

Lecture 25: Review of Program Analysis

17-355/17-665/17-819: Program Analysis

Rohan Padhye

December 4, 2025

* Course materials developed with Jonathan Aldrich and Claire Le Goues

What is was this course about?

- Program analysis is *the systematic examination of a program to determine its properties*.
- From 30,000 feet, this requires:
 - Precise program representations
 - Tractable, systematic ways to reason over those representations.
- ~~We will learn~~: What we learned:
 - How to unambiguously define the meaning of a program, and a programming language.
 - How to prove theorems about the behavior of particular programs.
 - How to use, build, and extend tools that do the above, automatically.

What is was this course about?

- Program analysis is *the systematic examination of a program to determine its properties*.
- Principal techniques:
 - **Dynamic:**
 - **Testing:** Direct execution of code on test data in a controlled environment.
 - **Analysis:** Tools extracting data from test runs.
 - **Static:**
 - **Inspection:** Human evaluation of code, design documents (specs and models), modifications.
 - **Analysis:** Tools reasoning about program behavior without executing it.
 - **Inference:** Statistical models of code (e.g., AI / ML)
 - ...and their combination.

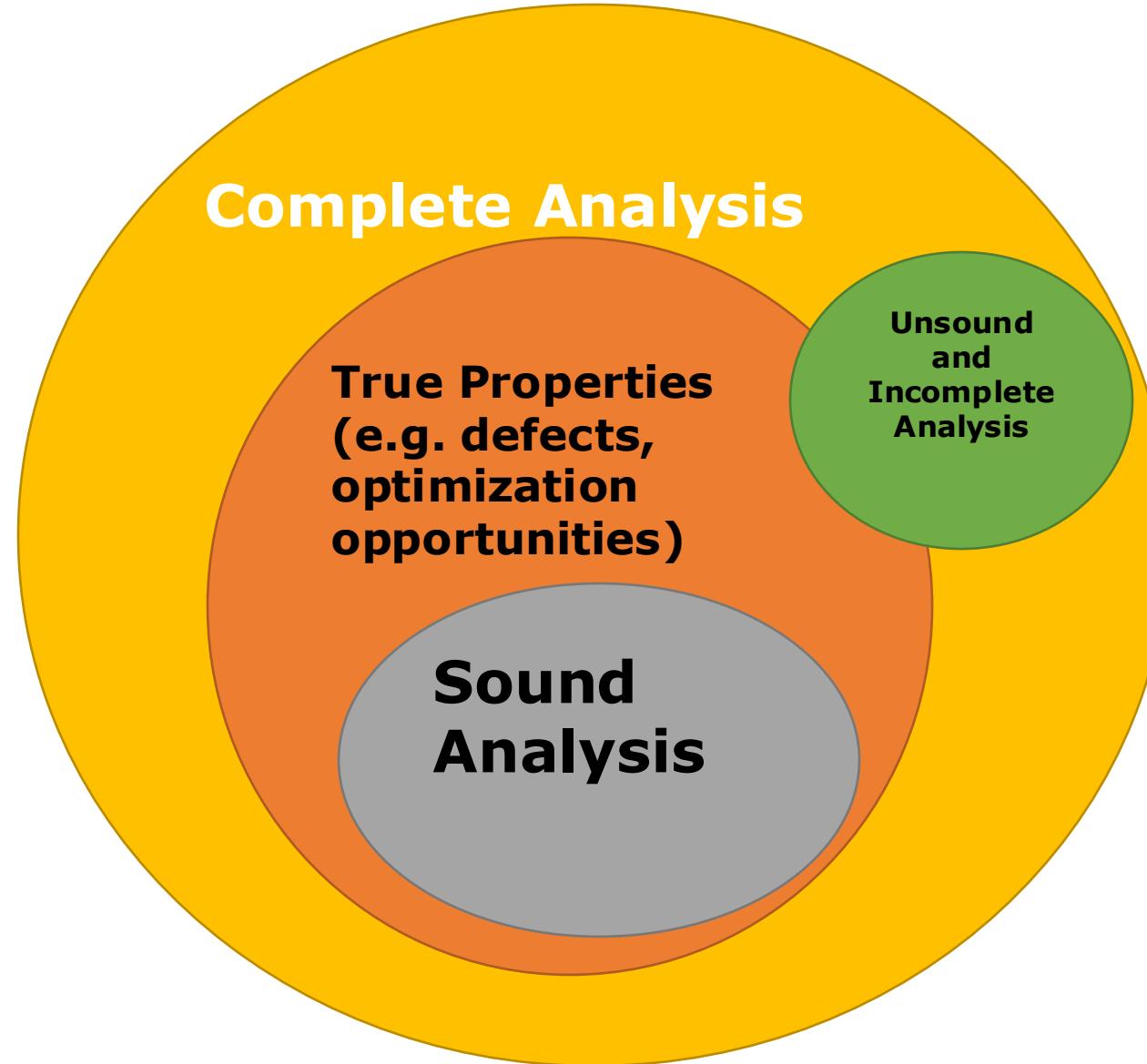
The Bad News: Rice's Theorem

"Any nontrivial property about the language recognized by a Turing machine is undecidable."

Henry Gordon Rice, 1953

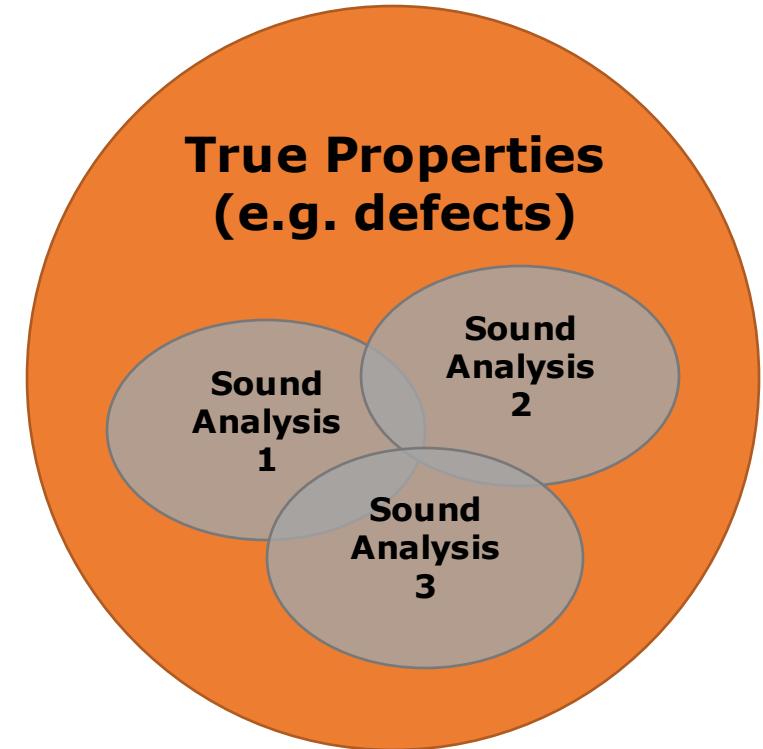
Soundness and Completeness

- An analysis is “sound” if every claim it makes is true
- An analysis is “complete” if it makes every true claim
- Soundness/Completeness correspond to under/over-approximation depending on context.
 - E.g. compilers and verification tools treat “soundness” as over-approximation since they make claims over all possible inputs
 - E.g. code quality tools often treat “sound” analyses as under-approximation because they make claims about existence of bugs



Soundness and Completeness Tradeoffs

- Sound + Complete is impossible in general (Rice's theorem)
- Most practical tools attempt to be either sound or complete for some specific application, using approximation
- Multiple classes of sound/complete techniques may exist, with trade-offs for accuracy and performance.
- Program analysis is a rich field because of the constant and never-ending battle to balance these trade-offs with ever-increasing software complexity



Fundamental concepts

- Abstraction
 - Elide details of a specific implementation.
 - Capture semantically relevant details; ignore the rest.
- The importance of semantics.
 - We prove things about analyses with respect to the semantics of the underlying language.
- Program proofs as inductive invariants.
- Implementation
 - You do not understand analysis until you have written several.

What you were supposed to get

- Beautiful and elegant theory
 - Mostly discrete mathematics, symbolic reasoning, inductive proofs
 - This is traditionally a “white-board” course [may not always use slides]
- Build awesome tools
 - Engineering of program analyses, compilers, and bug finding tools make great use of many fundamental ideas from computer science and software engineering
- New way to think about programs
 - Representations, control/data-flow, input state space
- Appreciate the limits and achievements in the space
 - What tools are *impossible* to build?
 - What tools are *impressive* that they exist at all?
 - When is it appropriate to use a particular analysis tool versus another?
 - How to interpret the results of a program analysis tool?

The WHILE language – Example program

```
y := x;  
z := 1;  
if y > 0 then  
    while y > 1 do  
        z := z * y;  
        y := y - 1  
else  
    skip
```

- Sample program computes $z = x!$ using y as a temp variable.
- WHILE uses assignment statements, if-then-else, while loops.
- All vars are integers.
- Expressions only arithmetic (for vars) or relational (for conditions).
- No I/O statements. Inputs and outputs are implicit.
 - Later on, we may use extensions with explicit `read x` and `print x`.

