

Advanced Static Analysis Technique

Woosuk Lee

CSE 6049 Program Analysis



Hanyang University, Korea

Goal of This Lecture

- Learn advanced techniques for more **precise** and **efficient** analysis

Combination of Multiple Abstractions

- Verifying a program often requires reason over multiple kinds of properties.
- We can construct abstract domains that can handle such cases.

Combination of Multiple Abstractions

- Suppose we have two abstract domains $\mathbb{A}_0, \mathbb{A}_1$ with concretization function $\gamma_0 : \mathbb{A}_0 \longrightarrow \mathcal{P}(\mathbb{M}), \gamma_1 : \mathbb{A}_1 \longrightarrow \mathcal{P}(\mathbb{M})$
- The *product abstraction* is defined by
 - $\mathbb{A}_\times = \mathbb{A}_0 \times \mathbb{A}_1$
 - **Concretization function:**
 $\forall (a_0, a_1) \in \mathbb{A}_\times, \gamma_\times(a_0, a_1) = \gamma_0(a_0) \cap \gamma_1(a_1)$

Combination of Multiple Abstractions

- Two abstract domains can characterize distinct kinds of properties.
 - E.g., \mathbb{A}_0 : interval domain, \mathbb{A}_1 : parity domain
 - Abstract value $([0, 100], \text{Even})$ describes even values between 0 and 100.

Synergistic Combination of Multiple Abstractions

- Multiple abstract domains may exchange information to derive stronger properties.

- e.g., Abstract value

$([1, 5], \text{Even})$

is concretized to $\{2, 4\}$.

The most precise abstract information is

$([2, 4], \text{Even})$

Reduced Product

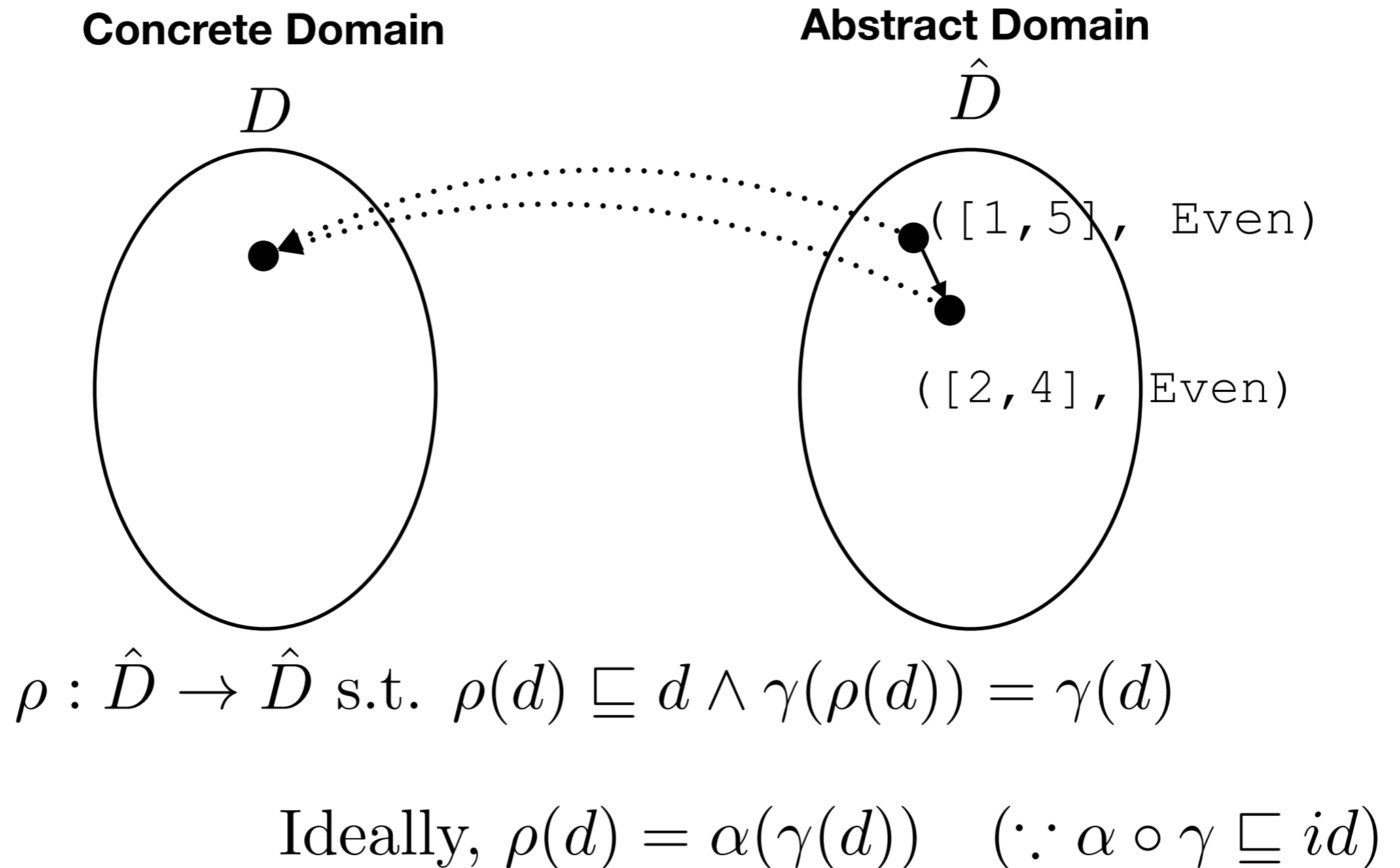
- The reduced product of abstract domains $\mathbb{A}_0, \mathbb{A}_1$ is the set \mathbb{A}_{\bowtie} of equivalence classes of $\mathbb{A}_0 \times \mathbb{A}_1$ for the relation \equiv defined by

$$(a_0, a_1) \equiv (a'_0, a'_1) \iff \gamma_{\times}(a_0, a_1) = \gamma_{\times}(a'_0, a'_1)$$

- e.g., $([1, 5], \text{Even}) \equiv ([2, 4], \text{Even})$

Reduction Operator

Reduction operator turns any abstract element into its optimal representation.



Reduction Operator

- Optimal reduction operator $(\alpha \circ \gamma)$ is often very expensive.
- In practice, a sound but approximate reduction operation $(\alpha \circ \gamma \sqsubseteq \rho)$ is often preferable.
- Reduction is typically performed only at specific points where it is known that it can provide a useful precision gain.

Advanced Iteration Techniques

- Loop invariant inference: sequences of abstract iterates
- Computes weaker and weaker abstract states until stabilization, hence a common source of imprecision (e.g., widening)
- There are various techniques to improve the precision of the static analysis of loops.

Review: Static Analysis of $\mathbf{while}(B)\{C\}$

- Concrete semantic function: $F = \llbracket C \rrbracket_{\mathcal{P}} \circ \mathcal{F}_B$
- Sound abstract semantic function: $F \circ \gamma \subseteq \gamma \circ F^\#$
- $M^\#$: abstract pre-condition for the loop
- Semantics of the loop:

$$M_{\text{loop}} = \bigcup_{i \geq 0} F^i(\gamma(M^\#)) = \mathbf{lfp} \ G$$

where $G(X) = \gamma(M^\#) \cup F(X)$

Review: Static Analysis of Loops

- To ensure termination of the analysis, we use widening. To compute the converging sequence

$$\begin{aligned}M_0^\# &= M^\# \\M_{k+1}^\# &= M_k^\# \nabla F^\#(M_k^\#)\end{aligned}$$

- $M_{\text{lim}}^\#$: the limit of the sequence (after finitely many iterates)
- The analysis returns the sound loop post-condition $\mathcal{F}_{\neg B}^\#(M_{\text{lim}}^\#)$

Techniques for Improving Precision of Analysis of Loops

- Loop unrolling
- Delayed widening
- Widening with thresholds

Loop Unrolling

- Motivation: hasty join

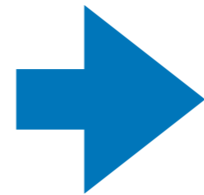
```
x = ?; // any value
i = 1;
while (i > 0) {
    if(x < 0 || x > 1000) {
        x = 0;
    } else {
        x = x + 1;
    }
    input(i);
}
// actually, x is in [0, 1001]
```

Initialization step

- The abstract value for x with a naive approach would be $[-\infty, +\infty]$
- Idea: **detach** the first iteration from the rest

Loop Unrolling

```
x = ?; // any value
i = 1;
while (i > 0) {
    if(x < 0 || x > 1000) {
        x = 0;
    } else {
        x = 1 + x;
    }
    input(i);
}
// actually, x is in [0, 1001]
```



```
x = ?; // any value
i = 1;
if(x < 0 || x > 1000) {
    x = 0;
} else {
    x = 1 + x;
}
input(i);
// x is in [0, 1001]
while (i > 0) {
    if(x < 0 || x > 1000) {
        x = 0;
    } else {
        x = 1 + x;
    }
    input(i);
}
// x is in [0, 1001]
```

} first iter.

} rest

Loop Unrolling

- The previous sequence for the loop

$$\begin{aligned}M_0^\# &= M^\# \\M_{k+1}^\# &= M_k^\# \nabla F^\#(M_k^\#)\end{aligned}$$

- With loop unrolling

$$\begin{aligned}M_0^\# &= M^\# \\M_{k+1}^\# &= \begin{cases} F^\#(M_k^\#) & \text{if } k < N \\ M_k^\# \nabla F^\#(M_k^\#) & \text{otherwise} \end{cases}\end{aligned}$$

and the analysis returns $\mathcal{F}_{\neg B}^\#(M_{\text{lim}}^\#)$

Delayed Widening

- Motivation: hasty widening

```
x = 0;
while (rand()) {
  if(rand()) {
    x = -1;
  } else {
    x = x + 2;
  }
}
// x >= -1
```

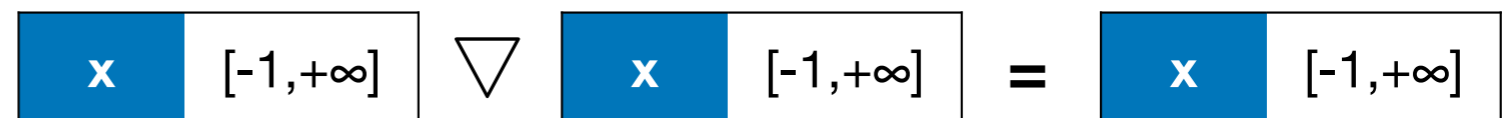


- The abstract value of x with a naive approach would be $[-\infty, +\infty]$
- Idea: **delay** the application of widening for the first N iterations

Delayed Widening

Delayed widening where $N = 1$

```
x = 0;  
while (rand()) {  
    if(rand()) {  
        x = -1;  
    } else {  
        x = x + 2;  
    }  
}  
// x >= -1
```



Fixed Point!

