bounded abstract representation

3-Valued Logic

- 1 = true
- 0 = false
- $1/2$ = unknown
- A join semi-lattice, $0 \sqcup 1 = 1/2$



3-Valued Logical Structures

- A set of individuals (nodes) U
- Relation meaning
 - Interpretation of relation symbols in \mathcal{P}
 - $\iota(p^0)(\cdot) \rightarrow \{0, 1, 1/2\}$
 - $\iota(p^1)(v) \rightarrow \{0, 1, 1/2\}$
 - $\iota(p^2)(u, v) \rightarrow \{0, 1, 1/2\}$
- A join semi-lattice: $0 \sqcup 1 = 1/2$

Boolean Connectives [Kleene]

\wedge	0	1/2	1
0	0	0	0
1/2	0	1/2	1/2
1	0	1/2	1

\vee	0	1/2	1
0	0	1/2	1
1/2	1/2	1/2	1
1	1	1	1

Property Space

- $3\text{-struct}[P]$ = the set of 3-valued logical structures over a vocabulary (set of predicates) P
- Abstract domain
 - $\wp(3\text{-Struct}[P])$
 - \sqsubseteq is \subseteq

Embedding Order

- Given two structures $S = \langle U, \iota \rangle$, $S' = \langle U', \iota' \rangle$ and an onto function $f : U \rightarrow U'$ mapping individuals in U to individuals in U'
- We say that f embeds S in S' (denoted by $S \sqsubseteq^f S'$) if
 - for every predicate symbol $p \in P$ of arity k : $u_1, \dots, u_k \in U$,
 $\iota(p)(u_1, \dots, u_k) \sqsubseteq \iota'(p)(f(u_1), \dots, f(u_k))$
 - and for all $u' \in U'$
 $(|\{u \mid f(u) = u'\}| > 1) \sqsubseteq \iota'(sm)(u')$
- We say that S can be embedded in S' (denoted by $S \sqsubseteq S'$) if there exists a function f such that $S \sqsubseteq^f S'$

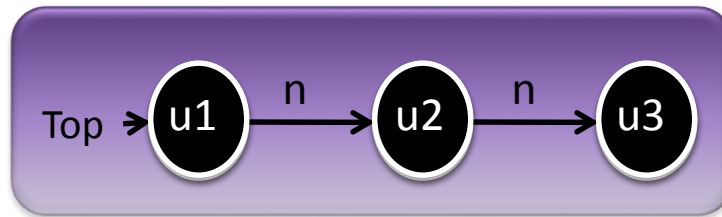
Tight Embedding

- $S' = \langle U', \iota' \rangle$ is a tight embedding of $S = \langle U, \iota \rangle$ with respect to a function f if:
 - S' does not lose unnecessary information

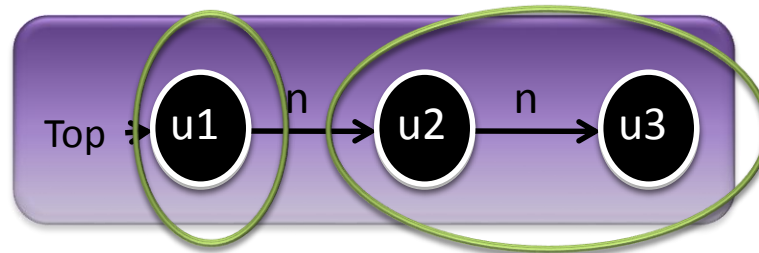
$$\iota'(u'_1, \dots, u'_k) = \sqcup \{ \iota(u_1, \dots, u_k) \mid f(u_1) = u'_1, \dots, f(u_k) = u'_k \}$$

- One way to get tight embedding is canonical abstraction

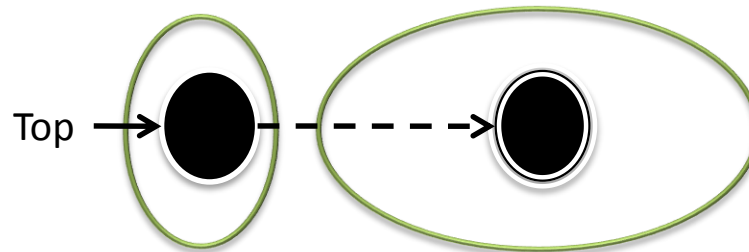
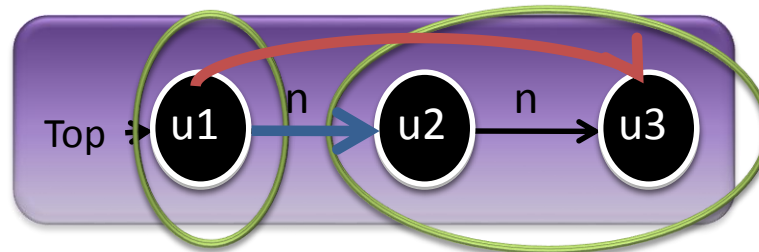
Canonical Abstraction



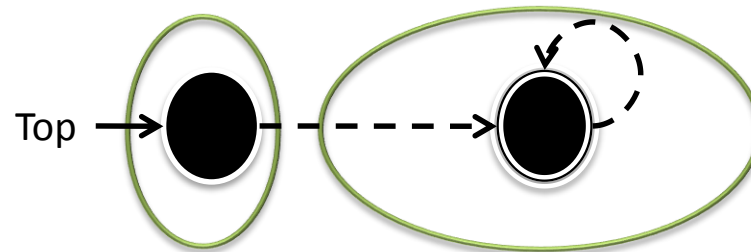
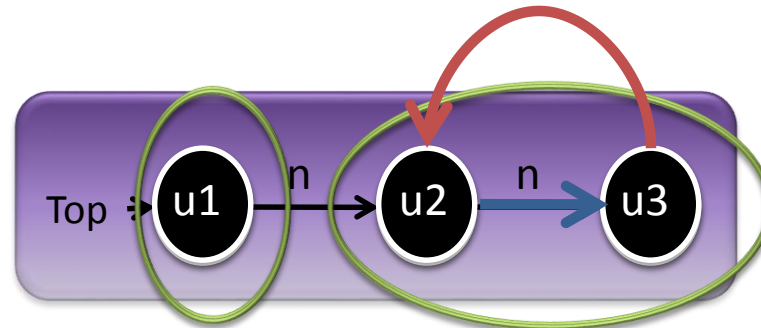
Canonical Abstraction



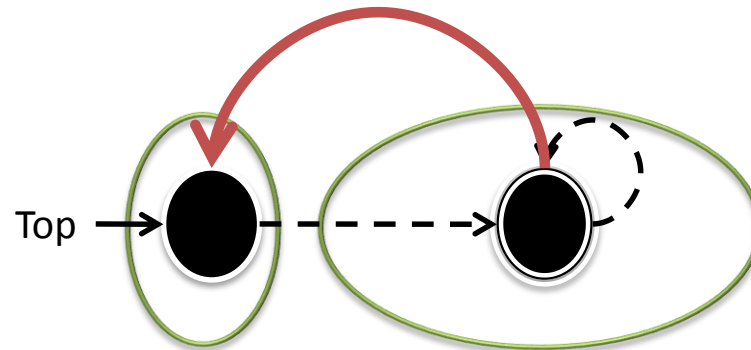
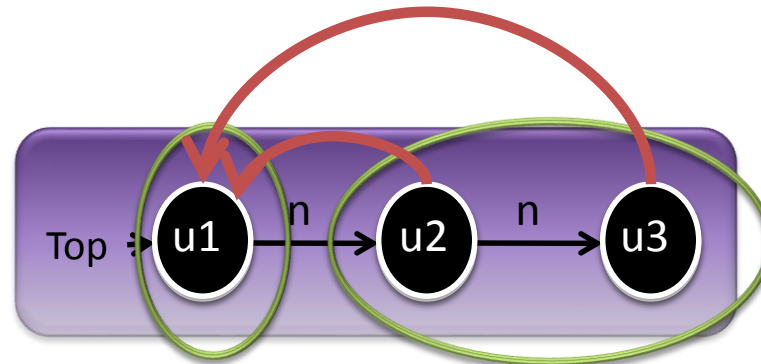
Canonical Abstraction



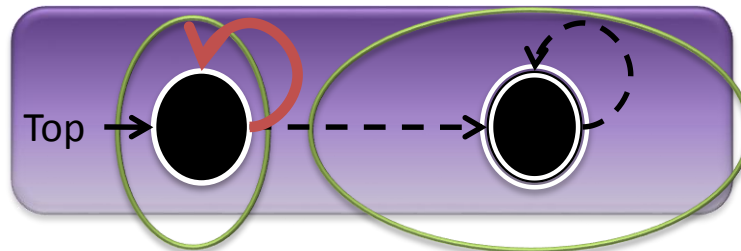
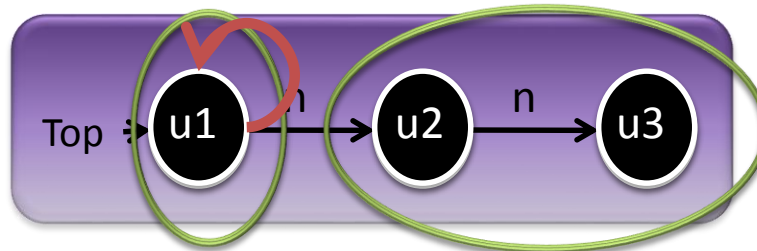
Canonical Abstraction



Canonical Abstraction



Canonical Abstraction



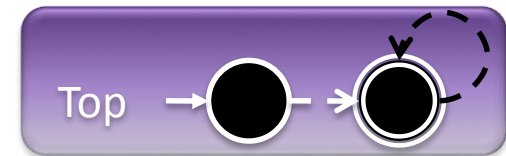
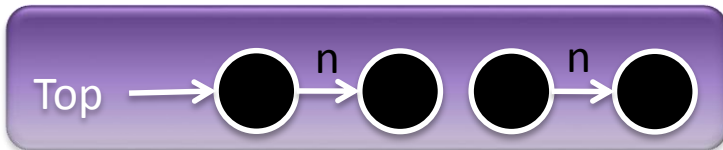
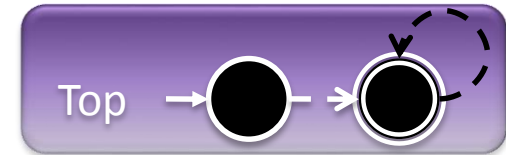
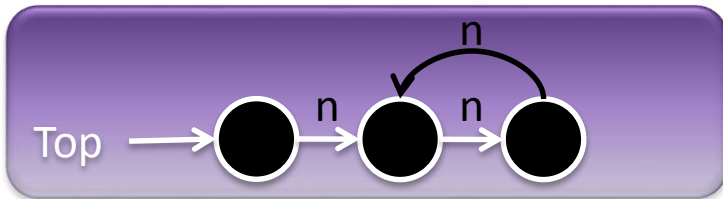
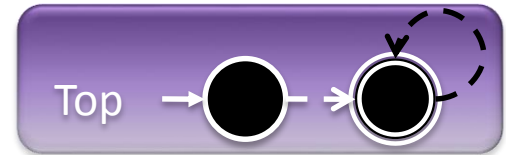
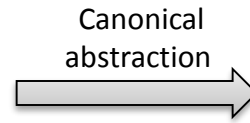
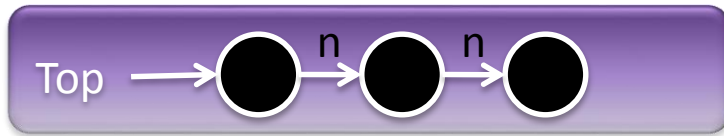
Canonical Abstraction (β)

- Merge all nodes with the **same unary predicate values** into a single summary node
- Join predicate values

$$\iota'(u'_1, \dots, u'_k) = \sqcup \{ \iota(u_1, \dots, u_k) \mid f(u_1) = u'_1, \dots, f(u_k) = u'_k \}$$

- Converts a state of **arbitrary** size into a 3-valued abstract state of **bounded** size
- $\alpha(C) = \sqcup \{ \beta(c) \mid c \in C \}$

Information Loss

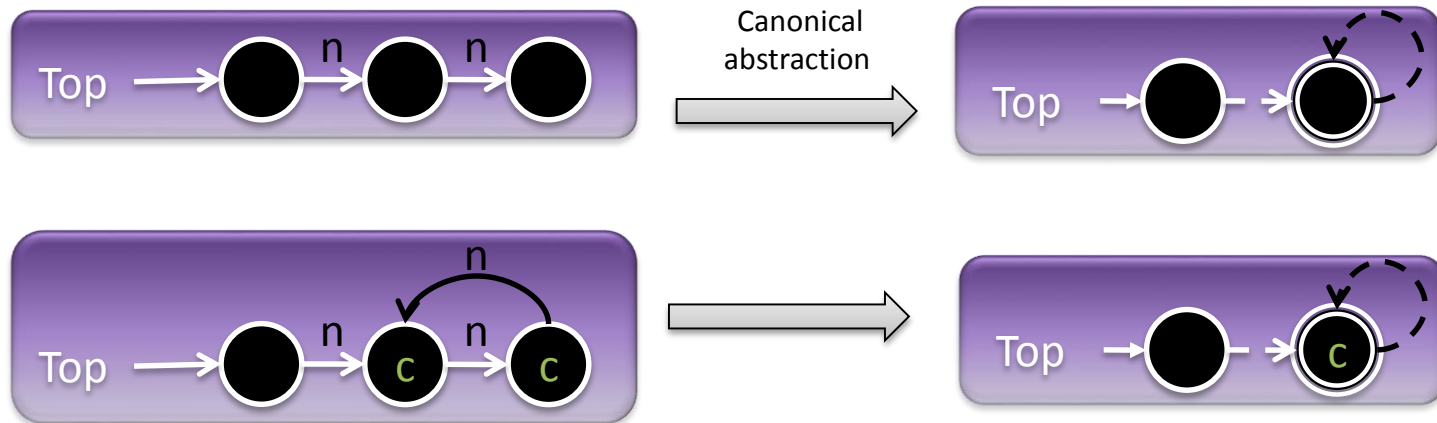


Instrumentation Predicates

- Record additional derived information via predicates

$$r_x(v) = \exists v1: x(v1) \wedge n^*(v1, v)$$

$$c(v) = \exists v1: n(v1, v) \wedge n^*(v, v1)$$



Embedding Theorem: Conservatively Observing Properties



No Cycles

$$\neg \exists v_1, v_2: n(v_1, v_2) \wedge n^*(v_2, v_1) \quad \mathbf{1/2}$$

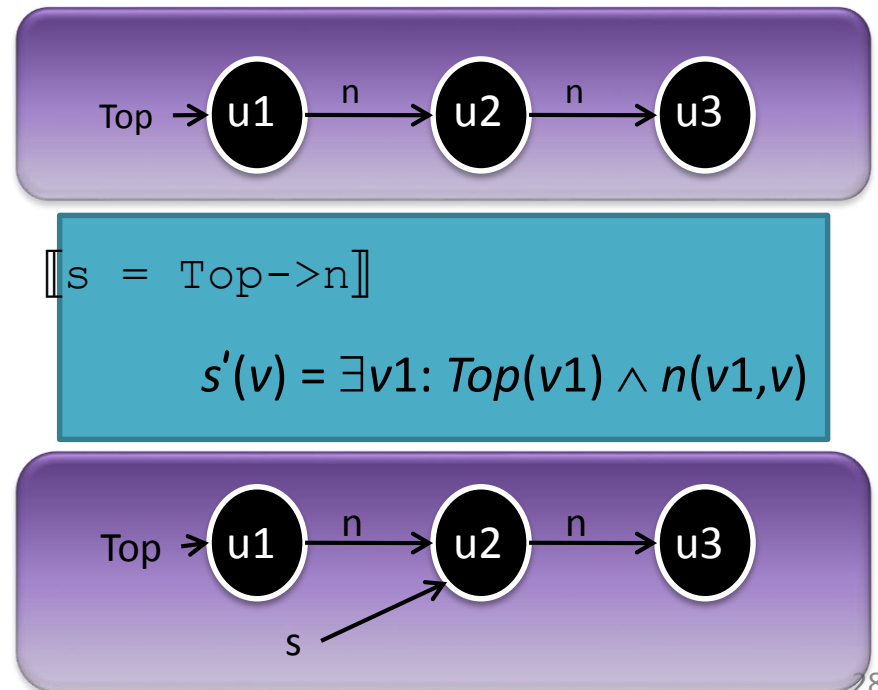
No cycles (derived)

$$\forall v: \neg c(v) \quad \mathbf{1}$$

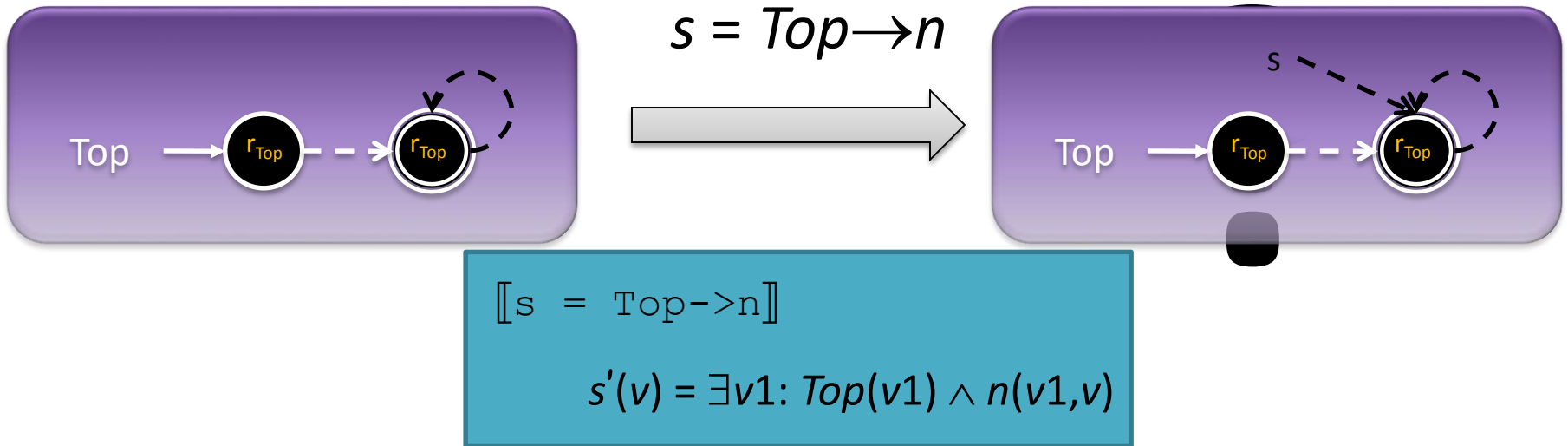
Operational Semantics

```
void push (int v) {  
  Node *x = malloc(sizeof(Node));  
  x->d = v;  
  x->n = Top;  
  Top = x;  
}
```

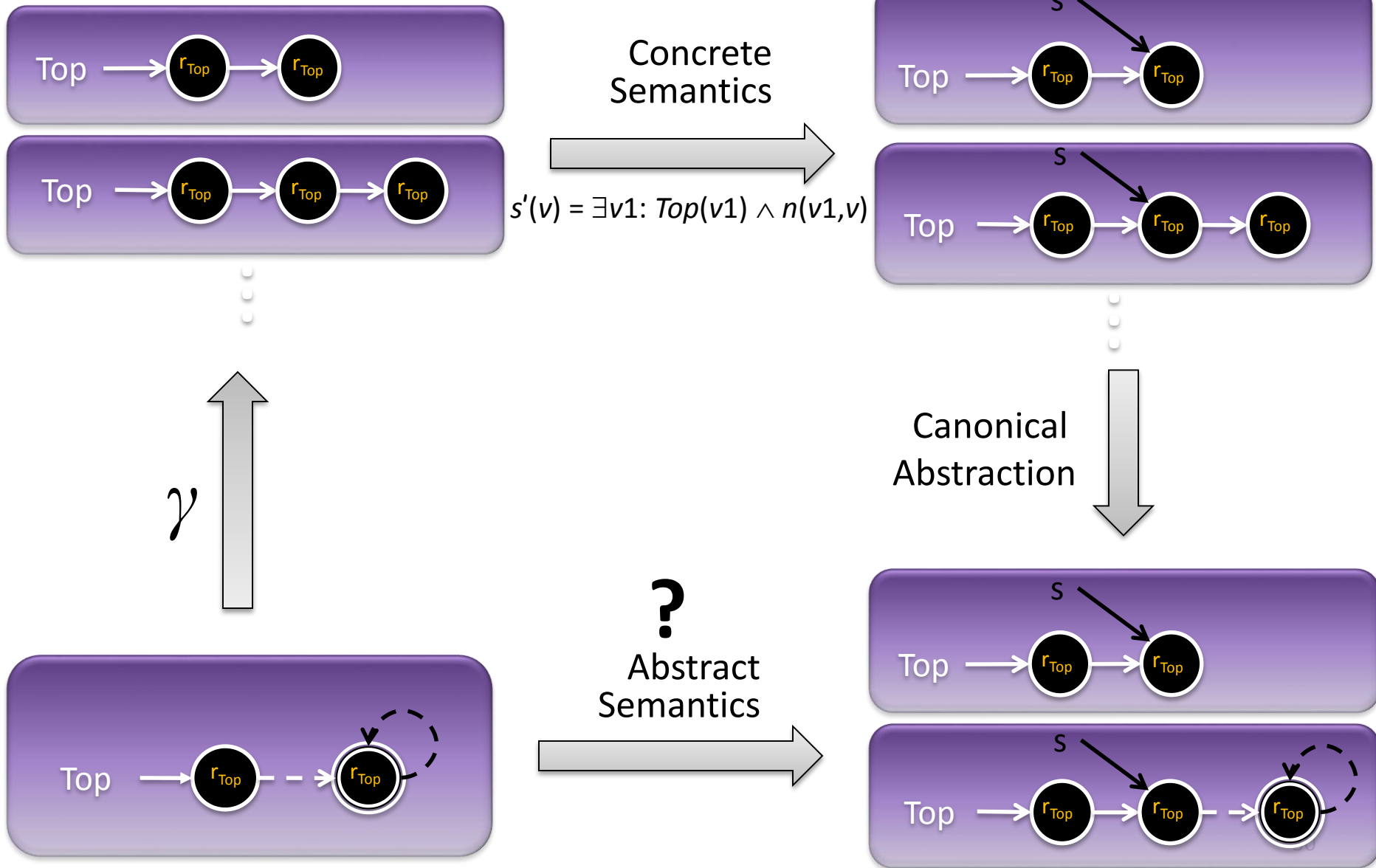
```
int pop () {  
  if (Top == NULL) return EMPTY;  
  Node *s = Top->n;  
  int r = Top->d;  
  Top = s;  
  return r;  
}
```



Abstract Semantics

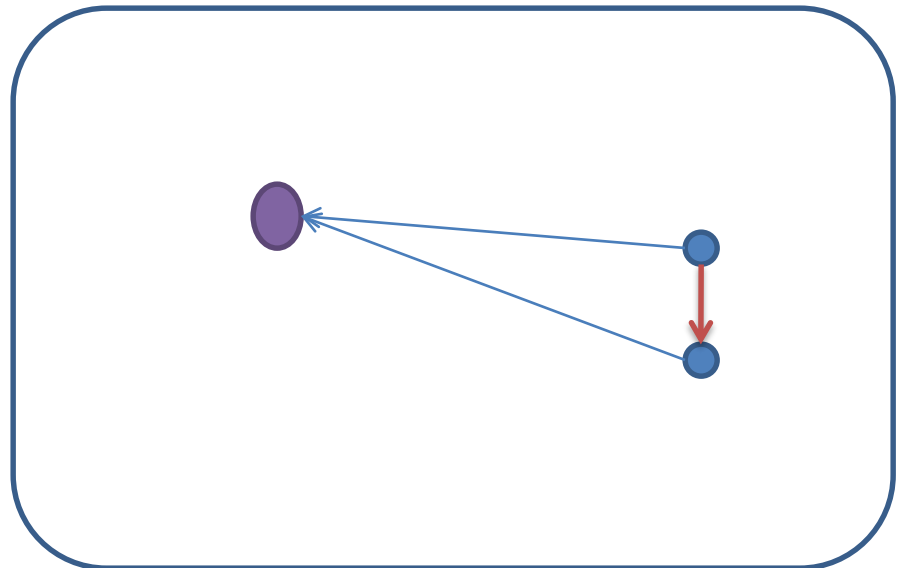


Best Transformer ($s = Top \rightarrow n$)



Semantic Reduction

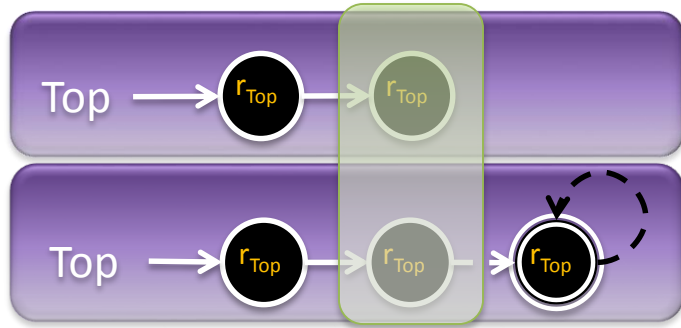
- Improve the precision of the analysis by recovering properties of the program semantics
- A Galois connection (C, α, γ, A)
- An operation $op:A \rightarrow A$ is a **semantic reduction** when
 - $\forall l \in L_2 \ op(l) \sqsubseteq l$ and
 - $\gamma(op(l)) = \gamma(l)$



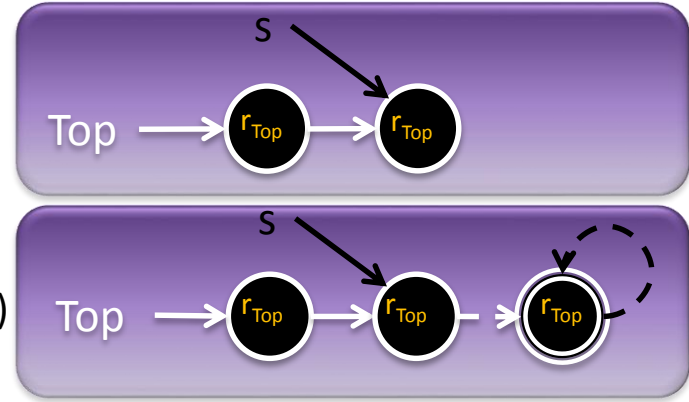
The Focus Operation

- Focus: $\text{Formula} \rightarrow (\wp(\mathcal{3}\text{-Struct}) \hookrightarrow \wp(\mathcal{3}\text{-Struct}))$
- Generalizes materialization
- For every formula φ
 - $\text{Focus}(\varphi)(X)$ yields structure in which φ evaluates to a definite values in all assignments
 - Only maximal in terms of embedding
 - $\text{Focus}(\varphi)$ is a semantic reduction
 - But $\text{Focus}(\varphi)(X)$ may be undefined for some X

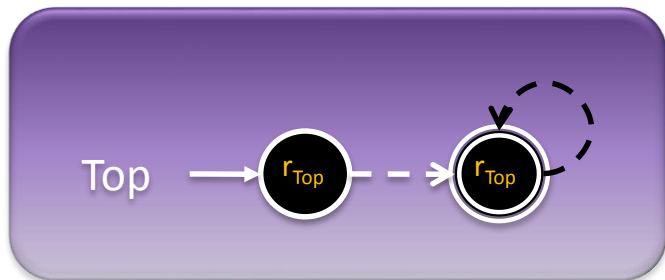
Partial Concretization Based on Transformer ($s=Top \rightarrow n$)



Abstract Semantics
 $s'(v) = \exists v1: Top(v1) \wedge n(v1, v)$

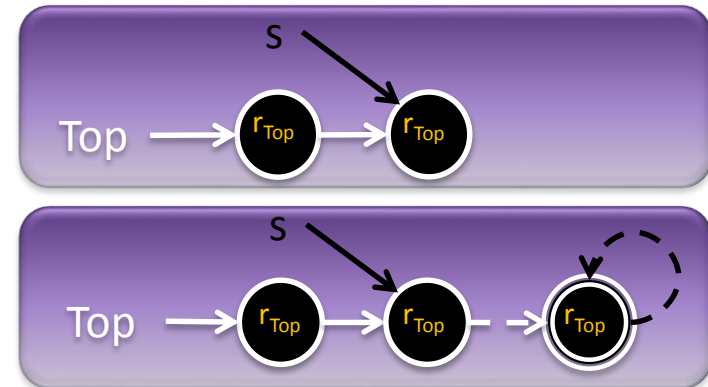


Partial Concretization



Abstract Semantics

Canonical Abstraction



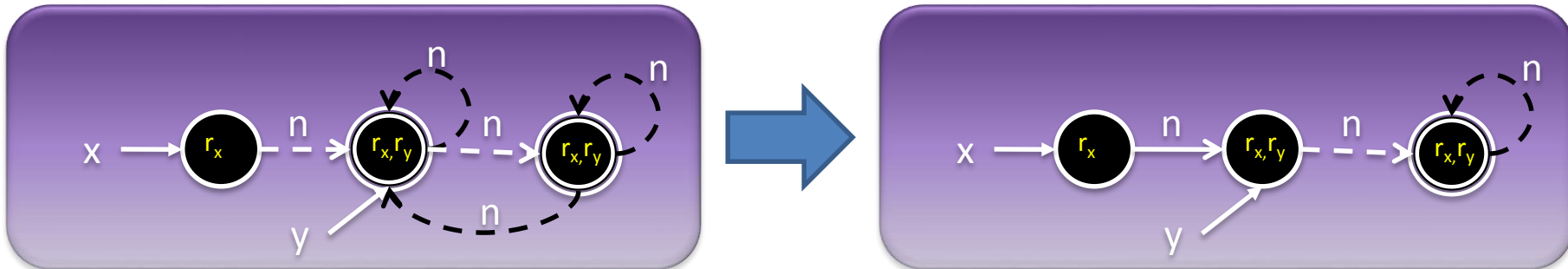
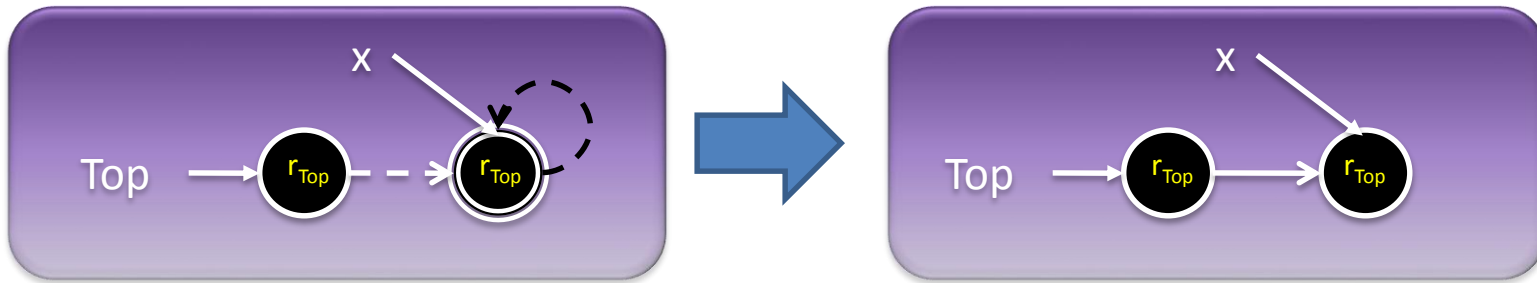
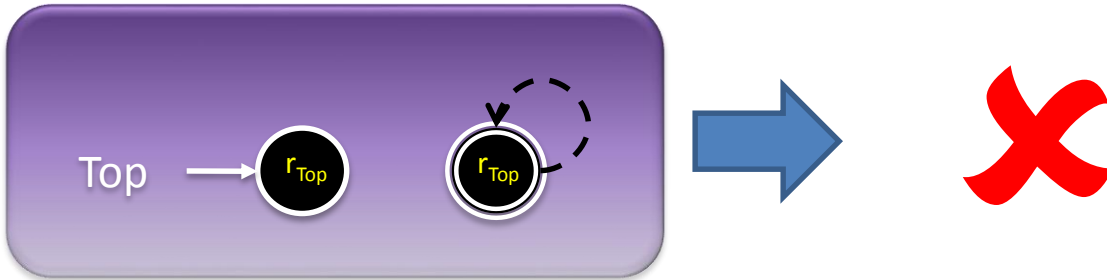
Partial Concretization

- Locally refine the abstract domain per statement
- Soundness is immediate
- Employed in other shape analysis algorithms
[Distefano et.al., TACAS'06, Evan et.al., SAS'07, POPL'08]

The Coercion Principle

- Another Semantic Reduction
- Can be applied after Focus or after Update or both
- Increase precision by exploiting some structural properties possessed by all stores (Global invariants)
- Structural properties captured by **constraints**
- Apply a constraint solver

Apply Constraint Solver



Sources of Constraints

- Properties of the operational semantics
- Domain specific knowledge
 - Instrumentation predicates
- User supplied

Example Constraints

$$x(v1) \wedge x(v2) \rightarrow eq(v1, v2)$$

$$n(v, v1) \wedge n(v, v2) \rightarrow eq(v1, v2)$$

$$n(v1, v) \wedge n(v2, v) \wedge \neg eq(v1, v2) \leftrightarrow is(v)$$

$$n^*(v1, v2) \leftrightarrow t[n](v1, v2)$$

Abstract Transformers: Summary

- Kleene evaluation yields sound solution
- Focus is a statement-specific partial concretization
- Coerce applies global constraints

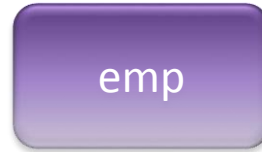
Abstract Semantics

$$\begin{aligned}
 SS[v] = & \left\{ \langle \emptyset, \emptyset \rangle \right\} && \text{if } v = \text{entry} \\
 & \bigcup \left\{ t_embed(\text{coerce}(\llbracket st(w) \rrbracket_3(\text{focus}_{F(w)}(SS[w]))) \right) \bigcup \\
 & (w,v) \in E(G), \\
 & w \in \text{Assignments}(G) \\
 & \bigcup \{ S \mid S \in SS[w] \} \bigcup && \text{otherwise} \\
 & (w,v) \in E(G), \\
 & w \in \text{Skip}(G) \\
 & \bigcup \{ t_embed(S) \mid S \in \text{coerce}(\llbracket st(w) \rrbracket_3(\text{focus}_{F(w)}(SS[w]))) \\
 & \text{and } S \models_3 \text{cond}(w) \} \bigcup \\
 & (w,v) \in \text{True-Branches}(G) \\
 & \bigcup \{ t_embed(S) \mid S \in \text{coerce}(\llbracket st(w) \rrbracket_3(\text{focus}_{F(w)}(SS[w]))) \\
 & \text{and } S \models_3 \neg \text{cond}(w) \} \bigcup \\
 & (w,v) \in \text{False-Branches}(G)
 \end{aligned}$$

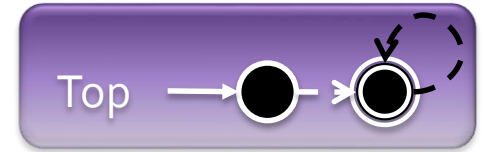
Recap

- Abstraction
 - canonical abstraction
 - recording derived information
- Transformers
 - partial concretization (focus)
 - constraint solver (coerce)
 - sound information extraction

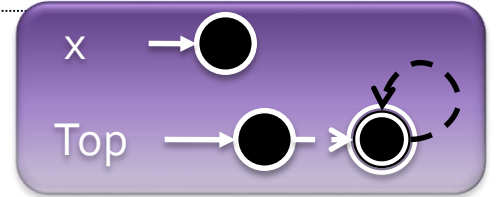
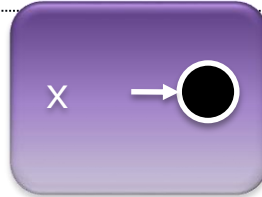
Stack Push



...



```
void push (int v) {
  Node *x =
    alloc(sizeof(Node));
```

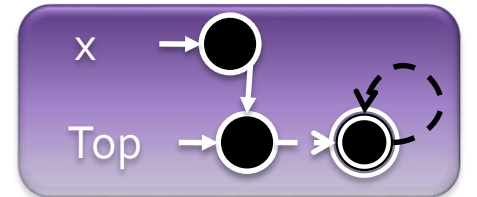
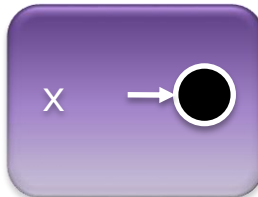


$$\exists v: x(v)$$

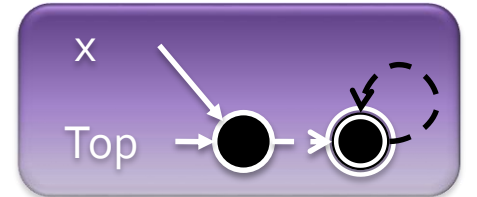
```
x->d = v;
```

$$\exists v: x(v)$$

```
x->n = Top;
```

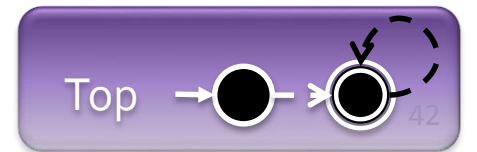
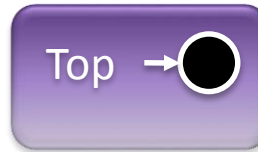


```
Top = x;
```



$$\neg \exists v_1, v_2: n(v_1, v_2) \wedge Top(v_2)$$

$$\forall v: \neg c(v)$$



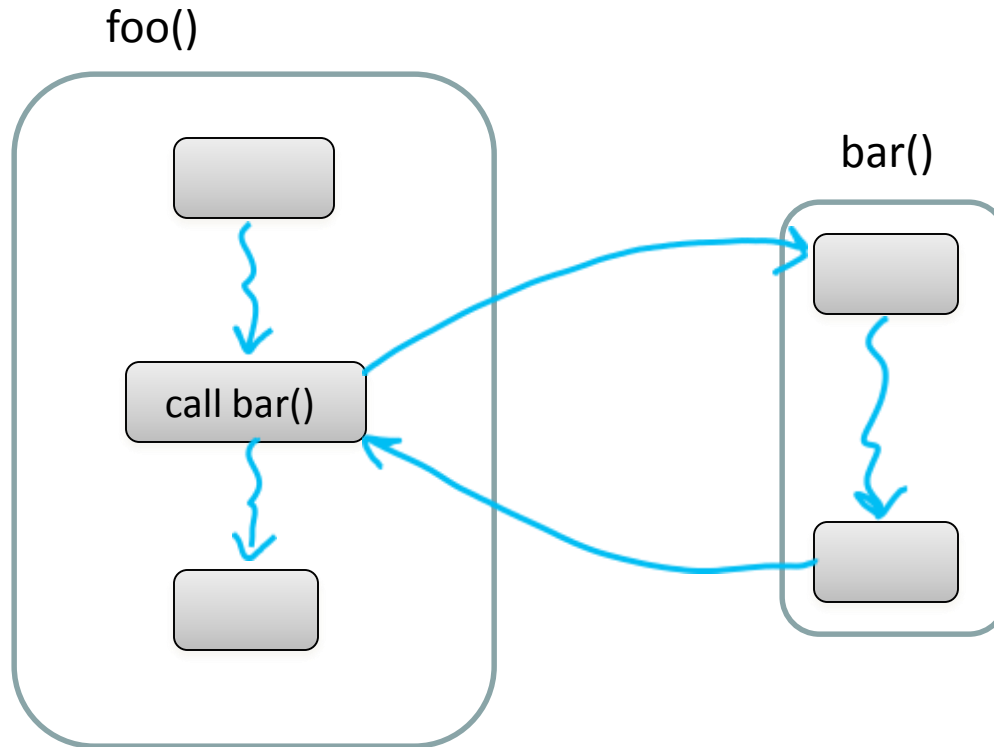
What about procedures?