WHILE abstract syntax

S statements
 a arithmetic expressions (AExp)
 x, y program variables (Vars)
 n number literals
 b boolean expressions (BExp)

We'll use these meta-variables frequently for ease of notation

$S ::= x := a$	$b ::= \text{true}$	$a ::= x$	$op_b ::= \text{and} \mid \text{or}$
$ \quad \text{skip}$	$ \quad \text{false}$	$ \quad n$	$op_r ::= < \mid \leq \mid =$
$ \quad S_1; S_2$	$ \quad \text{not } b$	$ \quad a_1 op_a a_2$	$ \quad > \mid \geq$
$ \quad \text{if } b \text{ then } S_1 \text{ else } S_2$	$ \quad b_1 op_b b_2$		$op_a ::= + \mid - \mid * \mid /$
$ \quad \text{while } b \text{ do } S$	$ \quad a_1 op_r a_2$		

Our first static analysis: AST walking

- One way to find “bugs” is to walk the AST, looking for particular patterns.
 - Traverse the AST, look for nodes of a particular type
 - Check the neighborhood of the node for the pattern in question.
 - Basically, a glorified “grep” that knows about the syntax but not semantics of a language.

CodeQL

- A language for querying code. Developed by GitHub.
- Supports many common languages.
- Library of common programming patterns and optimizations.

CodeQL queries 1.23

Dashboard / Java queries ...

Inefficient empty string test

Created by Documentation team, last modified on Mar 28, 2019

Name: Inefficient empty string test
Description: Checking a string for equality with an empty string is inefficient.
ID: java/inefficient-empty-string-test
Kind: problem
Severity: recommendation
Precision: high

Query: InefficientEmptyStringTest.ql Expand source

When checking whether a string `s` is empty, perhaps the most obvious solution is to write something like `s.equals("")` (or `"".equals(s)`). However, this actually carries a fairly significant overhead, because `String.equals` performs a number of type tests and conversions before starting to compare the content of the strings.

Recommendation

The preferred way of checking whether a string `s` is empty is to check if its length is equal to zero. Thus, the condition is `s.length() == 0`. The `length` method is implemented as a simple field access, and so should be noticeably faster than calling `equals`.

Note that in Java 6 and later, the `String` class has an `isEmpty` method that checks whether a string is empty. If the codebase does not need to support Java 5, it may be better to use that method instead.

Operational Semantics of WHILE

- The meaning of WHILE expressions depend on the values of variables
 - What does $x+5$ mean? It depends on x .
 - If $x = 8$ at some point, we expect $x+5$ to mean 13
- The value of integer variables at a given moment is abstracted as a function:

$$E : Var \rightarrow \mathbb{Z}$$

- We will augment our notation of big-step evaluation to include state:

$$\langle E, a \rangle \Downarrow n$$

- So, if $\{x \mapsto 8\} \in E$, then $\langle E, x + 5 \rangle \Downarrow 13$

Big-Step Semantics for WHILE expressions

$$\frac{}{\langle E, n \rangle \Downarrow n} \text{big-int} \quad \frac{}{\langle E, x \rangle \Downarrow E(x)} \text{big-var}$$

$$\frac{\langle E, a_1 \rangle \Downarrow n_1 \quad \langle E, a_2 \rangle \Downarrow n_2}{\langle E, a_1 + a_2 \rangle \Downarrow n_1 + n_2} \text{big-add}$$

- Similarly for other arithmetic and boolean expressions

Big-Step Semantics for WHILE statements

$$\frac{\langle E, b \rangle \Downarrow \text{false}}{\langle E, \text{while } b \text{ do } S \rangle \Downarrow E} \text{ big-whilefalse}$$

$$\frac{\langle E, b \rangle \Downarrow \text{true} \quad \langle E, S; \text{while } b \text{ do } S \rangle \Downarrow E'}{\langle E, \text{while } b \text{ then } S \rangle \Downarrow E'} \text{ big-whiletrue}$$

Small-Step Semantics for WHILE statements

$$\frac{\langle E, S_1 \rangle \rightarrow \langle E', S'_1 \rangle}{\langle E, S_1; S_2 \rangle \rightarrow \langle E', S'_1; S_2 \rangle} \text{ small-seq-congruence}$$

$$\frac{}{\langle E, \text{skip}; S_2 \rangle \rightarrow \langle E, S_2 \rangle} \text{ small-seq}$$

Small-Step Semantics for WHILE statements

$$\frac{\langle E, b \rangle \rightarrow_b b'}{\langle E, \text{if } b \text{ then } S_1 \text{ else } S_2 \rangle \rightarrow \langle E, \text{if } b' \text{ then } S_1 \text{ else } S_2 \rangle} \text{ small-if-congruence}$$
$$\frac{}{\langle E, \text{if true then } S_1 \text{ else } S_2 \rangle \rightarrow \langle E, S_1 \rangle} \text{ small-iftrue}$$

Proofs by Structural Induction

Example. Let $L(a)$ be the number of literals and variable occurrences in some expression a and $O(a)$ be the number of operators in a . Prove by induction on the structure of a that $\forall a \in \text{Aexp} . L(a) = O(a) + 1$:

Base cases:

- Case $a = n$. $L(a) = 1$ and $O(a) = 0$
- Case $a = x$. $L(a) = 1$ and $O(a) = 0$

Inductive case 1: Case $a = a_1 + a_2$

- By definition, $L(a) = L(a_1) + L(a_2)$ and $O(a) = O(a_1) + O(a_2) + 1$.
- By the induction hypothesis, $L(a_1) = O(a_1) + 1$ and $L(a_2) = O(a_2) + 1$.
- Thus, $L(a) = O(a_1) + O(a_2) + 2 = O(a) + 1$.

The other arithmetic operators follow the same logic.

Proofs by Structural Induction

- Prove that WHILE is *deterministic*. That is, if the program terminates, it evaluates to a unique value.

$$\forall a \in \mathbf{Aexp} . \quad \forall E . \quad \forall n, n' \in \mathbb{N} . \quad \langle E, a \rangle \Downarrow n \wedge \langle E, a \rangle \Downarrow n' \Rightarrow n = n'$$

$$\forall P \in \mathbf{Bexp} . \quad \forall E . \quad \forall b, b' \in \mathcal{B} . \quad \langle E, P \rangle \Downarrow b \wedge \langle E, P \rangle \Downarrow b' \Rightarrow b = b'$$

$$\forall S . \quad \forall E, E', E'' . \quad \langle E, S \rangle \Downarrow E' \wedge \langle E, S \rangle \Downarrow E'' \Rightarrow E' = E''$$

Rule for while is recursive;
doesn't depend only on
subexpressions

- Can prove for expressions via induction over syntax, but not for statements.
- But there's still a way.

To prove: $\boxed{\forall S . \quad \forall E, E', E'' . \quad \langle E, S \rangle \Downarrow E' \wedge \langle E, S \rangle \Downarrow E'' \Rightarrow E' = E''}$

Structural Induction over Derivations

Base case: the one rule with no premises, `skip`:

let $D ::= \langle E, \text{skip} \rangle \Downarrow E'$, and let $D' ::= \langle E, S \rangle \Downarrow E''$

$$D ::= \overline{\langle E, \text{skip} \rangle \Downarrow E}$$

By inversion, the last rule used in D' (which, again, produced E'') must also have been the rule for `skip`. By the structure of the `skip` rule, we know $E'' = E$.

Inductive cases: We need to show that the property holds when the last rule used in D was each of the possible non-skip WHILE commands. I will show you one representative case; the rest are left as an exercise. If the last rule used was the `while-true` statement:

$$D ::= \frac{D_1 ::= \langle E, b \rangle \Downarrow \text{true} \quad D_2 ::= \langle E, S \rangle \Downarrow E_1 \quad D_3 ::= \langle E_1, \text{while } b \text{ do } S \rangle \Downarrow E'}{\langle E, \text{while } b \text{ do } S \rangle \Downarrow E'}$$

Pick arbitrary E'' such that $D' ::= \langle E, \text{while } b \text{ do } S \rangle \Downarrow E''$

By inversion, D' must use either the `while-true` or the `while-false` rule. However, having proved that boolean expressions are deterministic (via induction on syntax), and given that D contains the judgment $\langle E, b \rangle \Downarrow \text{true}$, we know that D' cannot be the `while-false` rule, as otherwise it would have to contain a contradicting judgment $\langle E, b \rangle \Downarrow \text{false}$.

So, we know that D' is also using `while-true` rule. In its derivation, D' must also have subderivations $D'_2 ::= \langle E, S \rangle \Downarrow E'_1$ and $D'_3 ::= \langle E'_1, \text{while } b \text{ do } S \rangle \Downarrow E''$. By the induction hypothesis on D_2 with D'_2 , we know $E_1 = E'_1$. Using this result and the induction hypothesis on D_3 with D'_3 , we have $E'' = E'$.

Data-Flow Analysis

Computes universal properties about program state at specific program points. (e.g. will x be zero at line 7?)