Delayed Widening

- The sequence with delayed widening:

$$\begin{aligned}M_0^\# &= M^\# \\M_{k+1}^\# &= M_k^\# \sqcup^\# F^\#(M_k^\#) && \text{if } k \leq N \\M_{k+1}^\# &= M_k^\# \nabla F^\#(M_k^\#) && \text{if } k > N\end{aligned}$$

- Loop unrolling: postpones **join** first N iterations
- Delayed widening: postpones **widening** first N iterations

Widening with Thresholds

- Motivation: the standard widening is too conservative!

```
x = 0;
while (x <= 100) {
  if(x >= 50) {
    x = 10;
  } else {
    x = x + 1;
  }
}
// actually, x is in [0, 50]
```



- The abstract value of x with a naive approach is $[0, +\infty]$
- Idea: use a **slower and more precise** widening

Widening with Thresholds

- Take several small steps and stops at pre-defined threshold values
- For example, consider only one threshold B :

A naive widening operator

$$[n, p] \nabla [n, q] = \begin{cases} [n, p] & \text{if } p \geq q \\ [n, +\infty] & \text{if } p < q \end{cases}$$

A widening with thresholds

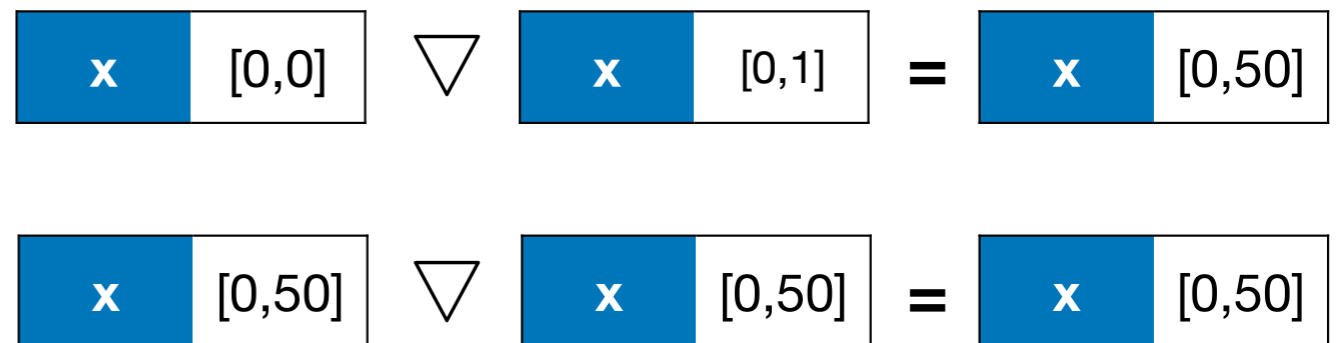
$$[n, p] \nabla [n, q] = \begin{cases} [n, p] & \text{if } p \geq q \\ [n, B] & \text{if } p < q \leq B \\ [n, +\infty] & \text{if } B < q \end{cases}$$

*only the right bounds, for brevity

Widening with Thresholds

```
x = 0;
while (x <= 100) {
  if(x >= 50) {
    x = 10;
  } else {
    x = x + 1;
  }
}
```

Thresholds = {50}



Fixed Point!

Scalability Challenge

The worklist algorithm often does not scale to large complex programs.

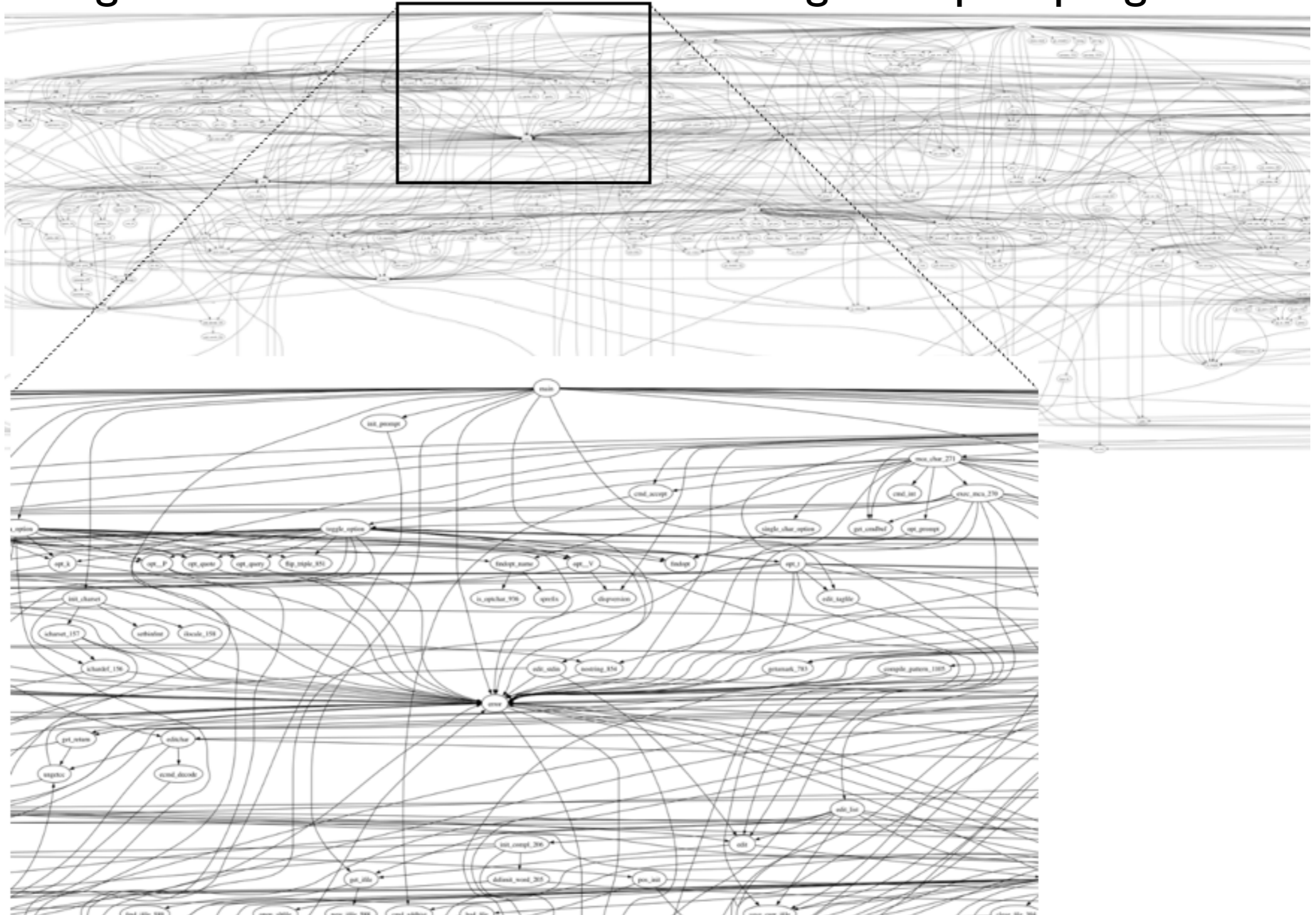


Figure: Call graph of less-382 (23,822 lines of code)

Sparse Analysis

- The worklist algorithm is often still not efficient enough to analyze large programs.
- Sparse analysis: exploit the semantic sparsity of the input program to analyze
- Spatial sparsity & temporal sparsity

Right part at right moment

Example Performance Gain by Sparse Analysis

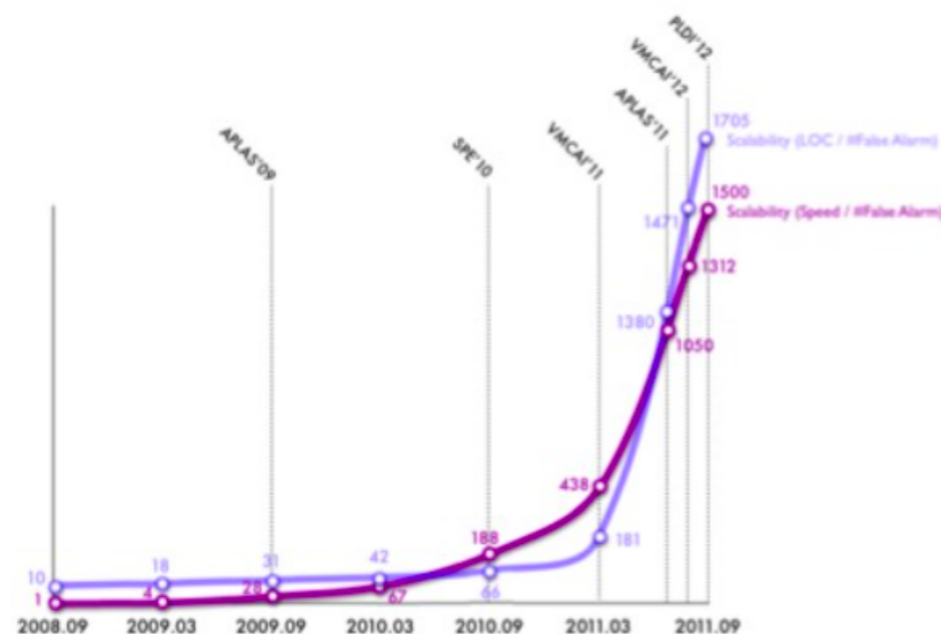
- Sparrow: a “sound”, global C analyzer for the memory safety property (no overrun, no null-pointer dereference, etc.)