Procedural program

```
void main() {  
    int x;  
    x = p(7);  
    x = p(9);  
}
```

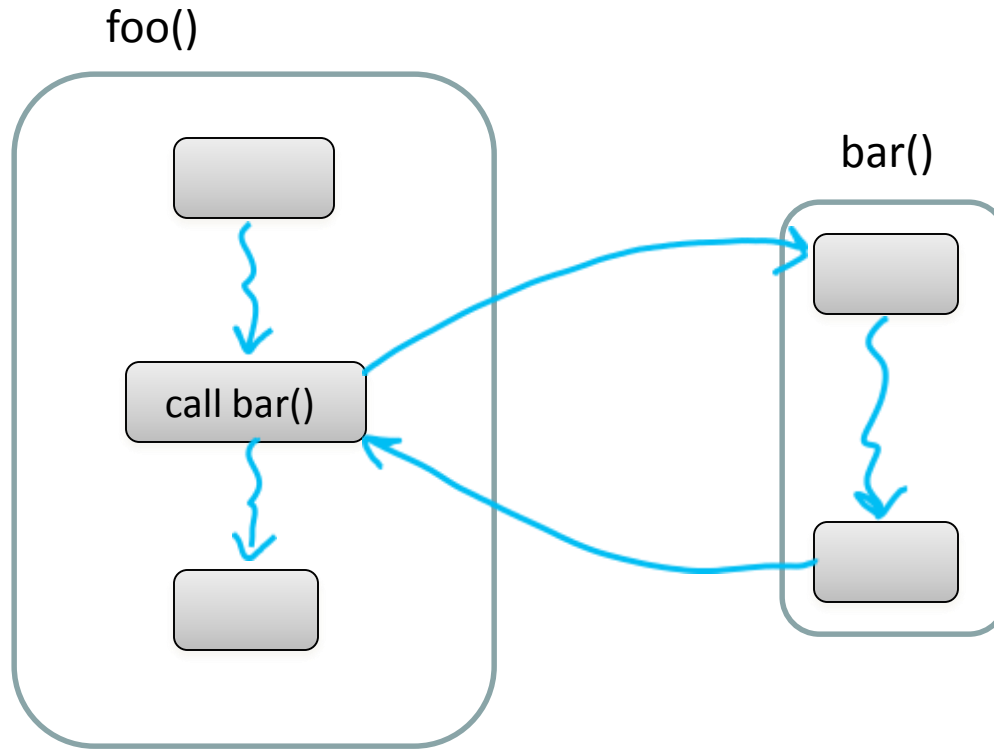
```
int p(int a) {  
    return a + 1;  
}
```

Effect of procedures



The effect of calling a procedure is the effect of executing its body

Interprocedural Analysis

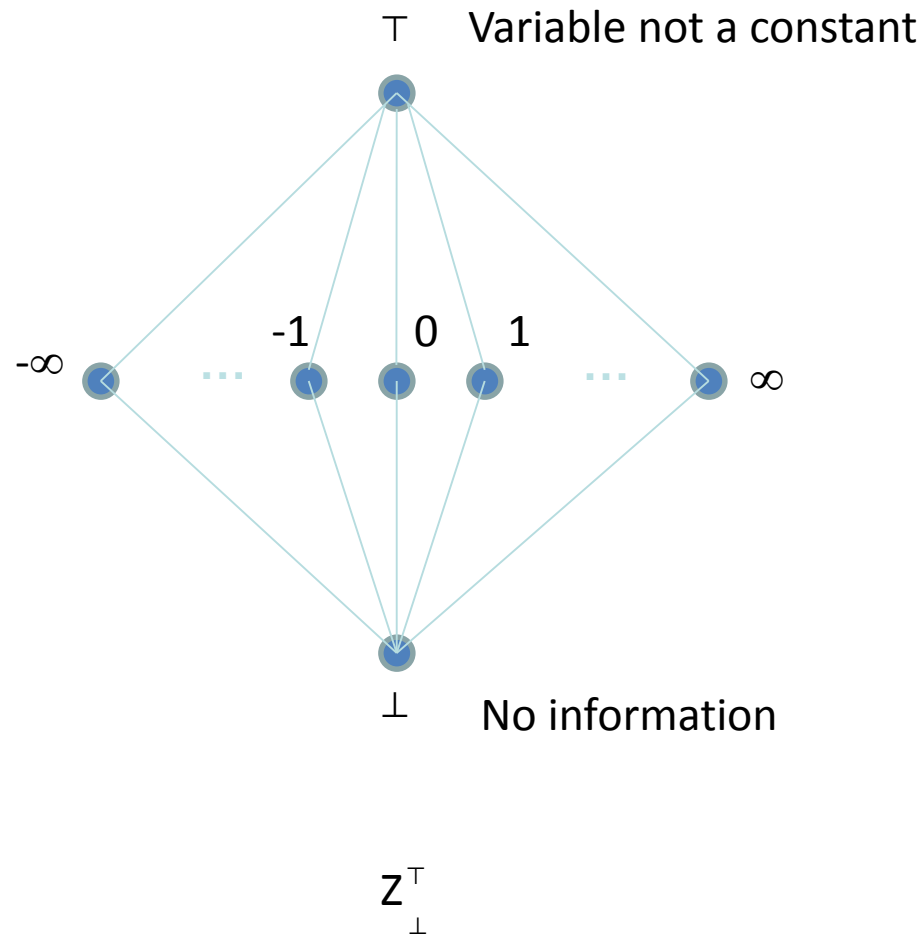


goal: compute the abstract effect of calling a procedure

Reduction to intraprocedural analysis

- Procedure inlining
- Naive solution: call-as-goto

Reminder: Constant Propagation



Reminder: Constant Propagation

- $L = (\text{Var} \rightarrow Z, \sqsubseteq)$
- $\sigma_1 \sqsubseteq \sigma_2$ iff $\forall x: \sigma_1(x) \sqsubseteq' \sigma_2(x)$
 - \sqsubseteq' ordering in the Z lattice
- Examples:
 - $[x \mapsto \perp, y \mapsto 42, z \mapsto \perp] \sqsubseteq [x \mapsto \perp, y \mapsto 42, z \mapsto 73]$
 - $[x \mapsto \perp, y \mapsto 42, z \mapsto 73] \sqsubseteq [x \mapsto \perp, y \mapsto 42, z \mapsto \top]$

Reminder: Constant Propagation

- Conservative Solution
 - Every detected constant is indeed constant
 - But may fail to identify some constants
 - Every potential impact is identified
 - Superfluous impacts

Procedure Inlining

```
void main() {  
    int x;  
    x = p(7);  
    x = p(9);  
}
```

```
int p(int a) {  
    return a + 1;  
}
```

Procedure Inlining

```
void main() {  
    int x;  
    x = p(7);  
    x = p(9);  
}
```

```
int p(int a) {  
    return a + 1;  
}
```

```
void main() {  
    int a, x, ret;  
    [a ↦ ⊥, x ↦ ⊥, ret ↦ ⊥]  
    a = 7; ret = a+1; x = ret;  
    [a ↦ 7, x ↦ 8, ret ↦ 8]  
    a = 9; ret = a+1; x = ret;  
    [a ↦ 9, x ↦ 10, ret ↦ 10]  
}
```

Procedure Inlining

- Pros
 - Simple
- Cons
 - Does not handle recursion
 - Exponential blow up
 - Reanalyzing the body of procedures

```
p1 {                p2 {                p3{
  call p2           call p3
  ...
  call p2           ...
}                  call p3
}                  }
```

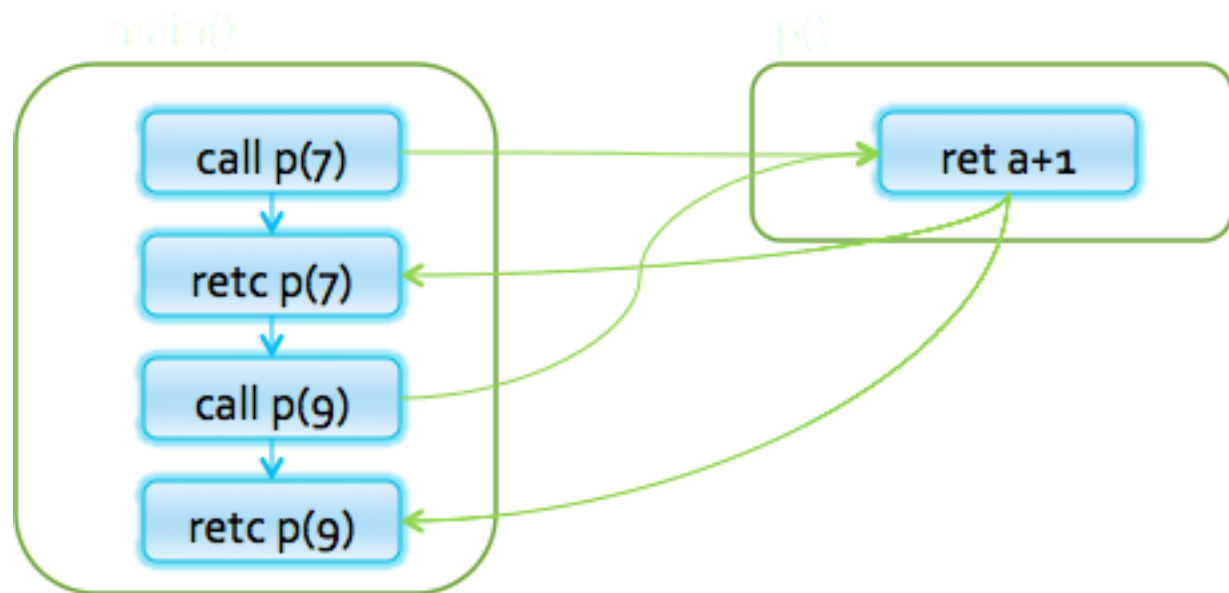
A Naive Interprocedural solution

- Treat procedure calls as gotos

Simple Example

```
void main() {  
    int x ;  
    → x = p(7);  
    x = p(9) ;  
}
```

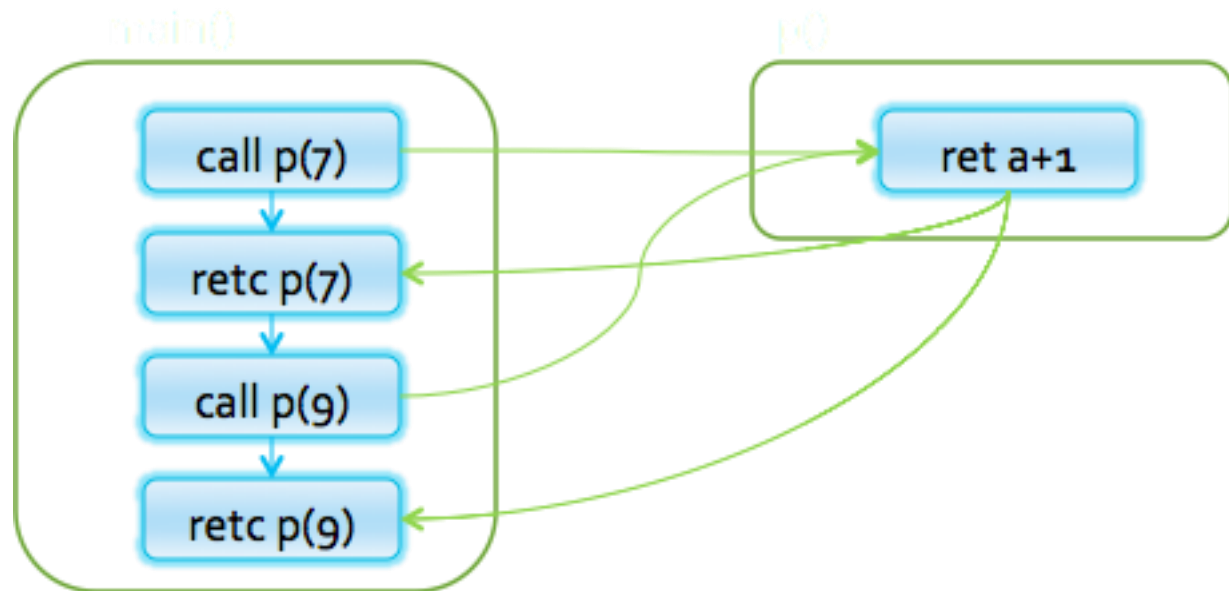
```
int p(int a) {  
    return a + 1;  
}
```



Simple Example

```
void main() {  
    int x ;  
    x = p(7);  
    x = p(9) ;  
}
```

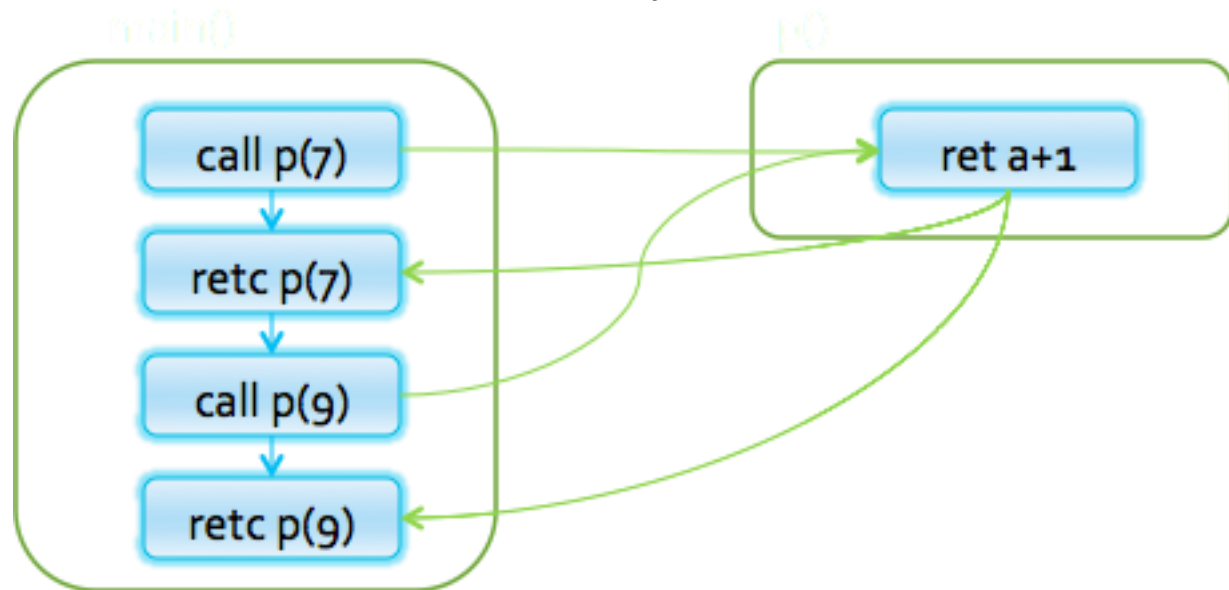
```
→ int p(int a) {  
    [a ↦ 7]  
    return a + 1;  
}
```



Simple Example

```
void main() {  
    int x ;  
    x = p(7);  
    x = p(9) ;  
}
```

```
int p(int a) {  
    [a ↦ 7]  
    → return a + 1;  
    [a ↦ 7, $$ ↦ 8]  
}
```



Simple Example

```
void main() {
```

```
  int x ;
```

```
  x = p(7); ←
```

```
  [x ↦ 8]
```

```
  x = p(9) ; ←
```

```
  [x ↦ 8]
```

```
}
```

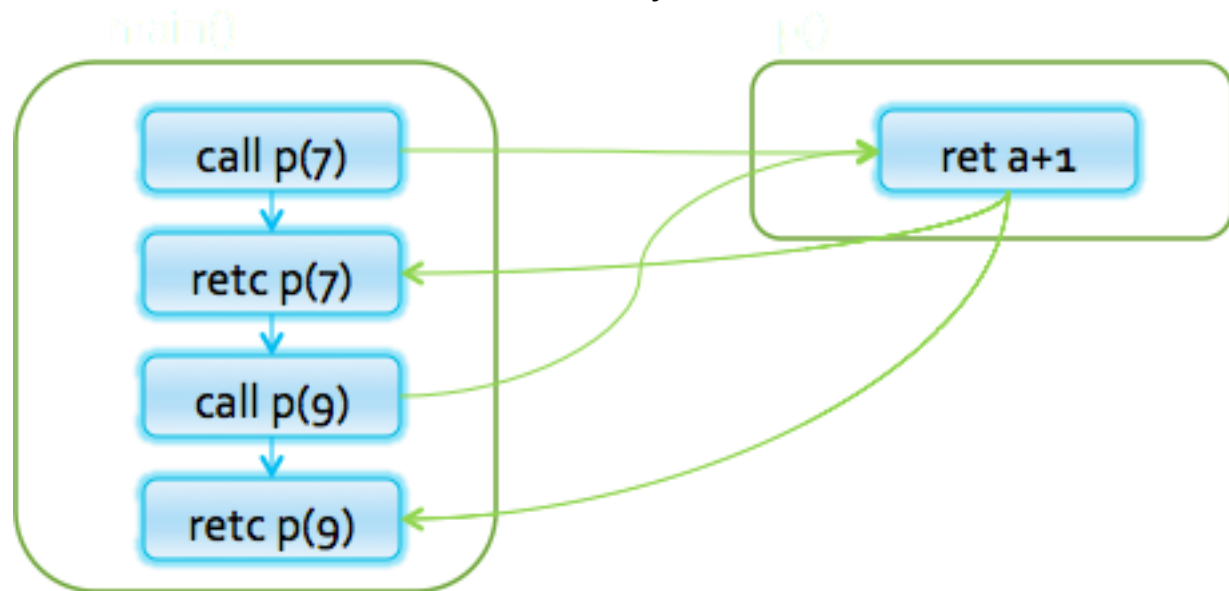
```
int p(int a) {
```

```
  [a ↦ 7]
```

```
  return a + 1;
```

```
  [a ↦ 7, $$ ↦ 8]
```

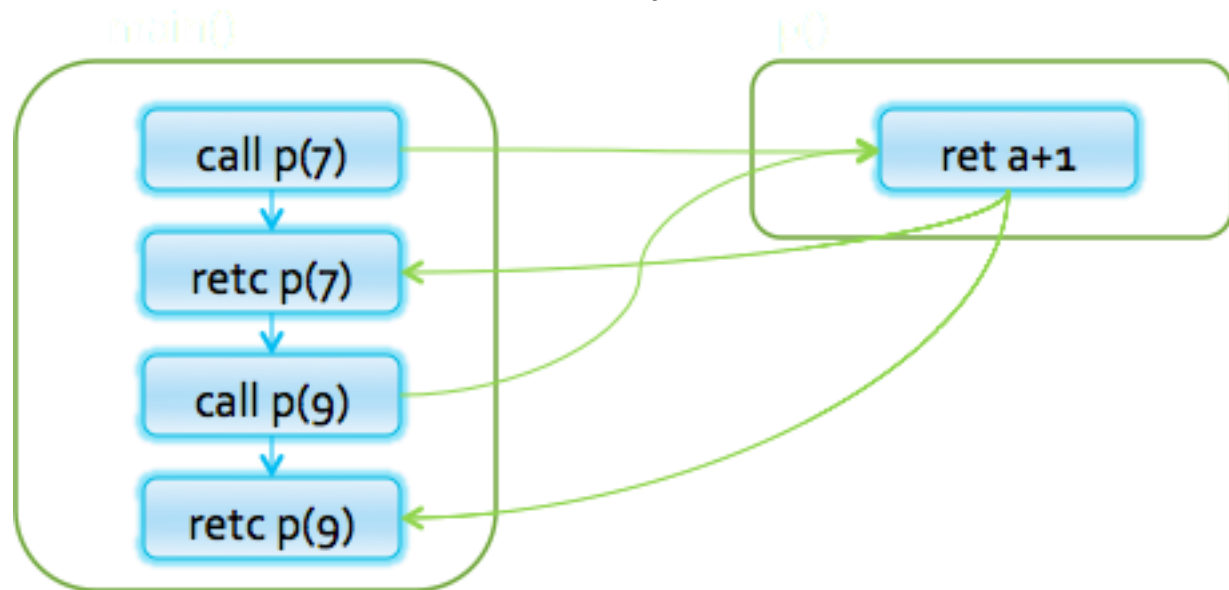
```
}
```



Simple Example

```
void main() {  
    int x ;  
    x = p(7);  
    [x ↦ 8]  
    → x = p(9) ;  
    [x ↦ 8]  
}
```

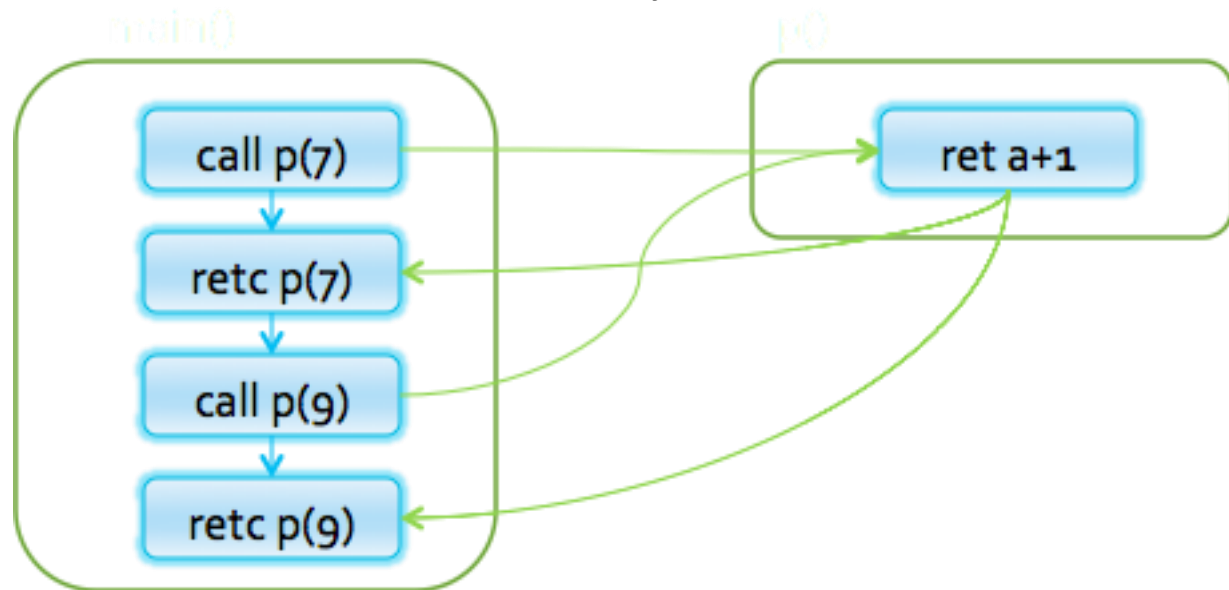
```
int p(int a) {  
    [a ↦ 7]  
    return a + 1;  
    [a ↦ 7, $$ ↦ 8]  
}
```



Simple Example

```
void main() {  
    int x ;  
    x = p(7);  
    [x ↦ 8]  
    → x = p(9) ;  
    [x ↦ 8]  
}
```

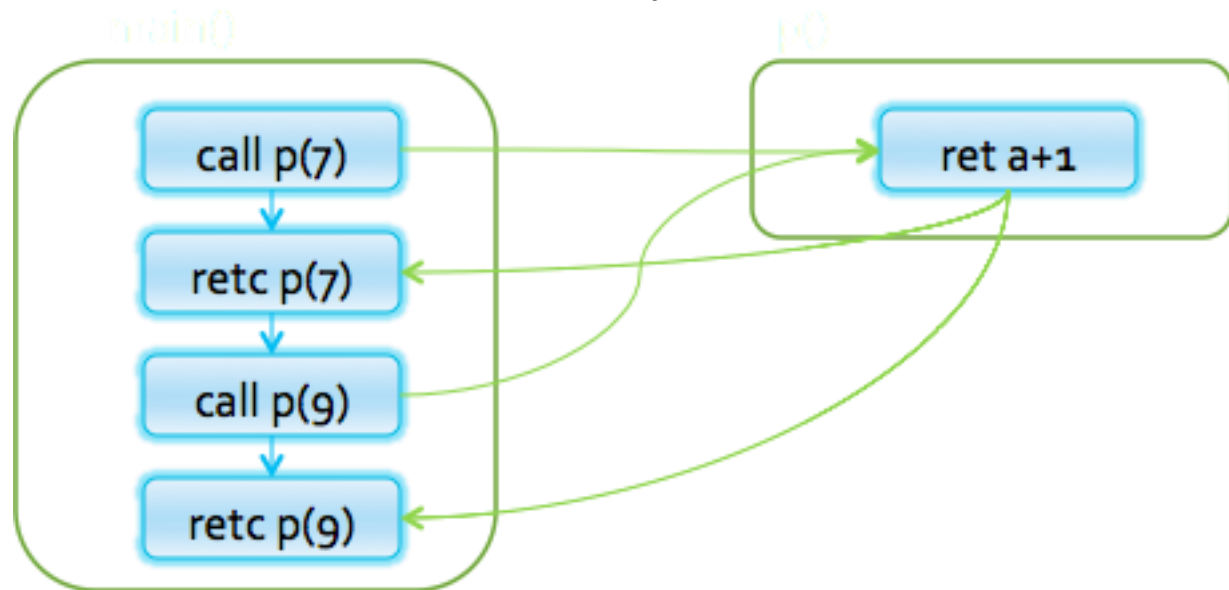
```
→ int p(int a) {  
    [a ↦ 7] [a ↦ 9]  
    return a + 1;  
    [a ↦ 7, $$ ↦ 8]  
}
```



Simple Example

```
void main() {  
    int x ;  
    x = p(7);  
    [x ↦ 8]  
    x = p(9) ;  
    [x ↦ 8]  
}
```

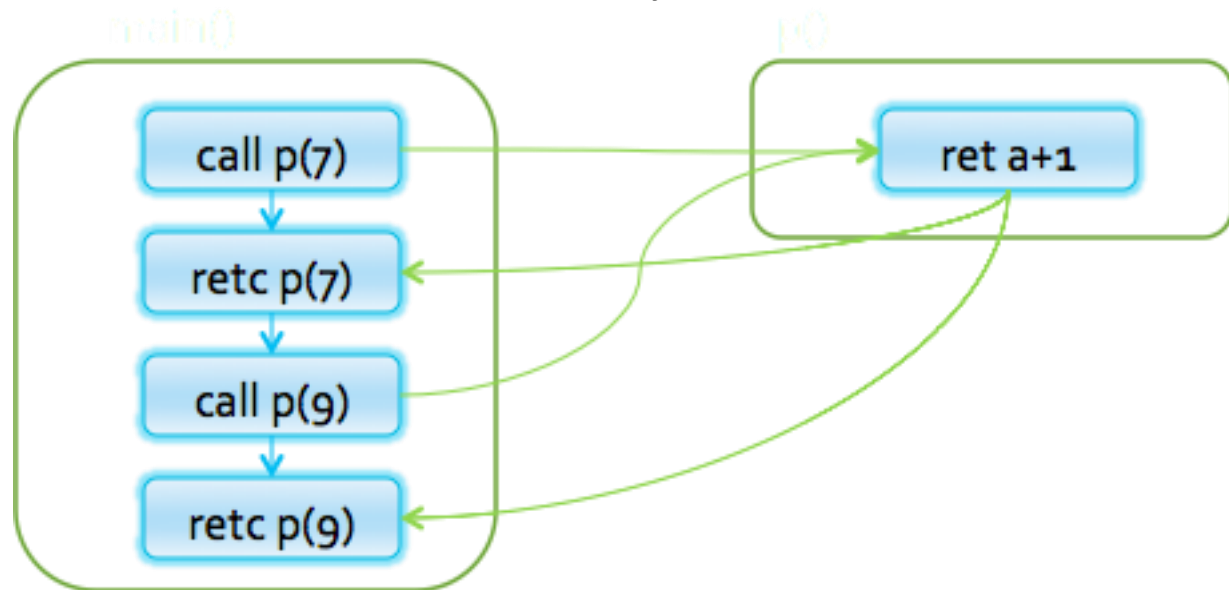
```
→ int p(int a) {  
    [a ↦ τ]  
    return a + 1;  
    [a ↦ 7, $$ ↦ 8]  
}
```



Simple Example

```
void main() {  
  int x ;  
  x = p(7);  
  [x ↦ 8]  
  x = p(9);  
  [x ↦ 8]  
}
```

```
int p(int a) {  
  [a ↦ T]  
  → return a + 1;  
  [a ↦ T, $$ ↦ T]  
}
```



Simple Example

```
void main() {
```

```
  int x ;
```

```
  x = p(7) ; ←
```

```
  [x ↦ T]
```

```
  x = p(9) ; ←
```

```
  [x ↦ T]
```

```
}
```

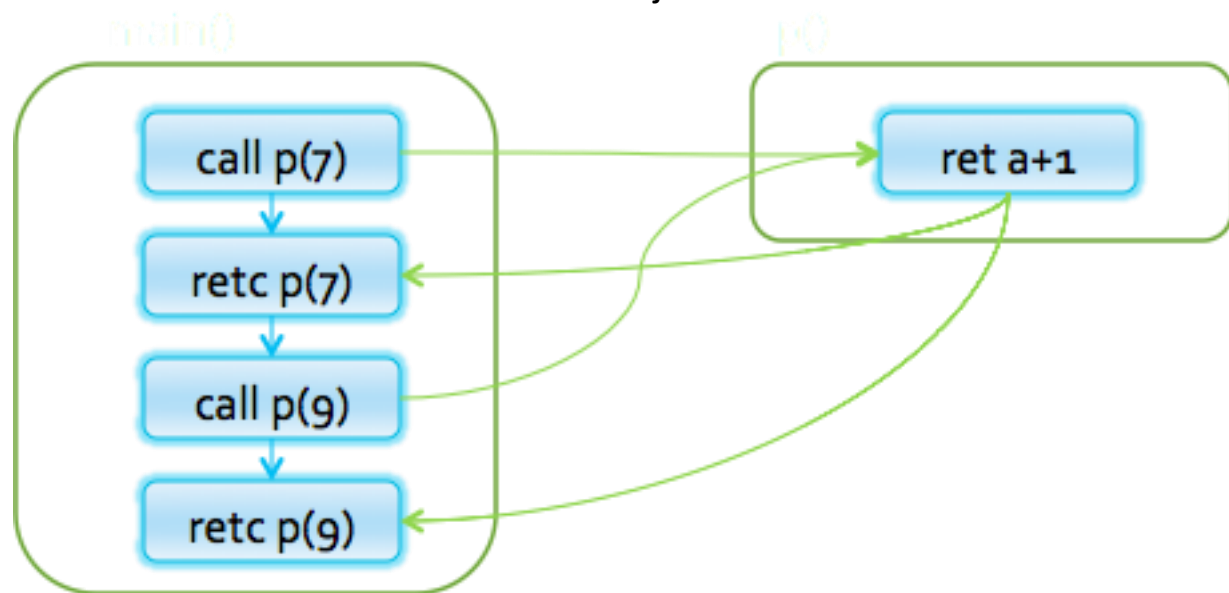
```
int p(int a) {
```

```
  [a ↦ T]
```

```
  return a + 1;
```

```
  [a ↦ T, $$ ↦ T]
```

```
}
```



A Naive Interprocedural solution

- Treat procedure calls as gotos
- Pros:
 - Simple
 - Usually fast
- Cons:
 - Abstract call/return correlations
 - Obtain a conservative solution

analysis by reduction

Call-as-goto

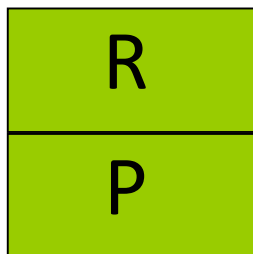
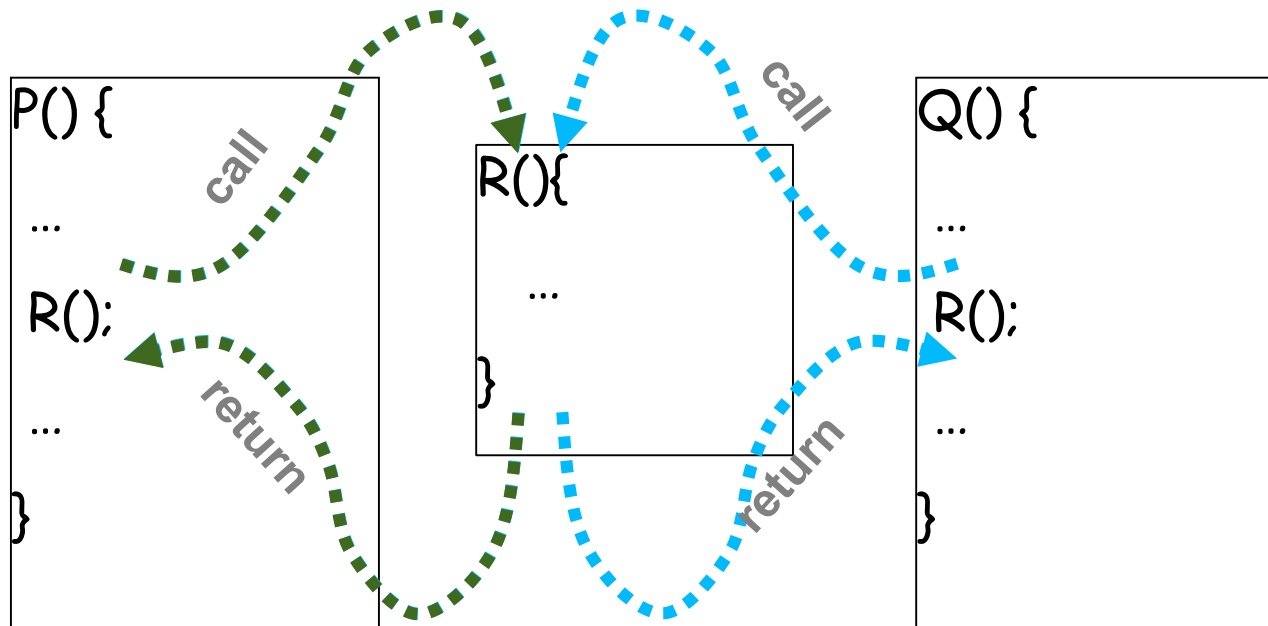
```
void main() {  
    int x;  
    x = p(7);  
    x = p(9);  
}  
  
int p(int a) {  
    return a + 1;  
}
```

Procedure inlining

```
void main() {  
    int a, x, ret;  
    a = 7; ret = a+1; x = ret;  
    a = 9; ret = a+1; x = ret;  
}
```

why was the naive solution less precise?