- About program state
 - About data store (e.g. variables, heap memory)
 - Not about control (e.g. termination, performance)
- At program points
 - Statically identifiable (e.g. line 7, or when `foo()` calls `bar()`)
 - Not dynamically computed (E.g. when x is 12 or when `foo()` is invoked 12 times)
- Universal
 - Reasons about all possible executions (always/never/maybe)
 - Not about specific program paths (see: symbolic execution, testing)

WHILE3ADDR: An Intermediate Representation

Simpler, more uniform than WHILE syntax

Categories:

$/ \in \text{Instruction}$	instructions
$x, y \in \text{Var}$	variables
$n \in \text{Num}$	number literals

Syntax:

$$\begin{aligned} / ::= & x := n \mid x := y \mid x := y \ op \ z \\ & \mid \text{goto } n \mid \text{if } x \ op_r \theta \text{ goto } n \\ op_a ::= & + \mid - \mid * \mid / \mid \dots \\ op_r ::= & < \mid \leq \mid = \mid > \mid \geq \mid \dots \\ P \in \text{Num} \rightarrow / \end{aligned}$$

$$\frac{P(n) = x := m}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E[x \mapsto m], n + 1 \rangle} \text{step-const}$$

$$\frac{P[n] = x := y}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E[x \mapsto E(y)], n + 1 \rangle} \text{step-copy}$$

$$\frac{P(n) = x := y \text{ op } z \quad E(y) \textbf{ op } E(z) = m}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E[x \mapsto m], n + 1 \rangle} \text{step-arith}$$

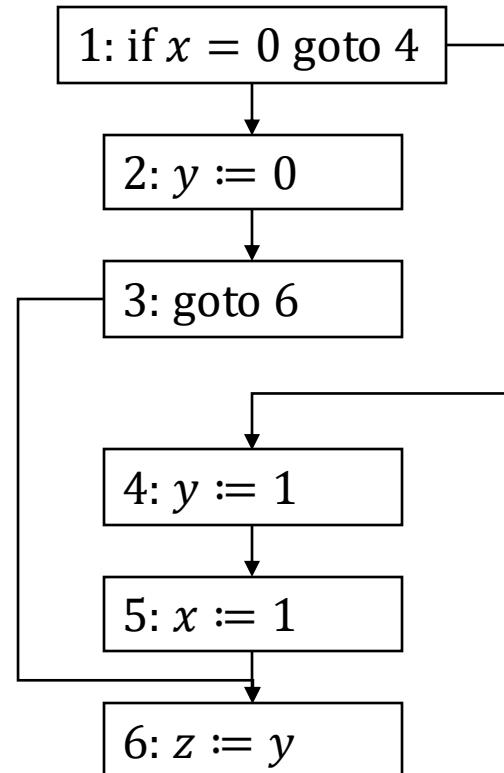
$$\frac{P(n) = \text{goto } m}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E, m \rangle} \text{step/goto}$$

$$\frac{P(n) = \text{if } x \text{ op}_r 0 \text{ goto } m \quad E(x) \textbf{ op_r } 0 = \text{true}}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E, m \rangle} \text{step-iftrue}$$

$$\frac{P(n) = \text{if } x \text{ op}_r 0 \text{ goto } m \quad E(x) \textbf{ op_r } 0 = \text{false}}{P \vdash \langle E, n \rangle \rightsquigarrow \langle E, n + 1 \rangle} \text{step-iffalse}$$

Control-flow Graphs

```
1 : if x = 0 goto 4
2 : y := 0
3 : goto 6
4 : y := 1
5 : x := 1
6 : z := y
```



Nodes = Statements

Edges = (s_1, s_2) is an edge iff s_1 and s_2 can be executed consecutively
aka "control flow"

Common properties of CFGs:

- Weakly connected
- Only one entry node
- Only one exit (terminal) node

Classic Data-Flow Analyses

- Zero Analysis
- Integer Sign Analysis
- Constant Propagation
- Reaching Definitions
- Live Variables Analysis

- Available Expressions
- Very Busy Expressions
- ...

Partial Order & Join on set L

$l_1 \sqsubseteq l_2$: l_1 is at least as precise as l_2

reflexive: $\forall l : l \sqsubseteq l$

transitive: $\forall l_1, l_2, l_3 : l_1 \sqsubseteq l_2 \wedge l_2 \sqsubseteq l_3 \Rightarrow l_1 \sqsubseteq l_3$

anti-symmetric: $\forall l_1, l_2 : l_1 \sqsubseteq l_2 \wedge l_2 \sqsubseteq l_1 \Rightarrow l_1 = l_2$

$l_1 \sqcup l_2$: **join** or *least-upper-bound*... “most precise generalization”

L is a *join-semilattice* iff: $l_1 \sqcup l_2$ always exists and is unique $\forall l_1, l_2 \in L$

\top (“top”) is the maximal element

Fixed point of Flow Functions

$$(\sigma_0, \sigma_1, \sigma_2, \dots, \sigma_n) \xrightarrow{f_z} (\sigma'_0, \sigma'_1, \sigma'_2, \dots, \sigma'_n)$$

Fixed point!

$$(\sigma_0, \sigma_1, \sigma_2, \dots, \sigma_n) = f_z(\sigma_0, \sigma_1, \sigma_2, \dots, \sigma_n)$$

Correctness theorem:

If data-flow analysis is well designed*, then
any fixed point of the analysis is sound.

$$\sigma'_0 = \sigma_0$$

$$\sigma'_1 = f_z[\![x := 10]\!](\sigma_0)$$

$$\sigma'_2 = f_z[\![y := 0]\!](\sigma_1)$$

$$\sigma'_3 = \sigma_2 \sqcup \sigma_7$$

$$\sigma'_4 = f_z[\![\text{if } x = 10 \text{ goto 7}]\!]_F(\sigma_3)$$

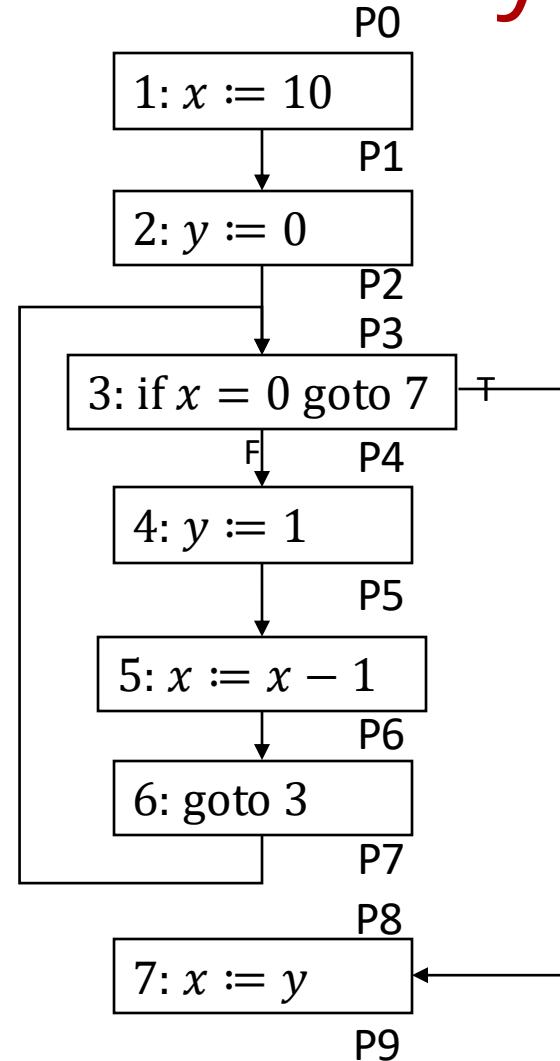
:

$$\sigma'_8 = f_z[\![\text{if } x = 10 \text{ goto 7}]\!]_T(\sigma_3)$$

$$\sigma'_9 = f_z[\![x := y]\!](\sigma_8)$$

Example of Zero Analysis: Looping Code

```
1 : x := 10
2 : y := 0
3 : if x = 0 goto 7
4 : y := 1
5 : x := x - 1
6 : goto 3
7 : x := y
```



	x	y	
P0	T	T	
P1	N	T	
P2	N	Z	
P3	T	T	<i>join</i>
P4	N_F	T	<i>updated</i>
P5	N	N	<i>already at fixed point</i>
P6	T	N	<i>already at fixed point</i>
P7	T	N	<i>already at fixed point</i>
P8	Z_T	T	<i>updated</i>
P9	T	T	<i>updated</i>

Kildall's Algorithm

```
worklist = Ø
for Node n in cfg
    input[n] = output[n] = ⊥
    add n to worklist
input[0] = initialDataflowInformation

while worklist is not empty
    take a Node n off the worklist
    output[n] = flow(n, input[n])
    for Node j in succs(n)
        newInput = input[j] ∪ output[n]
        if newInput ≠ input[j]
            input[j] = newInput
            add j to worklist
```

Worklist Algorithm Terminates at Fixed Point

At the fixed point, we therefore have the following equations satisfied:

$$\sigma_0 \sqsubseteq \sigma_1$$
$$\forall i \in P : \left(\bigsqcup_{j \in \text{preds}(i)} f[\![P[j]]\!](\sigma_j) \right) \sqsubseteq \sigma_i$$

The worklist algorithm shown above computes a fixed point when it terminates. We can prove this by showing that the following loop invariant is maintained:

$$\forall i . (\exists j \in \text{preds}(i) \text{ such that } f[\![P[j]]\!](\sigma_j) \not\sqsubseteq \sigma_i) \Rightarrow i \in \text{worklist}$$

Program Traces and DataFlow Soundness

A trace T of a program P is a potentially infinite sequence $\{c_0, c_1, \dots\}$ of program configurations, where $c_0 = E_0, 1$ is called the initial configuration, and for every $i \geq 0$ we have

$$P \vdash c_i \rightsquigarrow c_{i+1}$$

The result $\langle \sigma_n \mid n \in P \rangle$ of a program analysis running on program P is sound iff, for all traces T of P , for all i such that $0 \leq i < \text{length}(T)$, $\alpha(c_i) \sqsubseteq \sigma_{n_i}$

Fixed Point Theorem

Theorem 2 (A fixed point of a locally sound analysis is globally sound). *If a dataflow analysis's flow function f is monotonic and locally sound, and for all traces T we have $\alpha(c_0) \sqsubseteq \sigma_0$ where σ_0 is the initial analysis information, then any fixed point $\{\sigma_n \mid n \in P\}$ of the analysis is sound.*

Proof. To show that the analysis is sound, we must prove that for all program traces, every program configuration in that trace is correctly approximated by the analysis results. We consider an arbitrary program trace T and do the proof by induction on the program configurations $\{c_i\}$ in the trace.