<http://github.com/ropas/sparrow>

- ~ 10 hours in analyzing million lines of C



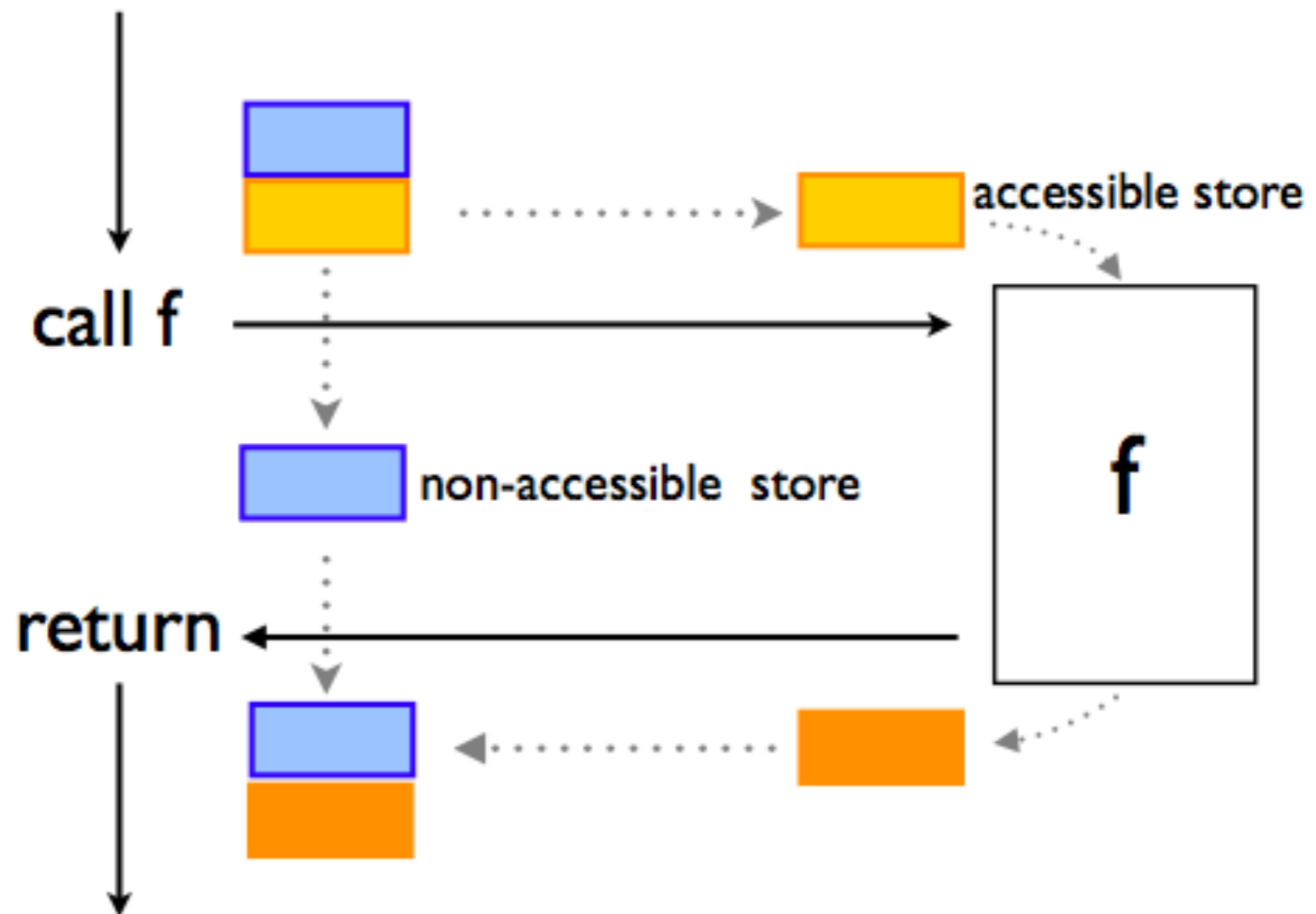
sound-&-global version



- < 1.4M in 10hrs with intervals
- < 0.14M in 20hrs with octagons

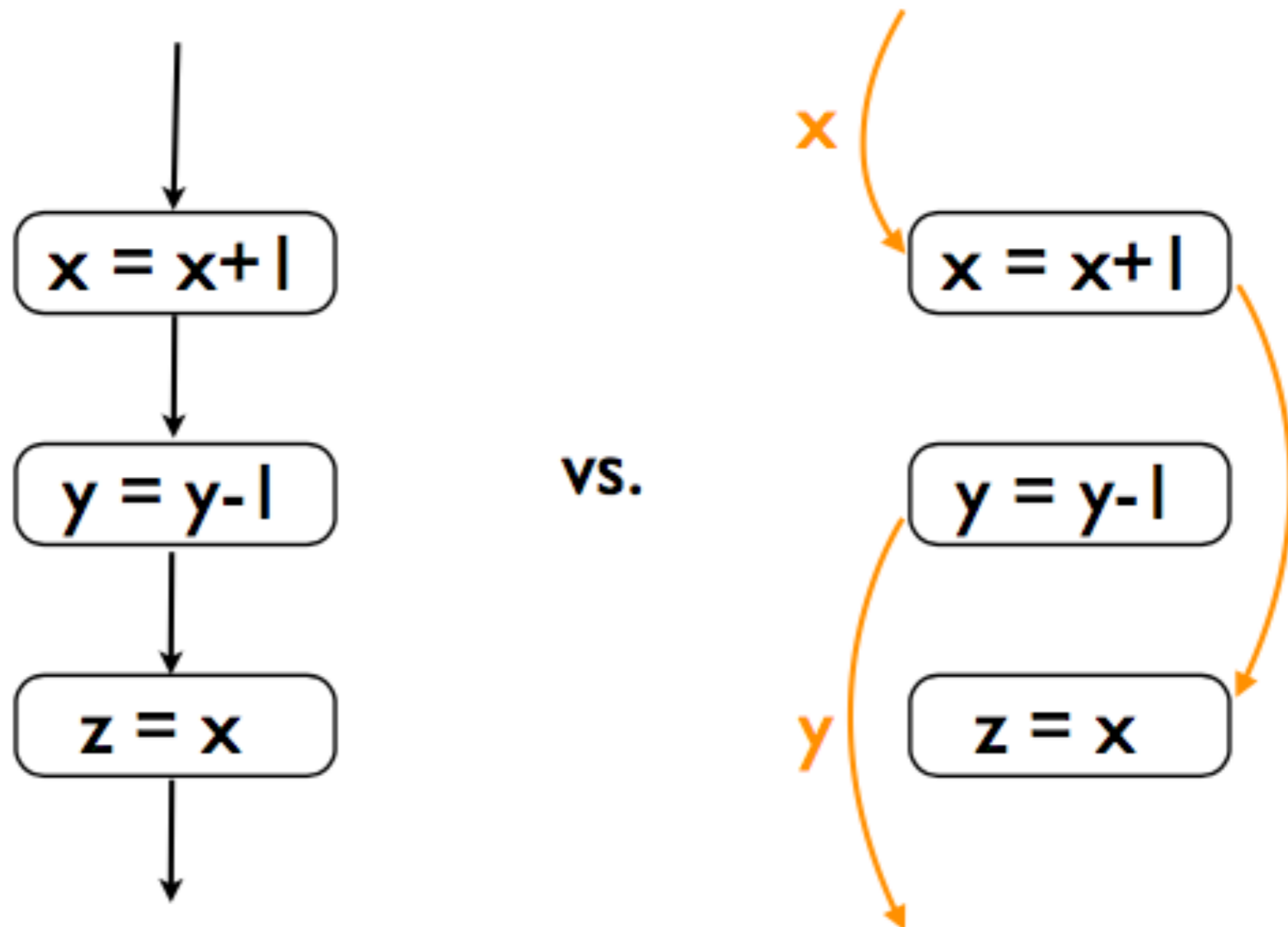
Spatial Sparcity

Each program portion accesses only a small part of the memory.



Temporal Sparsity

After the def of a memory, its use is far.