Stack regime



Guiding light

- Exploit stack regime
 - ➔ Precision
 - ➔ Efficiency



Simplifying Assumptions

- Parameter passed by value
 - No procedure nesting
 - No concurrency
- ✓ Recursion is supported

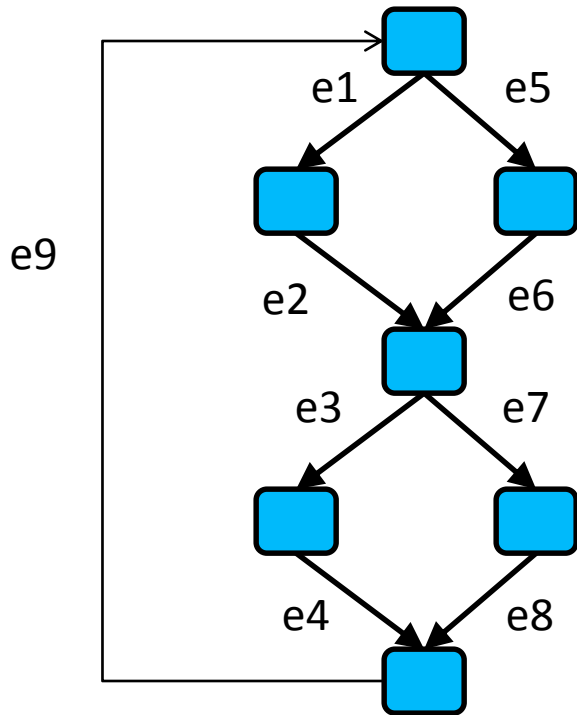
Topics Covered

- ✓ Procedure Inlining
- ✓ The naive approach
 - Valid paths
 - The callstring approach
 - The Functional Approach
 - IFDS: Interprocedural Analysis via Graph Reachability
 - IDE: Beyond graph reachability
- The trivial modular approach

Join-Over-All-Paths (JOP)

- Let $\text{paths}(v)$ denote the potentially infinite set paths from start to v (written as sequences of edges)
- For a sequence of edges $[e_1, e_2, \dots, e_n]$ define $f[e_1, e_2, \dots, e_n]: L \rightarrow L$ by composing the effects of basic blocks
$$f[e_1, e_2, \dots, e_n](l) = f(e_n) (\dots (f(e_2) (f(e_1) (l)) \dots))$$
- $\text{JOP}[v] = \sqcup \{f[e_1, e_2, \dots, e_n](l) \mid [e_1, e_2, \dots, e_n] \in \text{paths}(v)\}$

Join-Over-All-Paths (JOP)



Paths transformers:

$f[e1,e2,e3,e4]$

$f[e1,e2,e7,e8]$

$f[e5,e6,e7,e8]$

$f[e5,e6,e3,e4]$

$f[e1,e2,e3,e4,e9, e1,e2,e3,e4]$

$f[e1,e2,e7,e8,e9, e1,e2,e3,e4,e9,...]$

...

JOP:

$f[e1,e2,e3,e4](\text{initial}) \sqcup$

$f[e1,e2,e7,e8](\text{initial}) \sqcup$

$f[e5,e6,e7,e8](\text{initial}) \sqcup$

$f[e5,e6,e3,e4](\text{initial}) \sqcup \dots$

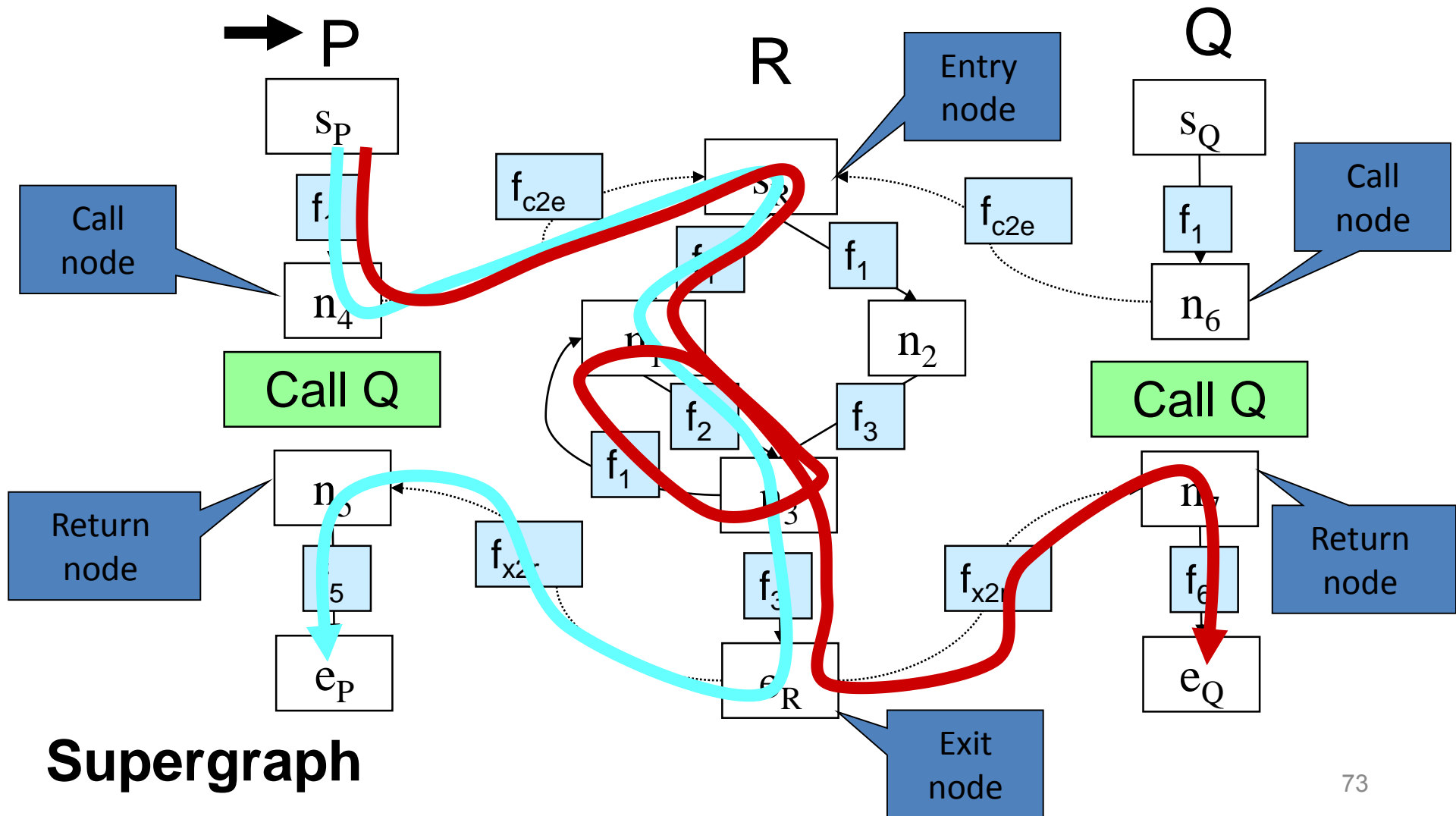
Number of program paths is unbounded due to loops

The lfp computation approximates JOP

- $JOP[v] = \sqcup \{ f [e_1, e_2, \dots, e_n](\perp) \mid [e_1, e_2, \dots, e_n] \in \text{paths}(v) \}$
- $LFP[v] = \sqcup \{ f [e](LFP[v']) \mid e = (v', v) \}$
 $LFP[v_0] = \perp$
- $JOP \sqsubseteq LFP$ - for a monotone function
 - $f(x \sqcup y) \supseteq f(x) \sqcup f(y)$
- $JOP = LFP$ - for a distributive function
 - $f(x \sqcup y) = f(x) \sqcup f(y)$

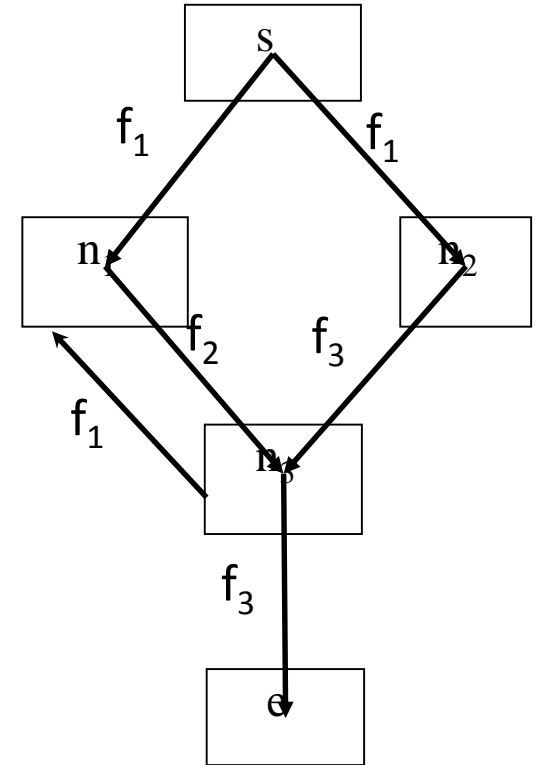
JOP may not be precise enough for interprocedural analysis!

Interprocedural analysis

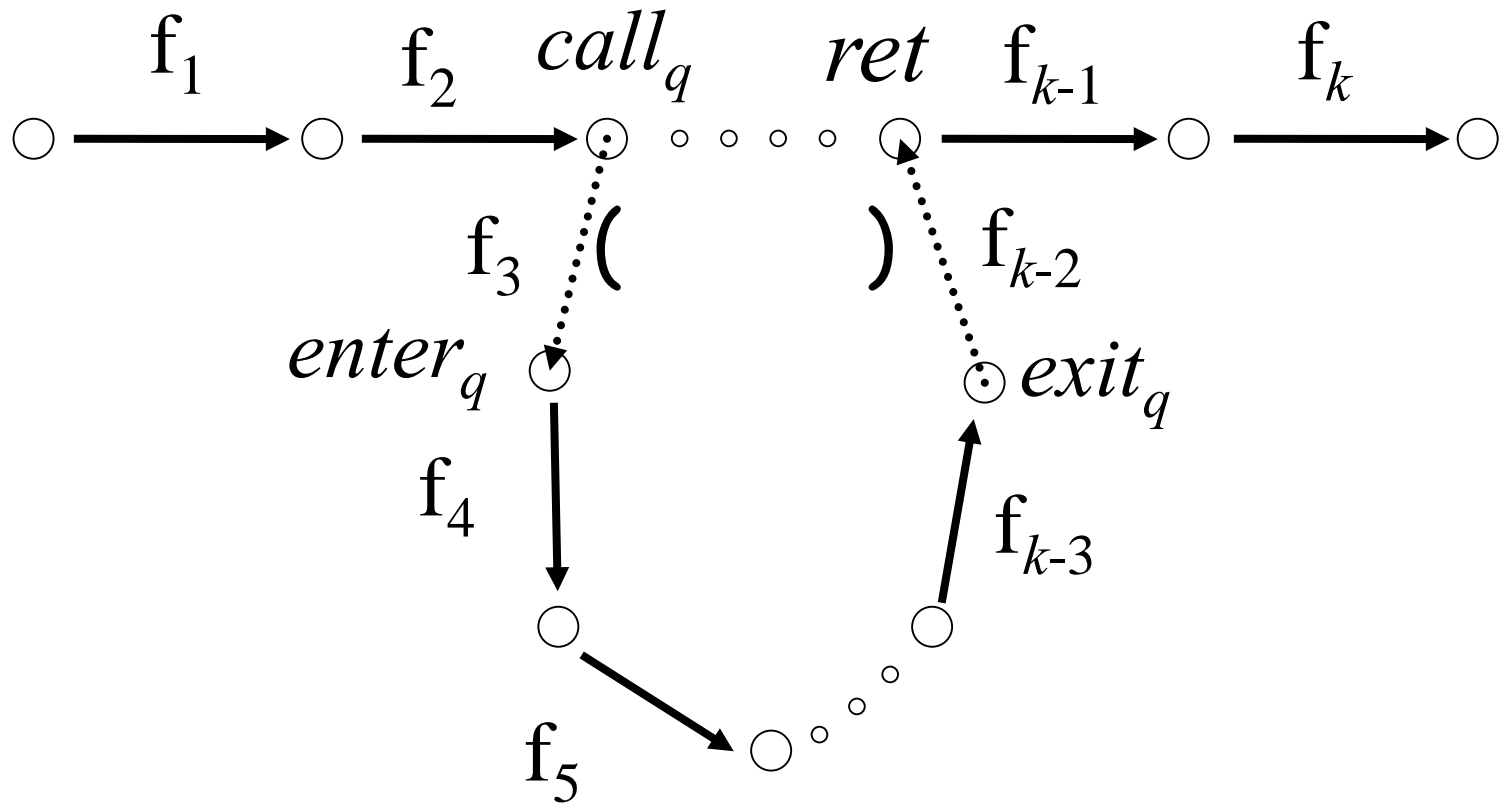


Paths

- **paths(n)** the set of paths from s to n
– $((s, n_1), (n_1, n_3), (n_3, n_1))$



Interprocedural Valid Paths



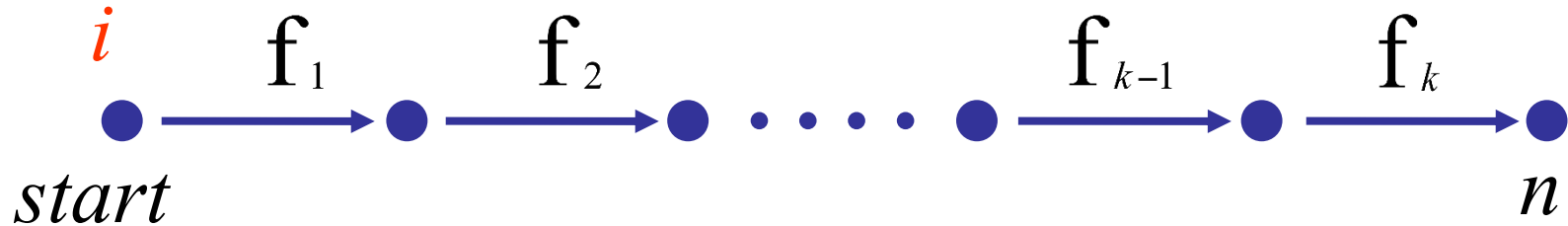
IVP: all paths with matching calls and returns •

And prefixes –

Interprocedural Valid Paths

- **IVP** set of paths
 - Start at program entry
- Only considers matching calls and returns
 - aka, **valid**
- Can be defined via context free grammar
 - $\text{matched} ::= \text{matched } (; \text{matched }) ; \mid \varepsilon$
 - $\text{valid} ::= \text{valid } (; \text{matched } \mid \text{matched}$
 - *paths* can be defined by a regular expression

Join Over All Paths (JOP)



$$\llbracket f_k \circ \dots \circ f_1 \rrbracket \in L \rightarrow L$$

- $JOP[v] = \sqcup \{ \llbracket [e_1, e_2, \dots, e_n] \rrbracket(i) \mid (e_1, \dots, e_n) \in \text{paths}(v) \}$
- $JOP \sqsubseteq LFP$
 - Sometimes $JOP = LFP$
 - precise up to “**symbolic execution**”
 - Distributive problem

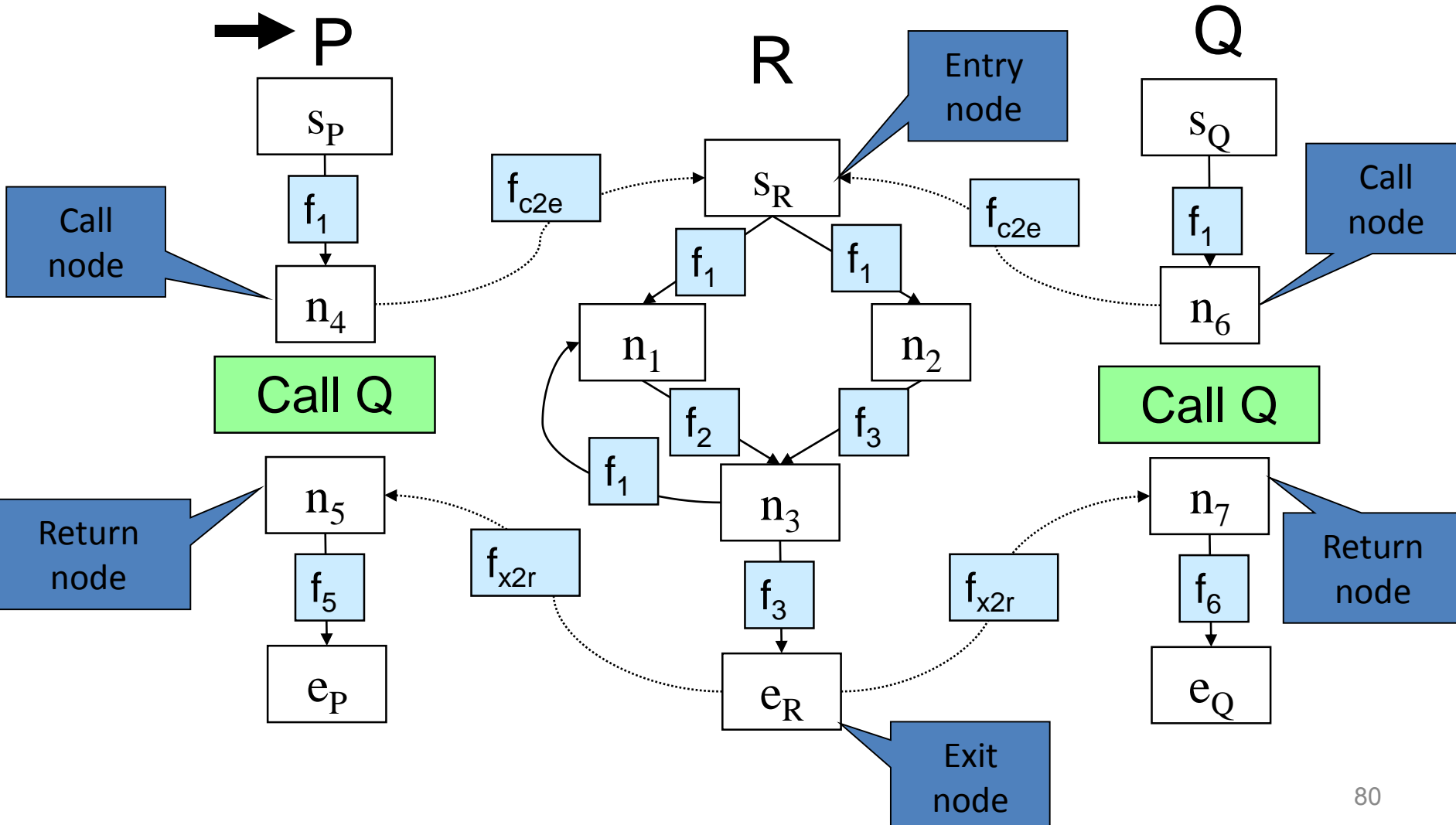
The Join-Over-Valid-Paths (JVP)

- $vpaths(n)$ all valid paths from program start to n
- $JVP[n] = \sqcup \{ [[e_1, e_2, \dots, e]](1) \mid (e_1, e_2, \dots, e) \in vpaths(n) \}$
- $JVP \sqsubseteq JOP$
 - In some cases the JVP can be computed
 - (Distributive problem)

The Call String Approach

- The data flow value is associated with sequences of calls (call string)
- Use Chaotic iterations over the supergraph

supergraph



Simple Example

```
void main() {  
    int x ;  
    → c1: x = p(7);  
    c2: x = p(9) ;  
}  
  
int p(int a) {  
    return a + 1;  
}
```

Simple Example

```
void main() {  
    int x ;  
    c1: x = p(7);  
    c2: x = p(9) ;  
  
}
```

```
→ int p(int a) {  
    c1: [a ↦ 7]  
    return a + 1;  
  
}
```

Simple Example

```
void main() {  
    int x ;  
    c1: x = p(7);  
    c2: x = p(9) ;  
  
}
```

```
int p(int a) {  
    c1: [a ↦ 7]  
    → return a + 1;  
    c1:[a ↦ 7, $$ ↦ 8]  
}
```

Simple Example

```
void main() {  
    int x ;  
    c1: x = p(7); ←  
     $\Sigma: x \mapsto 8$   
    c2: x = p(9) ;  
  
}
```

```
int p(int a) {  
    c1: [a  $\mapsto$  7]  
    return a + 1;  
    c1:[a  $\mapsto$  7, $$  $\mapsto$  8]  
}
```

Simple Example

```
void main() {  
    int x ;  
    c1: x = p(7);  
    ε: [x ↦ 8]  
    → c2: x = p(9) ;  
}
```

```
int p(int a) {  
    c1:[a ↦7]  
    return a + 1;  
    c1:[a ↦7, $$ ↦8]  
}
```

Simple Example

```
void main() {  
    int x ;  
    c1: x = p(7);  
    ε : [x ↦ 8]  
    c2: x = p(9) ;  
}
```

```
→ int p(int a) {  
    c1:[a ↦7]  
    c2:[a ↦9]  
    return a + 1;  
    c1:[a ↦7, $$ ↦8]  
}
```

Simple Example

```
void main() {  
    int x ;  
    c1: x = p(7);  
    ε : [x ↦ 8]  
    c2: x = p(9) ;  
}
```

```
int p(int a) {  
    c1:[a ↦7]  
    c2:[a ↦9]  
    → return a + 1;  
    c1:[a ↦7, $$ ↦8]  
    c2:[a ↦9, $$ ↦10]  
}
```

Simple Example

```
void main() {  
    int x ;  
    c1: x = p(7);  
    ε : [x ↦ 8]  
    c2: x = p(9) ; ←  
    ε : [x ↦ 10]  
}  
  
int p(int a) {  
    c1:[a ↦7]  
    c2:[a ↦9]  
    return a + 1;  
    c1:[a ↦7, $$ ↦8]  
    c2:[a ↦9, $$ ↦10]  
}
```


The Call String Approach

- The data flow value is associated with sequences of calls (call string)
- Use Chaotic iterations over the supergraph
- To guarantee termination limit the size of call string (typically 1 or 2)
 - Represents tails of calls
- Abstract inline

Another Example ($|cs|=2$)

```
void main() {  
    int x ;  
    c1: x = p(7);  
     $\epsilon : [x \mapsto 16]$   
    c2: x = p(9) ;  
     $\epsilon :: [x \mapsto 20]$   
}
```

```
int p(int a) {  
     $c1:[a \mapsto 7]$   
     $c2:[a \mapsto 9]$   
    return c3: p1(a + 1);  
     $c1:[a \mapsto 7, \$\$ \mapsto 16]$   
     $c2:[a \mapsto 9, \$\$ \mapsto 20]$   
}
```

```
int p1(int b) {  
     $c1.c3:[b \mapsto 8]$   
     $c2.c3:[b \mapsto 10]$   
    return 2 * b;  
     $c1.c3:[b \mapsto 8, \$\$ \mapsto 16]$   
     $c2.c3:[b \mapsto 10, \$\$ \mapsto 20]$   
}
```

Another Example ($|cs|=1$)

```
void main() {  
  int x ;  
  c1: x = p(7);  
   $\epsilon : [x \mapsto \tau]$   
  c2: x = p(9) ;  
   $\epsilon : [x \mapsto \tau]$   
}
```

```
int p(int a) {  
  c1:[a  $\mapsto$  7]  
  c2:[a  $\mapsto$  9]  
  return c3: p1(a + 1);  
  c1:[a  $\mapsto$  7, $$  $\mapsto$  T]  
  c2:[a  $\mapsto$  9, $$  $\mapsto$  T]  
}
```

```
int p1(int b) {  
  (c1 | c2)c3:[b  $\mapsto$  T]  
  return 2 * b;  
  (c1 | c2)c3:[b  $\mapsto$  T, $$  $\mapsto$  T]  
}
```

Handling Recursion

```
void main() {  
    c1: p(7);  
    ε : [x ↦ τ]  
}
```

```
int p(int a) {  
    c1: [a ↦ 7]  c1.c2+: [a ↦ τ]  
    if (...) {  
        c1: [a ↦ 7]  c1.c2+: [a ↦ τ]  
        a = a - 1 ;  
        c1: [a ↦ 6]  c1.c2+: [a ↦ τ]  
        c2: p (a);  
        c1.c2*: [a ↦ τ]  
        a = a + 1;  
        c1.c2*: [a ↦ τ]  
    }  
    c1.c2*: [a ↦ τ]  
    x = -2*a + 5;  
    c1.c2*: [a ↦ τ, x ↦ τ]  
}
```

Summary Call String

- Easy to implement
- Efficient for very small call strings
- Limited precision
 - Often loses precision for recursive programs
 - For finite domains can be precise even with recursion (with a bounded callstring)
- Order of calls can be abstracted
- Related method: procedure cloning

The Functional Approach

- The meaning of a procedure is mapping from states into states
- The abstract meaning of a procedure is function from an abstract state to abstract states
- Relation between input and output
- In certain cases can compute JVP

The Functional Approach

- Two phase algorithm
 - Compute the dataflow solution at the exit of a procedure as a function of the initial values at the procedure entry (functional values)
 - Compute the dataflow values at every point using the functional values

Phase 1

```
void main() {
```

```
    p(7);
```

```
}
```

```
int p(int a) {
```

```
    [a ↦ a0, x ↦ x0]
```

```
    if (...) {
```

```
        [a ↦ a0, x ↦ x0]
```

```
        a = a - 1;
```

```
        [a ↦ a0-1, x ↦ x0]
```

```
        p(a);
```

```
        [a ↦ a0-1, x ↦ -2a0+7]
```

```
        a = a + 1;
```

```
        [a ↦ a0, x ↦ -2a0+7]
```

```
    }
```

```
    [a ↦ a0, x ↦ x0] [a ↦ a0, x ↦ τ]
```

```
    x = -2*a + 5;
```

```
    [a ↦ a0, x ↦ -2*a0+5]
```

```
}
```

$p(a_0, x_0) = [a \mapsto a_0, x \mapsto -2a_0 + 5]$

Phase 2

```
void main() {  
    p(7);  
    [x ↦ -9]  
}
```

$p(a_0, x_0) = [a \mapsto a_0, x \mapsto -2a_0 + 5]$

```
int p(int a) {  
    [a ↦ 7, x ↦ 0]      [a ↦ T, x ↦ 0]  
    if (...) {  
        [a ↦ 7, x ↦ 0]      [a ↦ T, x ↦ 0]  
        a = a - 1;  
        [a ↦ 6, x ↦ 0]      [a ↦ T, x ↦ 0]  
        p(a);  
        [a ↦ 6, x ↦ -7]     [a ↦ T, x ↦ T]  
        a = a + 1;  
        [a ↦ 7, x ↦ -7]     [a ↦ T, x ↦ T]  
    }  
    [a ↦ 7, x ↦ 0]      [a ↦ T, x ↦ T]  
    x = -2*a + 5;  
    [a ↦ 7, x ↦ -9]      [a ↦ T, x ↦ T]  
}
```

Summary Functional approach

- Computes procedure abstraction
- Sharing between different contexts
- Rather precise
- Recursive procedures may be more precise/efficient than loops
- But requires more from the implementation
 - Representing (input/output) relations
 - Composing relations

Issues in Functional Approach

- How to guarantee that finite height for functional lattice?
 - It may happen that L has finite height and yet the lattice of monotonic function from L to L do not
- Efficiently represent functions
 - Functional join
 - Functional composition
 - Testing equality

Tabulation

- Special case: L is finite
- Data facts: $d \in L \times L$
- Initialization:
 - $f_{\text{start},\text{start}} = (\top, \top)$; otherwise (\perp, \perp)
 - $S[\text{start}, \top] = \top$
- Propagation of (x, y) over edge $e = (n, n')$
 - Maintain summary: $S[n', x] = S[n', x] \sqcup \llbracket n \rrbracket (y)$
 - n intra-node: $\rightarrow n' : (x, \llbracket n \rrbracket (y))$
 - n call-node:
 - $\rightarrow n' : (y, y)$ if $S[n', y] = \perp$ and $n' = \text{entry node}$
 - $\rightarrow n' : (x, z)$ if $S[\text{exit}(\text{call}(n), y)] = z$ and $n' = \text{ret-site-of } n$
 - n return-node: $\rightarrow n' : (u, y)$; $n_c = \text{call-site-of } n'$, $S[n_c, u] = x$

CFL-Graph reachability

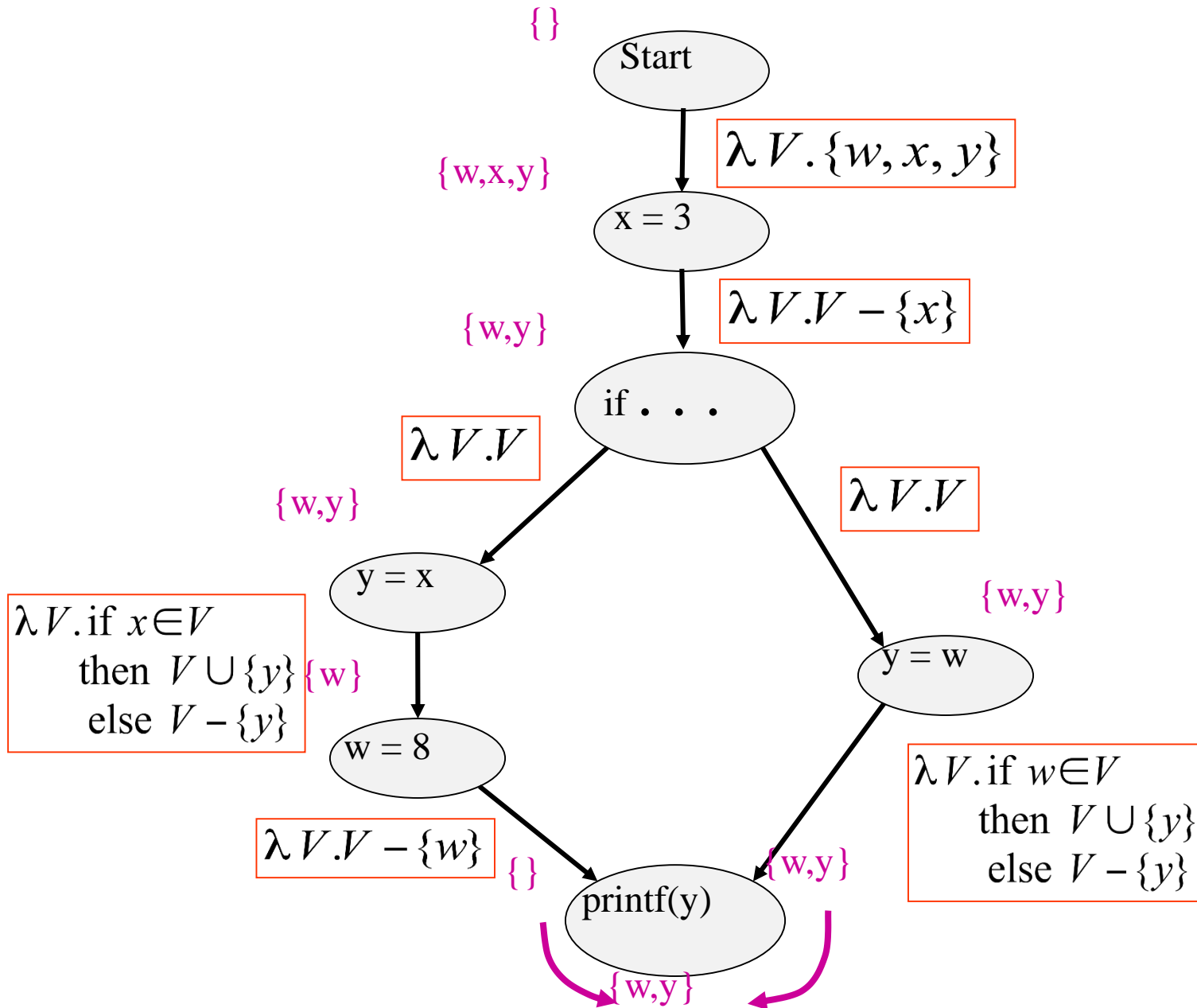
- Special cases of functional analysis
- Finite distributive lattices
- Provides more efficient analysis algorithms
- Reduce the interprocedural analysis problem to finding context free reachability



IDFS / IDE

- **IDFS** Interprocedural Distributive Finite Subset
Precise interprocedural dataflow analysis via graph reachability. *Reps, Horowitz, and Sagiv, POPL' 95*
- **IDE** Interprocedural Distributive Environment
Precise interprocedural dataflow analysis with applications to constant propagation. *Reps, Horowitz, and Sagiv, FASE' 95, TCS' 96*
 - *More general solutions exist*

Possibly Uninitialized Variables

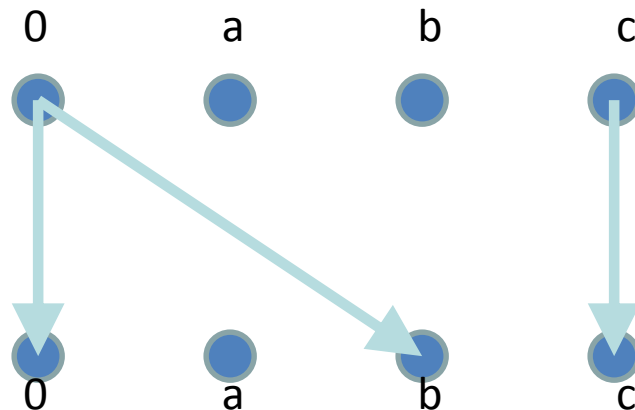


IFDS Problems

- Finite subset distributive
 - Lattice $L = \wp(D)$
 - \sqsubseteq is \subseteq
 - \sqcup is \cup
 - Transfer functions are distributive
- Efficient solution through formulation as CFL reachability

Encoding Transfer Functions

- Enumerate all input space and output space
- Represent functions as graphs with $2(D+1)$ nodes
- Special symbol “0” denotes empty sets (sometimes denoted Λ)
- Example: $D = \{ a, b, c \}$
 $f(S) = (S - \{a\}) \cup \{b\}$



Efficiently Representing Functions

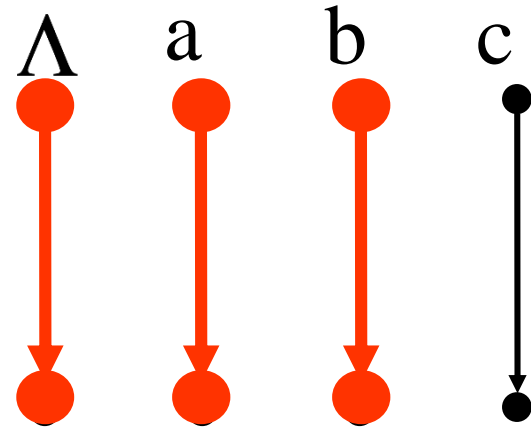
- Let $f:2^D \rightarrow 2^D$ be a distributive function
- Then:
 - $f(X) = f(\emptyset) \cup (\cup \{ f(\{z\}) \mid z \in X \})$
 - $f(X) = f(\emptyset) \cup (\cup \{ f(\{z\}) \setminus f(\emptyset) \mid z \in X \})$