Least Fixed Point (LFP)

The least fixed point solution of a composite flow function \mathcal{F} is the fixed-point result Σ^* such that $\mathcal{F}(\Sigma^*) = \Sigma^*$ and $\forall \Sigma : (\mathcal{F}(\Sigma) = \Sigma) \Rightarrow (\Sigma^* \sqsubseteq \Sigma)$.

Merge Over Paths (MOP)

We first enumerate all paths π of the form $\pi = n_1, n_2, \dots$ in the control-flow graph, where n_i are the instructions (nodes) in the path. For each such path π , we successively apply flow functions to form the sequence of tuples $\Pi = \langle \sigma_1, n_1 \rangle, \langle \sigma_2, n_2 \rangle, \dots$ such that $\sigma_{\Pi_j} = f[\![P[n_{\Pi_j}]]\!](\sigma_{\Pi_{j-1}})$, where Π_j is the j -th tuple in the sequence and σ_0 is the initial data flow information. We then join over all σ values computed for an instruction i to get the MOP:

$$\text{MOP}(i) = \bigsqcup \{\sigma \mid \langle \sigma, i \rangle \in \text{Some } \Pi \text{ for } P\}$$

The MOP solution is the most precise result if we consider all possible program paths through the CFG, even though it may be less precise than the optimal solution due to the consideration of infeasible paths. The MOP is computable when flow functions are *distributive* over join.

Distributivity

A function f is *distributive* iff $f(\sigma_1) \sqcup f(\sigma_2) = f(\sigma_1 \sqcup \sigma_2)$

Reaching Definitions

$$f_{RD}[\![I]\!](\sigma) = \sigma - KILL_{RD}[\![I]\!] \cup GEN_{RD}[\![I]\!]$$

$$KILL_{RD}[\![n: x := ...]\!] = \{x_m \mid x_m \in \text{DEFS}(x)\}$$

$$KILL_{RD}[\![I]\!] = \emptyset \quad \text{if } I \text{ is not an assignment}$$

$$GEN_{RD}[\![n: x := ...]\!] = \{x_n\}$$

$$GEN_{RD}[\![I]\!] = \emptyset \quad \text{if } I \text{ is not an assignment}$$

Live Variables

Flow functions map backward! (out --> in)

$$\mathbf{KILL}_{LV}[I] = \{x \mid I \text{ defines } x\}$$

$$\mathbf{GEN}_{LV}[I] = \{x \mid I \text{ uses } x\}$$

Constant Propagation

$$\begin{aligned}\sigma &\in \textit{Var} \rightarrow L_{CP} \\ \sigma_1 \sqsubseteq_{lift} \sigma_2 &\text{ iff } \forall x \in \textit{Var} : \sigma_1(x) \sqsubseteq \sigma_2(x) \\ \sigma_1 \sqcup_{lift} \sigma_2 &= \{x \mapsto \sigma_1(x) \sqcup \sigma_2(x) \mid x \in \textit{Var}\} \\ \top_{lift} &= \{x \mapsto \top \mid x \in \textit{Var}\} \\ \perp_{lift} &= \{x \mapsto \perp \mid x \in \textit{Var}\} \\ \alpha_{CP}(n) &= n \\ \alpha_{lift}(E) &= \{x \mapsto \alpha_{CP}(E(x)) \mid x \in \textit{Var}\} \\ \sigma_0 &= \top_{lift}\end{aligned}$$

Extend WHILE3ADDR with functions

$F ::= \text{fun } f(x) \{ \overline{n : I} \}$

$I ::= \dots \mid \text{return } x \mid y := f(x)$

1 : $\text{fun } double(x) : int$

2 : $y := 2 * x$

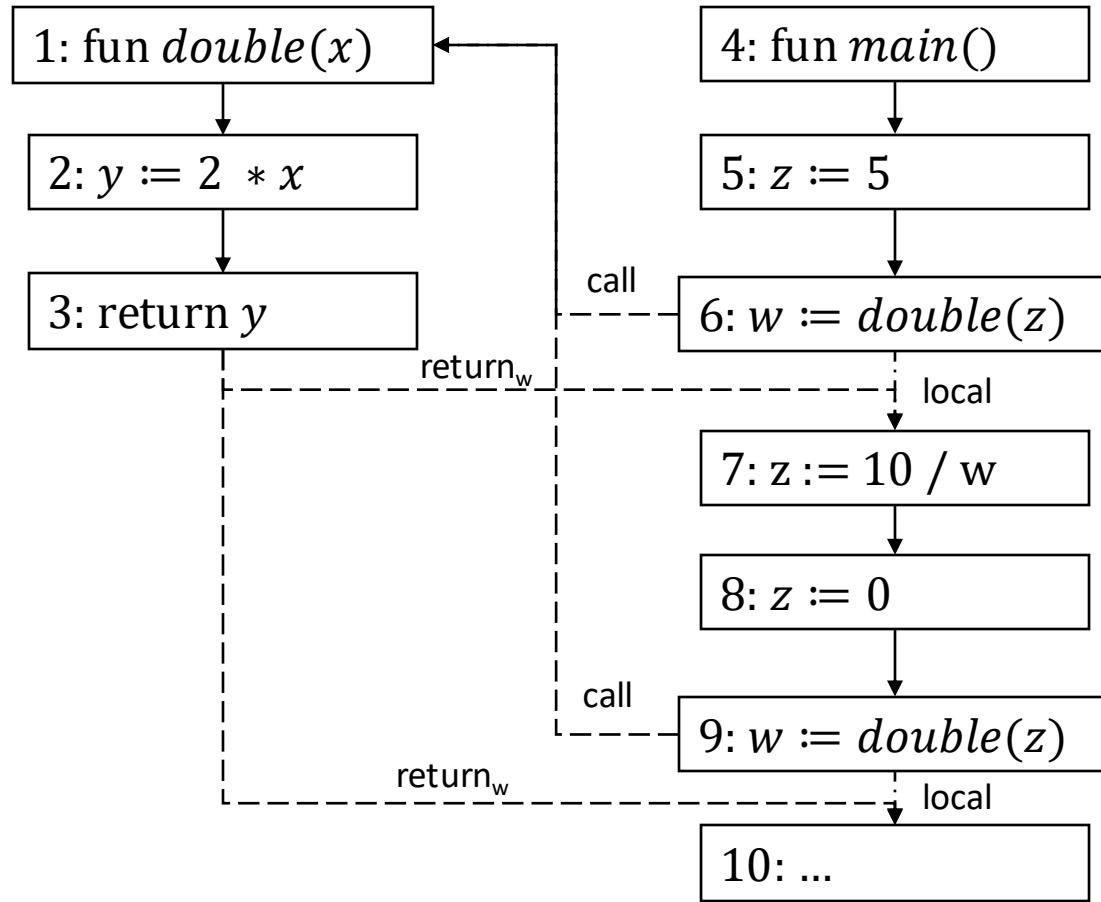
3 : $\text{return } y$

4 : $\text{fun } main() : void$

5 : $z := 0$

6 : $w := double(z)$

Interprocedural CFG



```

1 : fun double(x) : int
2 :   y := 2 * x
3 :   return y
4 : fun main()
5 :   z := 5
6 :   w := double(z)
7 :   z := 10 / w
8 :   z := 0
9 :   w := double(z)

```

$$\begin{aligned}
 f_Z[x := g(y)]_{local}(\sigma) &= \sigma \setminus (\{x\} \cup \text{Globals}) \\
 f_Z[x := g(y)]_{call}(\sigma) &= \{v \mapsto \sigma(v) \mid v \in \text{Globals}\} \cup \{\text{formal}(g) \mapsto \sigma(y)\} \\
 f_Z[\text{return } y]_{return_x}(\sigma) &= \{v \mapsto \sigma(v) \mid v \in \text{Globals}\} \cup \{x \mapsto \sigma(y)\}
 \end{aligned}$$

Context-Sensitive Analysis Example

```
1 : fun double(x) : int
2 :     y := 2 * x
3 :     return y
4 : fun main()
5 :     z := 5
6 :     w := double(z)
7 :     z := 10/w
8 :     z := 0
9 :     w := double(z)
```

Context	σ_{in}	σ_{out}
<main, T>	T	{w->Z, Z->Z}
<double, N>	{x->N}	{x->N, y->N}
<double, Z>	{x->Z}	{x->Z, y->Z}

Key Idea: Worklist of Contexts

```
val worklist : Set[Context]
val analyzing : Set[Context]
val results : Map[Context, Summary]
val callers : Map[Context, Set[Context]]
```

```
function ANALYZEPROGRAM
  initCtx ← GETCTX(main, nil, 0, ⊤)
  worklist ← {initCtx}
  results[initCtx] ← Summary(⊤, ⊥)
  while NOTEMPTY(worklist) do
    ctx ← REMOVE(worklist)
    ANALYZE(ctx, results[ctx].input)
  end while
end function
```

```
function ANALYZE(ctx, σin)
  σout ← results[ctx].output
  ADD(analyzing, ctx)
  σ'out ← INTRAPROCEDURAL(ctx, σin)
  REMOVE(analyzing, ctx)
  if σ'out ≠ σout then
    results[ctx] ← Summary(σin, σout ∪ σ'out)
    for c ∈ callers[ctx] do
      ADD(worklist, c)
    end for
  end if
  return σ'out
end function
```

```

function RESULTSFOR( $ctx, \sigma_{in}$ )
  if  $ctx \in \text{dom}(\text{results})$  then
    if  $\sigma_{in} \sqsubseteq \text{results}[ctx].\text{input}$  then
      return  $\text{results}[ctx].\text{output}$ 
    else
       $\text{results}[ctx].\text{input} \leftarrow \text{results}[ctx].\text{input} \sqcup \sigma_{in}$   $\triangleright$  keep track of more general input
    end if
  else
     $\text{results}[ctx] = \text{Summary}(\sigma_{in}, \perp)$   $\triangleright$  initially optimistic
  end if
  if  $ctx \in \text{analyzing}$  then
    return  $\text{results}[ctx].\text{output}$   $\triangleright \perp$  if it hasn't been analyzed yet; otherwise
  else
    return ANALYZE( $ctx, \text{results}[ctx].\text{input}$ )
  end if
end function