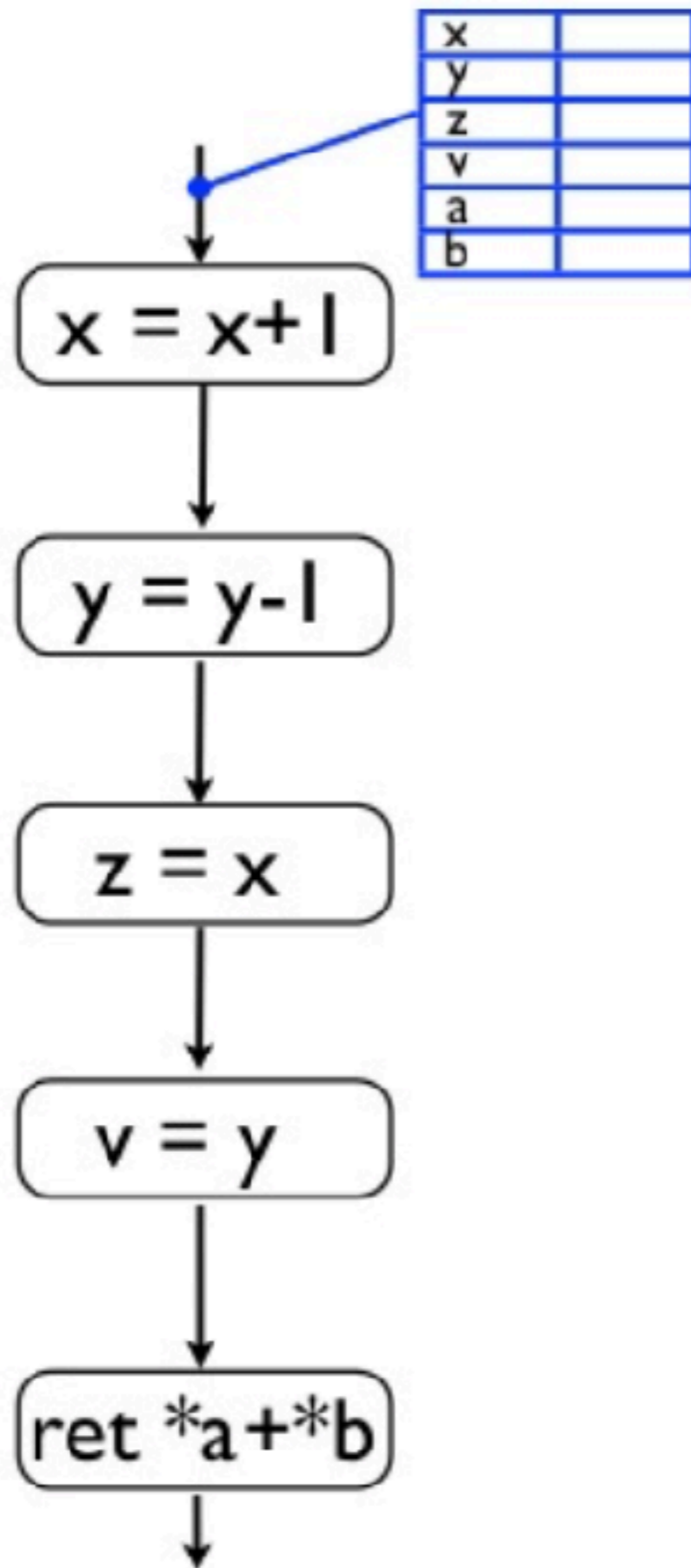


Example

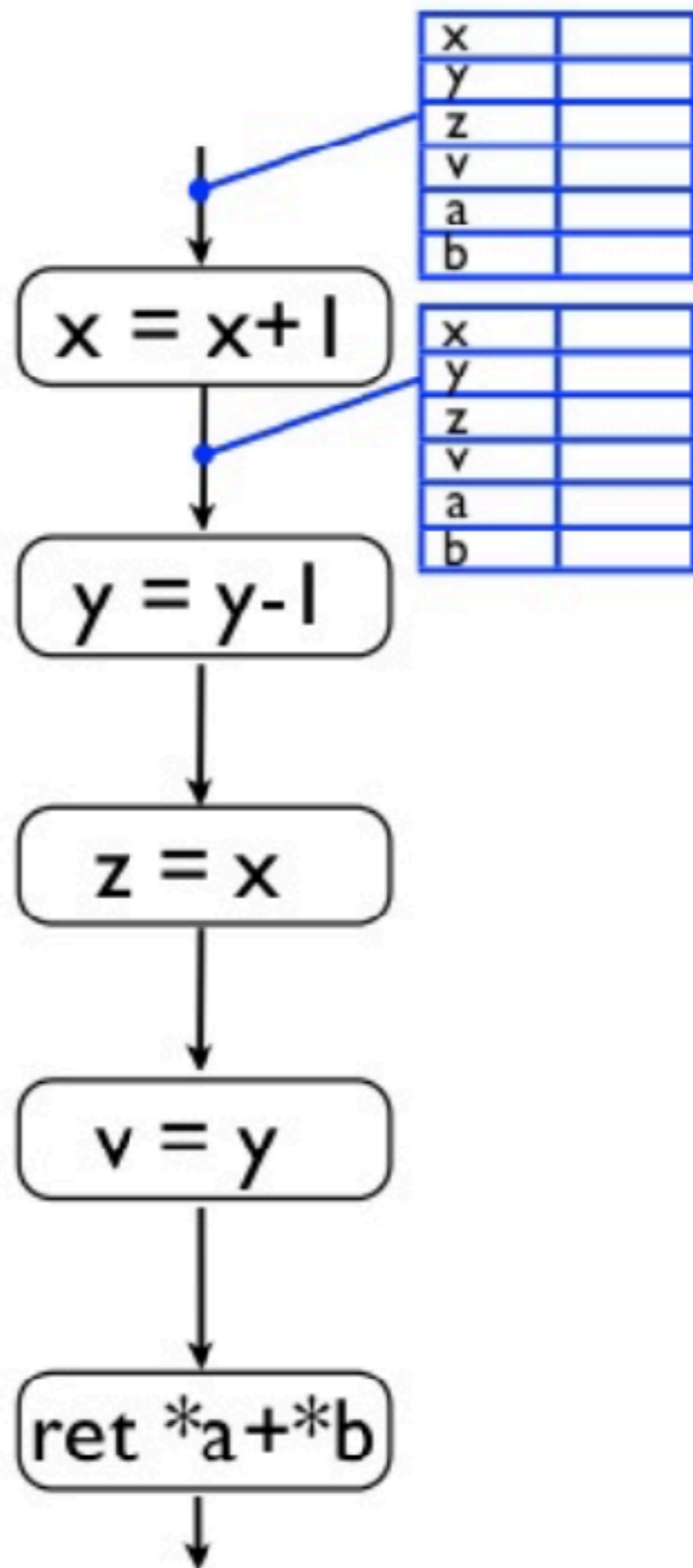
Example (Code fragment)

```
x = x + 1;  
y = y - 1;  
z = x;  
v = y;  
ret *a + *b
```

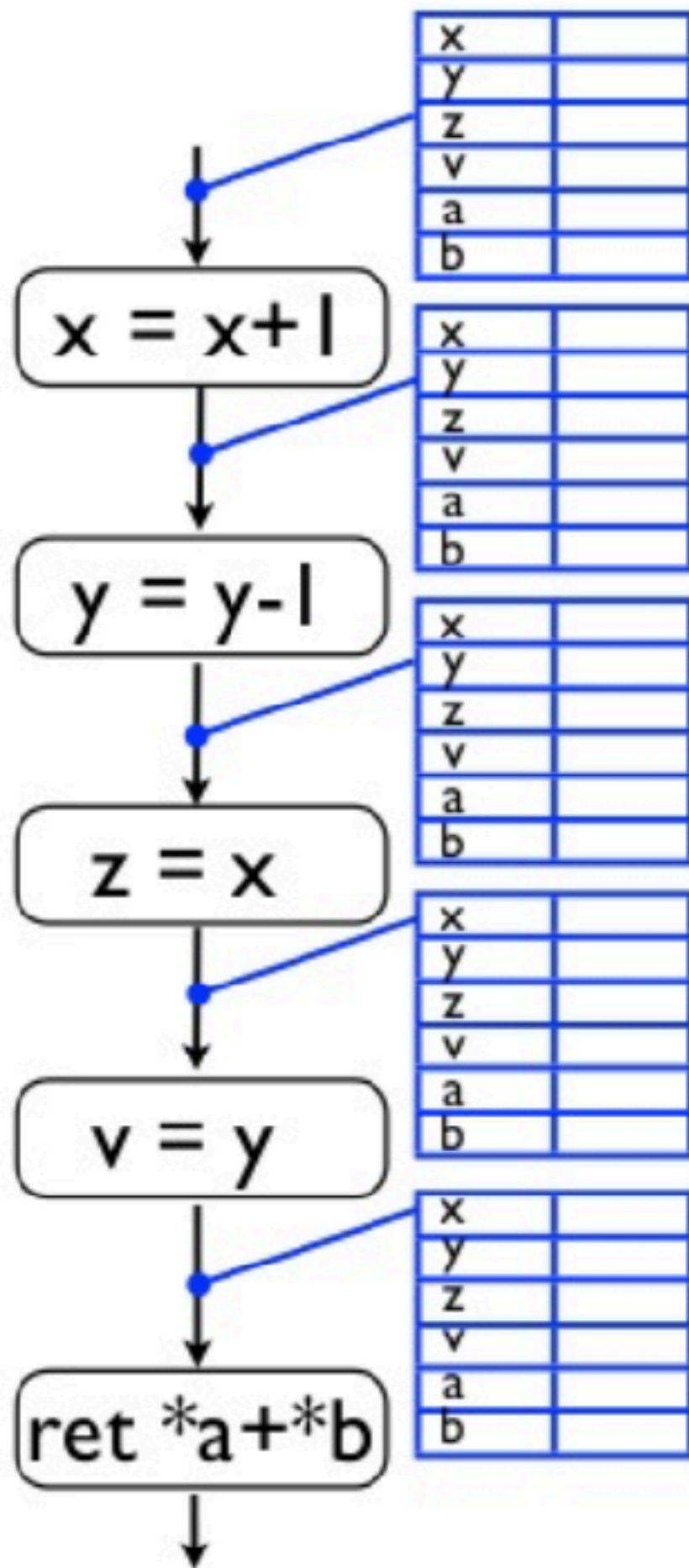
Assume that `a` points to `v` and `b` to `z`.



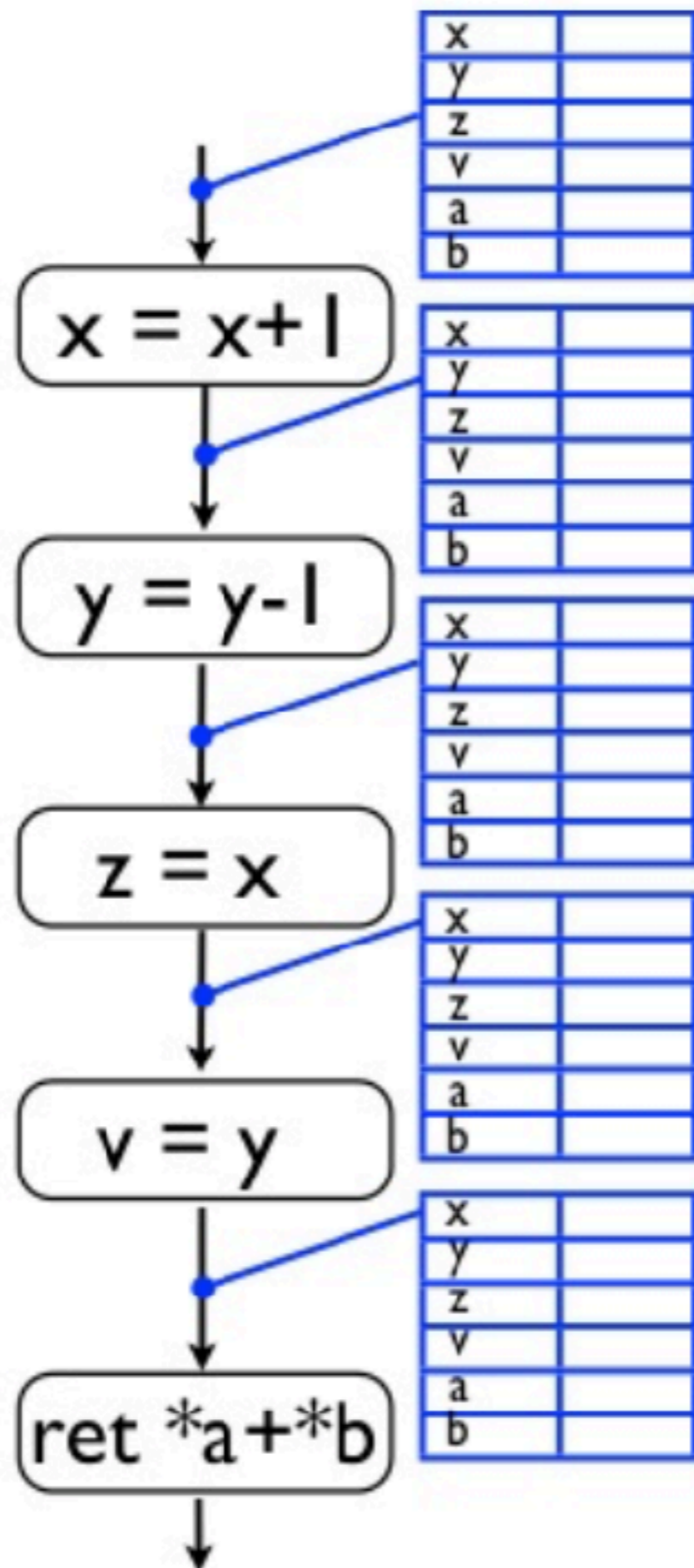
Consider this program and suppose that we analyze the program with the initial abstract state.