Representing Dataflow Functions

Identity Function

$$f = \lambda V.V$$

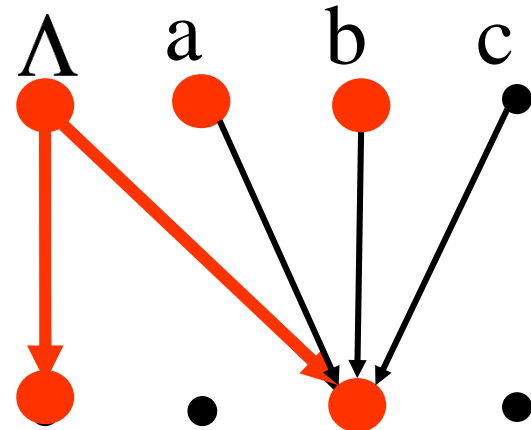
$$f(\{a, b\}) = \{a, b\}$$



Constant Function

$$f = \lambda V.\{b\}$$

$$f(\{a, b\}) = \{b\}$$

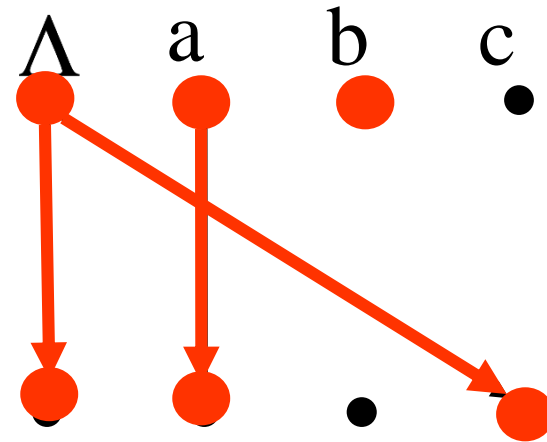


Representing Dataflow Functions

“Gen/Kill” Function

$$f = \lambda V. (V - \{b\}) \cup \{c\}$$

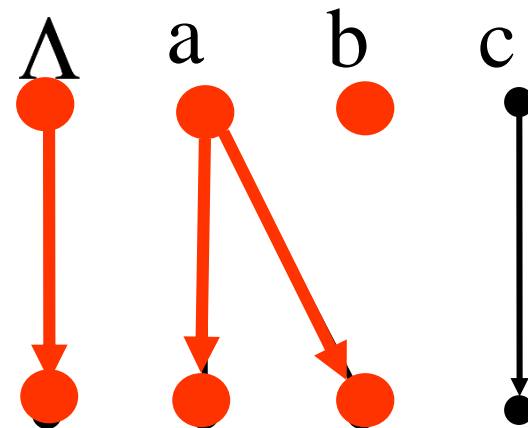
$$f(\{a, b\}) = \{a, c\}$$

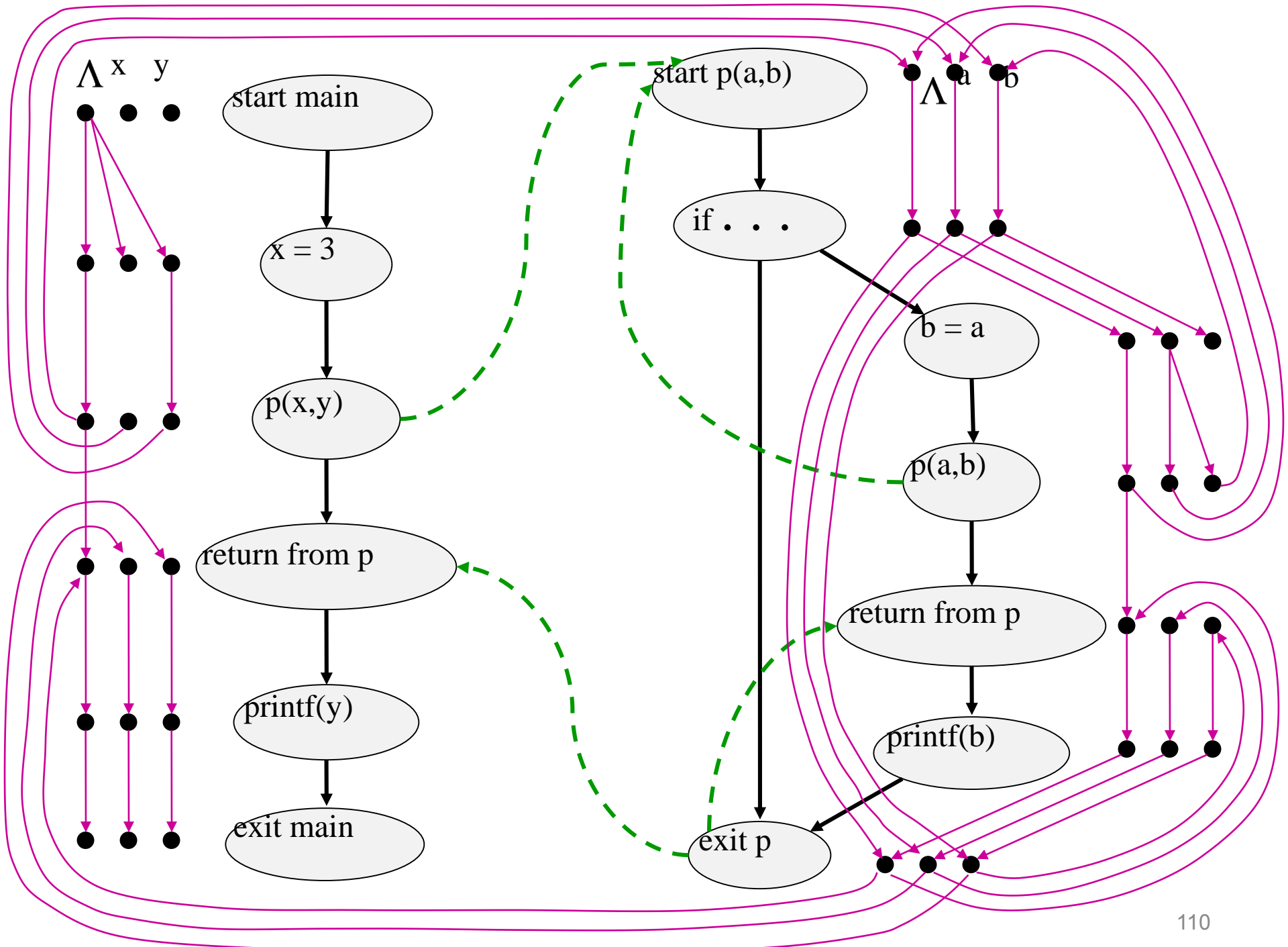


Non-“Gen/Kill” Function

$$f = \lambda V. \text{if } a \in V \\ \text{then } V \cup \{b\} \\ \text{else } V - \{b\}$$

$$f(\{a, b\}) = \{a, b\}$$

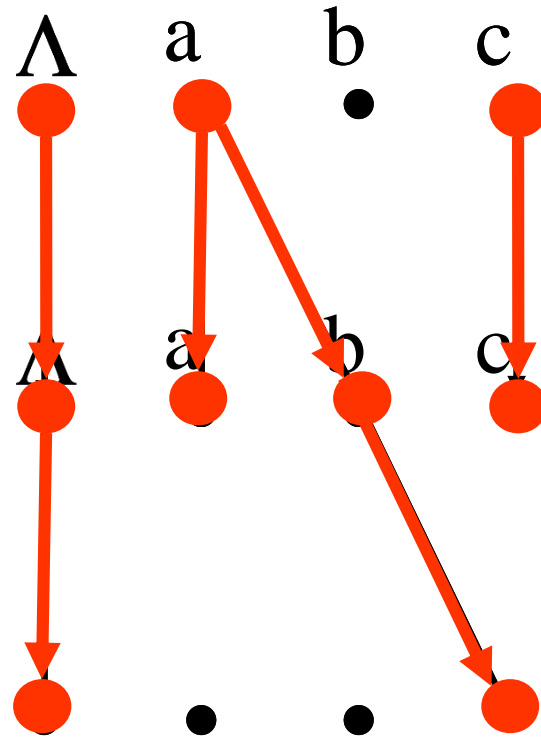




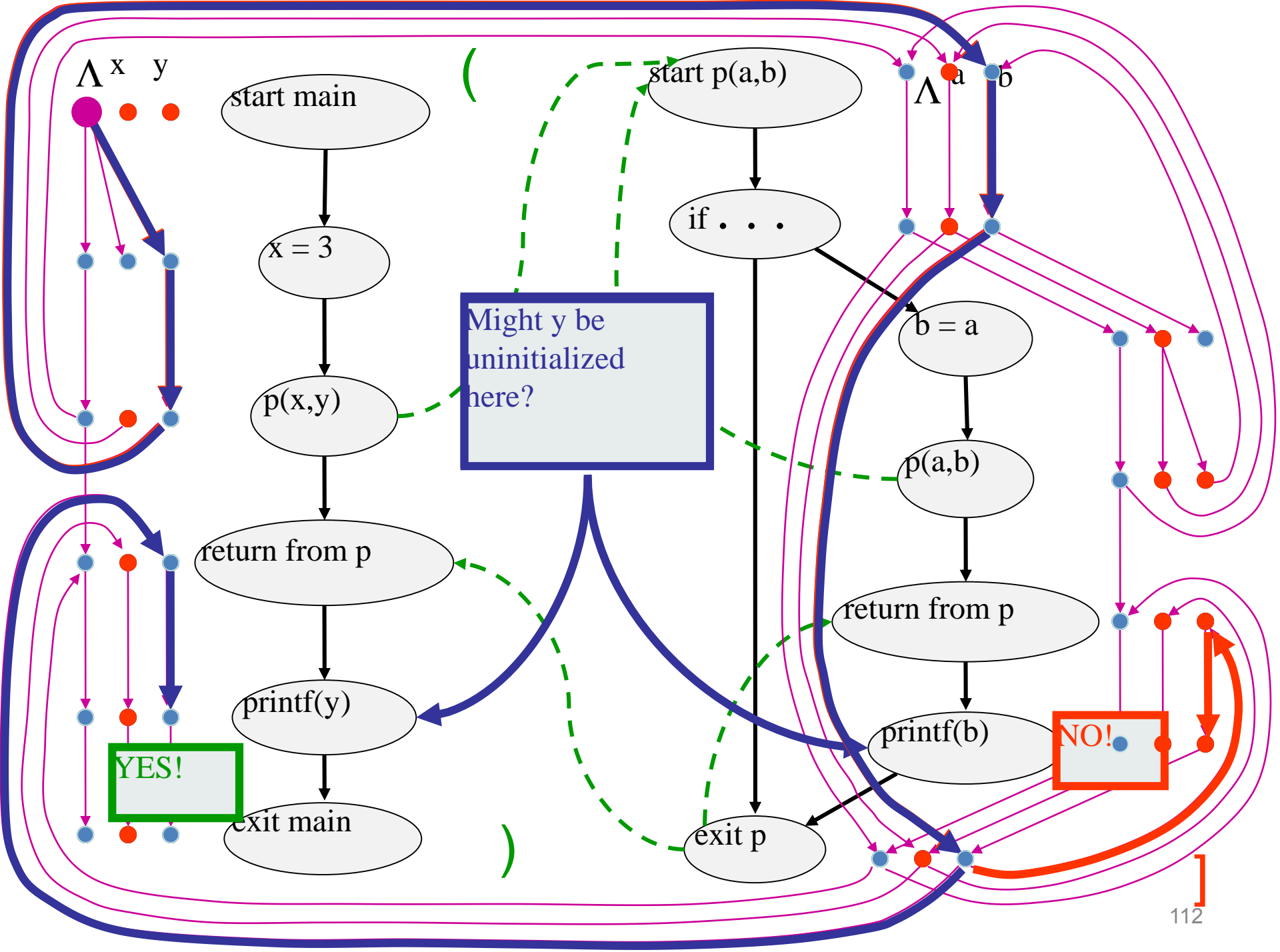
Composing Dataflow Functions

$f_1 = \lambda V. \text{if } a \in V$
 then $V \cup \{b\}$
 else $V - \{b\}$

$f_2 = \lambda V. \text{if } b \in V$
 then $\{c\}$
 else ϕ



$$f_2 \circ f_1(\{a, c\}) = \boxed{\{c\}}$$




The Tabulation Algorithm

- Worklist algorithm, start from entry of “main”
- Keep track of
 - Path edges: matched paren paths from procedure entry
 - Summary edges: matched paren call-return paths
- At each instruction
 - Propagate facts using transfer functions; **extend path edges**
- At each call
 - Propagate to procedure entry, start with an empty path
 - If a summary for that entry exists, use it
- At each exit
 - Store paths from corresponding call points as summary paths
 - When a new summary is added, propagate to the return node

Interprocedural Dataflow Analysis via CFL-Reachability

- Graph: Exploded control-flow graph
- $L: L(\text{unbalLeft})$
 - $\text{unbalLeft} = \text{valid}$
- Fact d holds at n iff there is an $L(\text{unbalLeft})$ -path from $\langle \text{start}_{\text{main}}, \Lambda \rangle$ to $\langle n, d \rangle$

Asymptotic Running Time

- CFL-reachability
 - Exploded control-flow graph: ND nodes
 - Running time: $O(N^3D^3)$
- Exploded control-flow graph  Special structure

Running time: $O(ED^3)$

Typically: $E \approx N$, hence $O(ED^3) \approx O(ND^3)$

“Gen/kill” problems: $O(ED)$

IDE

- Goes beyond IFDS problems
 - Can handle unbounded domains
- Requires special form of the domain
- Can be **much** more efficient than IFDS

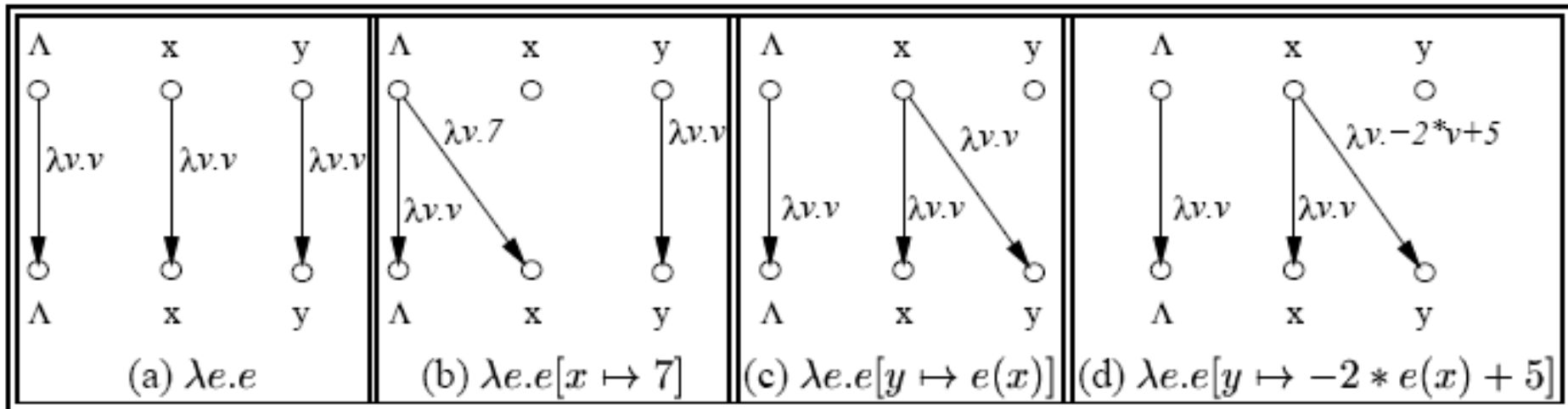
Example Linear Constant Propagation

- Consider the constant propagation lattice
- The value of every variable y at the program exit can be represented by:

$$y = \sqcup \{(a_x x + b_x) \mid x \in \text{Var}_*\} \sqcup c$$
$$a_x, c \in Z \cup \{\perp, \top\} \quad b_x \in Z$$

- Supports efficient composition and “functional” join
 - $[z := a * y + b]$
 - What about $[z:=x+y]$?

Linear constant propagation



Point-wise representation of environment transformers

IDE Analysis

- Point-wise representation closed under composition
- CFL-Reachability on the exploded graph
- Compose functions

```

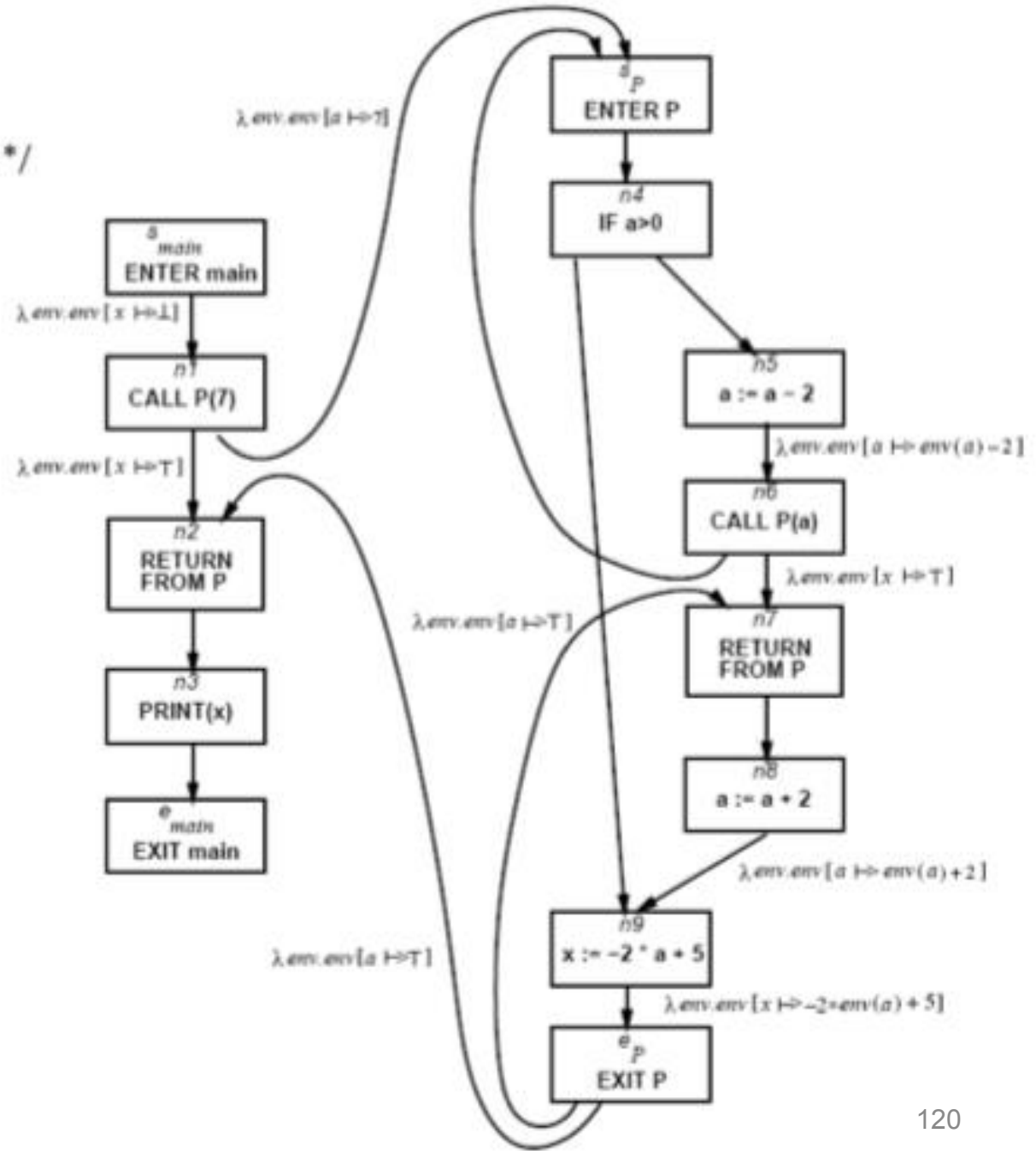
declare x: integer
program main
begin
    call P(7)
    print (x) /* x is a constant here */
end

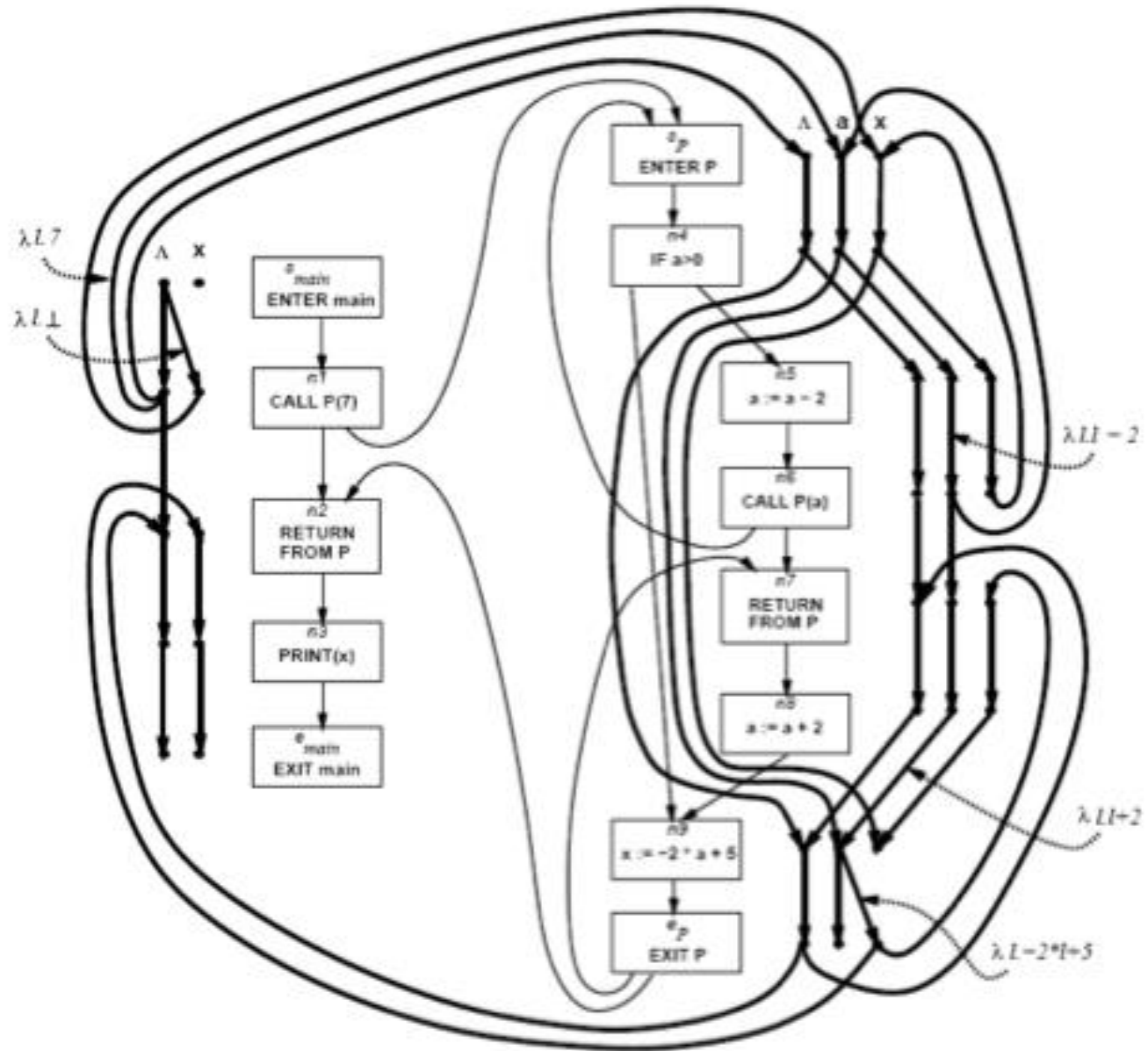
```

```

procedure P (value a : integer)
begin /* a is not a constant here */
    if a > 0 then
        a := a - 2
        call P (a)
        a := a + 2
    fi
    x := -2 * a + 5
    /* x is not a constant here */
end

```





Costs

- $O(ED^3)$
- Class of value transformers $F \subseteq L \rightarrow L$
 - $\text{id} \in F$
 - Finite height
- Representation scheme with (efficient)
 - Application
 - Composition
 - Join
 - Equality
 - Storage

Conclusion

- Handling functions is crucial for abstract interpretation
- Virtual functions and exceptions complicate things
- But scalability is an issue
 - Small call strings
 - Small functional domains
 - Demand analysis

Challenges in Interprocedural Analysis

- Respect call-return mechanism
- Handling recursion
- Local variables
- Parameter passing mechanisms
- The called procedure is not always known
- The source code of the called procedure is not always available

A trivial treatment of procedure

- Analyze a single procedure
- After every call continue with conservative information
 - Global variables and local variables which “may be modified by the call” have unknown values
- Can be easily implemented
- Procedures can be written in different languages
- Procedure inline can help

Disadvantages of the trivial solution

- Modular (object oriented and functional) programming encourages small frequently called procedures
- Almost all information is lost

Bibliography

- **Textbook 2.5**
- Patrick Cousot & Radhia Cousot. Static determination of dynamic properties of recursive procedures In *IFIP Conference on Formal Description of Programming Concepts*, E.J. Neuhold, (Ed.), pages 237-277, St-Andrews, N.B., Canada, 1977. North-Holland Publishing Company (1978).
- **Two Approaches to interprocedural analysis by Micha Sharir and Amir Pnueli**
- **IDFS** Interprocedural Distributive Finite Subset Precise interprocedural dataflow analysis via graph reachability. *Reps, Horowitz, and Sagiv, POPL' 95*
- **IDE** Interprocedural Distributive Environment Precise interprocedural dataflow analysis with applications to constant propagation. *Sagiv, Reps, Horowitz, and TCS' 96*

A Semantics for Procedure Local Heaps and its Abstractions

Noam Rinetzky Tel Aviv University

Jörg Bauer Universität des Saarlandes

Thomas Reps University of Wisconsin

Mooly Sagiv Tel Aviv University

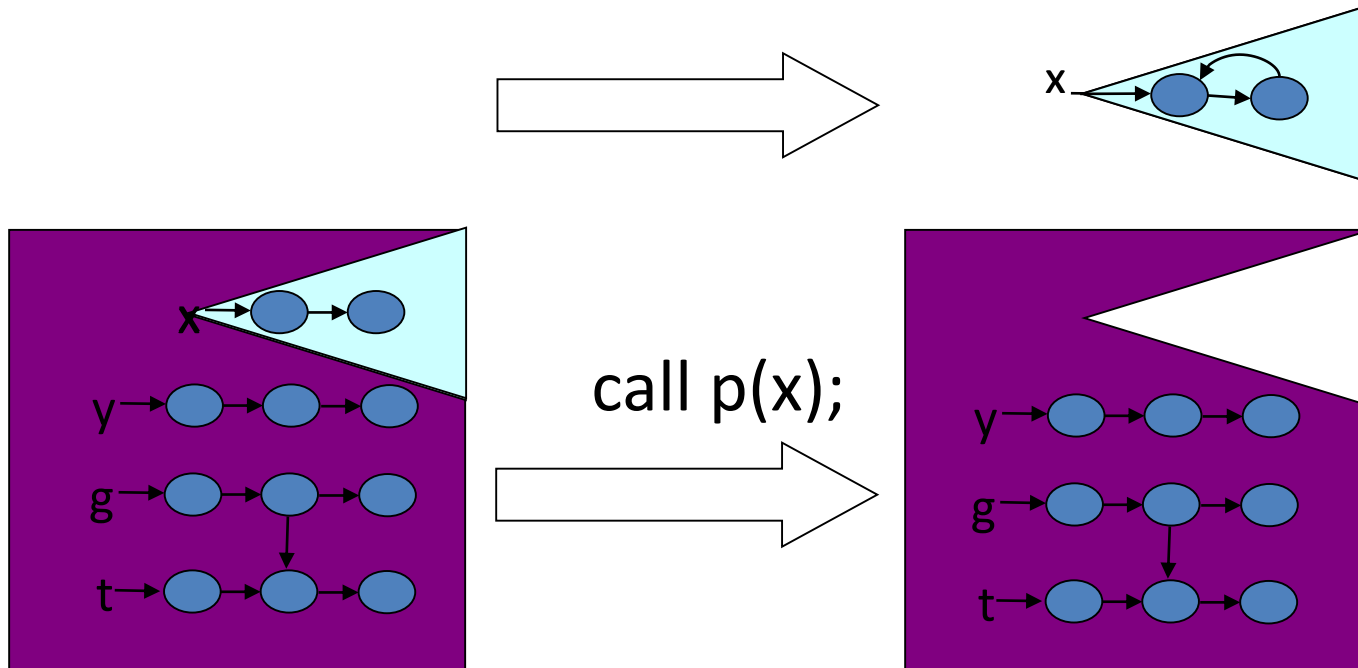
Reinhard Wilhelm Universität des Saarlandes

Motivation

- Interprocedural shape analysis
 - Conservative static pointer analysis
 - Heap intensive programs
 - Imperative programs with procedures
 - Recursive data structures
- Challenge
 - Destructive update
 - Localized effect of procedures

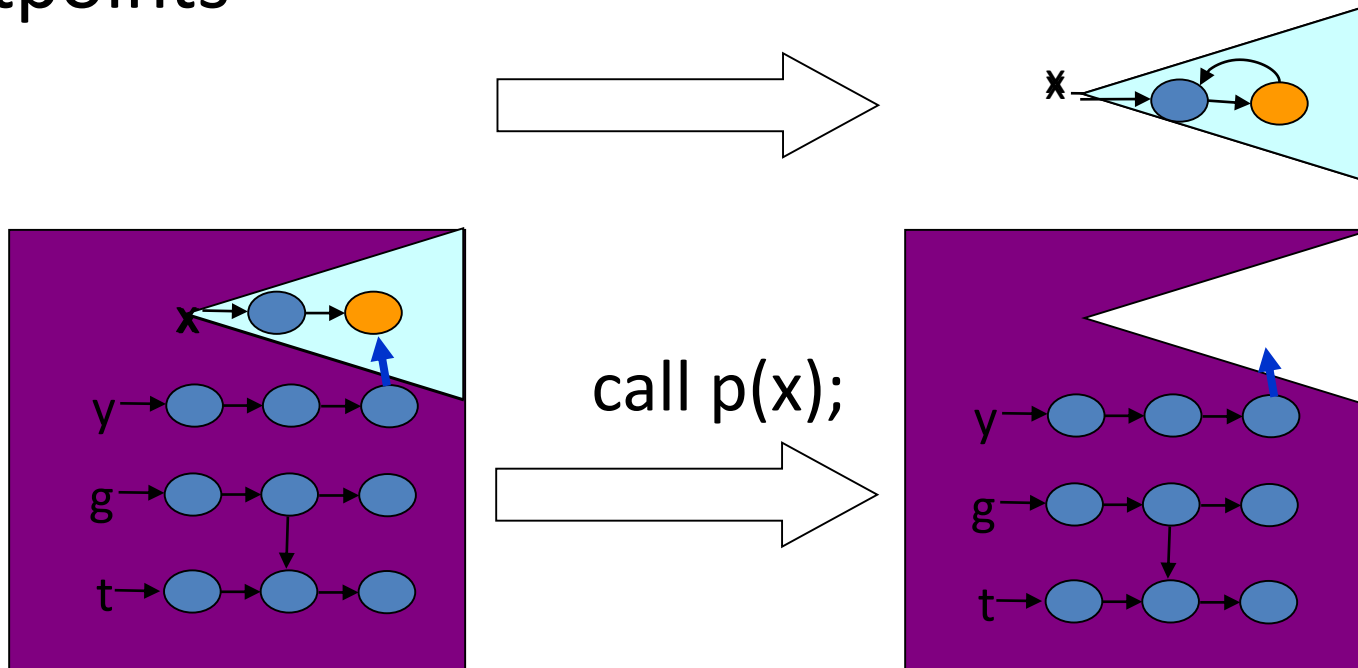
Main idea

- Local heaps



Main idea

- Local heaps
- Cutpoints



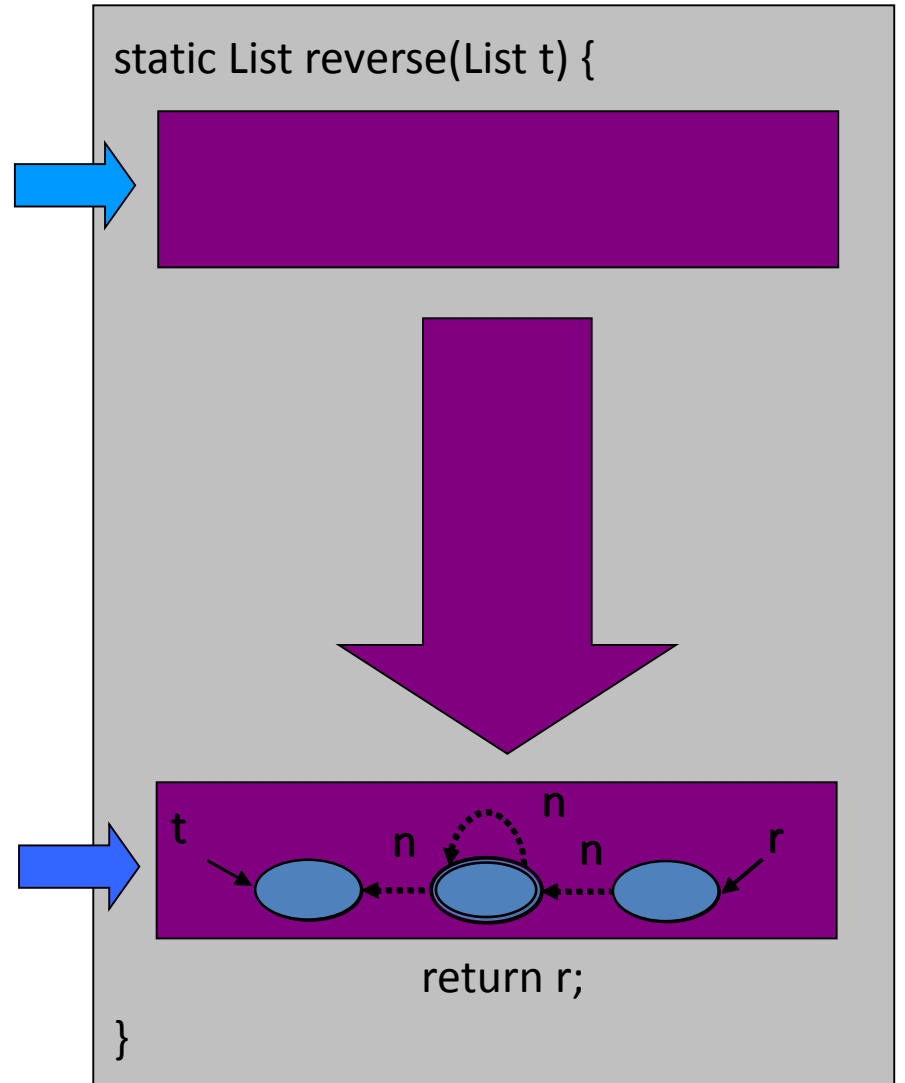
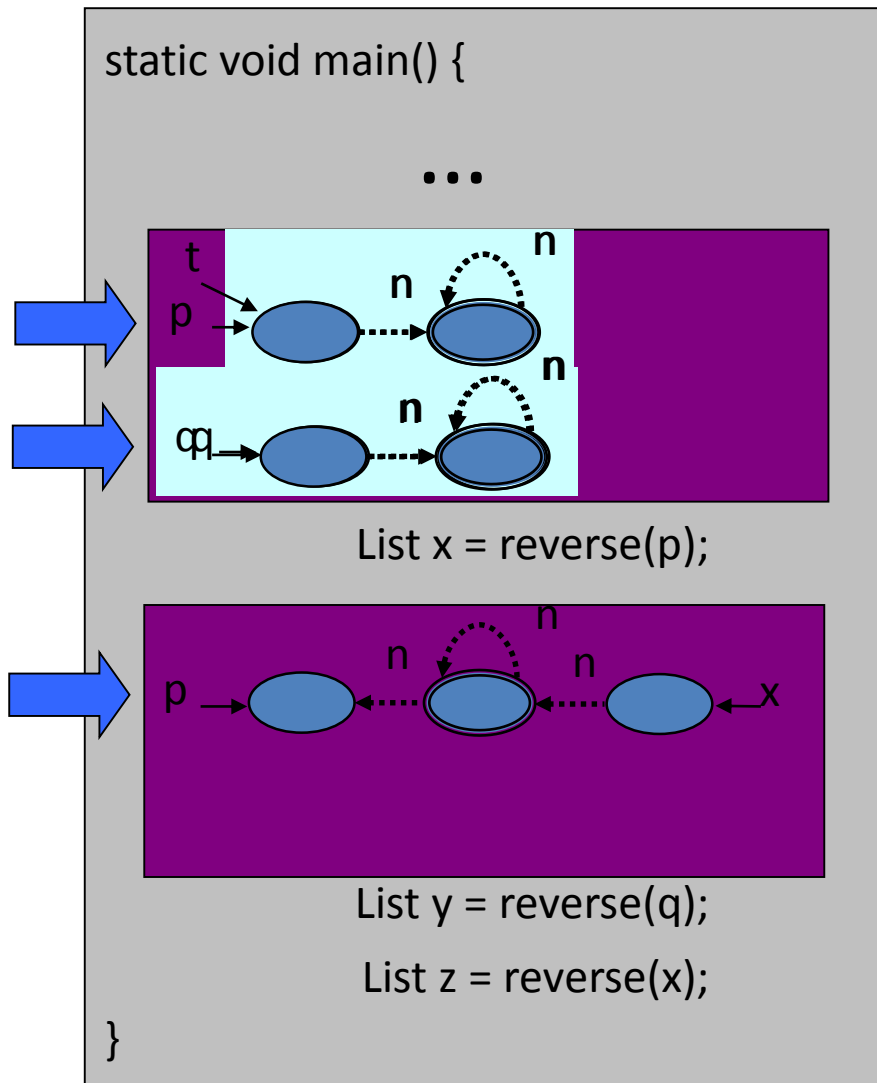
Main Results

- Concrete operational semantics
 - Large step
 - Functional analysis
 - Storeless
 - Shape abstractions
 - Local heap
 - Observationally equivalent to “standard” semantics
 - Java and “clean” C
- Abstractions
 - Shape analysis [Sagiv, Reps, Wilhelm, TOPLAS ‘02]
 - May-alias [Deutsch, PLDI ‘94]
 - ...