```

```

function FLOW( $\llbracket n : x := f(y) \rrbracket, ctx, \sigma_n$ )
   $\sigma_{in} \leftarrow [formal(f) \mapsto \sigma_n(y)]$   $\triangleright$ 
   $calleeCtx \leftarrow \text{GETCTX}(f, ctx, n, \sigma_{in})$ 
   $\sigma_{out} \leftarrow \text{RESULTSFOR}(calleeCtx, \sigma_{in})$ 
  ADD( $\text{callers}[calleeCtx]$ ,  $ctx$ )
  return  $\sigma_n[x \mapsto \sigma_{out}[\text{result}]]$ 

```

```

function ANALYZE( $ctx, \sigma_{in}$ )
   $\sigma_{out} \leftarrow \text{results}[ctx].\text{output}$ 
  ADD( $\text{analyzing}, ctx$ )
   $\sigma'_{out} \leftarrow \text{INTRAPROCEDURAL}(ctx, \sigma_{in})$ 
  REMOVE( $\text{analyzing}, ctx$ )
  if  $\sigma'_{out} \not\sqsubseteq \sigma_{out}$  then
     $\text{results}[ctx] \leftarrow \text{Summary}(\sigma_{in}, \sigma_{out} \sqcup \sigma'_{out})$ 
    for  $c \in \text{callers}[ctx]$  do
      ADD( $\text{worklist}, c$ )
    end for
  end if
  return  $\sigma'_{out}$ 
end function

```

Extending WHILE3ADDR with Pointers

$I ::= \dots$

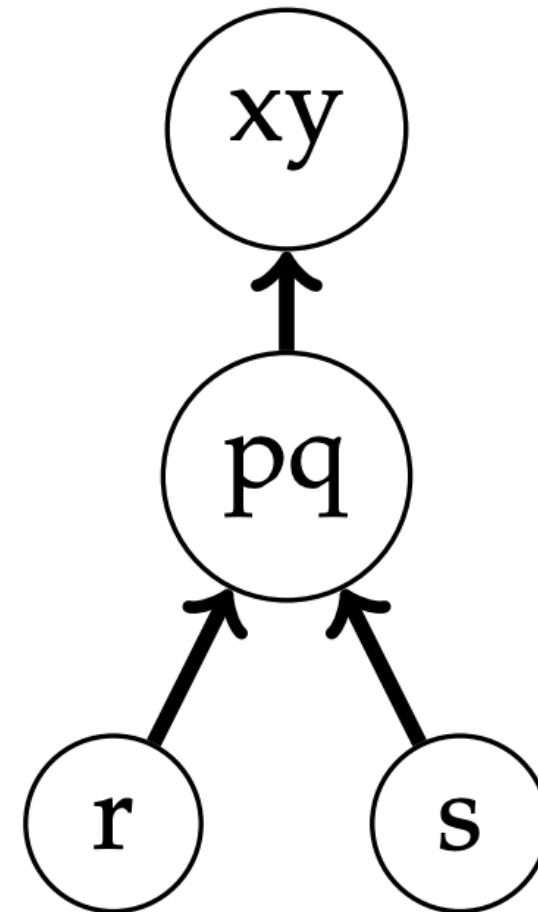
- | $p := \&x$ taking the address of a variable
- | $p := q$ copying a pointer from one variable to another
- | $*p := q$ assigning through a pointer
- | $p := *q$ dereferencing a pointer

Andersen's Analysis

$$\frac{\overline{[\![p := \&x]\!] \hookrightarrow l_x \in p} \quad \text{address-of}}{[\![p := q]\!] \hookrightarrow p \supseteq q} \text{ copy}$$
$$\frac{l_x \in q}{\exists p} \text{ copy}$$
$$\frac{\overline{[\![*p := q]\!] \hookrightarrow *p \supseteq q} \quad \text{assign}}{\vdash p \quad l_x \in q} \text{ assign}$$
$$\frac{\overline{[\![p := *q]\!] \hookrightarrow p \supseteq *q} \quad \text{derefence}}{l_x \in r} \text{ dereference}$$
$$\frac{}{[\![p := *q]\!] \hookrightarrow p \supseteq *q} \text{ dereference}$$

Steengaard's Analysis - Example

```
1 : p := &x  
2 : r := &p  
3 : q := &y  
4 : s := &q  
5 : r := s
```



Steensgaard's Analysis

	join(ℓ_1, ℓ_2)	
$\boxed{[p := q]} \hookrightarrow \text{join}(*p, *q)$	<i>copy</i>	if (find(ℓ_1) == find(ℓ_2)) return
$\boxed{[p := \&x]} \hookrightarrow \text{join}(*p, x)$	<i>address-of</i>	$n_1 \leftarrow * \ell_1$ $n_2 \leftarrow * \ell_2$ union(ℓ_1, ℓ_2)
$\boxed{[p := *q]} \hookrightarrow \text{join}(*p, **q)$	<i>dereference</i>	join(n_1, n_2)
$\boxed{[*p := q]} \hookrightarrow \text{join}(**p, *q)$	<i>assign</i>	

Hoare Triple

$$\{ P \} S \{ Q \}$$

- P is the precondition
- Q is the postcondition
- S is any statement (in WHILE, at least for our class)
- Semantics: if P holds in some state E and if $\langle S; E \rangle \Downarrow E'$, then Q holds in E'
 - This is *partial correctness*: termination of S is not guaranteed
 - *Total correctness* additionally implies termination, and is written $[P] S [Q]$

Semantics of Hoare Triples

- A partial correctness assertion $\models \{P\} S \{Q\}$ is defined formally to mean:

$$\forall E. \forall E'. (E \models P \wedge \langle E, S \rangle \Downarrow E') \Rightarrow E' \models Q$$

- How would we define total correctness $[P] S [Q]$?
- This is a good formal definition—but it doesn't help us prove many assertions because we have to reason about all environments. How can we do better?

Derivation Rules for Hoare Logic

- Judgment form $\vdash \{P\} S \{Q\}$ means “we can prove the Hoare triple $\{P\} S \{Q\}$ ”

$$\frac{}{\vdash \{P\} \text{skip } \{P\}} \textit{skip} \quad \frac{}{\vdash \{[a/x]P\} x:=a \{P\}} \textit{assign}$$

$$\frac{\vdash \{P\} S_1 \{P'\} \quad \vdash \{P'\} S_2 \{Q\}}{\vdash \{P\} S_1; S_2 \{Q\}} \textit{seq} \quad \frac{\vdash \{P \wedge b\} S_1 \{Q\} \quad \vdash \{P \wedge \neg b\} S_2 \{Q\}}{\vdash \{P\} \text{if } b \text{ then } S_1 \text{ else } S_2 \{Q\}} \textit{if}$$

$$\frac{\vdash P' \Rightarrow P \quad \vdash \{P\} S \{Q\} \quad \vdash Q \Rightarrow Q'}{\vdash \{P'\} S \{Q'\}} \textit{consq}$$

Hoare Triples and Weakest Preconditions

- Theorem: $\{P\} S \{Q\}$ holds if and only if $P \Rightarrow wp(S, Q)$
 - In other words, a Hoare Triple is still valid if the precondition is stronger than necessary, but not if it is too weak
 - Can use this to prove $\{P\} S \{Q\}$ by computing $wp(S, Q)$ and checking implication.
- Question: Could we state a similar theorem for a strongest postcondition function?
 - e.g. $\{P\} S \{Q\}$ holds if and only if $sp(S, P) \Rightarrow Q$
 - A: Yes, but it's harder to compute (see text for why)

Proving loops correct

- First consider *partial correctness*
 - The loop may not terminate, but if it does, the postcondition will hold
- $\{P\} \text{ while } B \text{ do } S \{Q\}$
 - Find an invariant Inv such that:
 - $P \Rightarrow \text{Inv}$
 - The invariant is initially true
 - $\{ \text{Inv} \&& B \} S \{\text{Inv}\}$
 - Each execution of the loop preserves the invariant
 - $(\text{Inv} \&& \neg B) \Rightarrow Q$
 - The invariant and the loop exit condition imply the postcondition

What if we just went forwards?