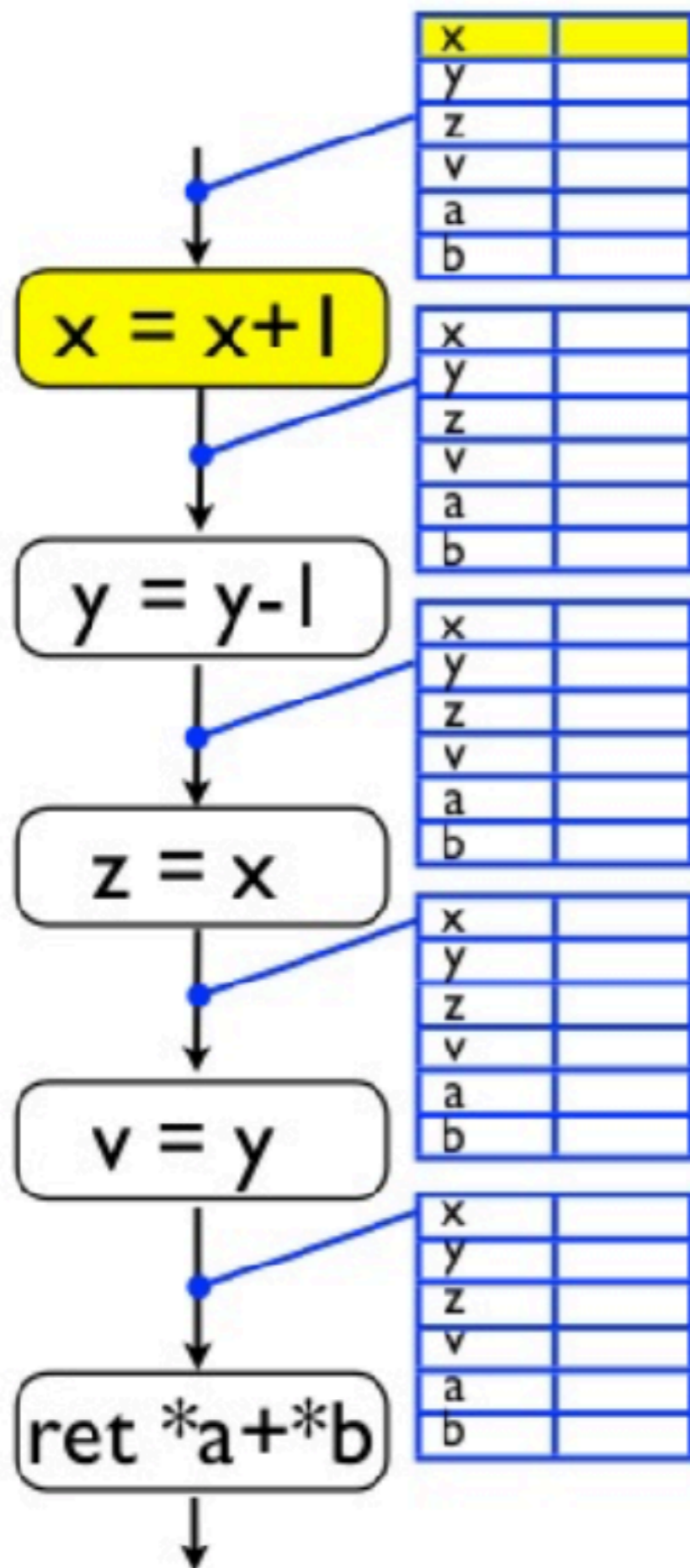
Conventional static analysis propagates the entire abstract states along control flows of the program.



So, after the analysis is terminated, each program point is associated with the entire abstract states.

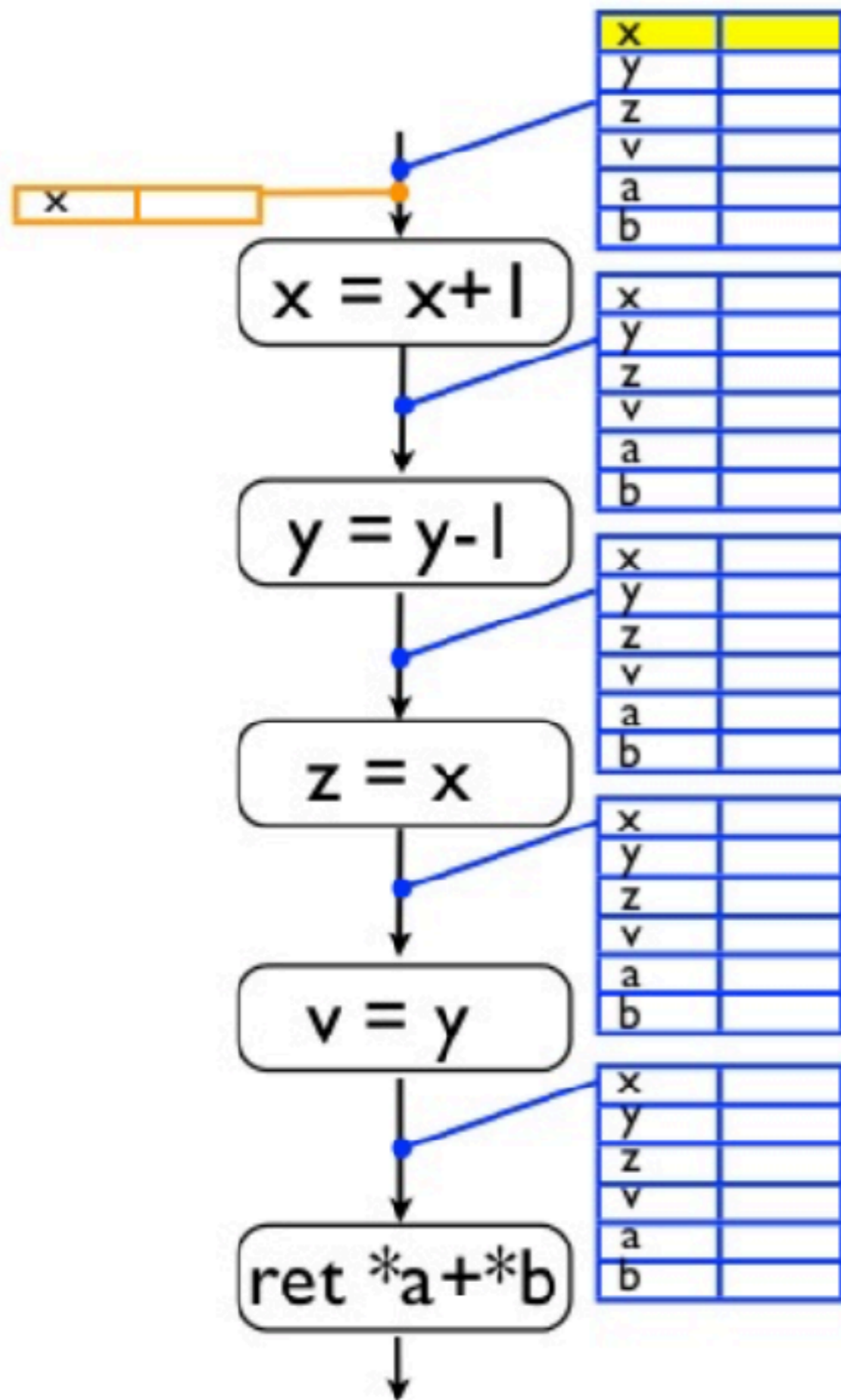


Sparse analysis aims to optimize this conventional static analysis based on two observations.