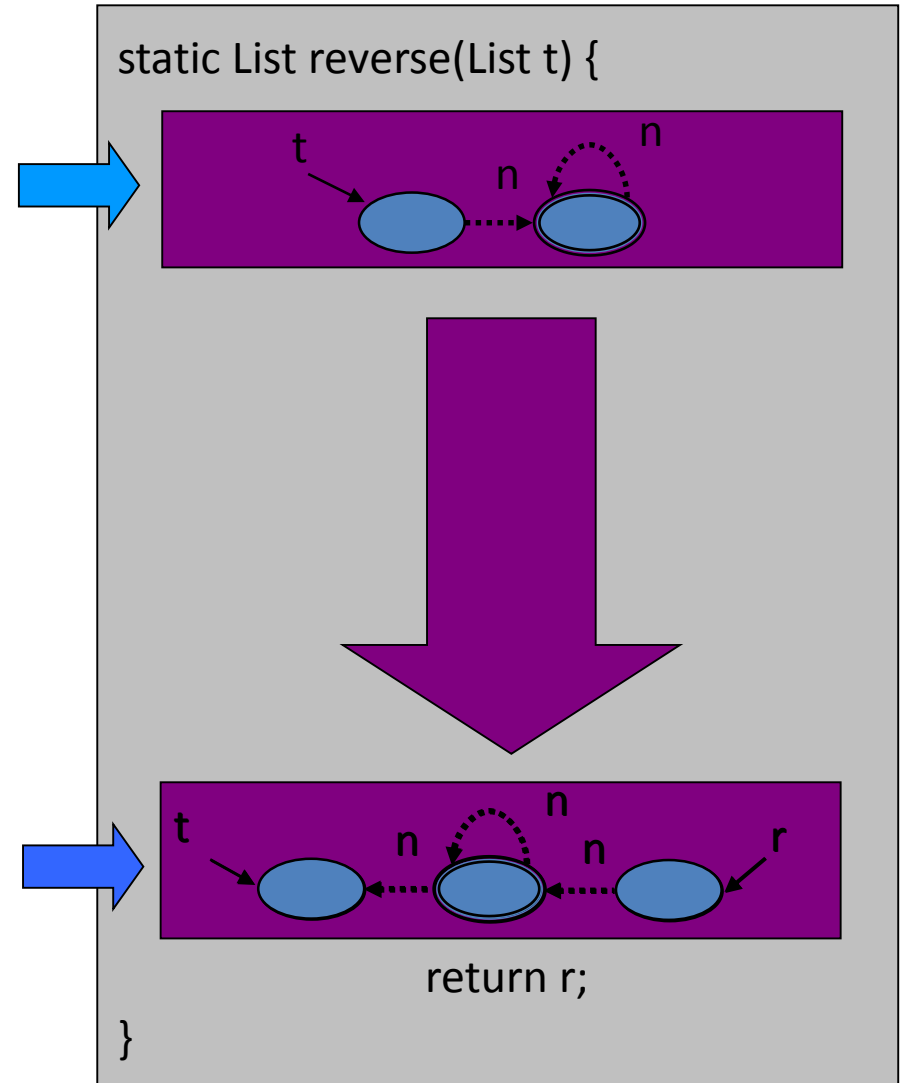
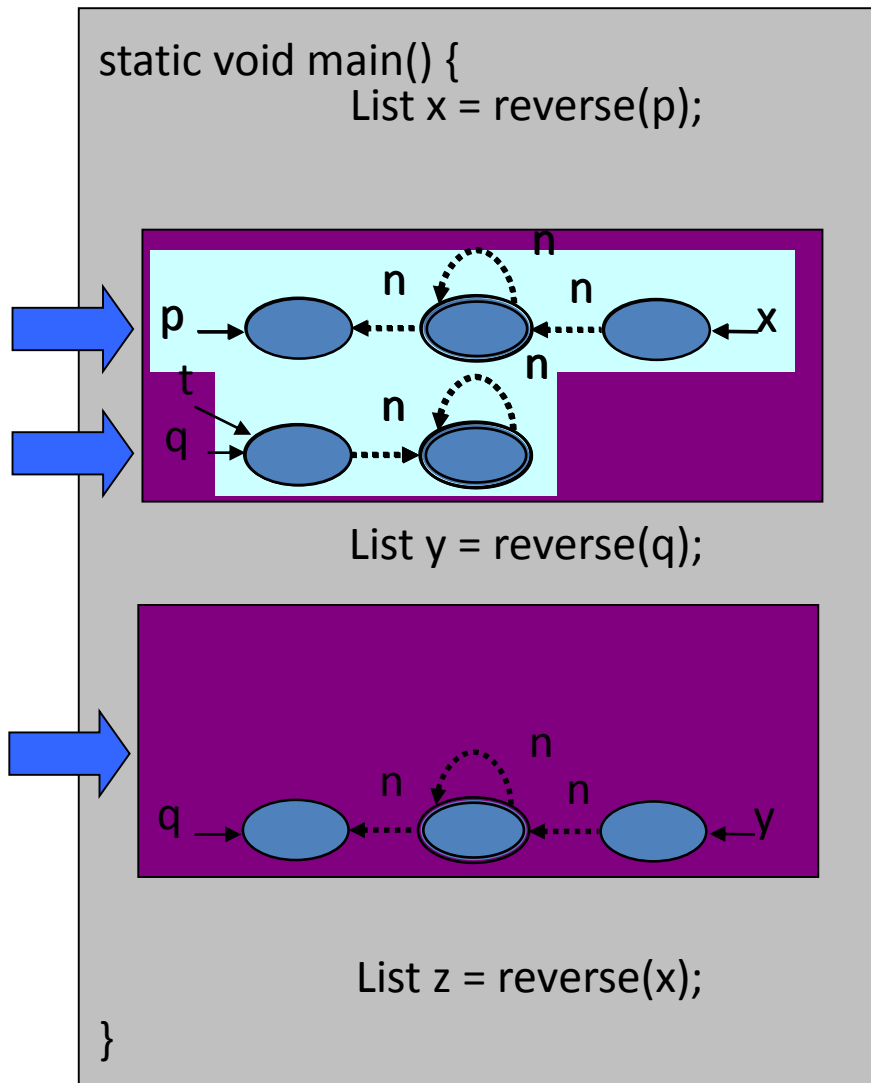
Outline

- Motivating example
 - Local heaps
 - Cutpoints
- Why semantics
- Local heap storeless semantics
- Shape abstraction

Example

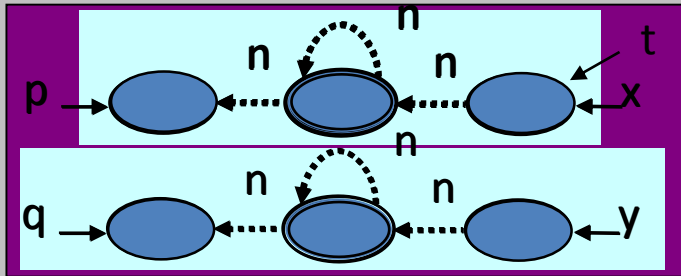


Example

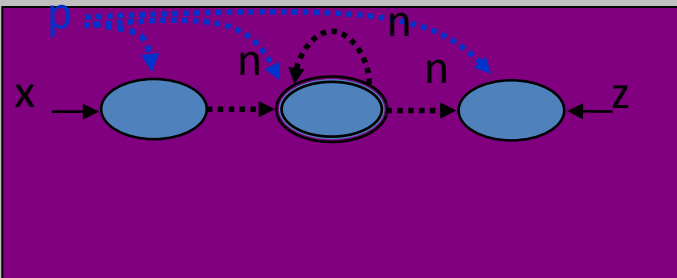


Example

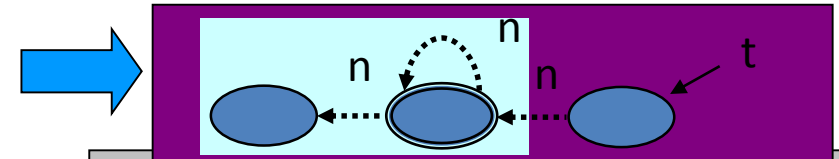
```
static void main() {  
    List x = reverse(p);  
    List y = reverse(q);
```



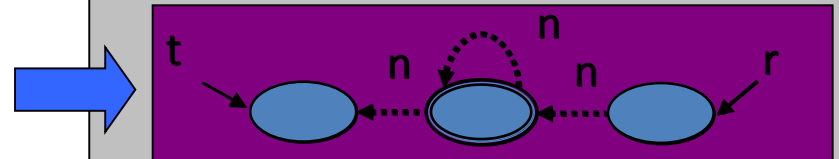
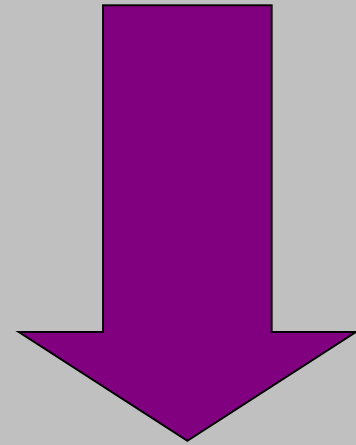
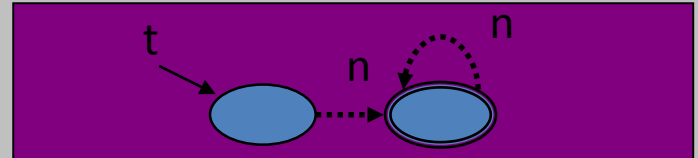
```
    List z = reverse(x);
```



```
}
```



```
static List reverse(List t) {
```



```
    return r;
```

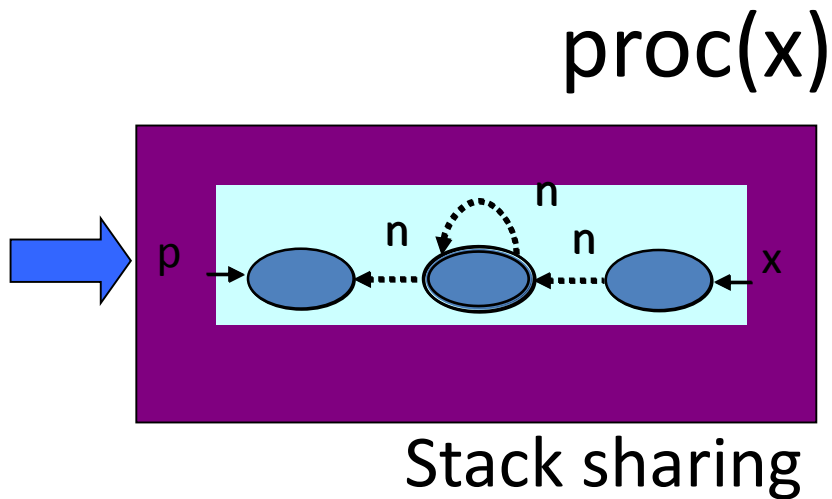
```
}
```

Cutpoints

- **Separating** objects
 - Not pointed-to by a parameter

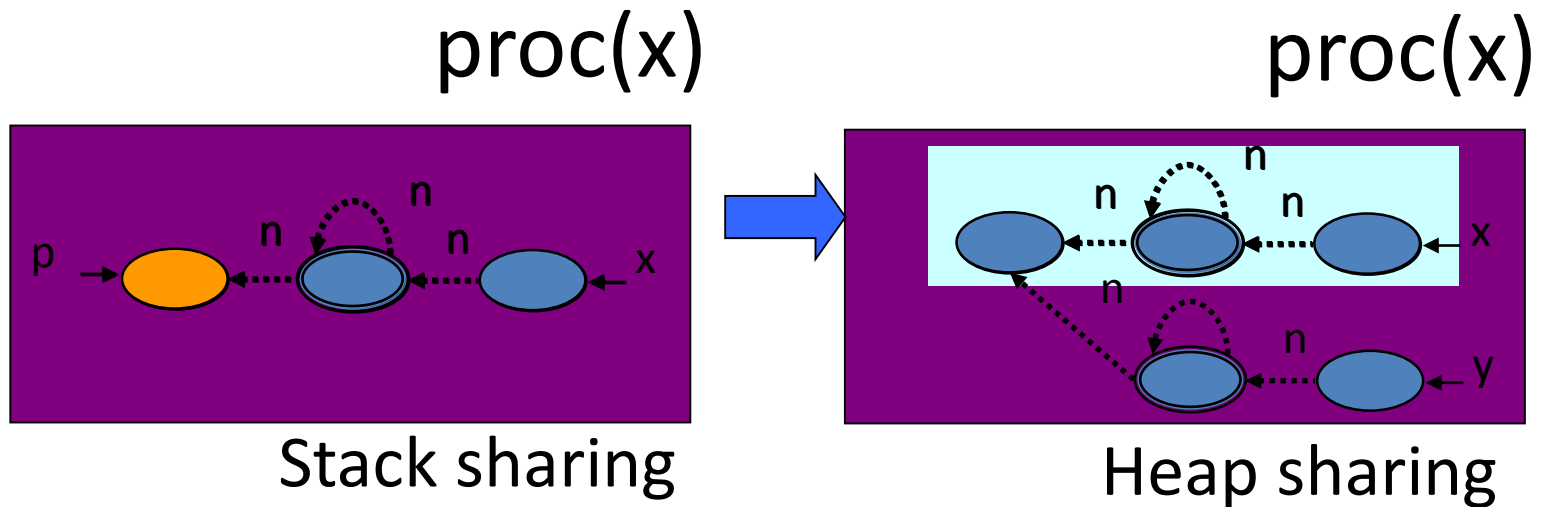
Cutpoints

- Separating objects
 - Not pointed-to by a parameter



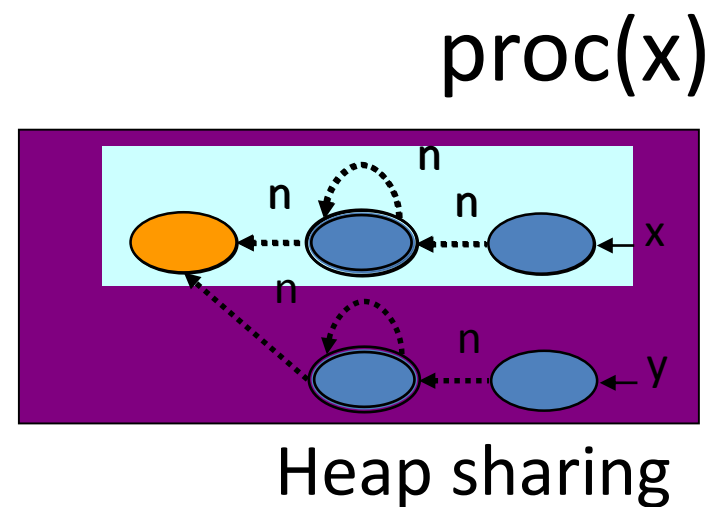
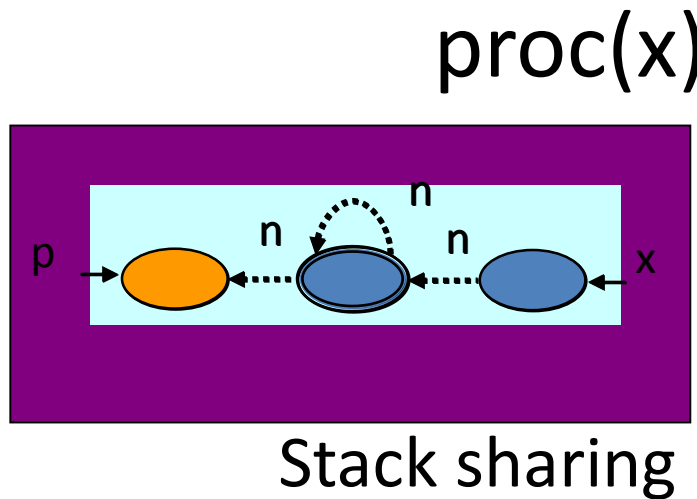
Cutpoints

- Separating objects
 - Not pointed-to by a parameter



Cutpoints

- Separating objects
 - Not pointed-to by a parameter
- Capture external sharing patterns

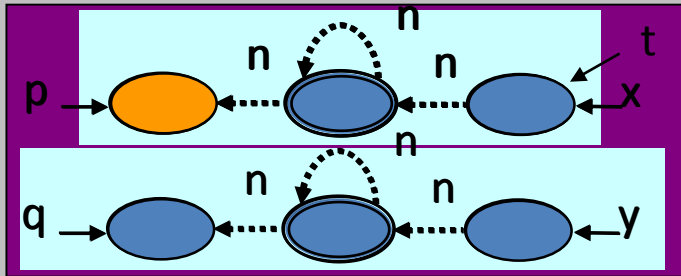


Example

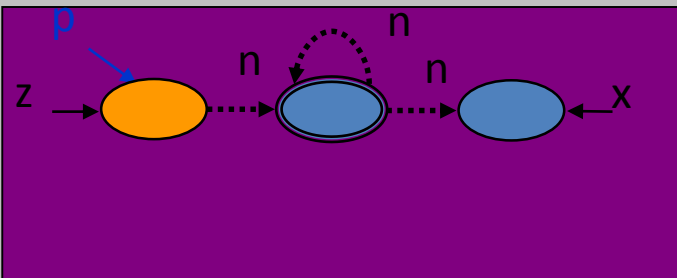
```
static void main() {
```

```
    List x = reverse(p);
```

```
    List y = reverse(q);
```

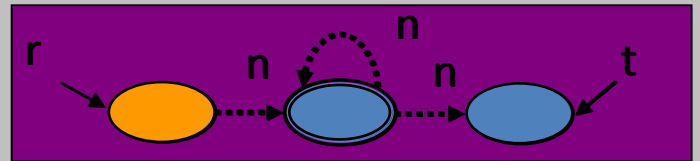
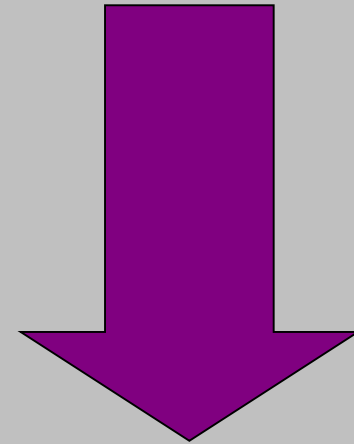
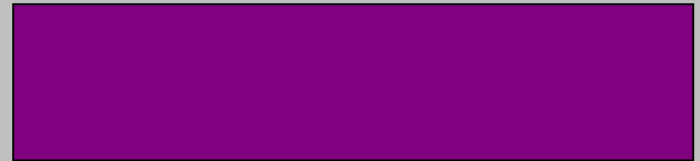


```
    List z = reverse(x);
```



```
}
```

```
static List reverse(List t) {
```



```
    return r;
```

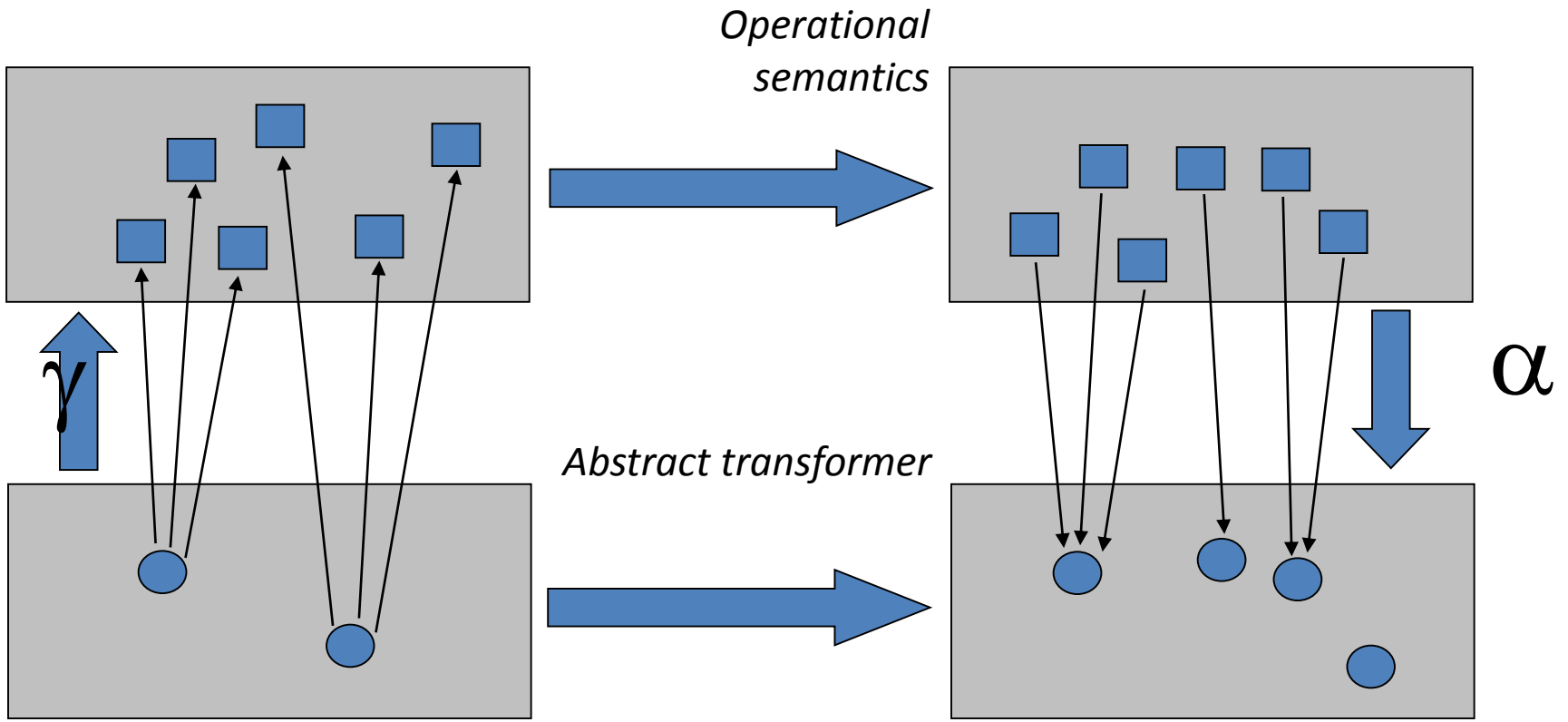
```
}
```

Outline

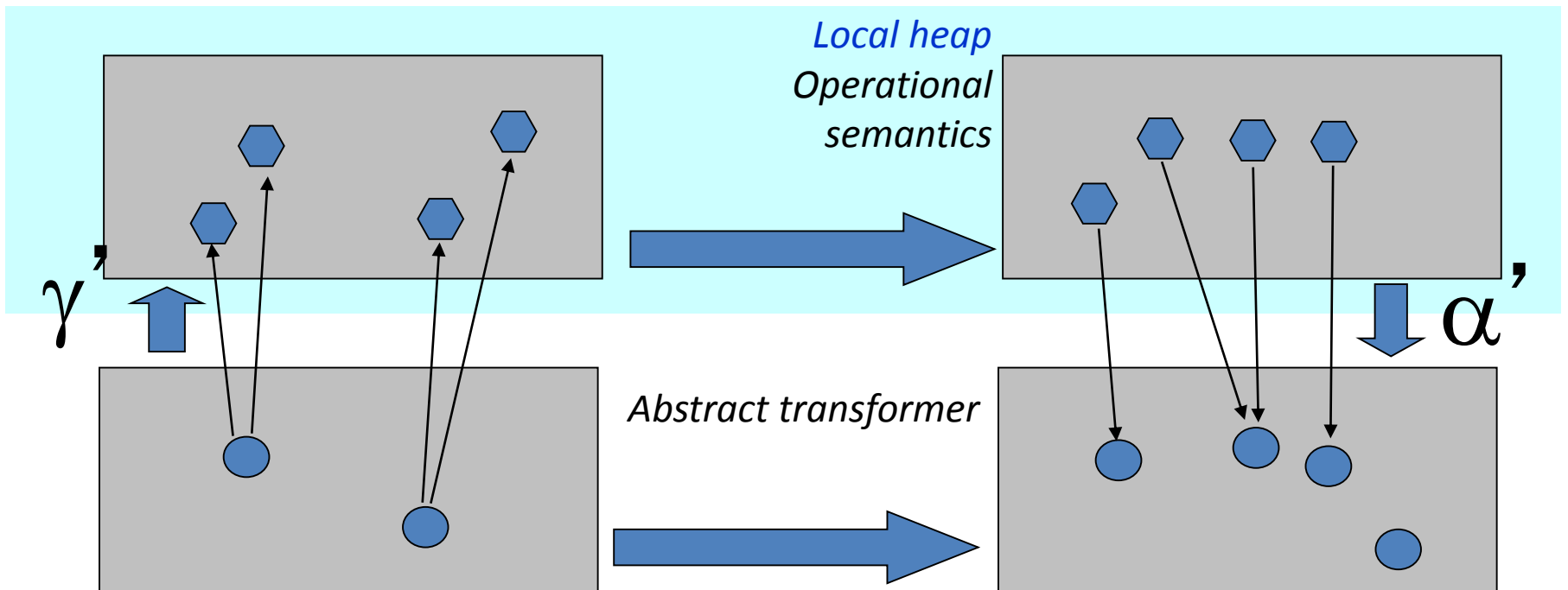
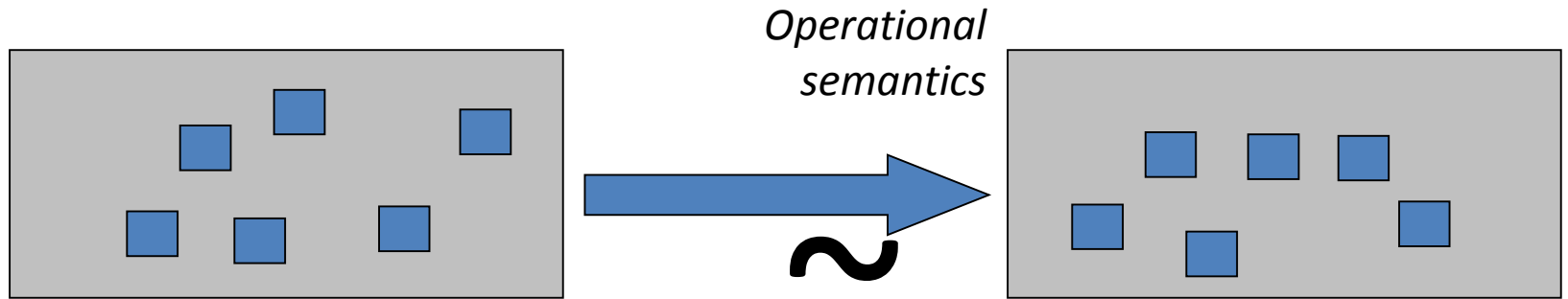
- ✓ Motivating example
 - Why semantics
 - Local heap storeless semantics
 - Shape abstraction

Abstract Interpretation

[Cousot and Cousot, POPL '77]



Introducing local heap semantics



Outline

- ✓ Motivating example
- ✓ Why semantics
- Local heap storeless semantics
- Shape abstraction

Programming model

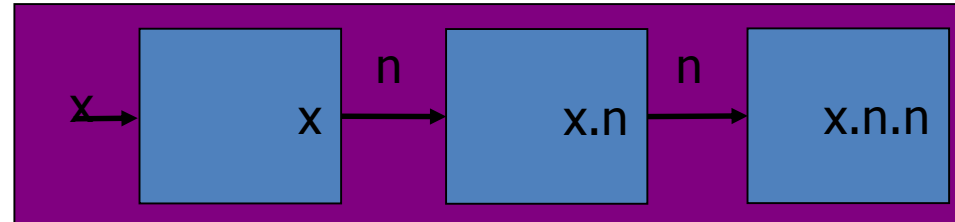
- Single threaded
- Procedures
 - ✓ Value parameters
 - ✓ Recursion
- Heap
 - ✓ Recursive data structures
 - ✓ Destructive update
 - ✗ No explicit addressing (&)
 - ✗ No pointer arithmetic

Simplifying assumptions

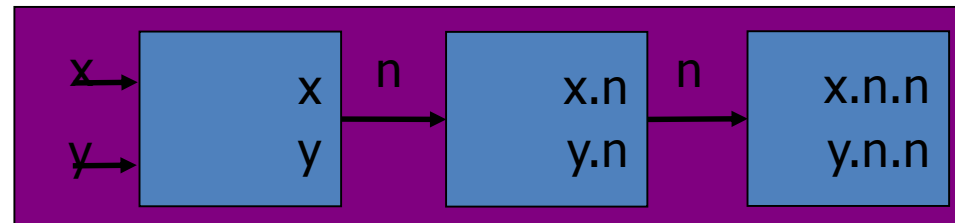
- No primitive values (only references)
- No globals
- Formals not modified

Storeless semantics

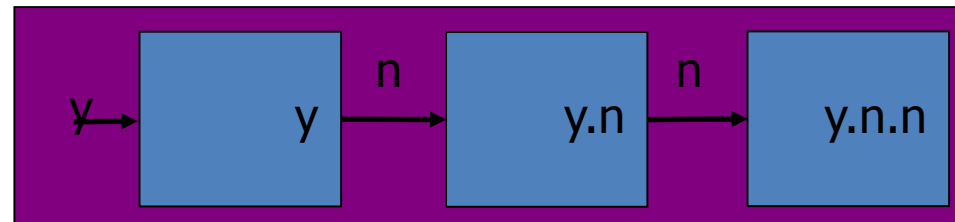
- No addresses
- Memory state:
 - Object: $2^{\text{Access paths}}$
 - Heap: 2^{Object}
- Alias analysis



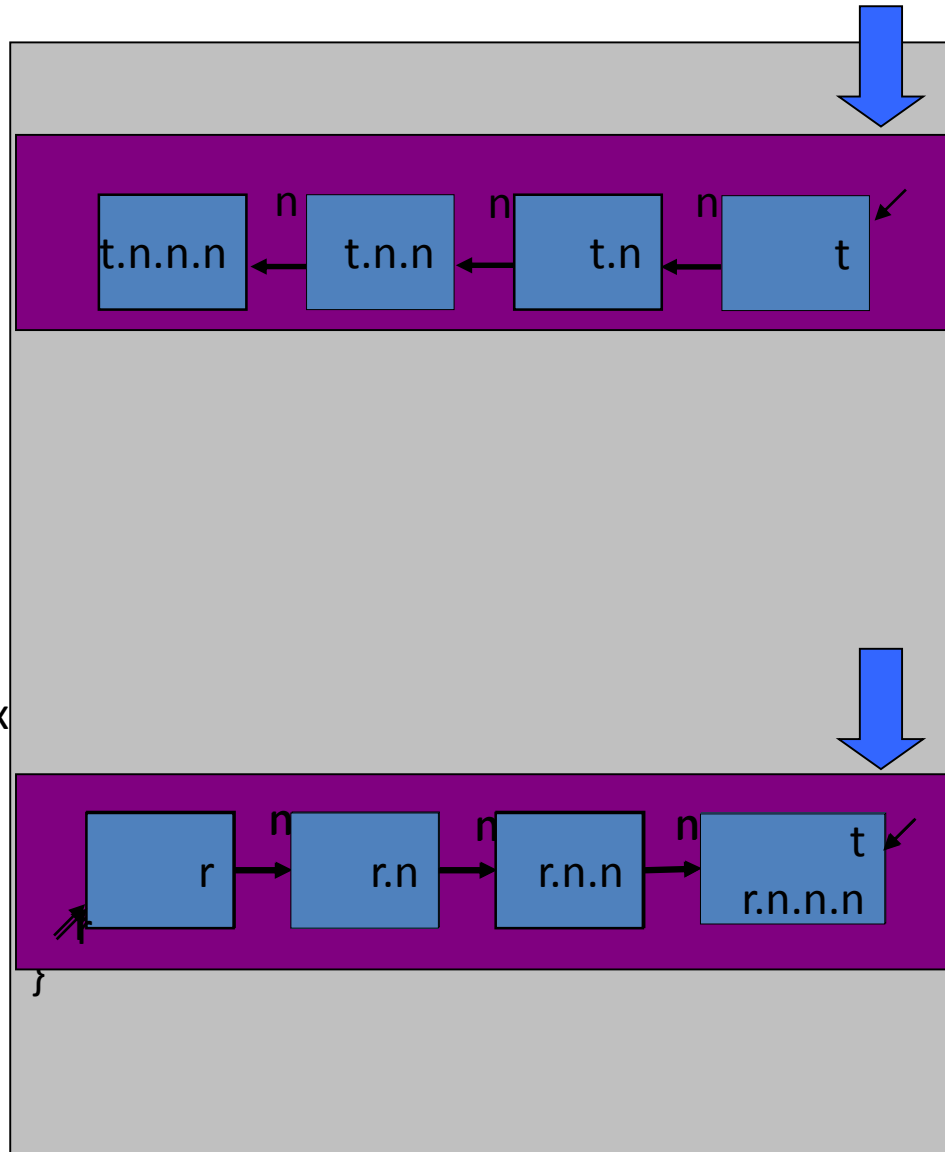
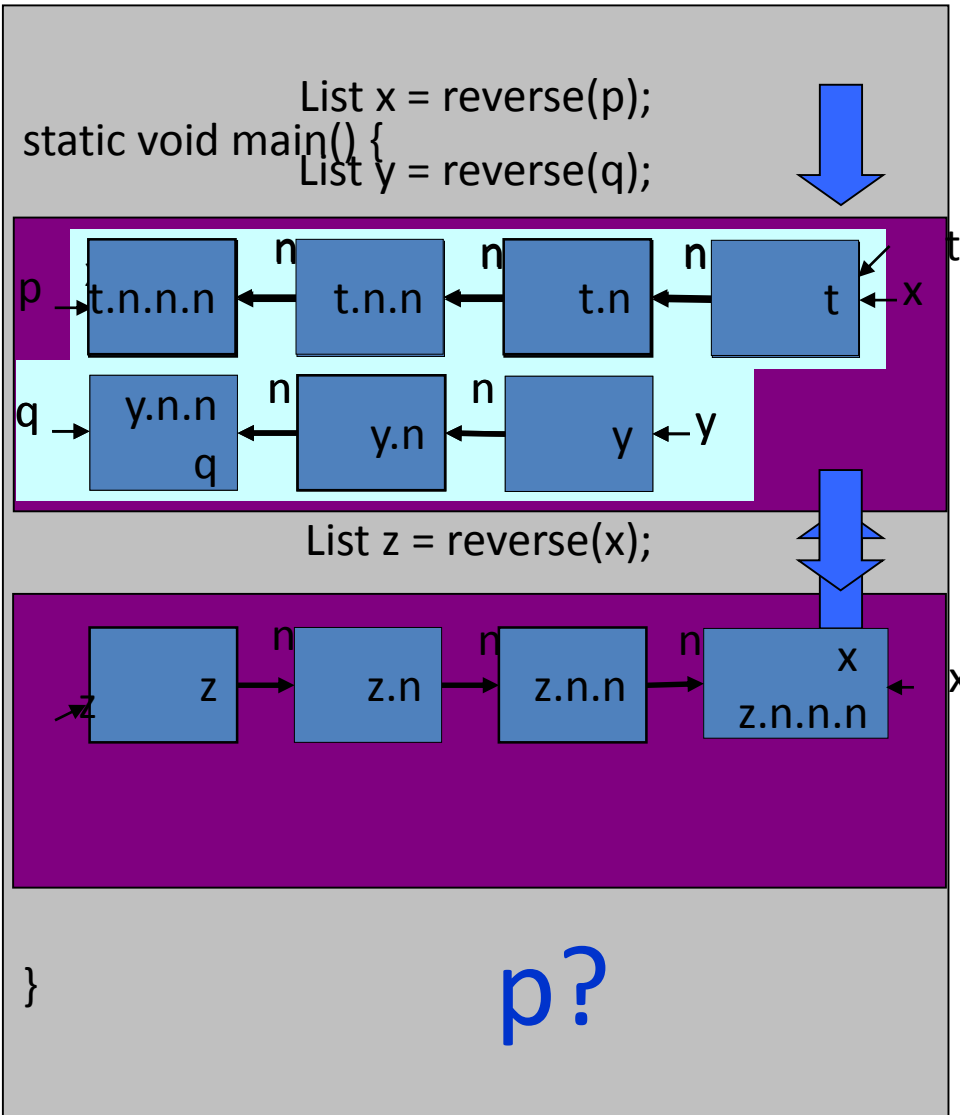
$y = x$