$\{P\}$

$x := e_1$

$x := e_2$

$\{Q\}$

Generate “fresh” math variables
for every mutable program
variable

Proof Obligation:

$$\forall x_n : ([x_0/x]P \wedge x_1 = [x_0/x]e_1 \wedge x_2 = ([x_1/x]e_2)) \Rightarrow [x_2/x]Q$$

Dealing with conditional paths

$\{true\}$

if ($x < 0$) :

$y := -x$

Dynamic Symbolic Execution:

$$\forall x_0, y_0 \in \mathbb{Z} : (x_0 < 0 \wedge y_0 = -x_0) \Rightarrow y_0 \geq 0$$

else :

$y := x$

$\{y \geq 0\}$

$$\forall x_0, y_0 \in \mathbb{Z} : (x_0 \geq 0 \wedge y_0 = x_0) \Rightarrow y_0 \geq 0$$

Static Symbolic Execution:

$$\forall x_0, y_0 \in \mathbb{Z} : ((x_0 < 0 \Rightarrow y_0 = -x_0) \wedge (x_0 \geq 0 \Rightarrow y_0 = x_0)) \Rightarrow y_0 \geq 0$$

Symbolic Execution of Statements (DSE)

$$\frac{}{\langle g, \Sigma, \text{skip} \rangle \Downarrow \langle g, \Sigma \rangle} \textit{big-skip}$$

$$\frac{\langle g, \Sigma, s_1 \rangle \Downarrow \langle g', \Sigma' \rangle \quad \langle g', \Sigma', s_2 \rangle \Downarrow \langle g'', \Sigma'' \rangle}{\langle g, \Sigma, s_1; s_2 \rangle \Downarrow \langle g'', \Sigma'' \rangle} \textit{big-seq}$$

$$\frac{\langle \Sigma, a \rangle \Downarrow a_s}{\langle g, \Sigma, x := a \rangle \Downarrow \langle g, \Sigma[x \mapsto a_s] \rangle} \textit{big-assign}$$

Symbolic Execution with Branching (DSE)

$$\frac{\langle \Sigma, b \rangle \Downarrow g' \quad g \wedge g' \text{SAT} \quad \langle g \wedge g', \Sigma, s_1 \rangle \Downarrow \langle g'', \Sigma' \rangle}{\langle g, \Sigma, \text{if } b \text{ then } s_1 \text{ else } s_2, \rangle \Downarrow \langle g'', \Sigma' \rangle} \text{ big-iftrue}$$

$$\frac{\langle \Sigma, b \rangle \Downarrow g' \quad g \wedge \neg g' \text{SAT} \quad \langle g \wedge \neg g', \Sigma, s_2 \rangle \Downarrow \langle g'', \Sigma' \rangle}{\langle g, \Sigma, \text{if } b \text{ then } s_1 \text{ else } s_2, \rangle \Downarrow \langle g'', \Sigma' \rangle} \text{ big-iffalse}$$

Symbolic Execution of Loops

Bounded exploration (k -limited)

$$\frac{k > 0 \quad \langle \Sigma, b \rangle \Downarrow g' \quad g \wedge g' \text{SAT} \quad \langle g \wedge g', \Sigma, s; \text{while}_{k-1} b \text{ do } s \rangle \Downarrow \langle g'', \Sigma' \rangle}{\langle g, \Sigma, \text{while}_k b \text{ do } s, \rangle \Downarrow \langle g'', \Sigma' \rangle} \text{ big-whiletrue}$$

$$\frac{\langle \Sigma, b \rangle \Downarrow g' \quad g \wedge \neg g' \text{SAT}}{\langle g, \Sigma, \text{while}_k b \text{ do } s, \rangle \Downarrow \langle g \wedge \neg g', \Sigma \rangle} \text{ big-whilefalse}$$

Satisfiability (SAT) solving

- Let's start by considering Boolean formulas:
variables connected with $\wedge \vee \neg$
- First step: convert to conjunctive normal form (CNF)
 - A conjunction of disjunctions of (possibly negated) variables
$$(a \vee \neg b) \wedge (\neg a \vee c) \wedge (b \vee c)$$
- If formula is not in CNF, we transform it: use De Morgan's laws, the double negative law, and the distributive laws:

$$\neg(P \vee Q) \iff \neg P \wedge \neg Q$$

$$\neg(P \wedge Q) \iff \neg P \vee \neg Q$$

$$\neg\neg P \iff P$$

$$(P \wedge (Q \vee R)) \iff ((P \wedge Q) \vee (P \wedge R))$$

$$(P \vee (Q \wedge R)) \iff ((P \vee Q) \wedge (P \vee R))$$

Concolic Execution

```
1 int double (int v) {  
2     return 2*v;  
3 }  
4  
5 void bar(int x, int y) {  
6     z = double (y);  
7     if (z == x) {  
8         if (x > y+10) {  
9             ERROR;  
10        }  
11    }  
12 }
```

1. Input: $x=0, y=1$
 - Path: $(2*y \neq x)$
 - Next: $(2*y == x) :: \text{SAT}$
2. Input: $x=2, y=1$
 - Path: $(2*y == x) \&& (x \leq y+10)$
 - Next: $(2*y == x) \&& (x > y+10) :: \text{SAT}$
3. Input: $x=22, y=11$
 - Path: $(2*y == x) \&& (x > y+10)$
 - **Bug found!!**

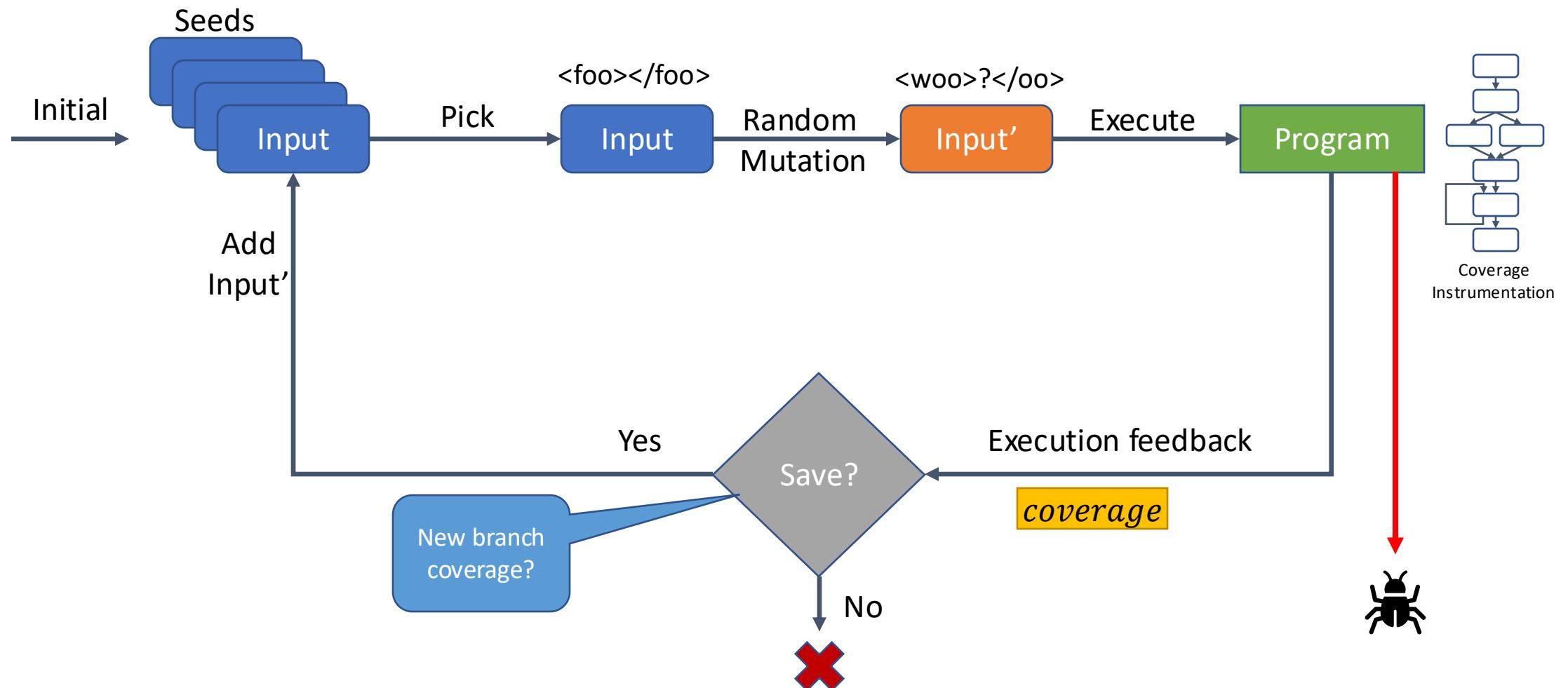
Dynamic analysis

- Observe program behavior during *execution* on *one or more inputs*.
- Examples:
 - Code coverage (→ Greybox fuzzing, fault localization)
 - Performance Profiling
 - Code profiling, memory profiling, algorithmic profiling
 - Invariant Generation
 - Concolic Execution
 - Data structure analysis
 - Concurrency analysis: Race detection
 - Concurrency analysis: Deadlock detection
 - Taint Analysis (→ Security & Privacy)
 - ... (many many more)

Collecting execution info

- What to collect? *Only what's necessary*
- Key idea (again): *Abstraction*
- Examples:
 - Code coverage → Log branches
 - Profiling → Log loops, function calls, allocations, frees, etc.
 - Invariant generation → Log predicates over vars in scope
 - Concolic execution → Track symbolic values; log branch constraints
 - Race detection → Track locks, vector clocks; log accesses

Coverage-Guided Fuzzing with AFL



Mutation testing measures test effectiveness on artificial bugs

```
bool is_even(int x) {  
    int a = x / 2;  
    int b = a * 2;  
    return x == b;  
}
```

```
void test_even() {  
    assert(is_even(4) == true);  
}
```

Kills mutant 1 & 3 but not 2

```
bool is_even(int x) {  
    int a = x / 2;  
    int b = a * 2;  
    return x != b;  
}
```

```
bool is_even(int x) {  
    int a = x / 2;  
    int b = a * 2;  
    return b == b;  
}
```