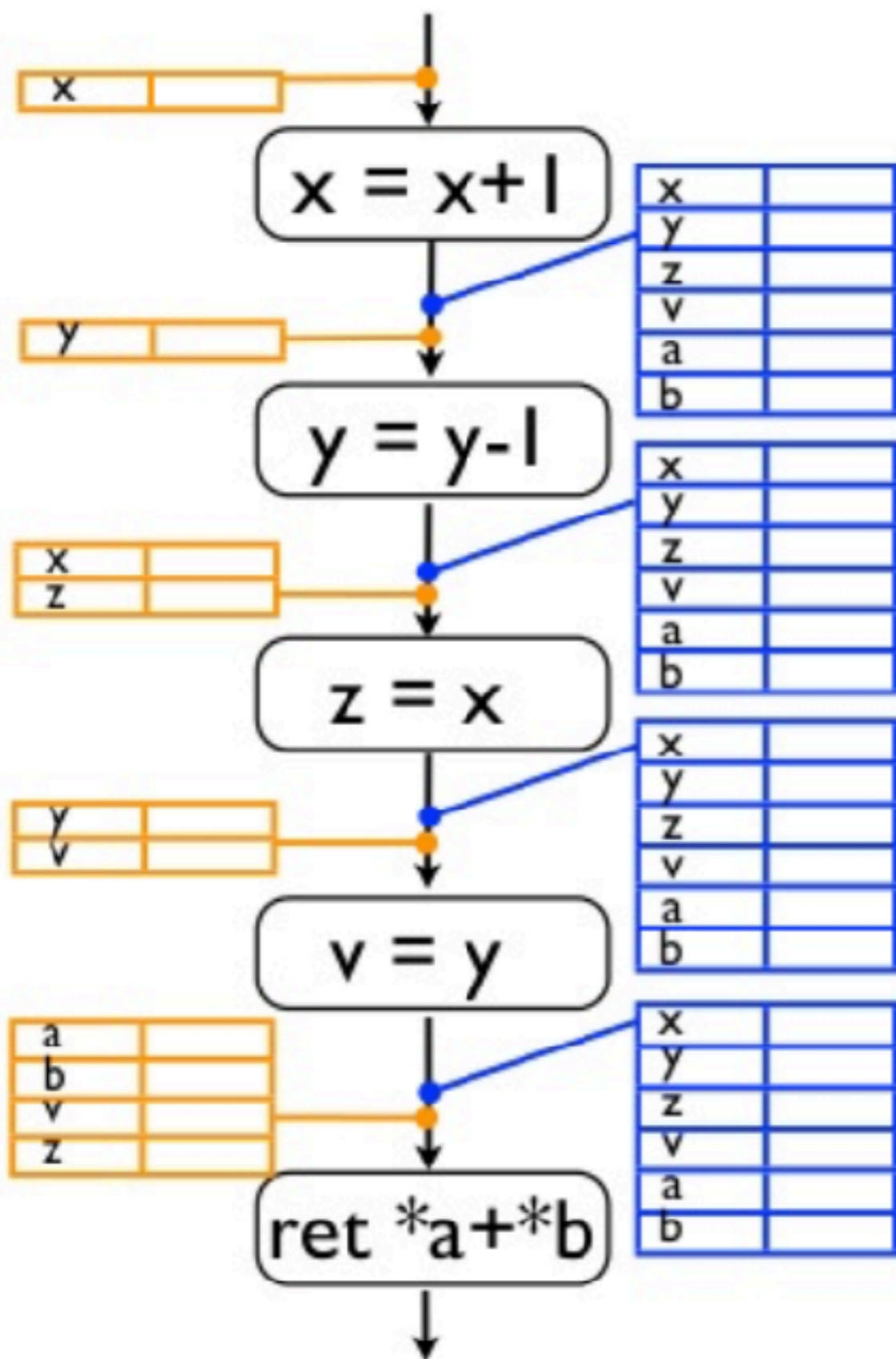


First, in the analysis of each statement, only a small subset of the state is actually used.

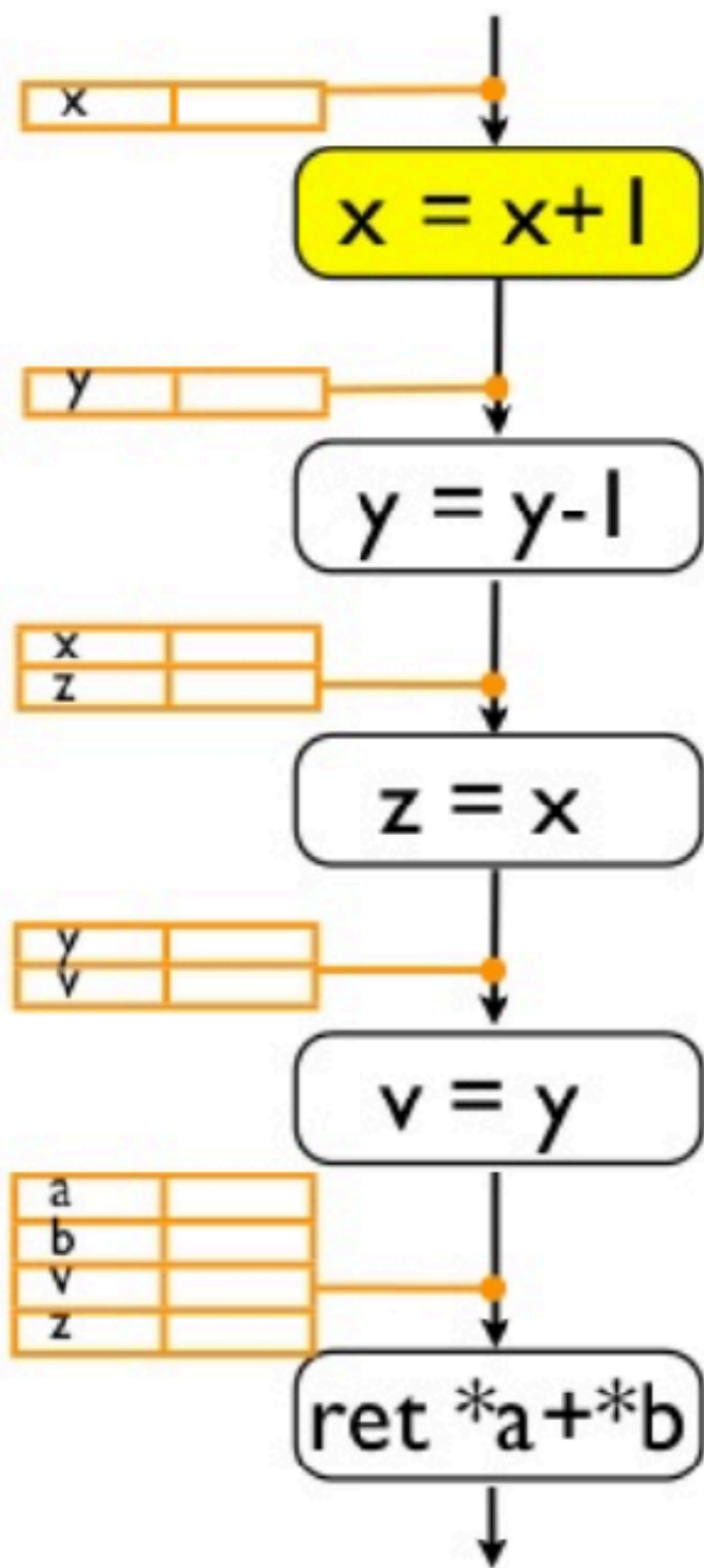
For example, only the value of x is necessary to analyze the first statement.



So, in sparse analysis, we keep `x` here and remove other values.

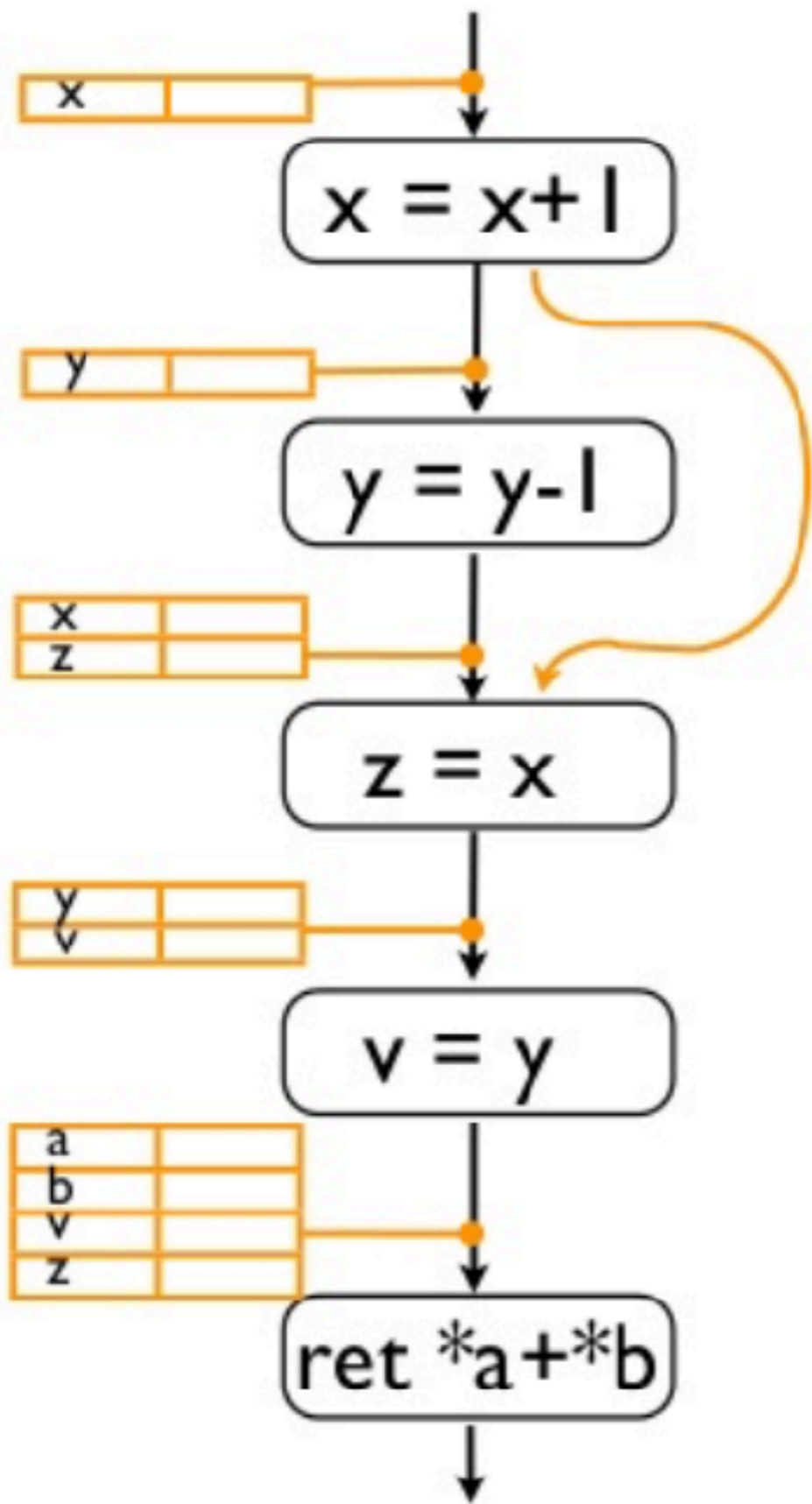


In this way, at each program point, sparse analysis stores only the values that will be used in the analysis.