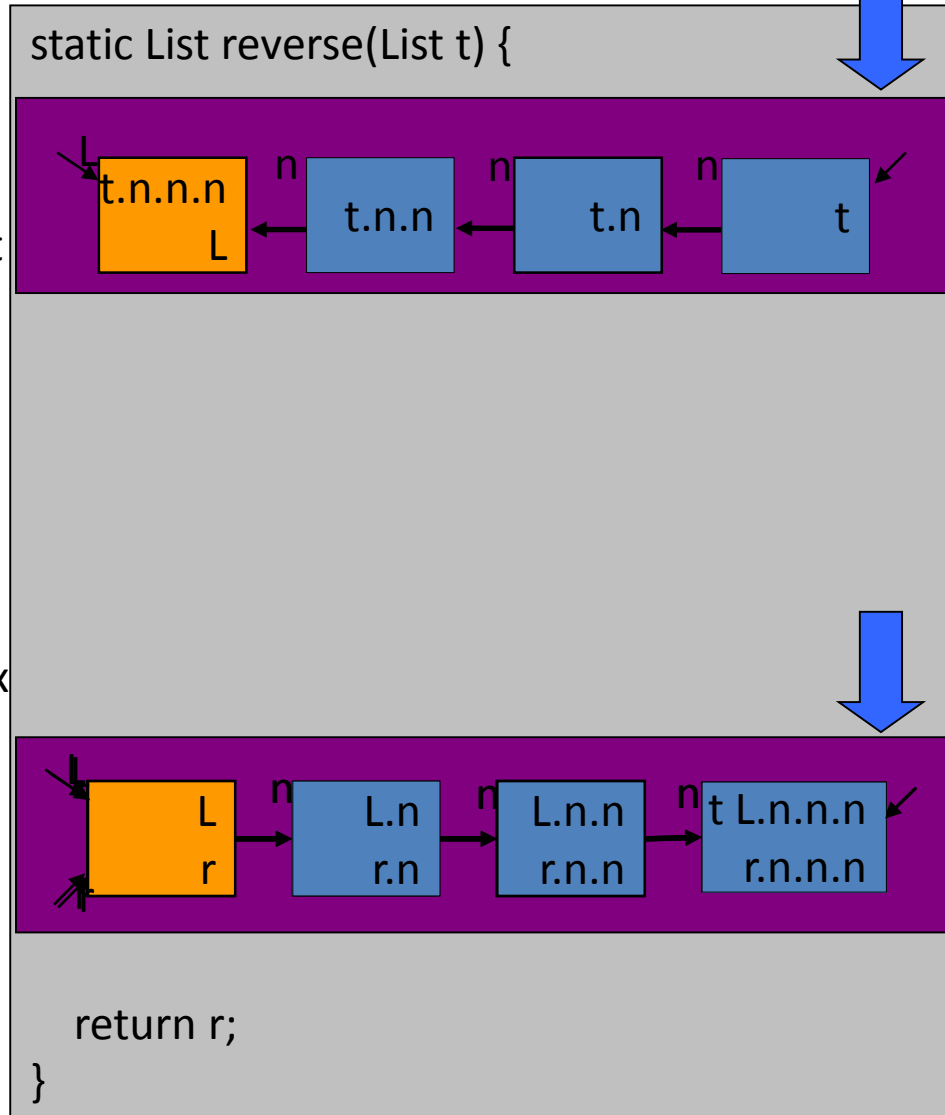
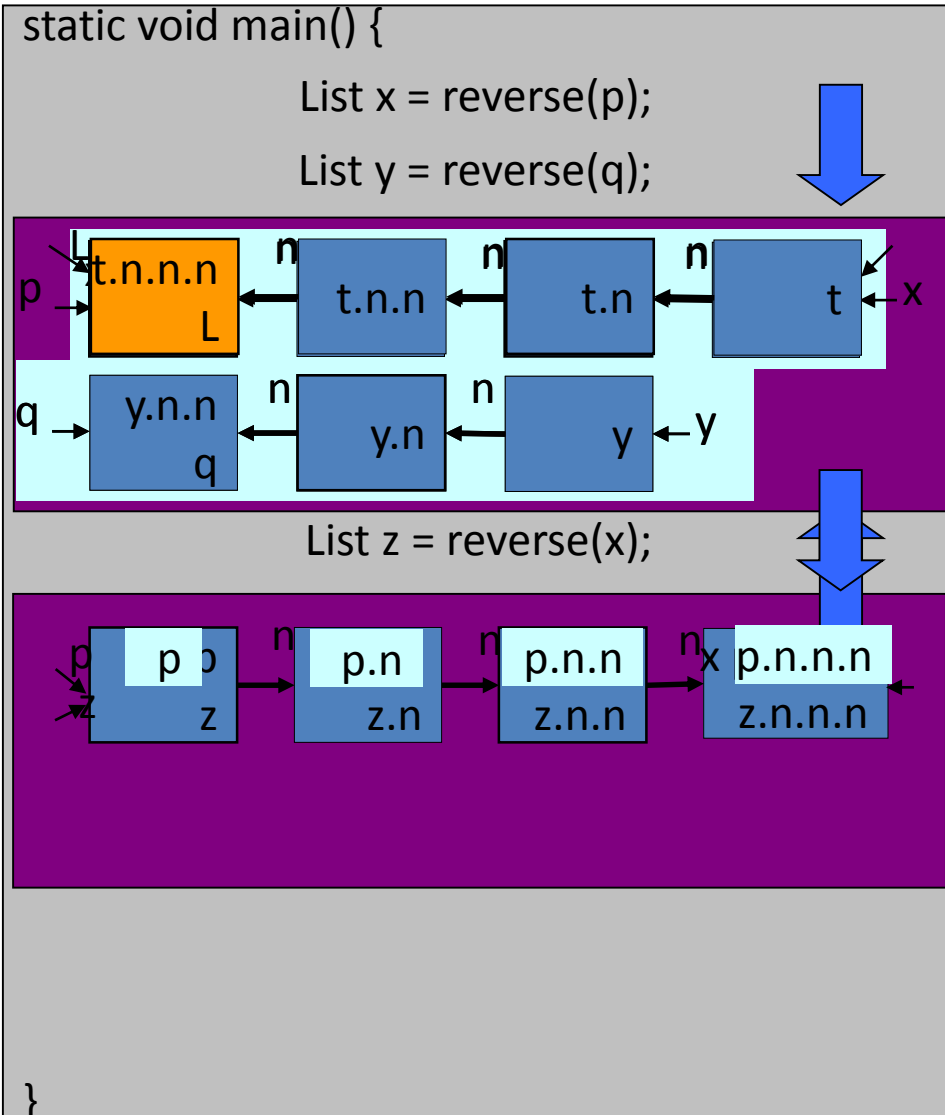
$x = \text{null}$



Example



Example

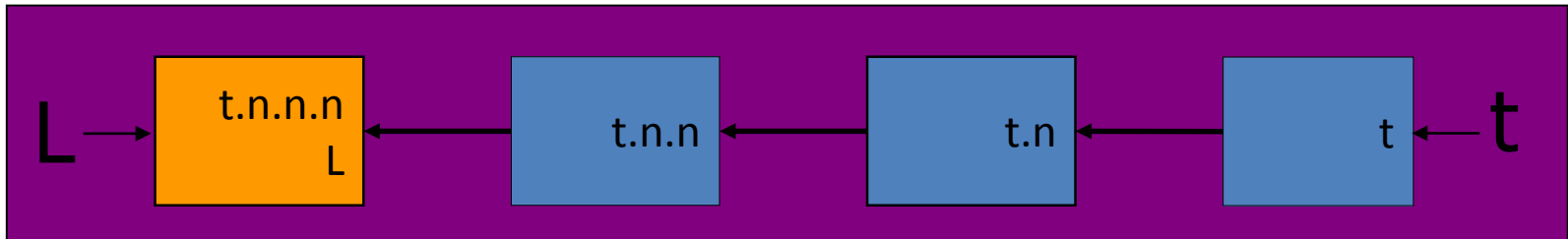


Cutpoint labels

- Relate pre-state with post-state
- Additional roots
- Mark cutpoints at and **throughout** an invocation

Cutpoint labels

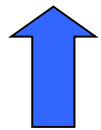
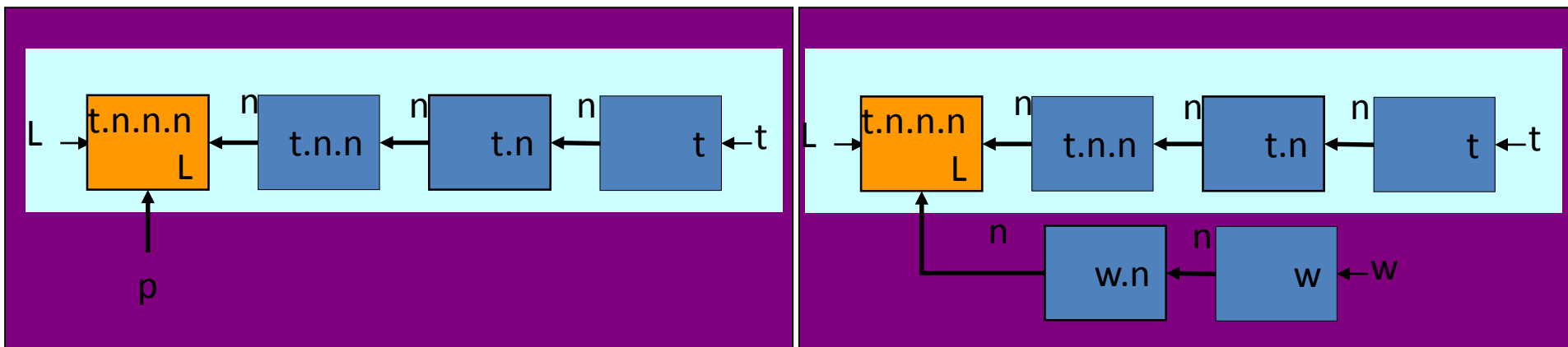
- **Cutpoint label:** the set of access paths that point to a cutpoint
 - when the invoked procedure starts



$$L \equiv \{t.n.n.n\}$$

Sharing patterns

- Cutpoint labels encode **sharing patterns**



Stack sharing



Heap sharing

$$L \equiv \{t.n.n.n\}$$

Observational equivalence

- $\sigma_L \in \Sigma_L$ (Local-heap Storeless Semantics)
- $\sigma_G \in \Sigma_G$ (Global-heap Store-based Semantics)

σ_L and σ_G observationally equivalent

when for every access paths AP_1, AP_2

$$\llbracket AP_1 = AP_2 \rrbracket(\sigma_L) \Leftrightarrow \llbracket AP_1 = AP_2 \rrbracket(\sigma_G)$$

Main theorem: semantic equivalence

- $\sigma_L \in \Sigma_L$ (Local-heap Storeless Semantics)
- $\sigma_G \in \Sigma_G$ (Global-heap Store-based Semantics)
- σ_L and σ_G observationally equivalent

$$\langle st, \sigma_L \rangle \xrightarrow{\text{LSL}} \sigma'_L \Leftrightarrow \langle st, \sigma_G \rangle \xrightarrow{\text{GSB}} \sigma'_G$$

σ'_L and σ'_G are observationally equivalent

Corollaries

- Preservation of invariants
 - Assertions: $AP_1 = AP_2$
- Detection of memory leaks

Applications

- Develop new static analyses
 - Shape analysis
- Justify soundness of existing analyses

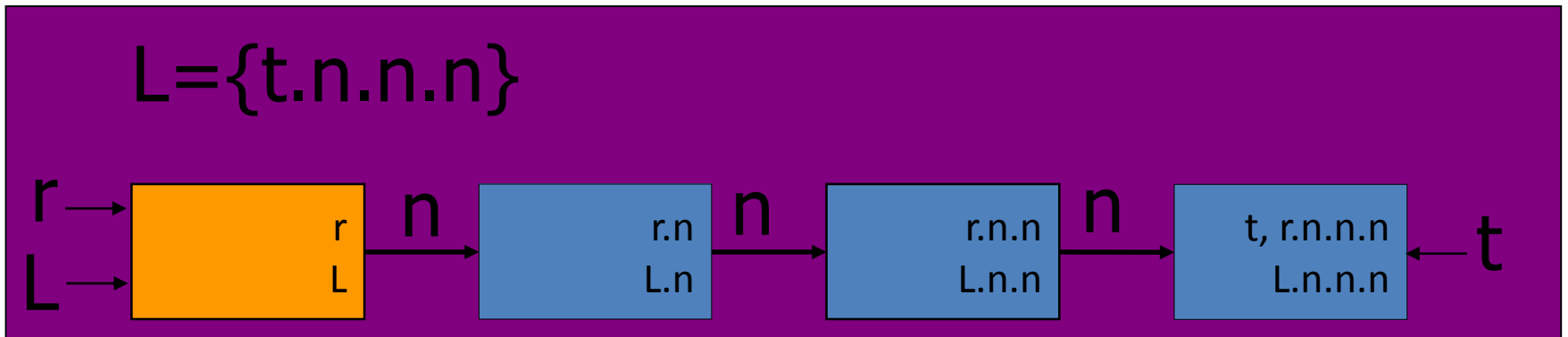
Related work

- **Storeless semantics**
 - Jonkers, Algorithmic Languages '81
 - Deutsch, ICCL '92

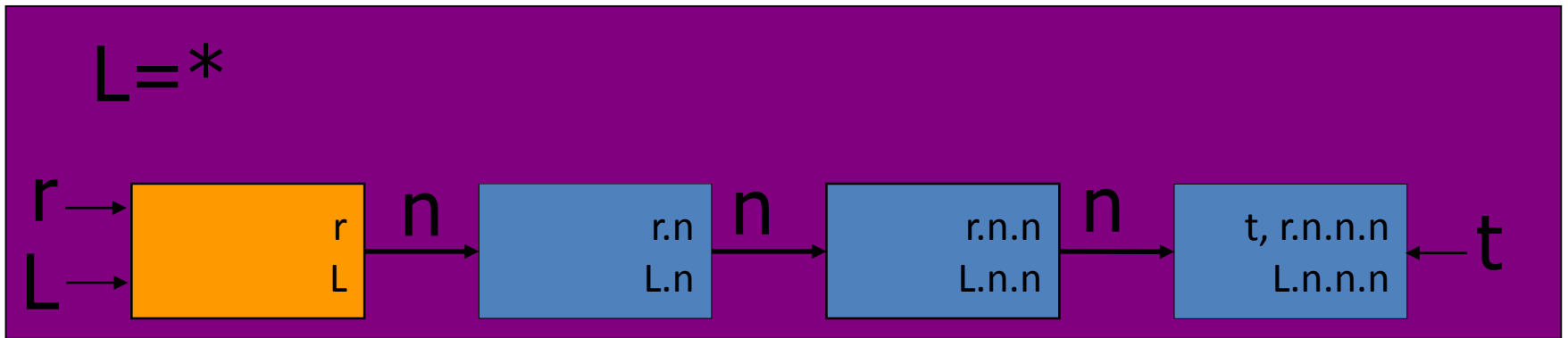
Shape abstraction

- Shape descriptors represent **unbounded** memory states
 - Conservatively
 - In a bounded way
- Two dimensions
 - Local heap (objects)
 - Sharing pattern (cutpoint labels)

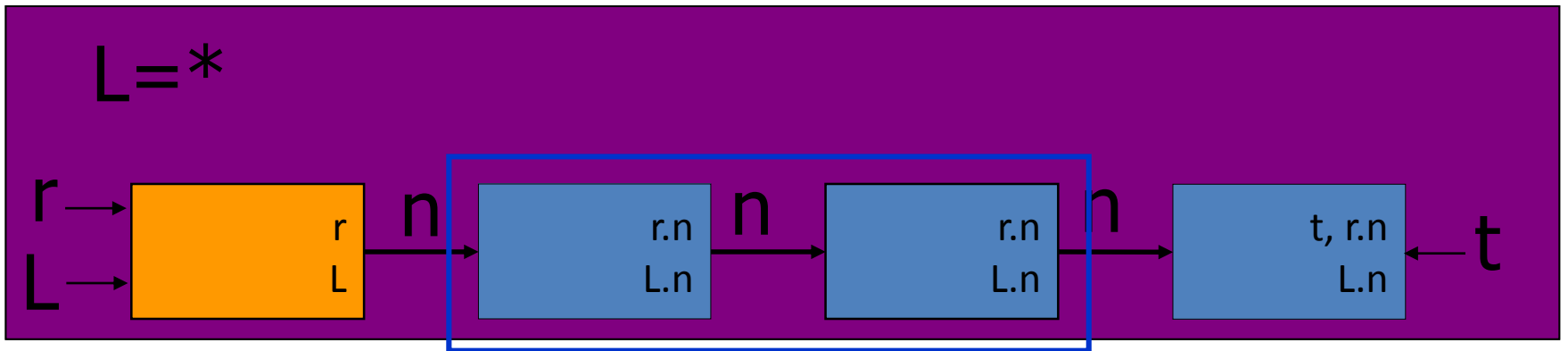
A Shape abstraction



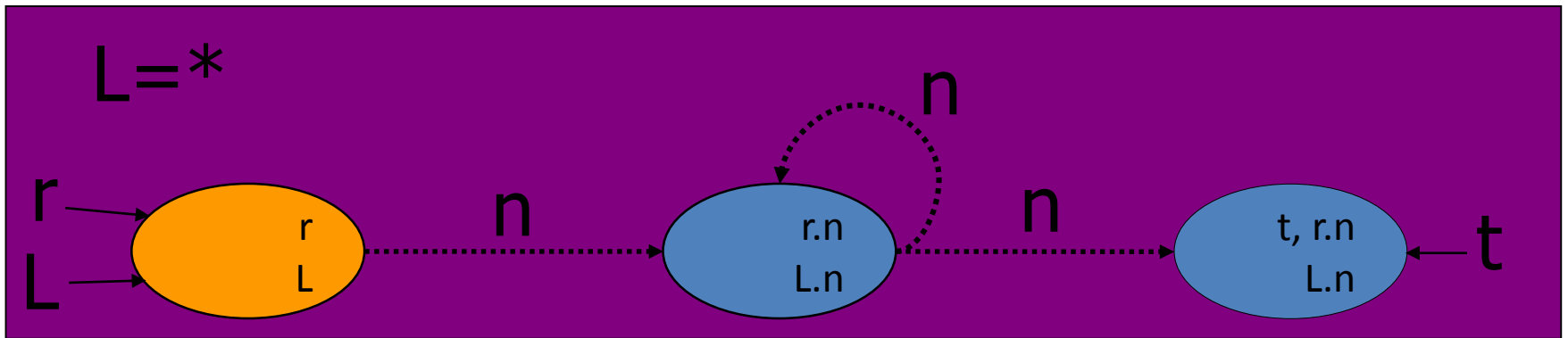
A Shape abstraction



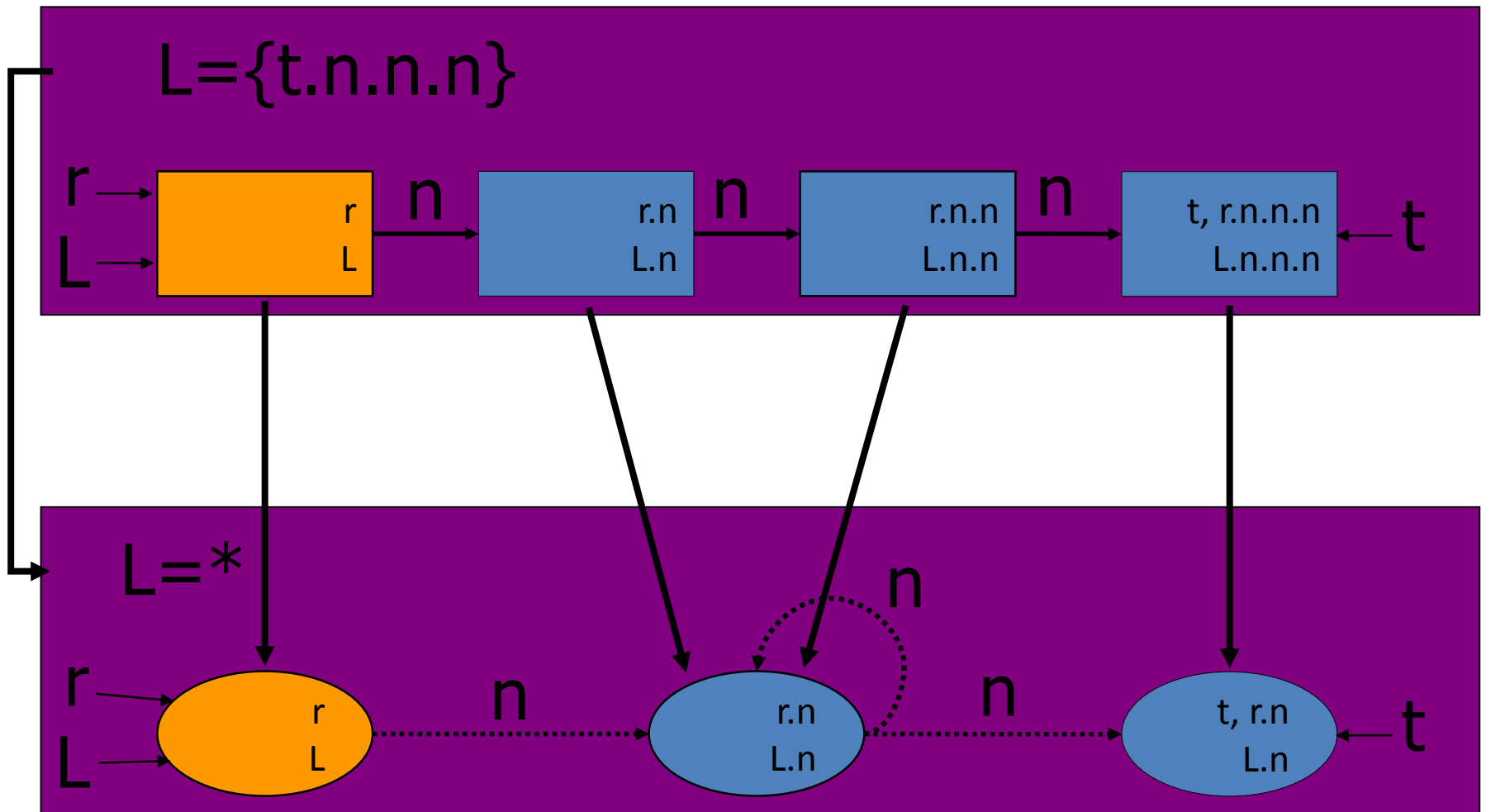
A Shape abstraction



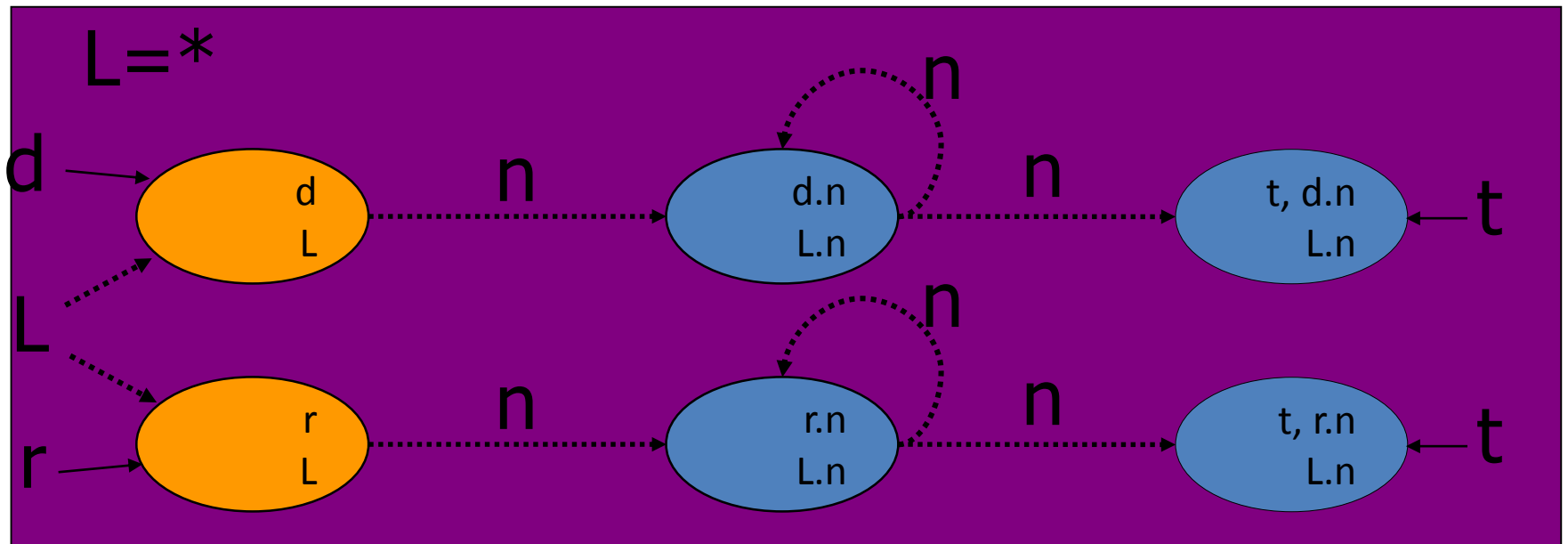
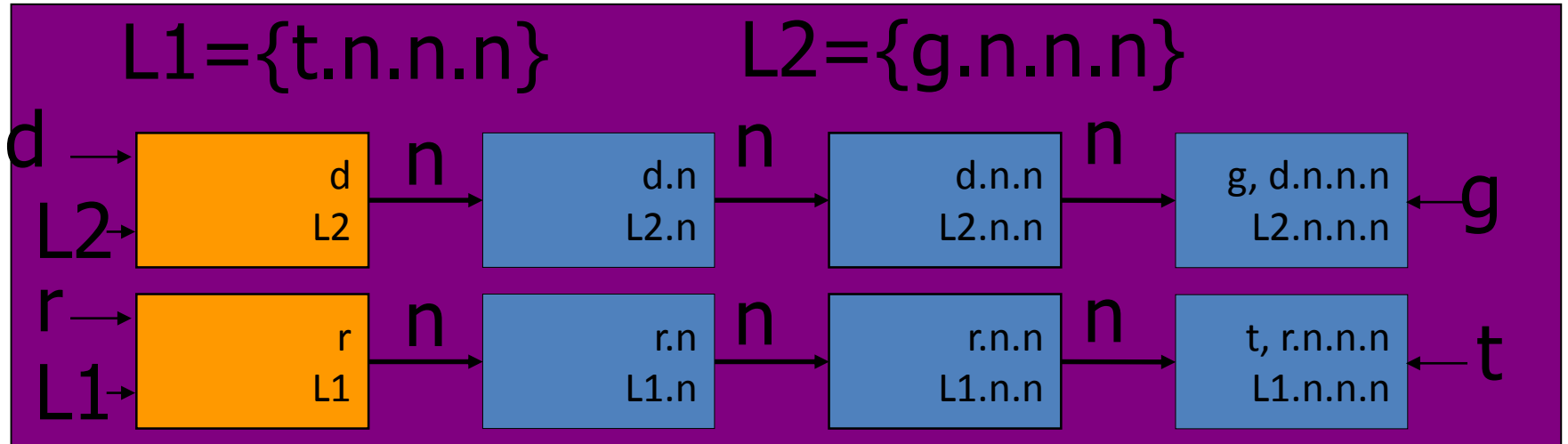
A Shape abstraction



A Shape abstraction



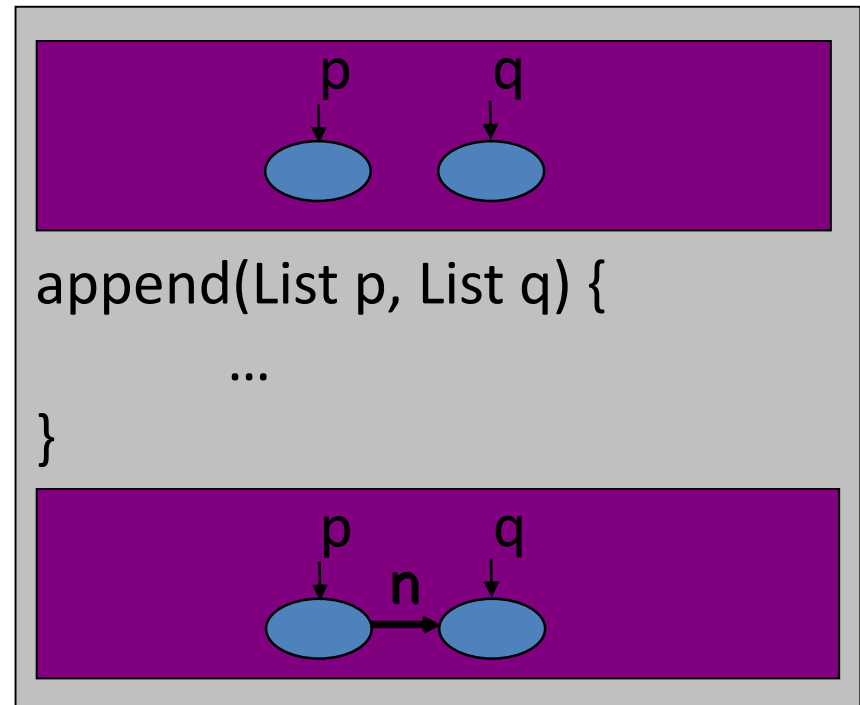
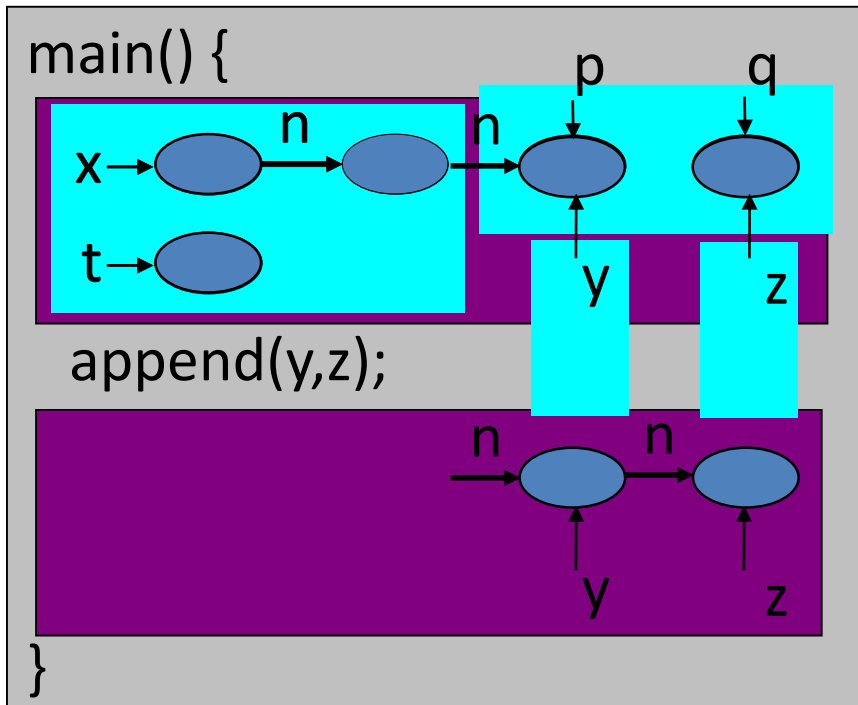
A Shape abstraction



Cutpoint-Freedom

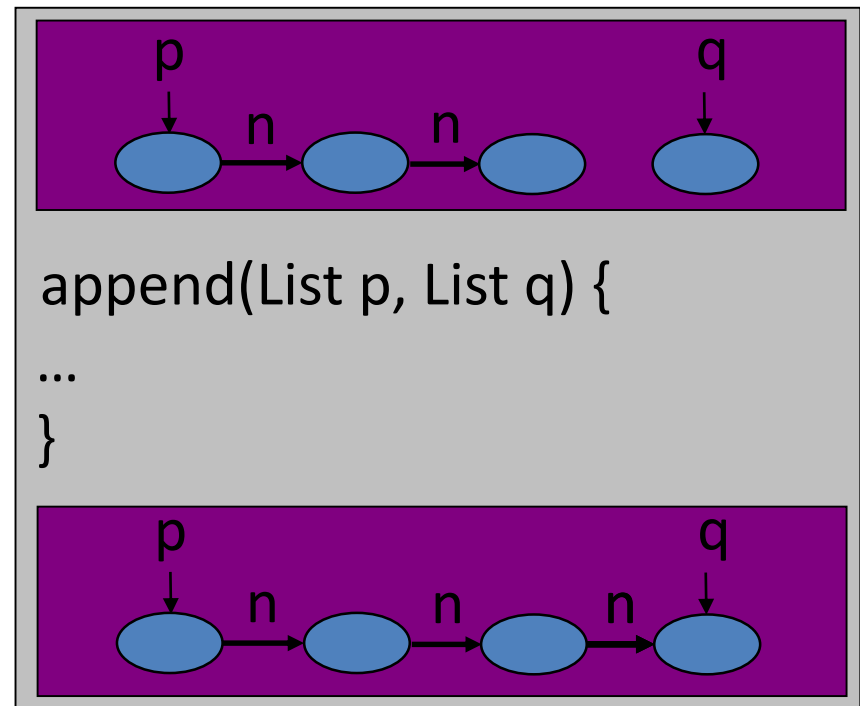
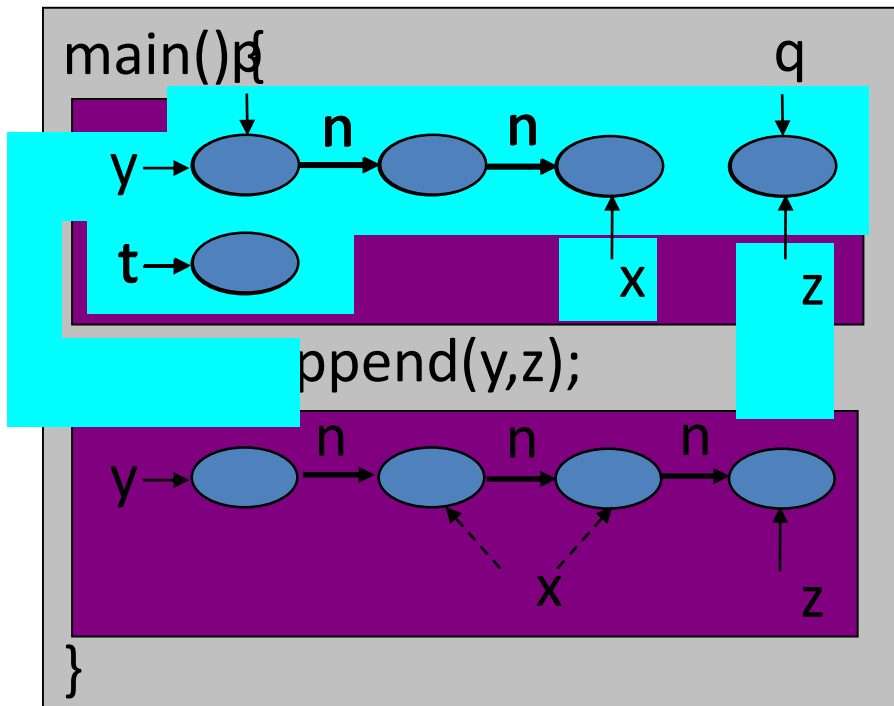
How to tabulate procedures?

- Procedure \equiv input/output relation
 - Not reachable \rightarrow Not effected
 - proc: local (\equiv reachable) heap \rightarrow local heap



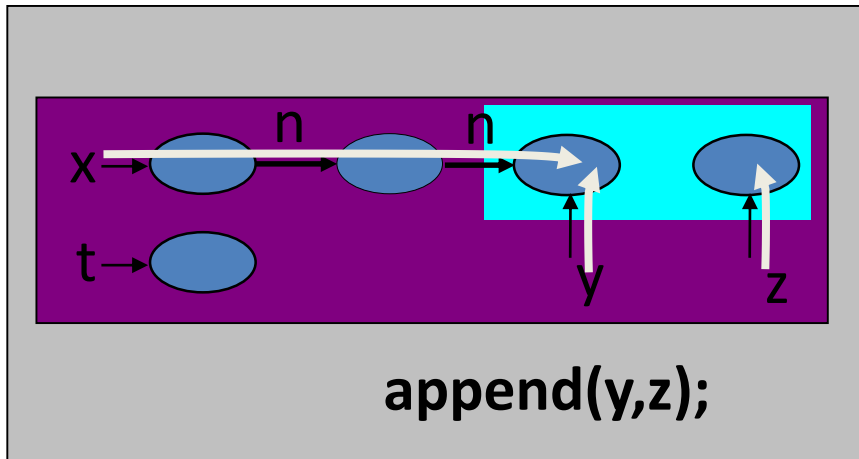
How to handle sharing?