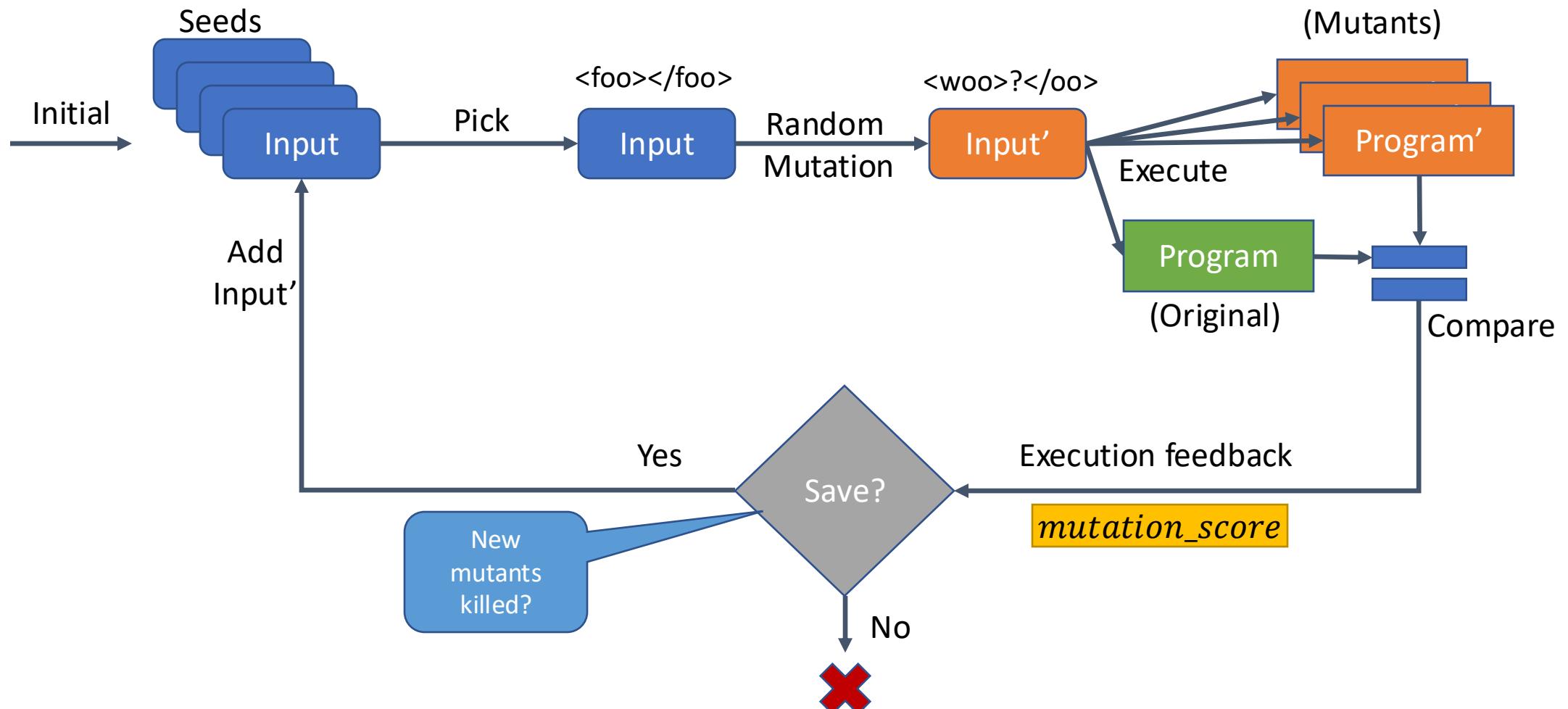
```
bool is_even(int x) {  
    int a = x / 2;  
    int b = a * 2;  
    return x == a;  
}
```