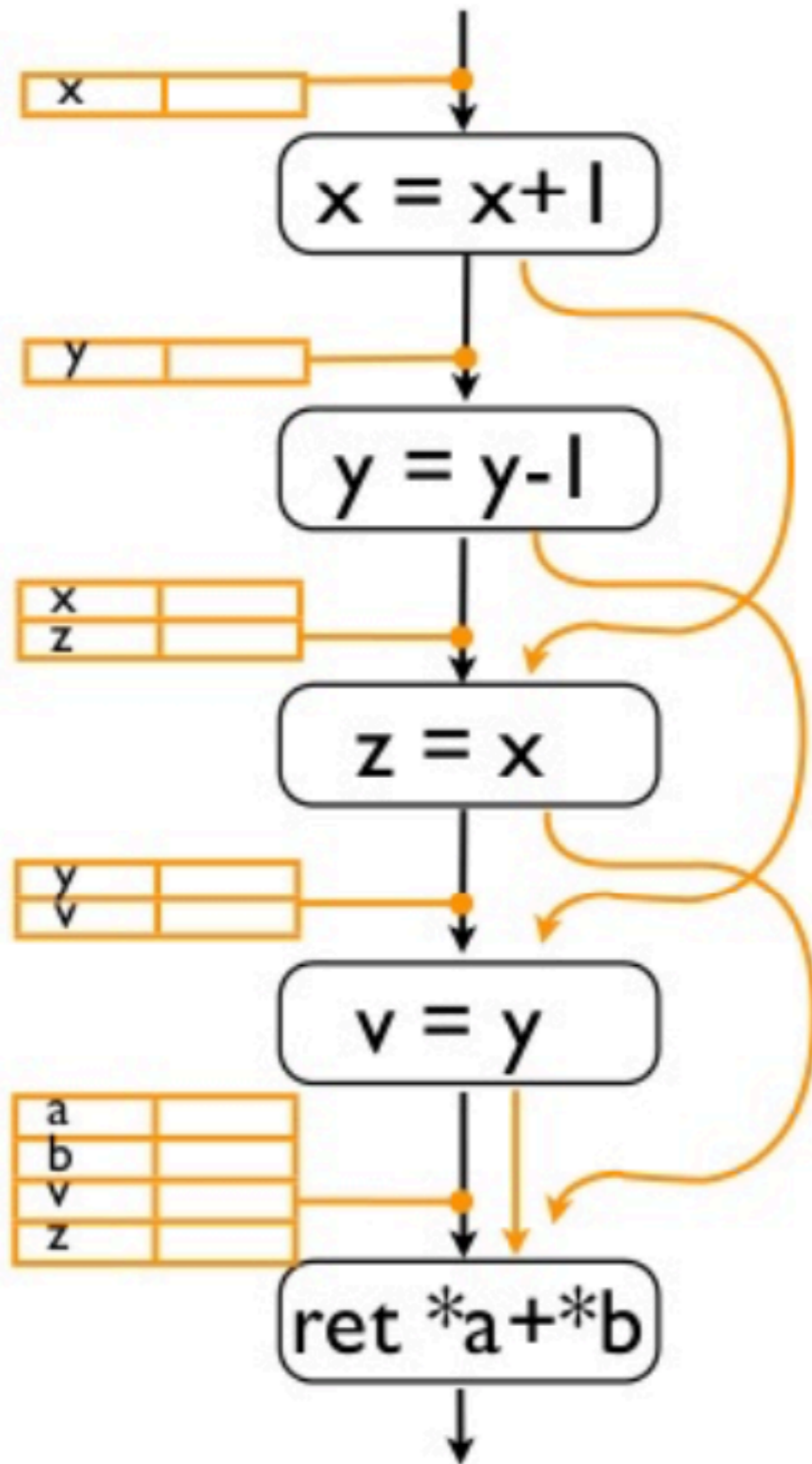


The second observation is that the semantic dependencies among statements are usually sparse.

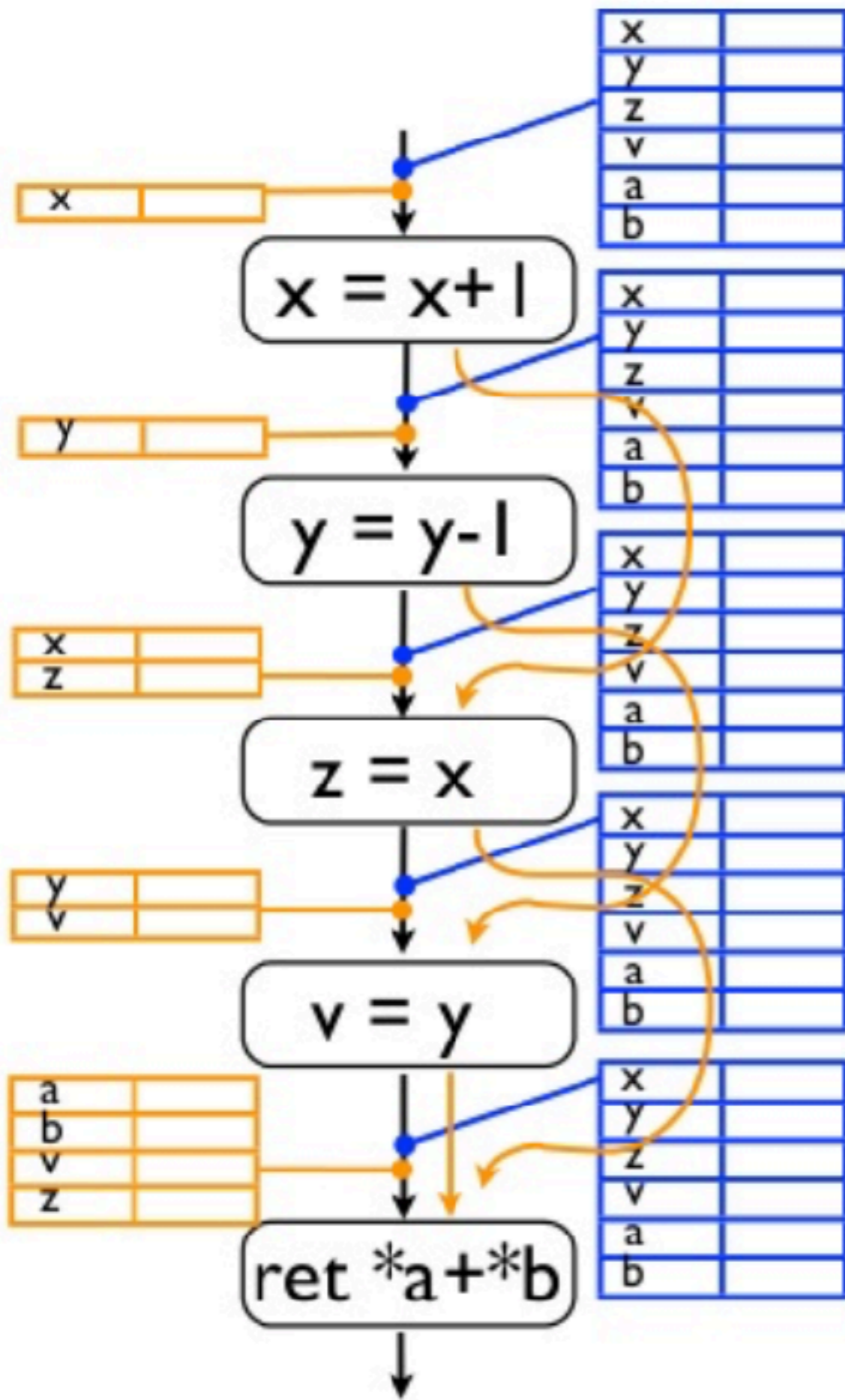
For example, the value of x in the first statement is not used at the next statement but used at the third statement.



So, sparse analysis propagates the value directly to the use point.



Similarly, other values are also propagated along semantic dependencies of the program.



These two “localizations” significantly improves the scalability of the non-sparse analysis.

Exploiting Sparsity

$$F^\# : (\mathbb{L} \rightarrow \mathbb{M}^\#) \rightarrow (\mathbb{L} \rightarrow \mathbb{M}^\#)$$

becomes

$$F^\#_{sparse} : (\mathbb{L} \rightarrow \mathbb{M}^\#_{sparse}) \rightarrow (\mathbb{L} \rightarrow \mathbb{M}^\#_{sparse})$$

where

$$\mathbb{M}^\#_{sparse} = \{M^\# \in \mathbb{M}^\# \mid \text{dom}(M^\#) = \text{Access}^\#(l), l \in \mathbb{L}\} \cup \{\perp\}.$$



**Variables used at label l
of the input program**

Exploiting Temporal Sparsity

Need the def-use chain information as follows.

- we streamline the abstract one-step relation

$$(l, M^\#) \hookrightarrow^\# (l', M'^\#) \quad \text{for } l' \in \text{next}^\#(l, M^\#).$$

so that the link $\hookrightarrow^\#$ should follow the **def-use chain**:

- ▶ from (def) a label where a location is defined
- ▶ to (use) a label where the defined location is read

Precision-Preserving Sparse Analysis

Goal

$$F^\# : D^\# \rightarrow D^\# \xrightarrow{\text{sparsify}} F_{\text{sparse}}^\# : D^\# \rightarrow D^\#$$

$$\text{lfp} F^\# \stackrel{\text{still}}{=} \text{lfp} F_{\text{sparse}}^\#$$

Step 1: Estimating Accessed Variables at Each Label

Need to safely estimate

$$Access^\#(l).$$

Use yet another sound static analysis, a further abstraction:

$$(\mathbb{L} \rightarrow \mathbb{M}^\#, \sqsubseteq) \xrightleftharpoons[\alpha]{\gamma} (\mathbb{M}^\#, \sqsubseteq_M)$$

(a “flow-insensitive” version of the “flow-sensitive” analysis design)



Flow-insensitive pointer analysis will be introduced in the upcoming lectures

Step 2: Computing Def-Use Information

- Let

$$D^\# : \mathbb{L} \rightarrow \wp(\mathbb{X}) \text{ and } U^\# : \mathbb{L} \rightarrow \wp(\mathbb{X})$$

be the def and use sets from the original analysis.