- External sharing may break the functional view

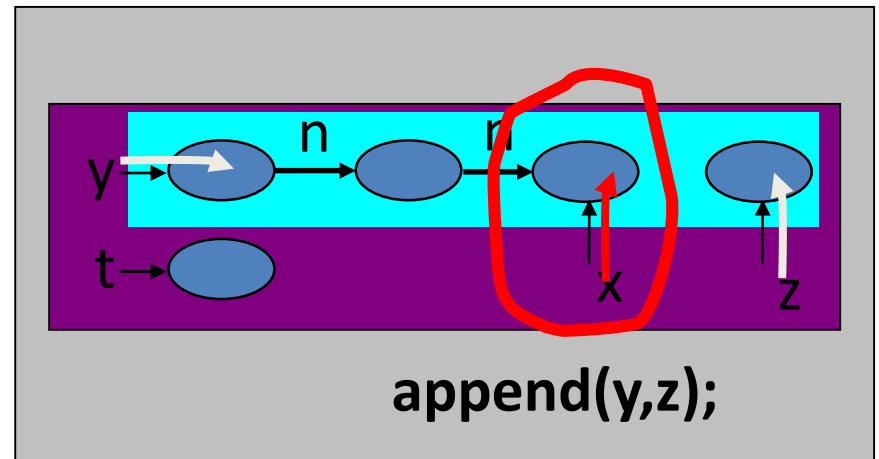


What's the difference?

1st Example



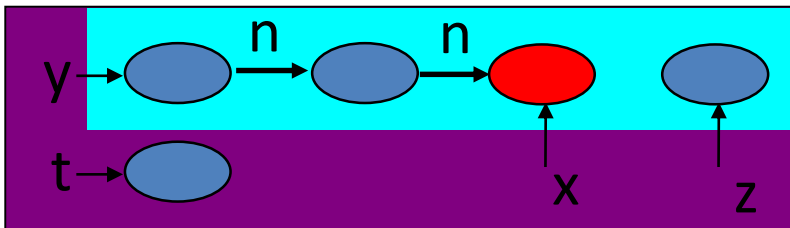
2nd Example



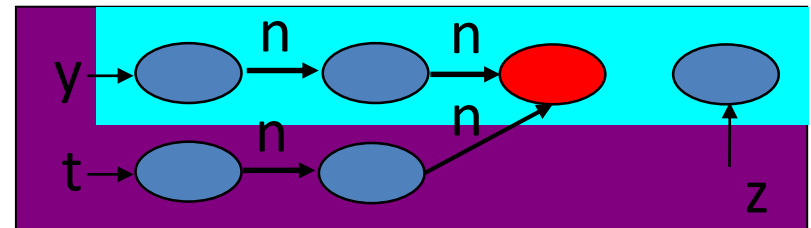
Cutpoints

- An object is a **cutpoint** for an invocation
 - Reachable from actual parameters
 - Not pointed to by an actual parameter
 - Reachable without going through a parameter

append(y,z)



append(y,z)

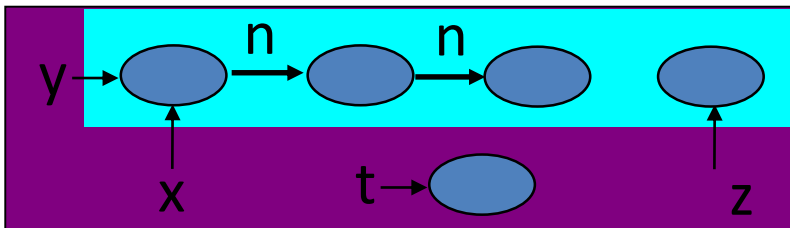


Cutpoint freedom

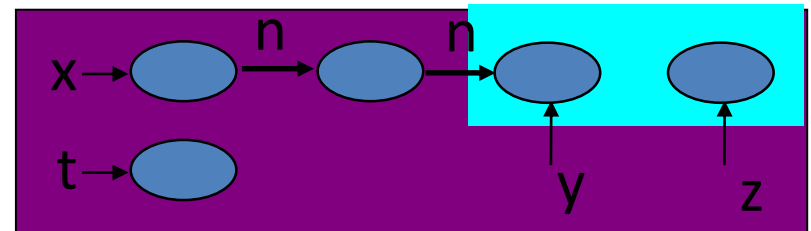
- **Cutpoint-free**

- Invocation: has no cutpoints
- Execution: every invocation is cutpoint-free
- Program: every execution is cutpoint-free

append(y,z)



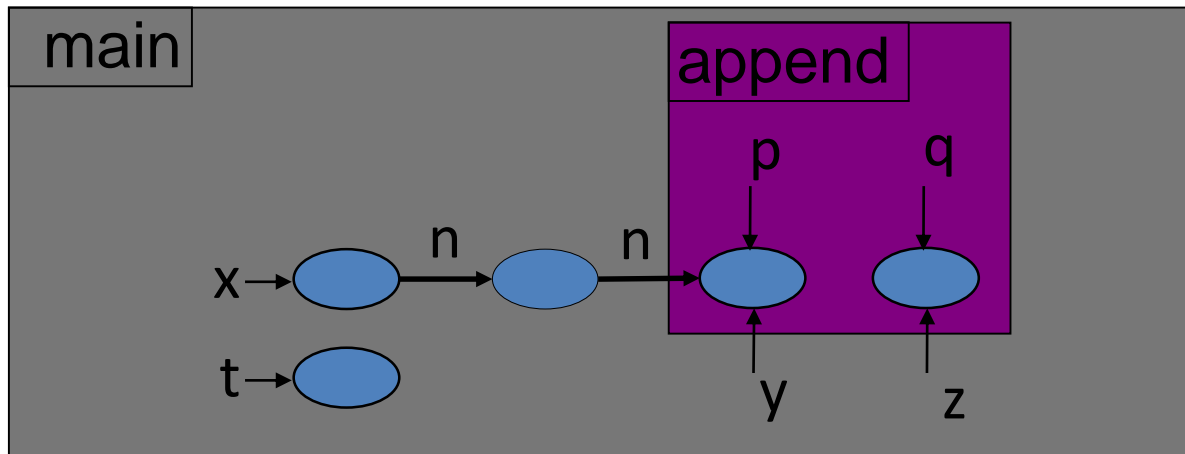
append(y,z)



**Interprocedural shape analysis
for cutpoint-free programs
using 3-Valued Shape Analysis**

Memory states: 2-Valued Logical Structure

- A memory state encodes a **local heap**
 - Local variables of the **current procedure invocation**
 - Relevant part of the heap
 - Relevant \equiv Reachable



Memory states

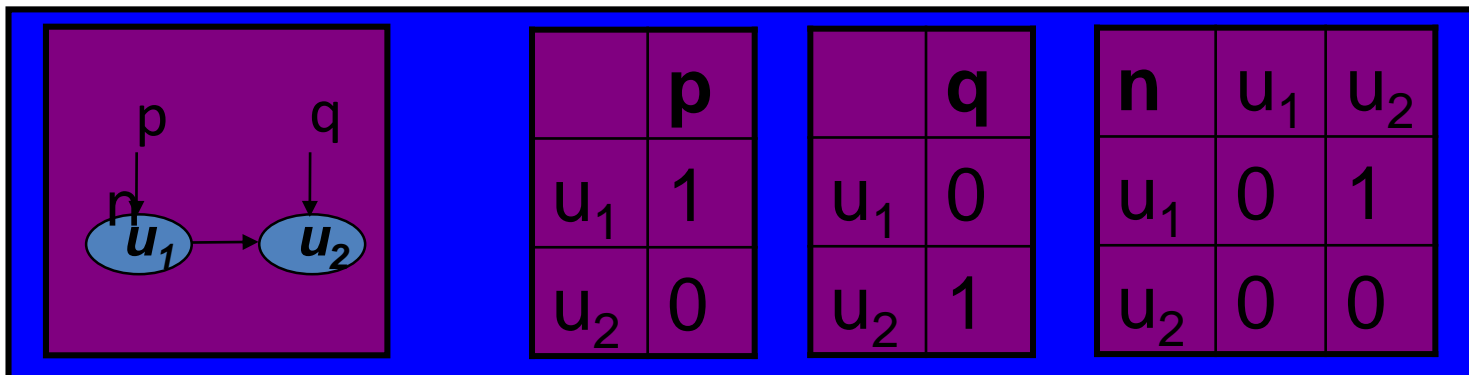
- Represented by first-order logical structures

Predicate	Meaning
$\mathbf{x}(v)$	Variable \mathbf{x} points to v
$\mathbf{n}(v_1, v_2)$	Field \mathbf{n} of object v_1 points to v_2

Memory states

- Represented by first-order logical structures

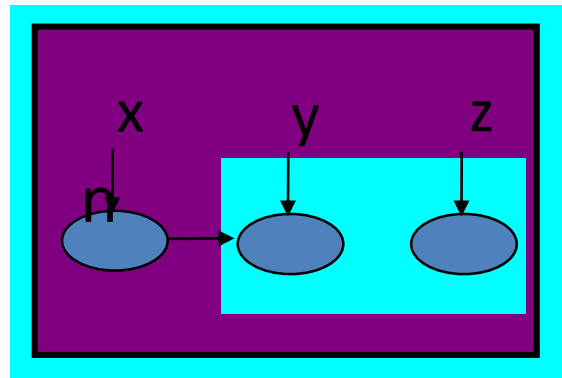
Predicate	Meaning
$\mathbf{x}(v)$	Variable \mathbf{x} points to v
$\mathbf{n}(v_1, v_2)$	Field \mathbf{n} of object v_1 points to v_2



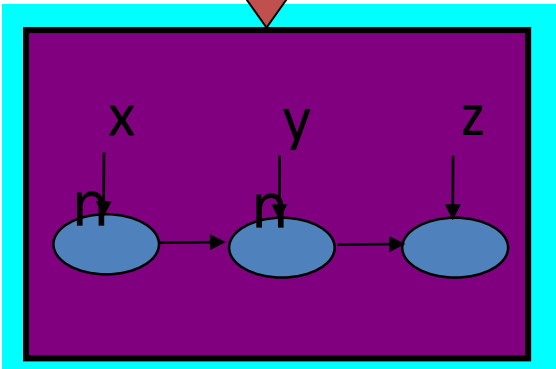
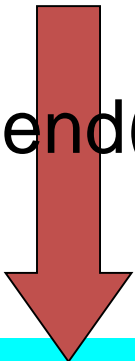
Operational semantics

- Statements modify values of predicates
- Specified by predicate-update formulae
 - Formulae in FO-TC

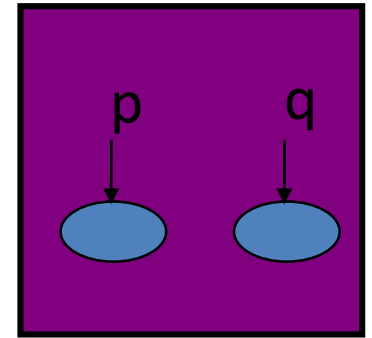
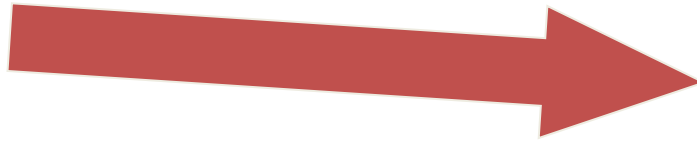
Procedure calls



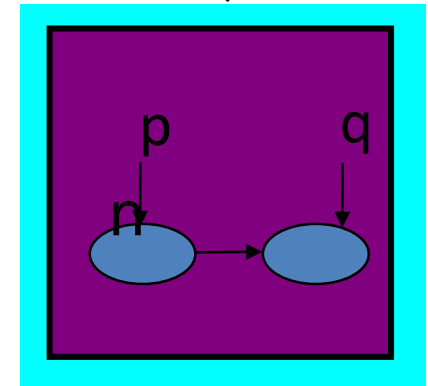
append(y,z)



append(p,q)



append body



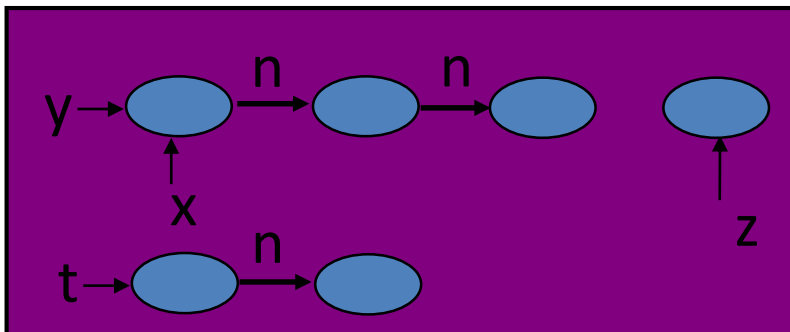
1. Verify cutpoint freedom
 2. Compute input
 3. Combine output
- ... Execute callee ...

Procedure call:

1. Verifying cutpoint-freedom

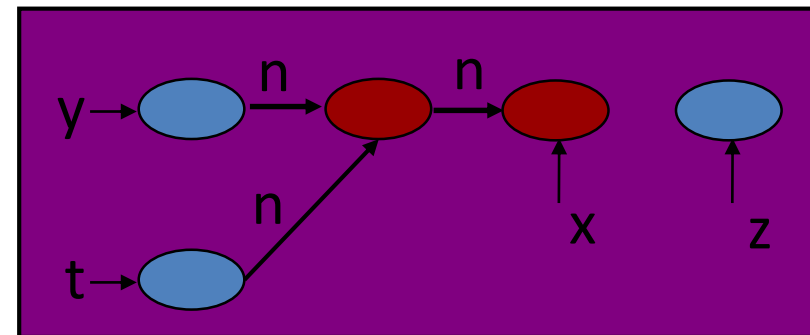
- An object is a **cutpoint** for an invocation
 - Reachable from actual parameters
 - Not pointed to by an actual parameter
 - Reachable without going through a parameter

append(y,z)



Cutpoint free

append(y,z)



Not Cutpoint free

Procedure call:

1. Verifying cutpoint-freedom

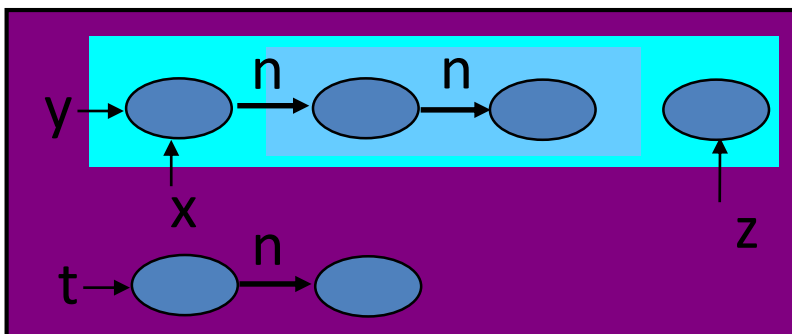
- Invoking `append(y,z)` in `main`

- $R_{\{y,z\}}(v) = \exists v_1: y(v_1) \wedge n^*(v_1, v) \vee \exists v_1: z(v_1) \wedge n^*(v_1, v)$

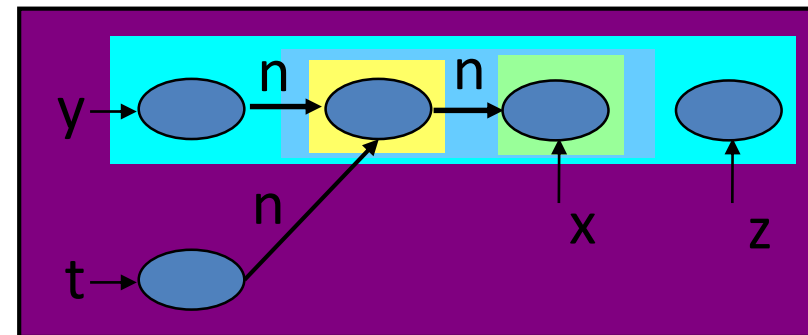
- $\text{isCP}_{\text{main},\{y,z\}}(v) = R_{\{y,z\}}(v) \wedge (\neg y(v) \wedge \neg z(v_1)) \wedge$

- $(x(v) \vee t(v) \vee \exists v_1: \neg R_{\{y,z\}}(v_1) \wedge n(v_1, v))$

(main's locals: x,y,z,t)



Cutpoint free



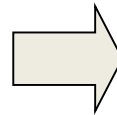
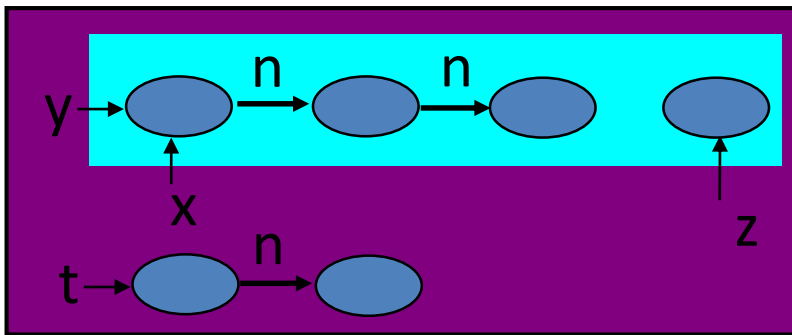
Not Cutpoint free

Procedure call:

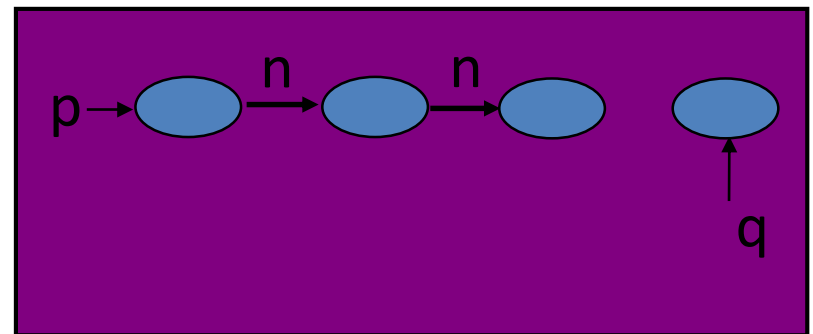
2. Computing the input local heap

- Retain only reachable objects
- Bind formal parameters

Call state

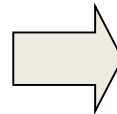
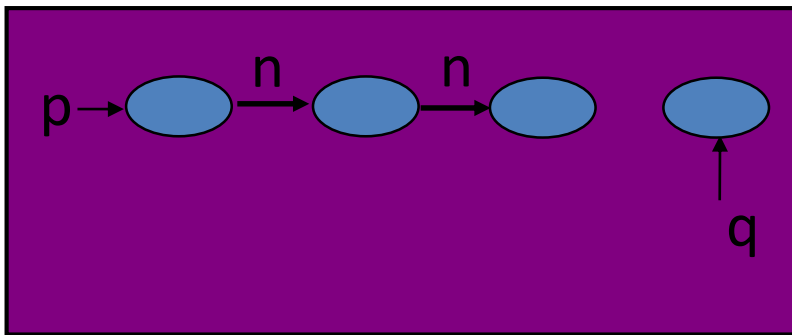


Input state

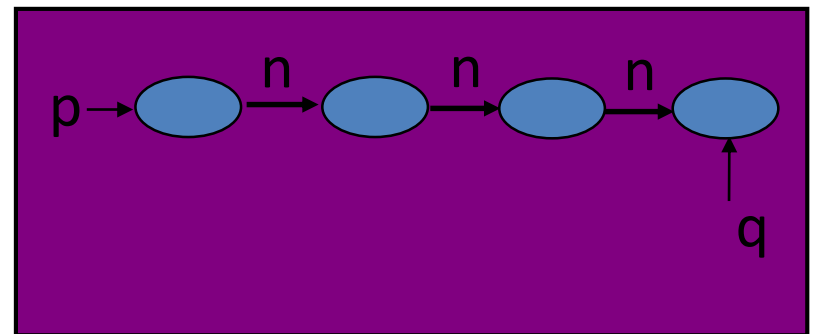


Procedure body: `append(p,q)`

Input state



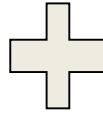
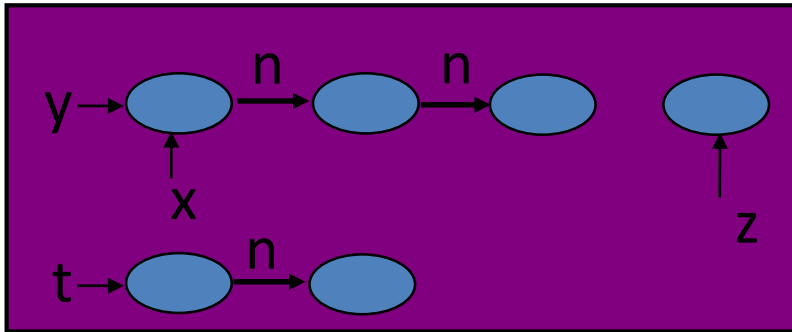
Output state



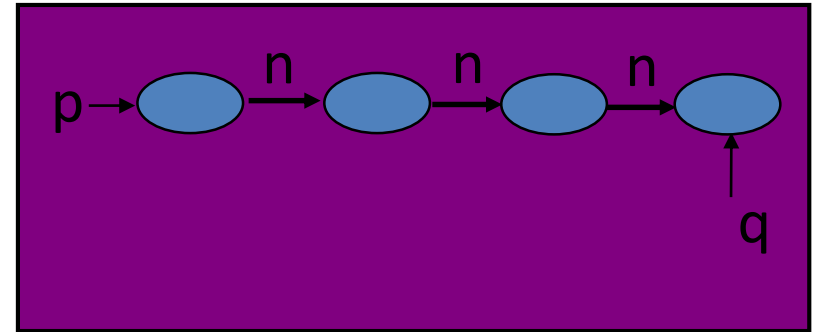
Procedure call:

3. Combine output

Call state

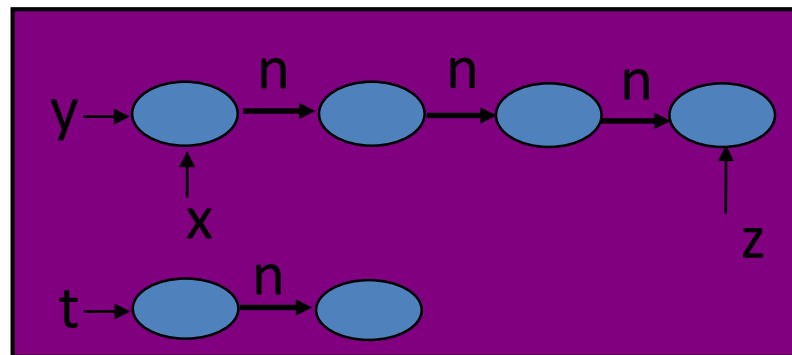
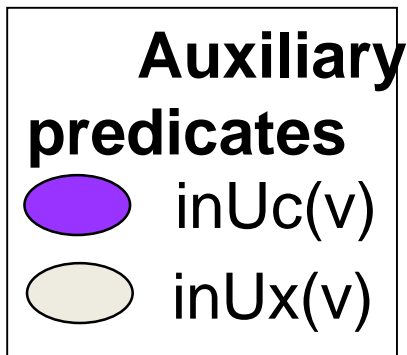
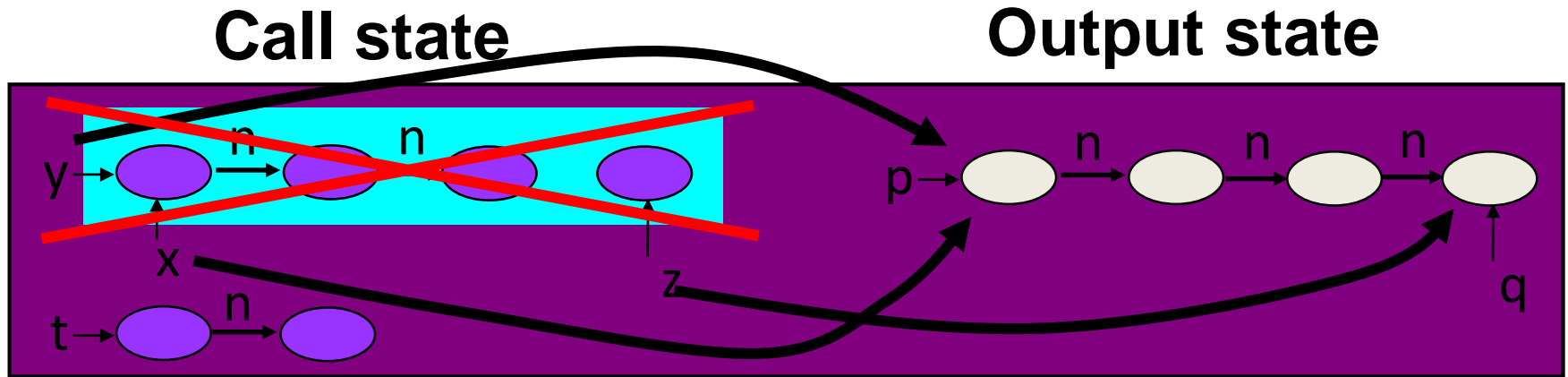


Output state



Procedure call:

3. Combine output



Observational equivalence

- $\sigma_{\text{CPF}} \in \Sigma_{\text{CPF}}$ (Cutpoint free semantics)
- $\sigma_{\text{GSB}} \in \Sigma_{\text{GSB}}$ (Standard semantics)

σ_{CPF} and σ_{GSB} **observationally equivalent**

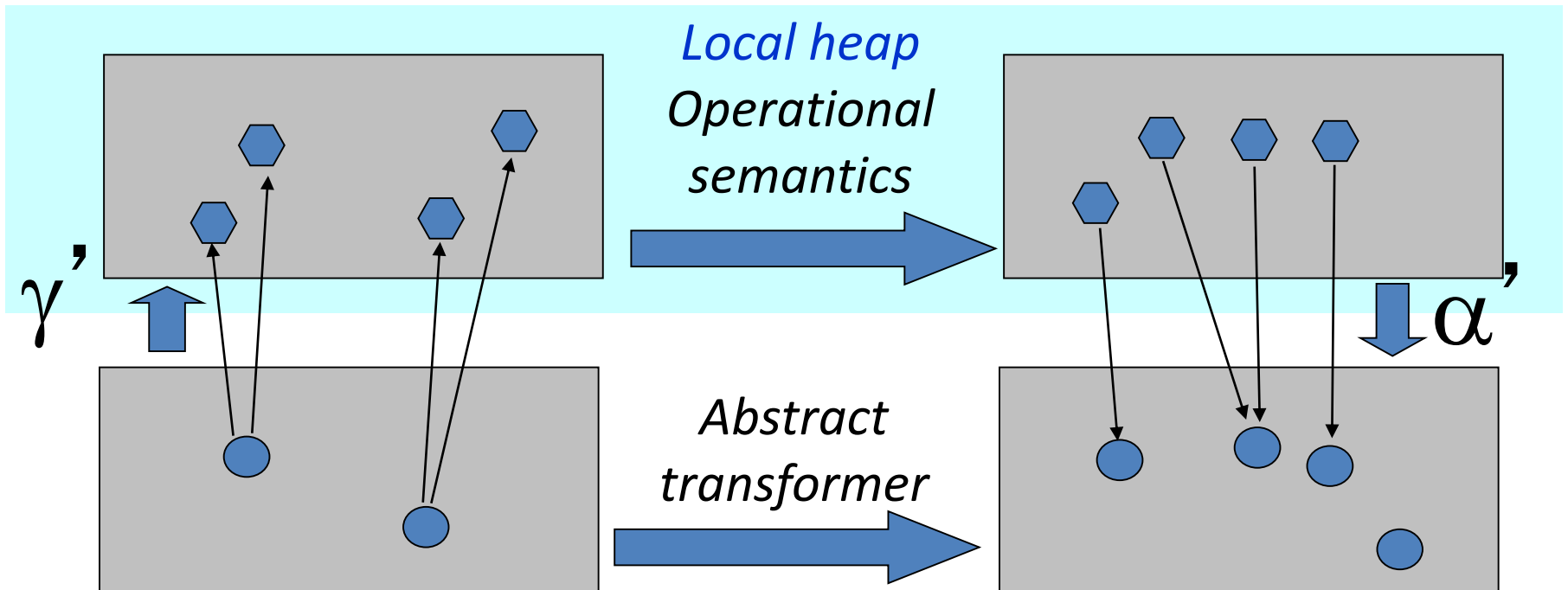
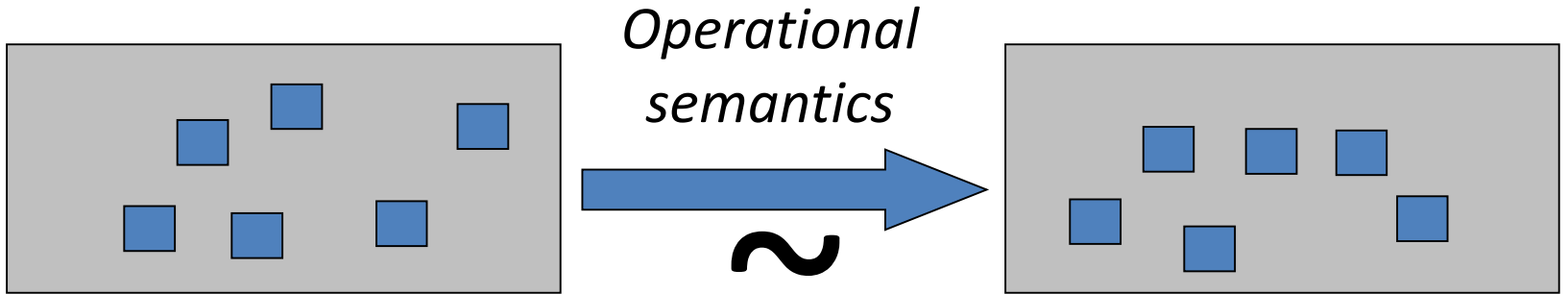
when for every access paths AP_1, AP_2

$$\llbracket AP_1 = AP_2 \rrbracket(\sigma_{\text{CPF}}) \Leftrightarrow \llbracket AP_1 = AP_2 \rrbracket(\sigma_{\text{GSB}})$$

Observational equivalence

- For cutpoint free programs:
 - $\sigma_{\text{CPF}} \in \Sigma_{\text{CPF}}$ (Cutpoint free semantics)
 - $\sigma_{\text{GSB}} \in \Sigma_{\text{GSB}}$ (Standard semantics)
 - σ_{CPF} and σ_{GSB} observationally equivalent
- It holds that
 - $\langle st, \sigma_{\text{CPF}} \rangle \rightsquigarrow \sigma'_{\text{CPF}} \Leftrightarrow \langle st, \sigma_{\text{GSB}} \rangle \rightsquigarrow \sigma'_{\text{GSB}}$
 - σ'_{CPF} and σ'_{GSB} are observationally equivalent

Introducing local heap semantics



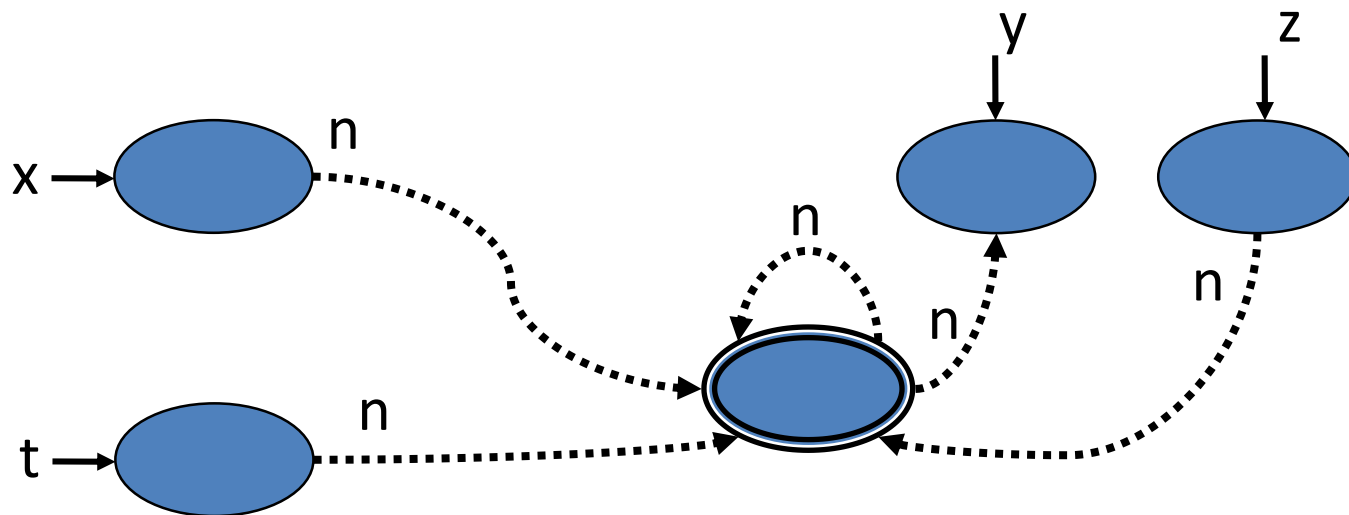
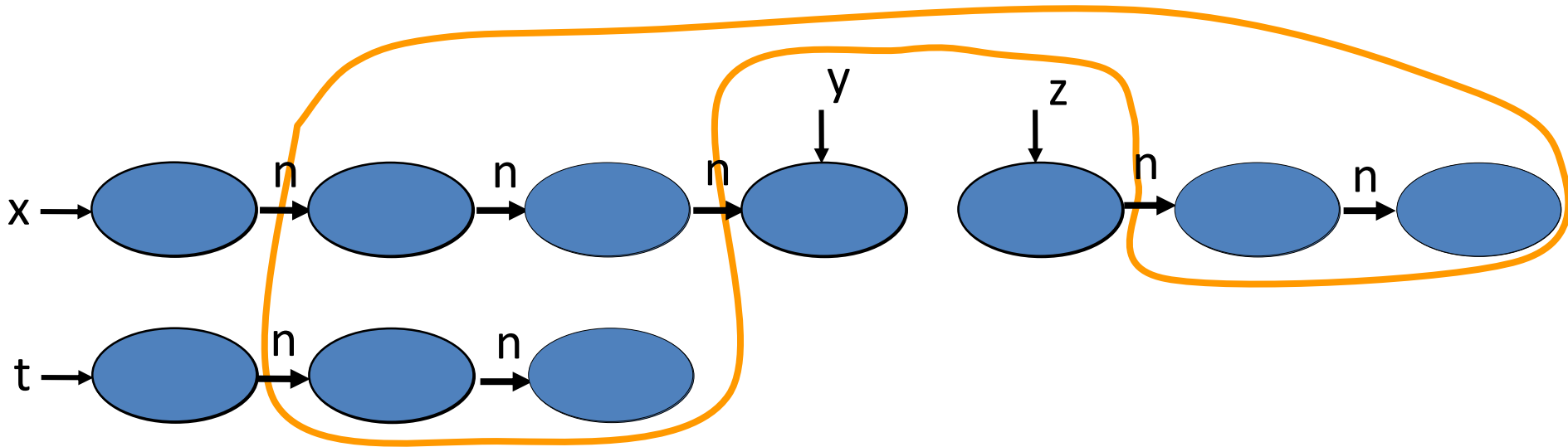
Shape abstraction

- Abstract memory states represent **unbounded** concrete memory states
 - Conservatively
 - In a bounded way
 - Using 3-valued logical structures

3-Valued logic

- $1 = \text{true}$
- $0 = \text{false}$
- $1/2 = \text{unknown}$
- A join semi-lattice, $0 \sqcup 1 = 1/2$