Mutant 1

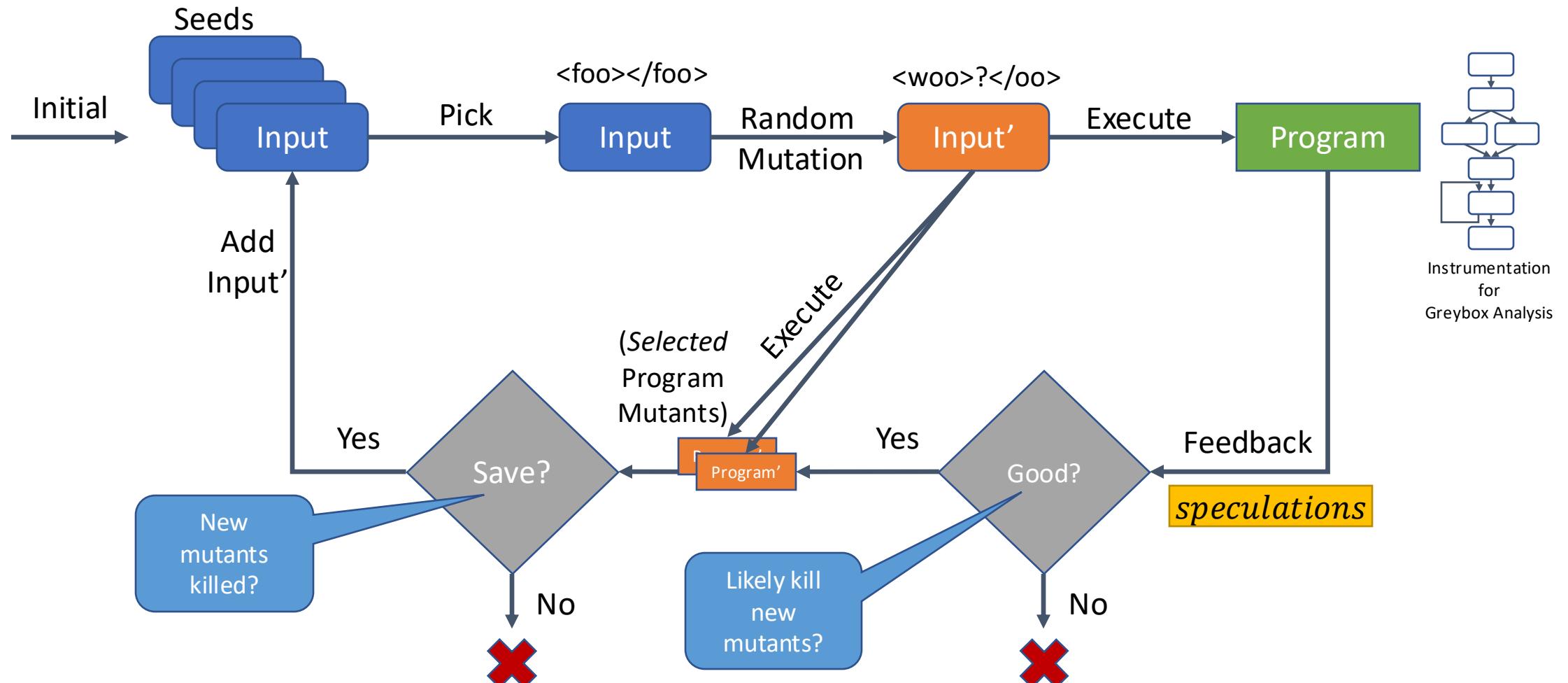
Mutant 2

Mutant 3

Mutation-analysis guided fuzzing

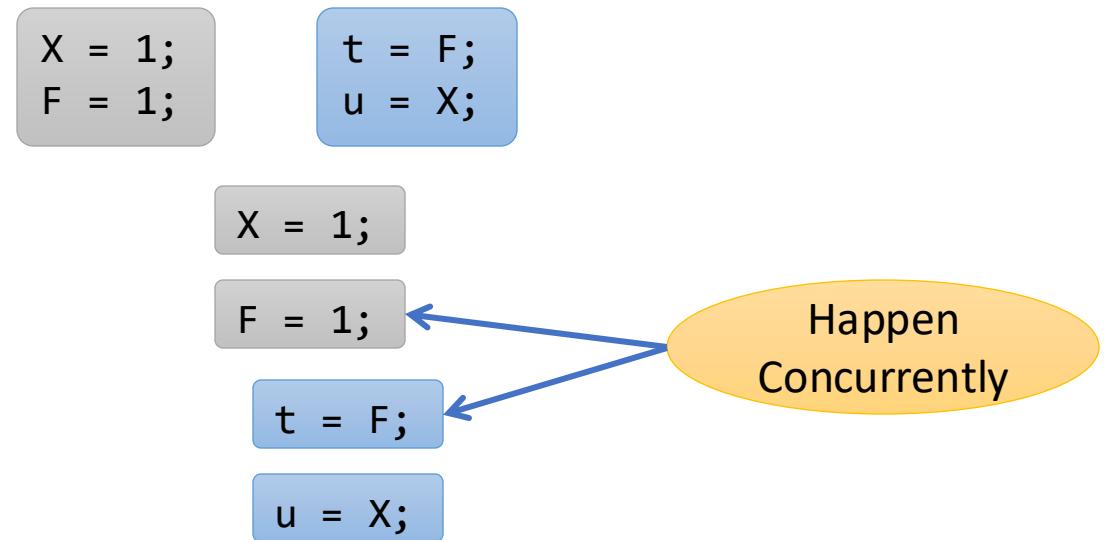


Speculative mutation analysis

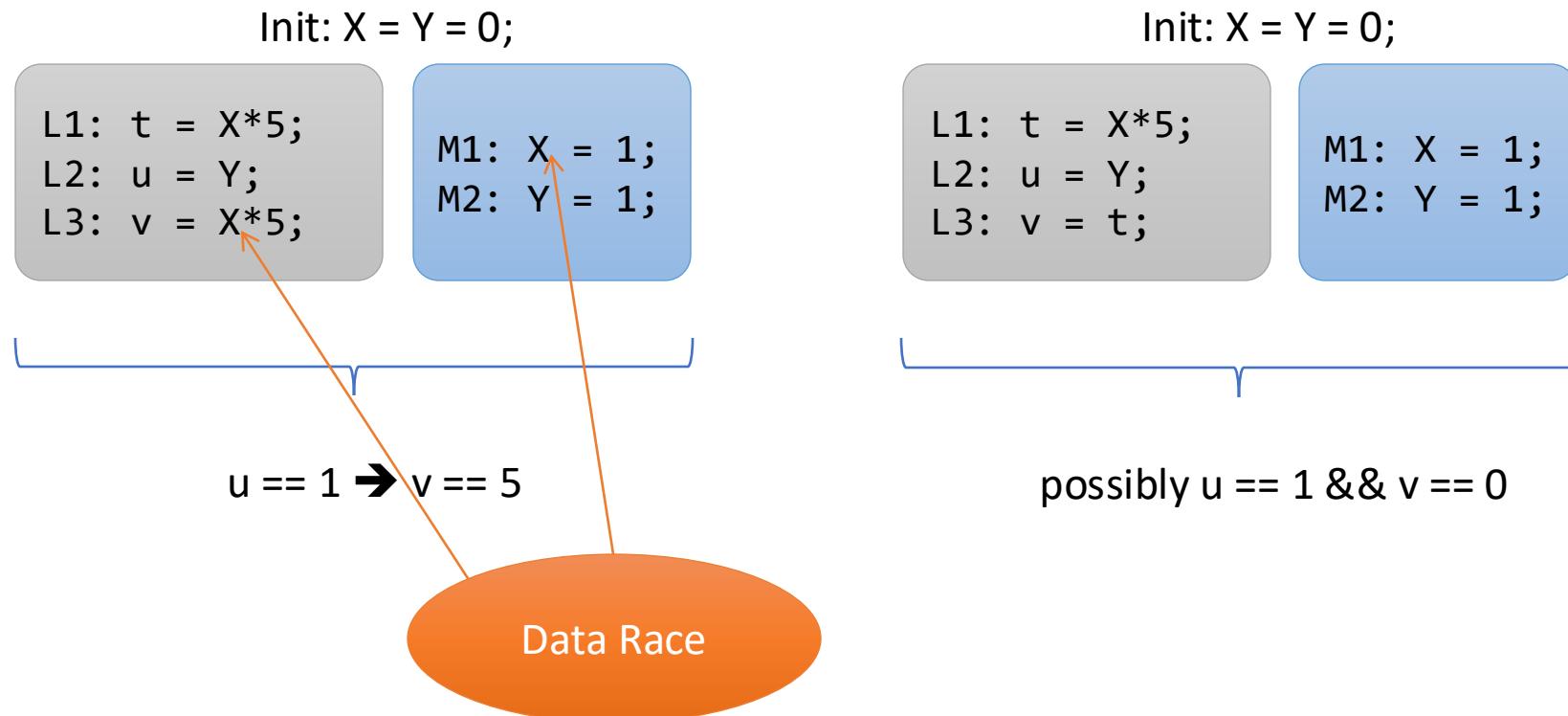


Data Races

- A data race is a pair of conflicting accesses that happen concurrently



Data Races Can Break Sequentially Consistent Semantics



Lockset Algorithm Overview

- Checks a sufficient condition for data-race-freedom
- Consistent locking discipline
 - Every data structure is protected by a single lock
 - All accesses to the data structure made while holding the lock
- Example:

```
// Remove a received packet  
AcquireLock( RecvQueueLk );  
pkt = RecvQueue.RemoveTop();  
ReleaseLock( RecvQueueLk );  
  
... // process pkt  
  
// Insert into processed  
AcquireLock( ProcQueueLk );  
ProcQueue.Insert(pkt);  
ReleaseLock( ProcQueueLk );
```

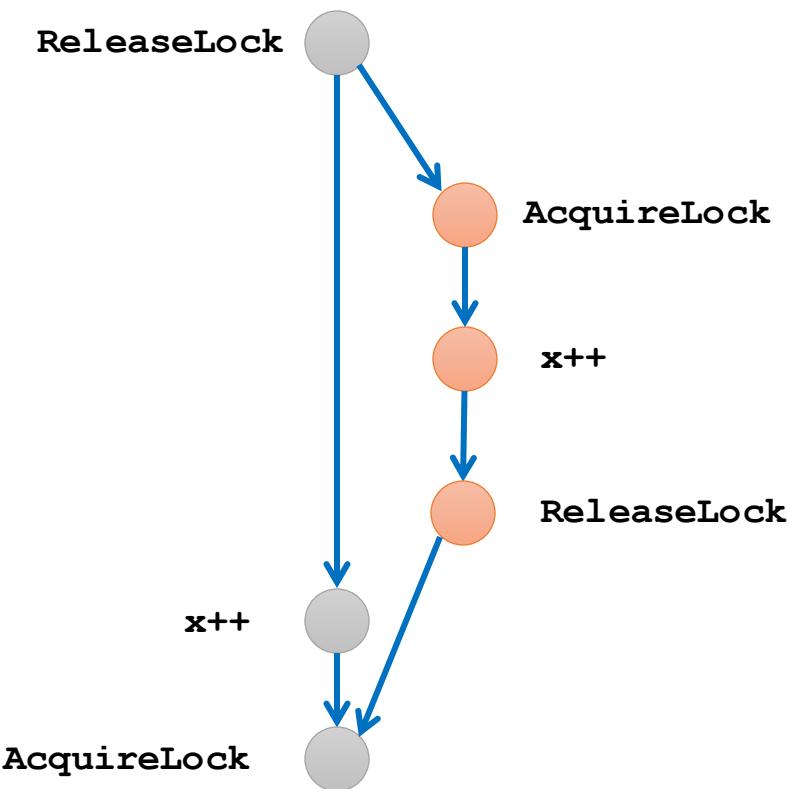
RecvQueue is
consistently protected
by RecvQueueLk

ProcQueue is
consistently protected
by ProcQueueLk

Happens-Before Relation And Data Races

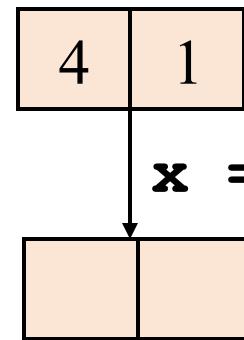
- If all conflicting accesses are ordered by happens-before
- → data-race-free execution
- → All linearizations of partial-order are valid program executions

- If there exists conflicting accesses not ordered
- → a data race



Vector Clocks

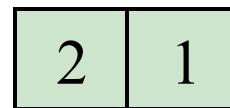
VC_A



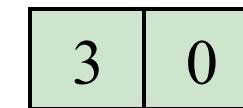
VC_B



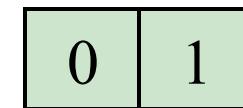
L_m



W_x



R_x



Write-Write Check: $W_x \sqsubseteq \text{VC}_A$?

$$\begin{array}{|c|c|} \hline 3 & 0 \\ \hline \end{array} \sqsubseteq \begin{array}{|c|c|} \hline 4 & 1 \\ \hline \end{array} ? \text{ Yes}$$

Read-Write Check: $R_x \sqsubseteq \text{VC}_A$?

$$\begin{array}{|c|c|} \hline 0 & 1 \\ \hline \end{array} \sqsubseteq \begin{array}{|c|c|} \hline 4 & 1 \\ \hline \end{array} ? \text{ Yes}$$

$O(n)$ time

Guest Lectures



Thread 1:

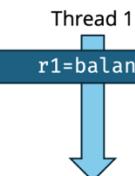
```
🔒1[b].lock()  
r1=AtomicRead(b)  
🔓1[b].unlock();  
if(r1>=amount)  
    🔒1[b].lock()  
    r1=AtomicRead(b)  
    🔓1[b].unlock();
```

Thread 2:

```
🔒2[b].lock()  
r1=AtomicRead(b)  
🔓2[b].unlock();  
if(r1>=amount)  
    🔒2[b].lock()  
    r1=AtomicRead(b)  
    🔓2[b].unlock();
```

Lock Status:

🔒₁[b]: Held by Fray
🔒₂[b]: Held by Fray



How useful is uReview at Uber?

We establish usefulness using a few different metrics:

- **Retrospective Detection Rate:** uReview was able to detect a good fraction of historic bugs that lead to past production incidents
- **Preventable Incident Count:** uReview detected many bugs before they reached production*
- **Developer Satisfaction Rate:** Median developer feedback was very positive*
- **Comment Addressal Rate:** ~65% of the posted comments were resolved

Considerations: AI-based Dynamic Analysis in practice

Now, we can automatically (and easily) do things that required substantial manual effort. **This substantially increases the bug-finding surface!**

Things to consider:

- Why are we able to find more bugs? Is it the program analysis technique? Are the LLMs just really capable?
- How would the next LLM impact the dynamic analysis system you've built?
- How do we properly report these bugs? See Google vs. ffmpeg battle

Static vs Dynamic Analysis

- Over-approximation vs Under-approximation
- When is one better than other? Tradeoffs!
 - Soundness/Completeness
 - Static analysis often “sound” for over-approximate reasoning (e.g. verification)
 - Dynamic Analysis can be “sound” for under-approximate reasoning (e.g. hot spots or bugs).
 - Neither technique is complete in general.
 - Scalability
 - Static analysis often scales super-linearly with program size
 - Dynamic analysis *tries* to scale linearly with execution length
 - Feasibility
 - Static analysis may be impossible with incomplete information (e.g. dynamically loaded code, dependency injection, multi-language code, hardware interaction)
 - Dynamic analysis is only useful when appropriate program inputs are available

Course Evaluations – cmu.smartevals.com



17-355 (Undergrad)

??



17-665 (Masters)

17-819 (PhD)

Next Steps – Course Project

- Project discussions today (Dec 4)
- Recitation tomorrow (Dec 5) reserved for project discussions
- No Office Hours next week
- Project Presentations (Dec 11)
 - In-person 1-4pm at GHC 4215
 - Bring your laptops and adapters, if any
 - 6 min talks + 2min Q&A (= 8 min firm cut-off)
- Project Slides + Report due Dec 11 at midnight on Gradescope