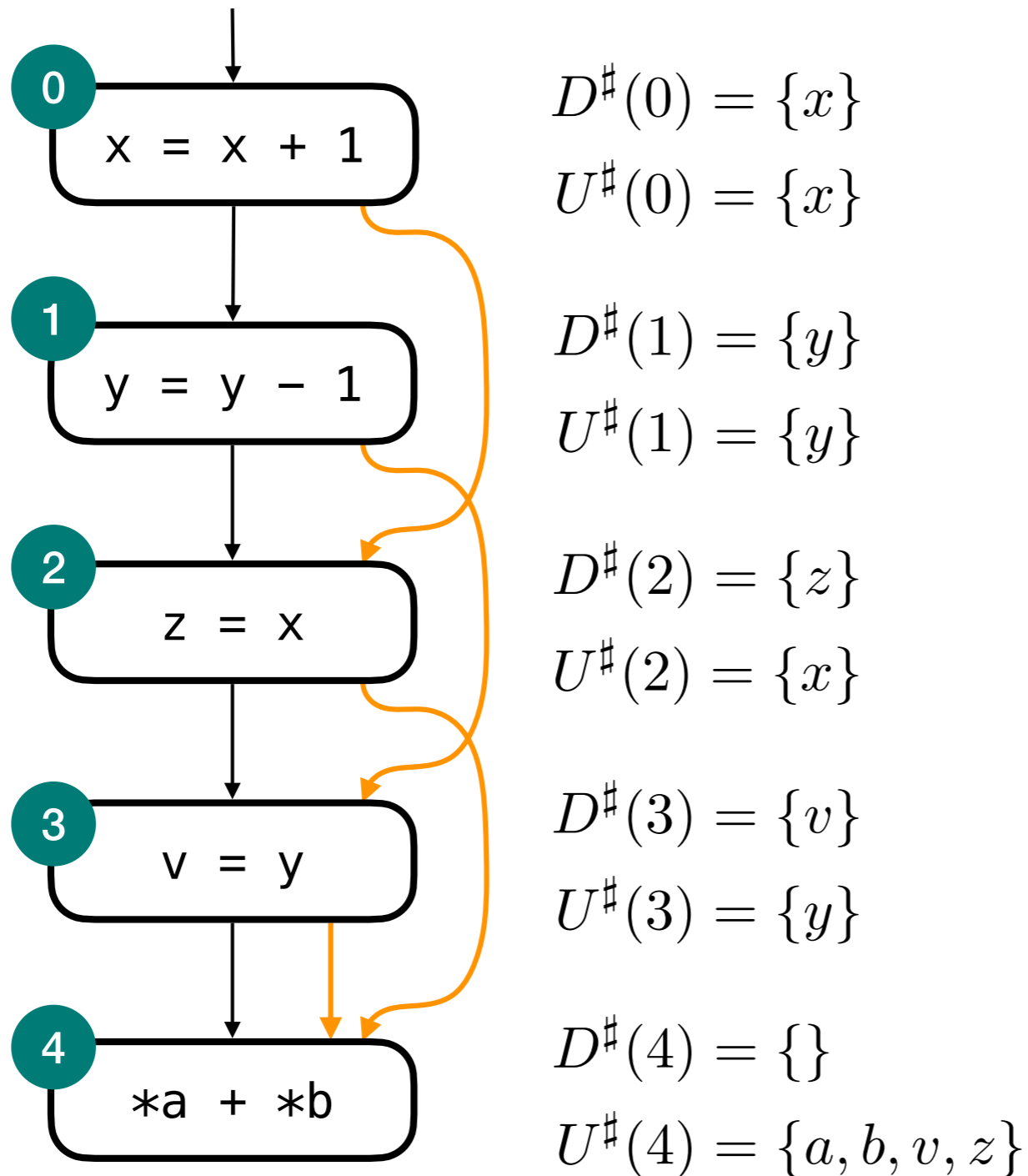
- Need to safely estimate $D^\#$ and $U^\#$.
- Use yet another sound static analysis to compute

$$D_{pre}^\# \text{ and } U_{pre}^\#$$

such that

- ▶ $\forall l \in \mathbb{L} : D_{pre}^\#(l) \supseteq D^\#(l) \text{ and } U_{pre}^\#(l) \supseteq U^\#(l).$
- ▶ $\forall l \in \mathbb{L} : U_{pre}^\#(l) \supseteq D_{pre}^\#(l) \setminus D^\#(l).$

Def and Use Sets



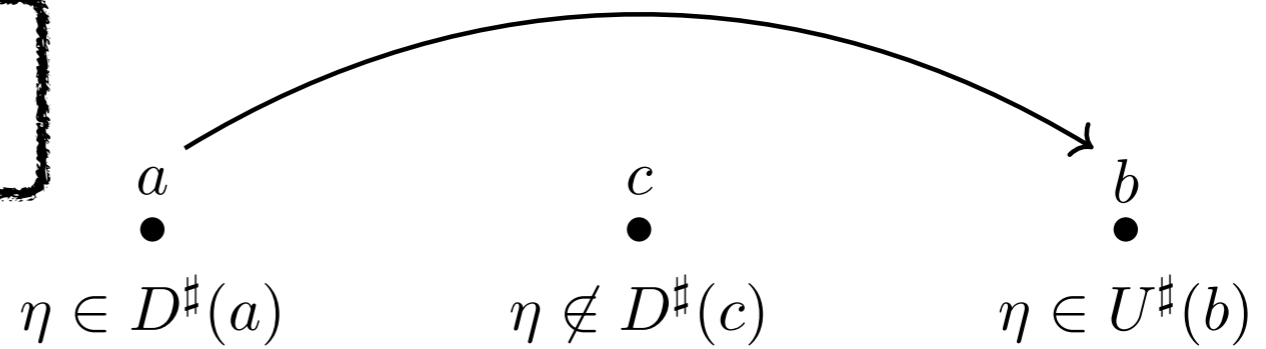
*Assume a and b point to z and v

Step 3: Analysis with Def-Use Chains

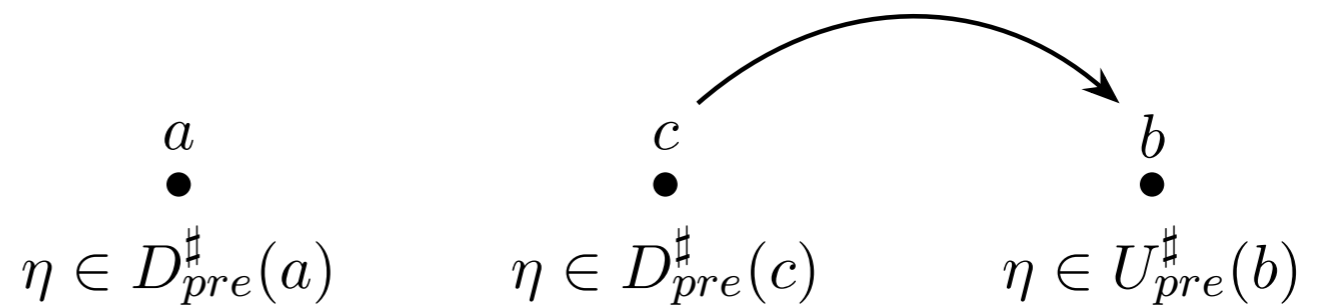
- **Def-use chain:** Label a to label b is a def-use chain for a variable x whenever x is defined at a and used at b , and x is not re-defined in-between a and b .
- The resulting sparse analysis with safe def-use chains has the same precision as the original non-sparse analysis.

Need for the Second Condition for $D_{pre}^\#$ and $U_{pre}^\#$

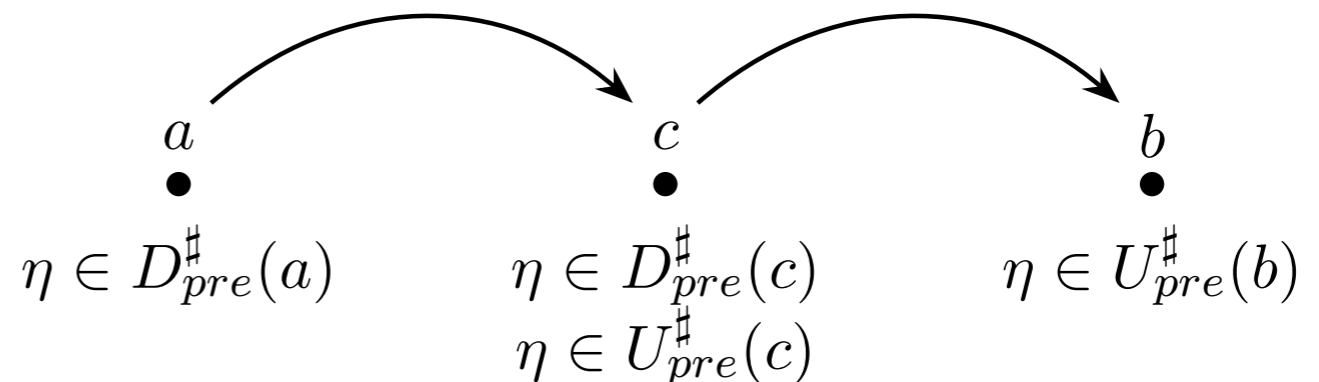
$$\forall l \in \mathbb{L} : U_{pre}^\#(l) \supseteq D_{pre}^\#(l) \setminus D^\#(l).$$



(d) Original analysis def-use edge for η



(e) Missing def-use edge (a to b) for η because of over-approximate $D_{pre}^\#(c)$



(f) Recovered def-use edge (a to b via c) for η by safe $U_{pre}^\#(c)$