Canonical abstraction

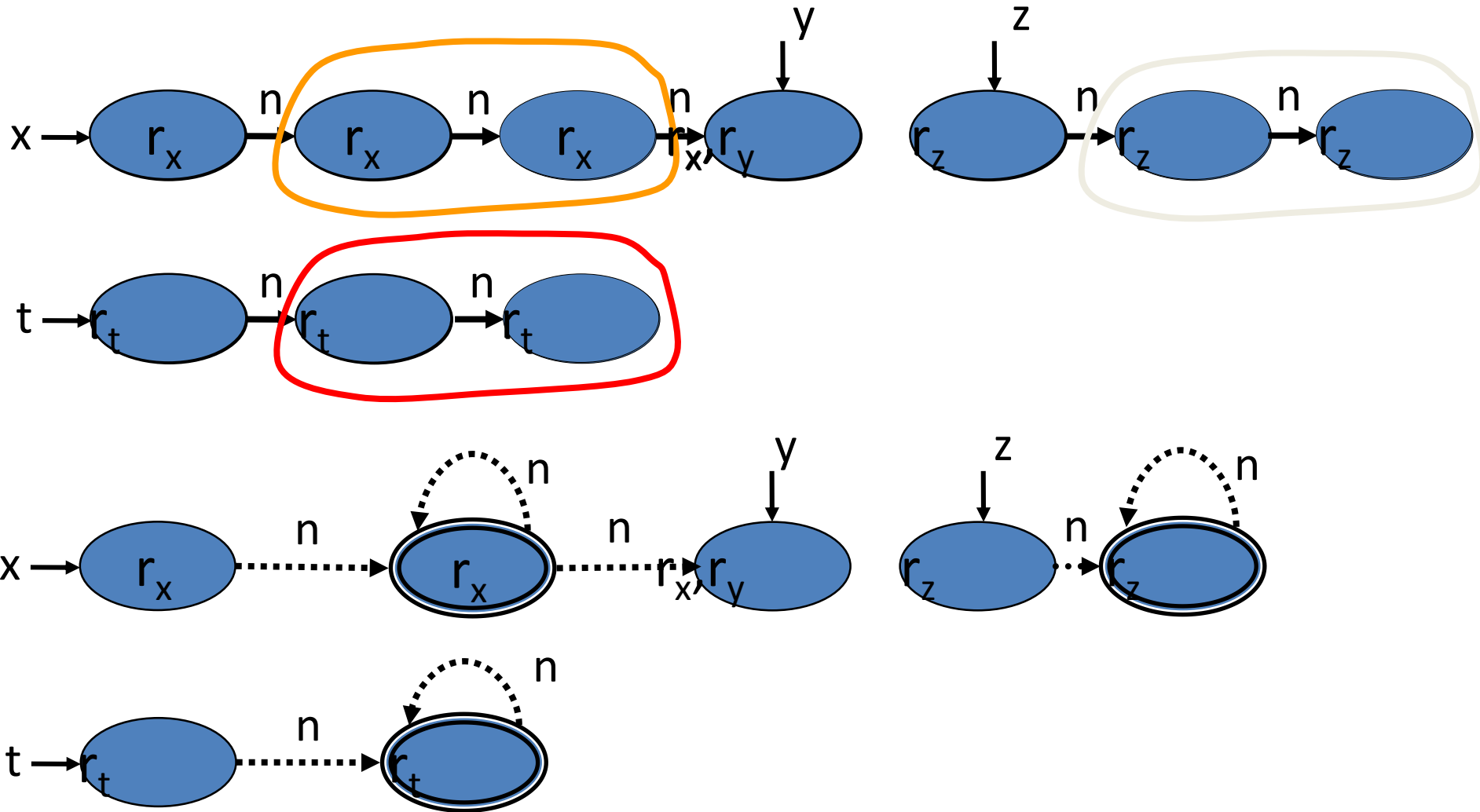


Instrumentation predicates

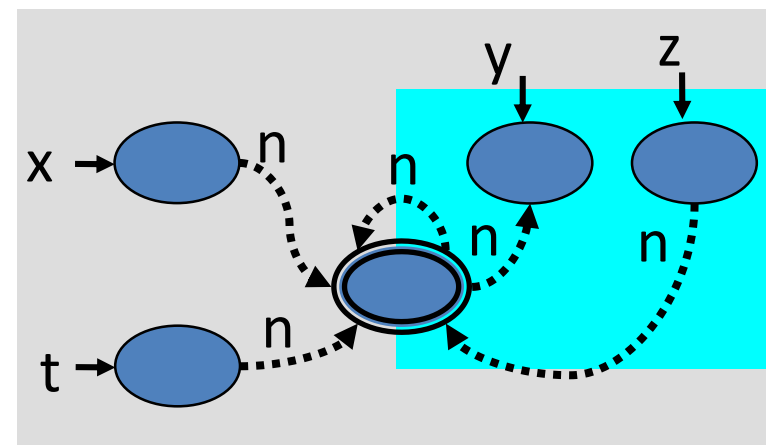
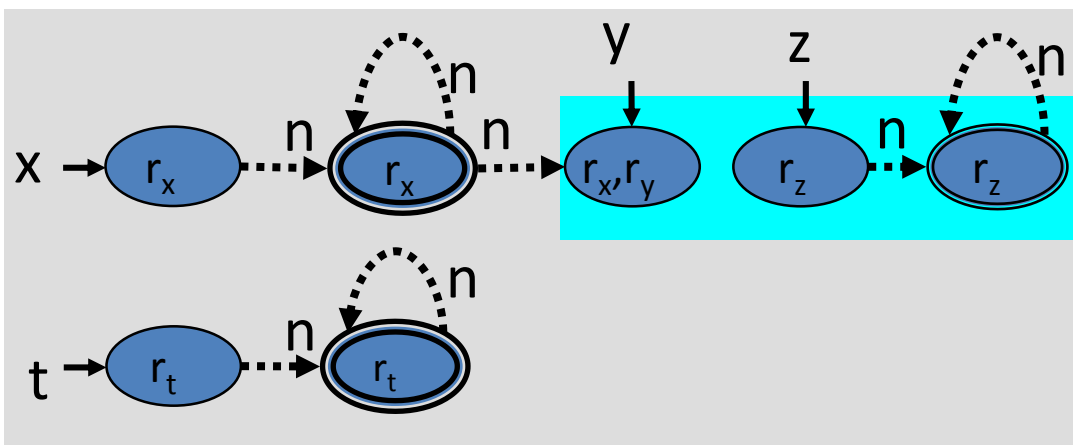
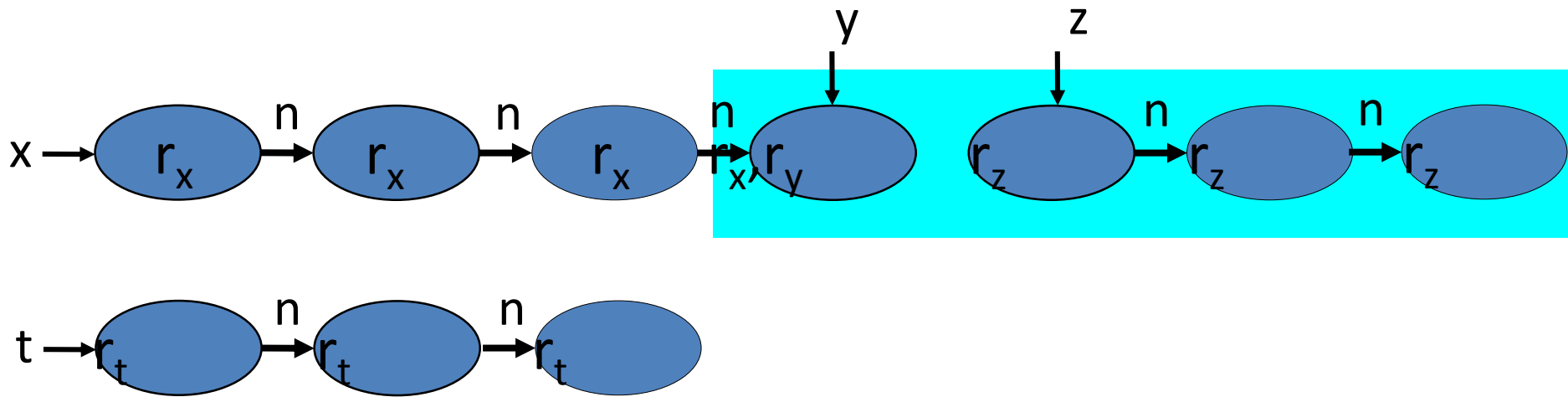
- Record derived properties
- Refine the abstraction
 - Instrumentation principle [SRW, TOPLAS'02]
- Reachability is central!

Predicate	Meaning
$r_x(v)$	v is reachable from variable x
$r_{\text{obj}}(v_1, v_2)$	v_2 is reachable from v_1
$\text{ils}(v)$	v is heap-shared
$c(v)$	v resides on a cycle

Abstract memory states (with reachability)



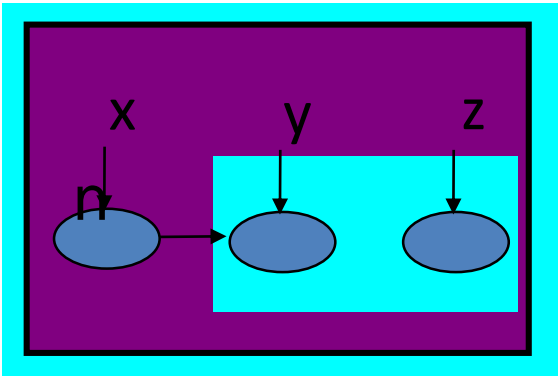
The importance of reachability: Call append(y,z)



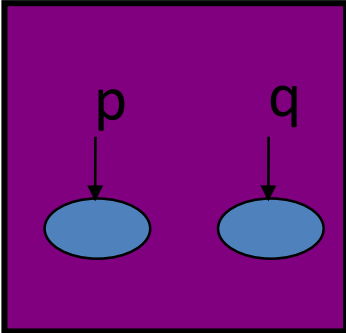
Abstract semantics

- Conservatively apply statements on abstract memory states
 - Same formulae as in concrete semantics
 - Soundness guaranteed [SRW, TOPLAS'02]

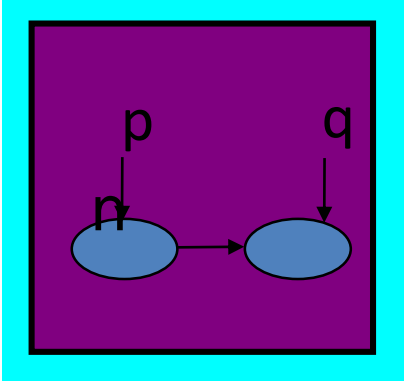
Procedure calls



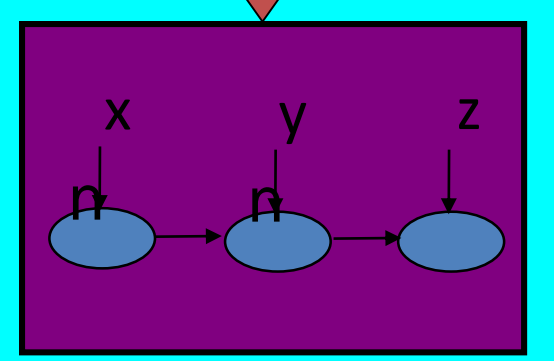
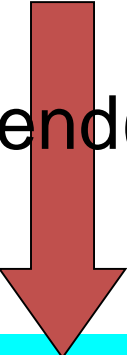
append(p,q)



append body



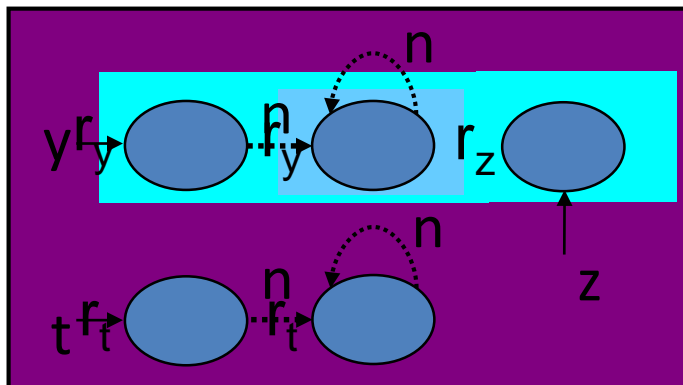
append(y,z)



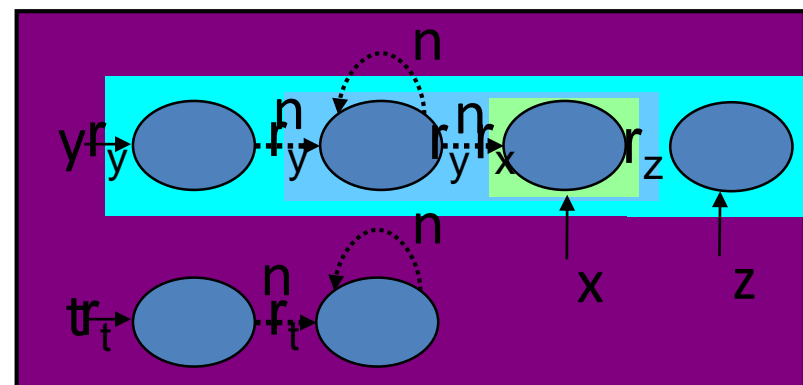
- 1. Verify cutpoint freedom
Compute input 2
- 2 ... Execute callee ...
- 3 Combine output

Conservative verification of cutpoint-freeness

- Invoking `append(y,z)` in main
 - $R_{\{y,z\}}(v) = \exists v_1: y(v_1) \wedge n^*(v_1, v) \vee \exists v_1: z(v_1) \wedge n^*(v_1, v)$
 - $\text{isCP}_{\text{main}, \{y,z\}}(v) = R_{\{y,z\}}(v) \wedge (\neg y(v) \wedge \neg z(v_1)) \wedge (x(v) \vee t(v) \vee \exists v_1: \neg R_{\{y,z\}}(v_1) \wedge n(v_1, v))$

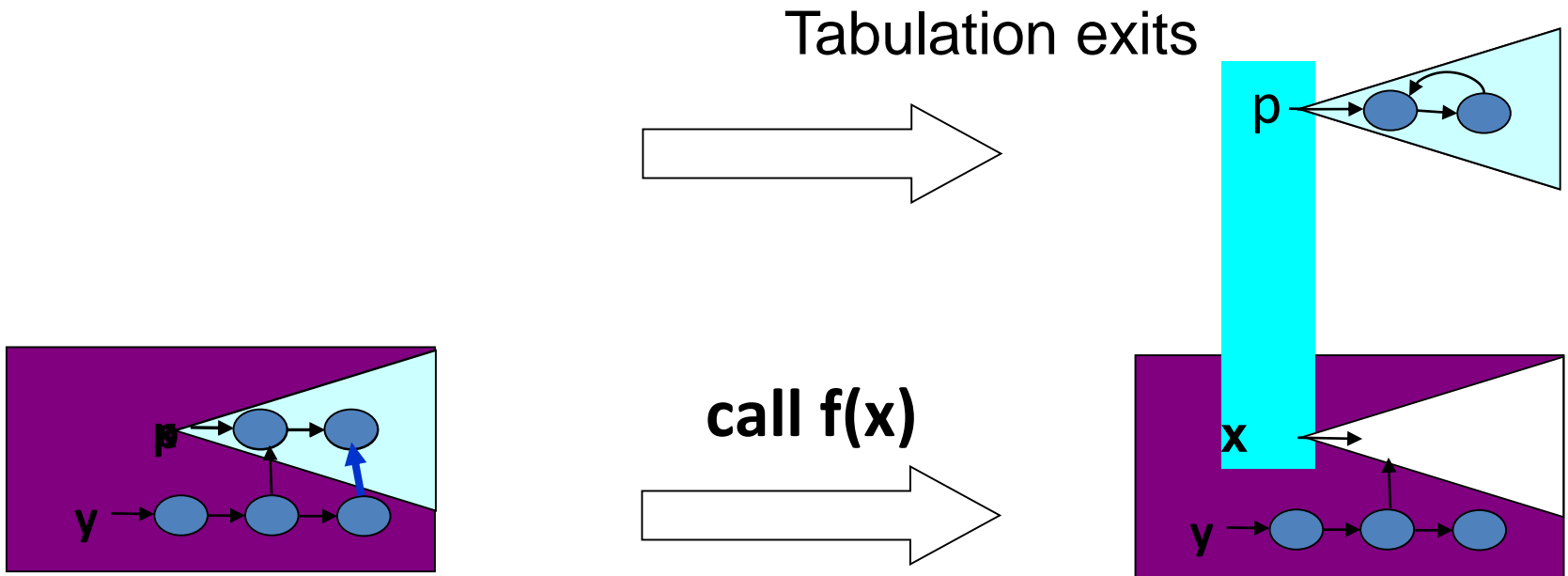


Cutpoint free



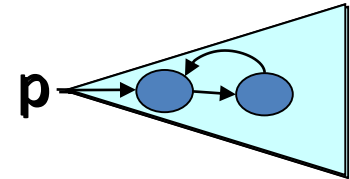
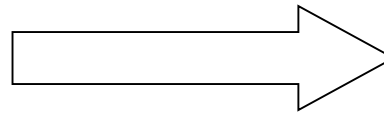
Not Cutpoint free

Interprocedural shape analysis

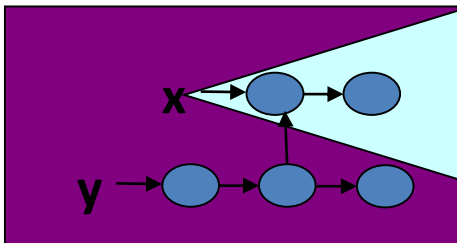
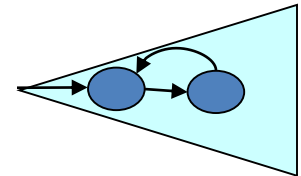
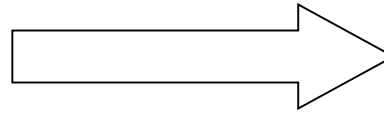


Interprocedural shape analysis

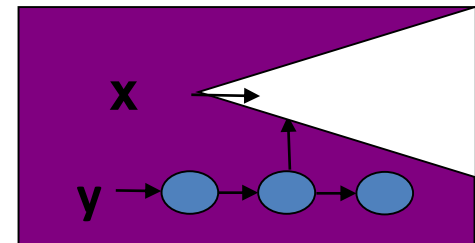
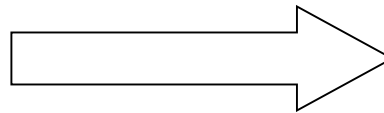
Analyze f



Tabulation exits

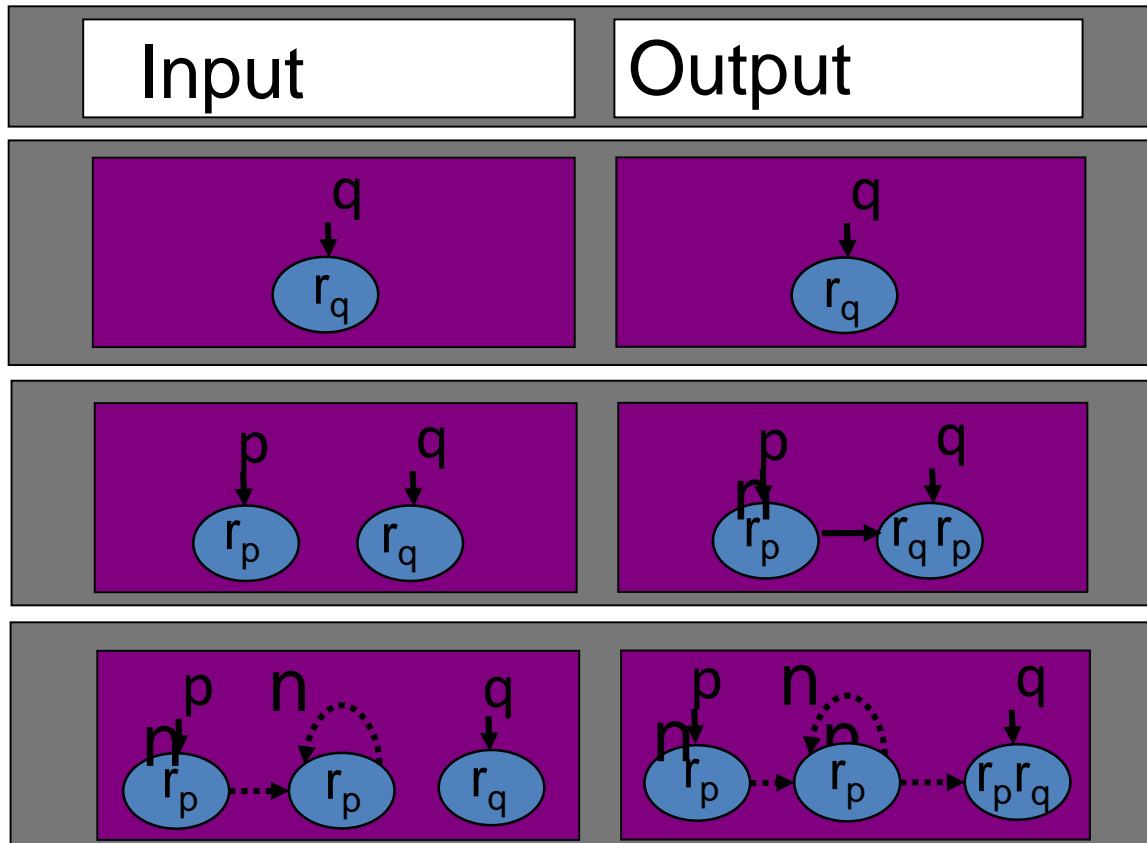


call f(x)



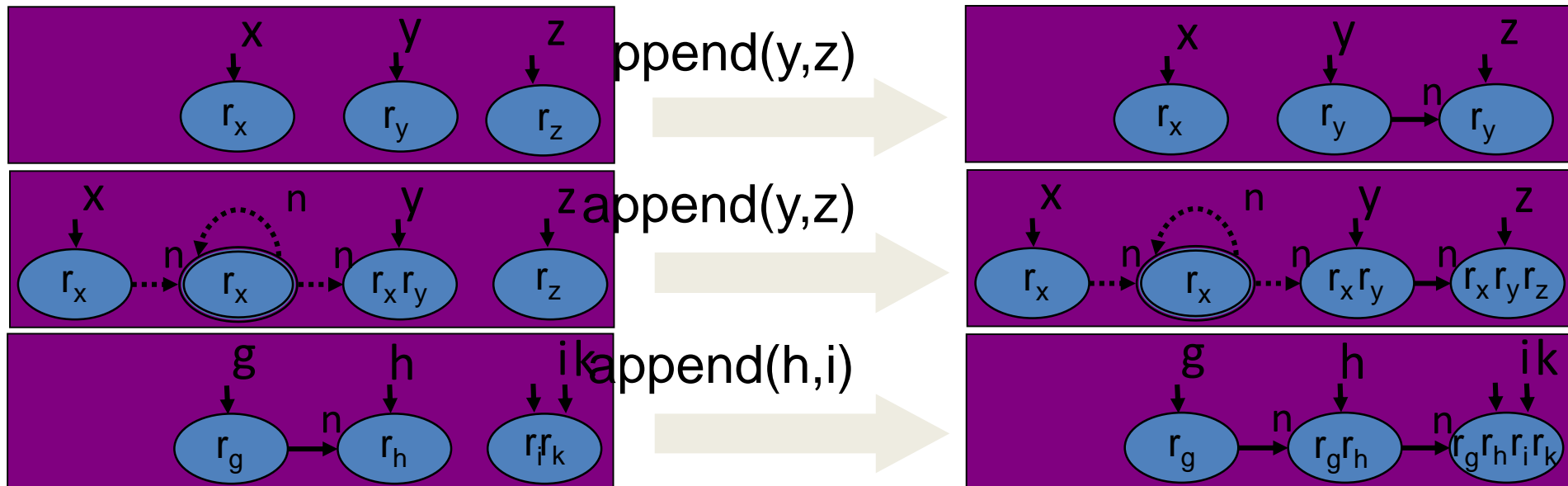
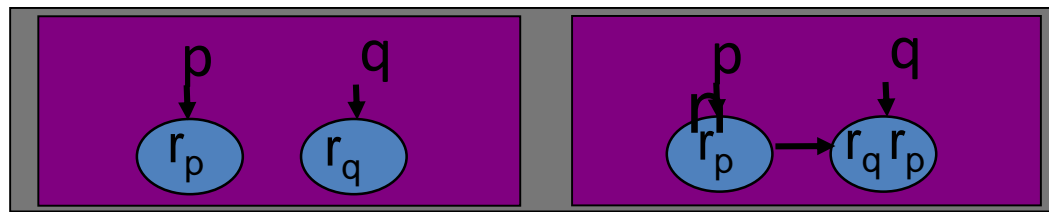
Interprocedural shape analysis

- Procedure \equiv input/output relation



Interprocedural shape analysis

- Reusable procedure summaries
 - Heap modularity



Plan

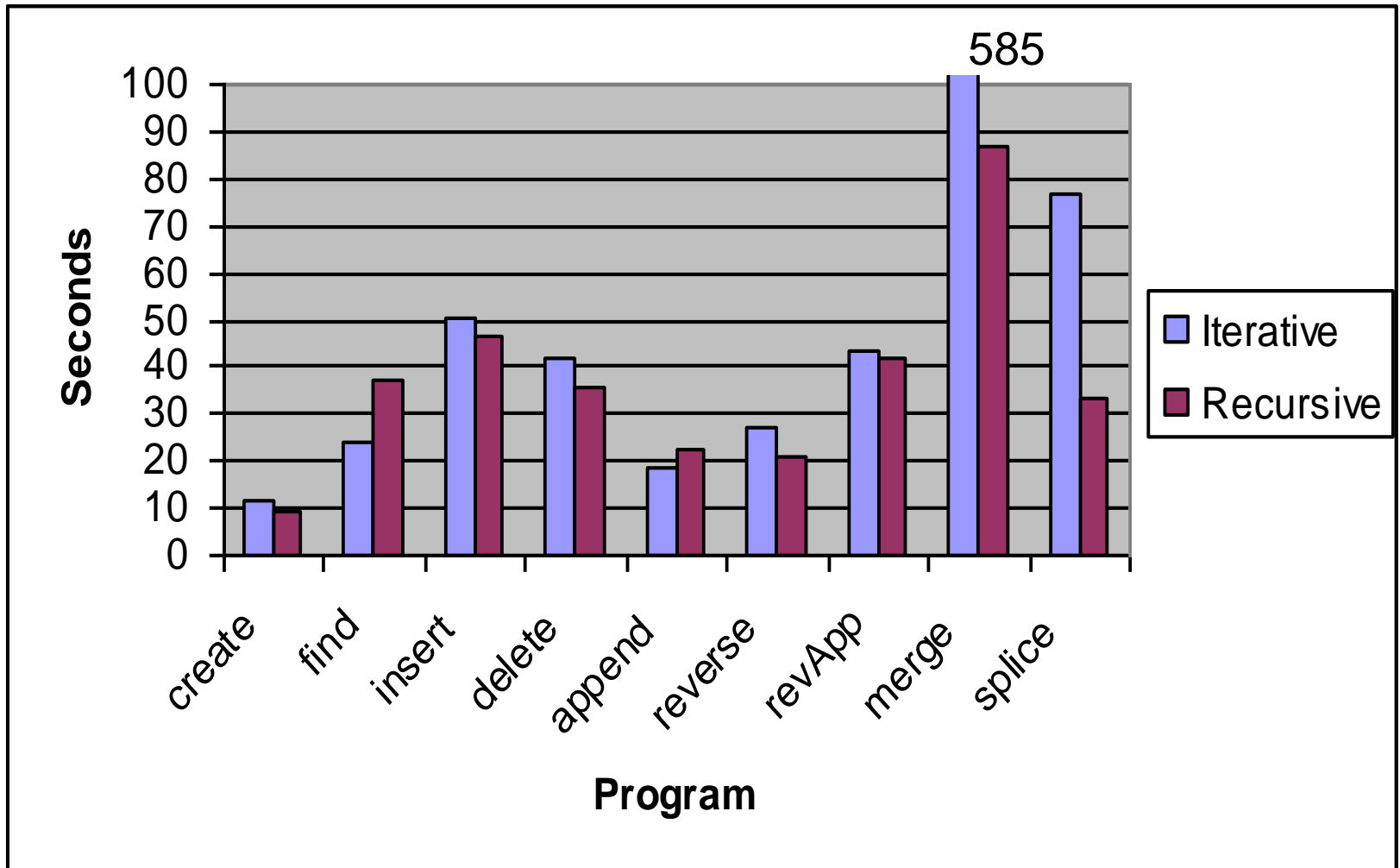
- ✓ Cutpoint freedom
- ✓ Non-standard concrete semantics
- ✓ Interprocedural shape analysis
- Prototype implementation

Prototype implementation

- TVLA based analyzer
- Soot-based Java front-end
- Parametric abstraction

Data structure	Verified properties
Singly linked list	Cleanness, acyclicity
Sorting (of SLL)	+ Sortedness
Unshared binary trees	Cleanness, tree-ness

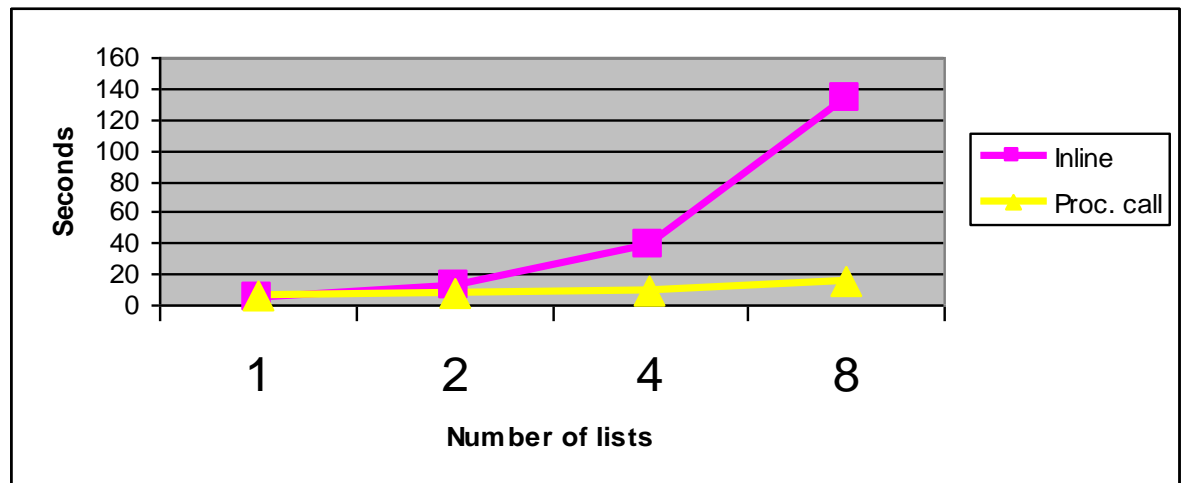
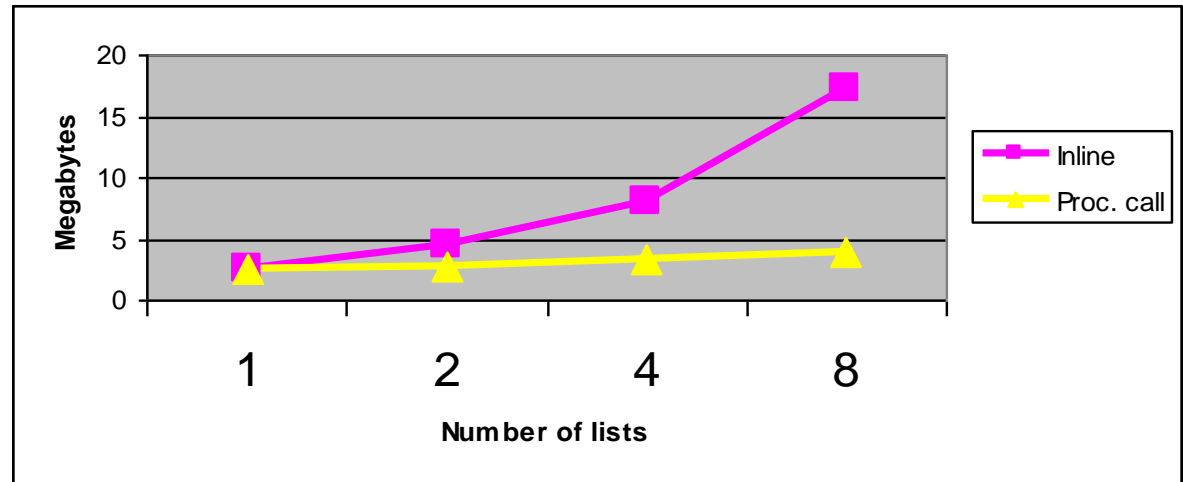
Iterative vs. Recursive (SLL)



Inline vs. Procedural abstraction

```
// Allocates a list of
// length 3
List create3(){
    ...
}

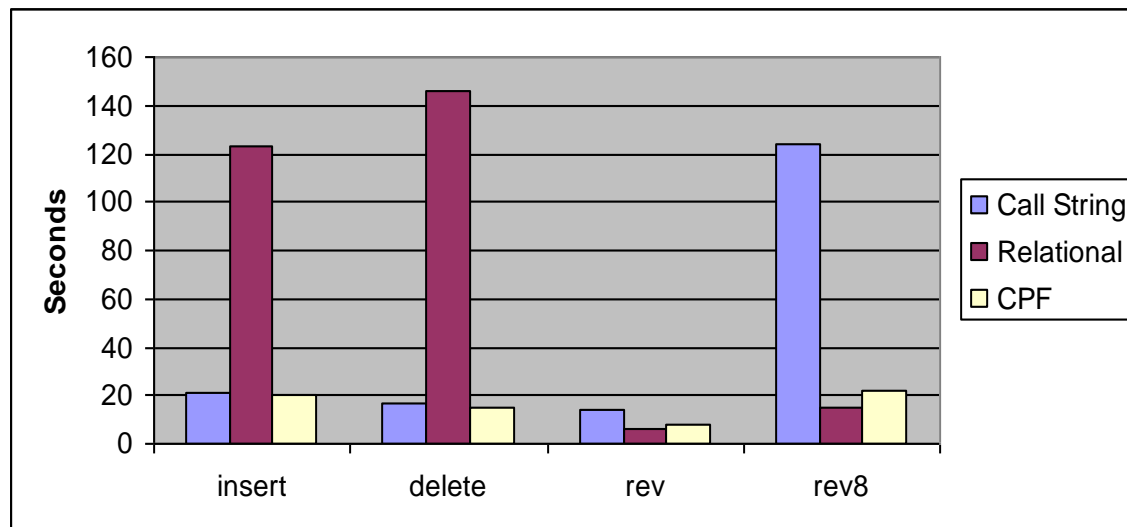
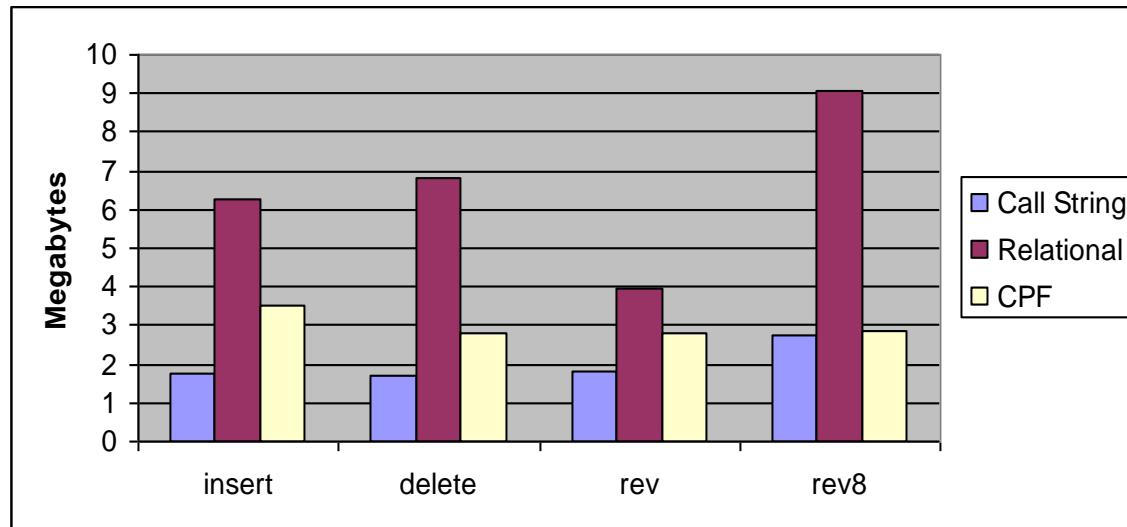
main() {
    List x1 = create3();
    List x2 = create3();
    List x3 = create3();
    List x4 = create3();
    ...
}
```



Call string vs. Relational vs. CPF

[Rinetzky and Sagiv, CC' 01]

[Jeannet et al., SAS' 04]



Summary

- Cutpoint freedom
- Non-standard operational semantics
- Interprocedural shape analysis
 - Partial correctness of quicksort
- Prototype implementation

Application

- Properties proved
 - Absence of null dereferences
 - Listness preservation
 - API conformance
- Recursive \approx Iterative
- Procedural abstraction

Related Work

- **Interprocedural shape analysis**
 - Rinetzky and Sagiv, CC '01
 - Chong and Rugina, SAS '03
 - Jeannet et al., SAS '04
 - Hackett and Rugina, POPL '05
 - Rinetzky et al., POPL '05
- **Local Reasoning**
 - Ishtiaq and O'Hearn, POPL '01
 - Reynolds, LICS '02
- **Encapsulation**
 - Noble et al. IWACO '03
 - ...

Summary

- Operational semantics
 - Storeless
 - Local heap
 - **Cutpoints**
 - Equivalence theorem
- Applications
 - Shape analysis
 - May-alias analysis

Project

- 1-2 Students in a group
 - 3-4: Bigger projects
- Theoretical + Practical
- Your choice of topic
 - Contact me in 2 weeks
- Submission – 15/Sep
 - Code + Examples
 - Document
 - 20 minutes presentation

Past projects

- JavaScript Dominator Analysis
- Attribute Analysis for JavaScript
- Simple Pointer Analysis for C
- Adding program counters to Past Abstraction
- Verification of Asynchronous programs
- Verifying SDNs using TVLA
- Verifying independent accesses to arrays in GO

Past projects

- Detecting index out of bound errors in C programs
- Lattice-Based Semantics for Combinatorial Models Evolution