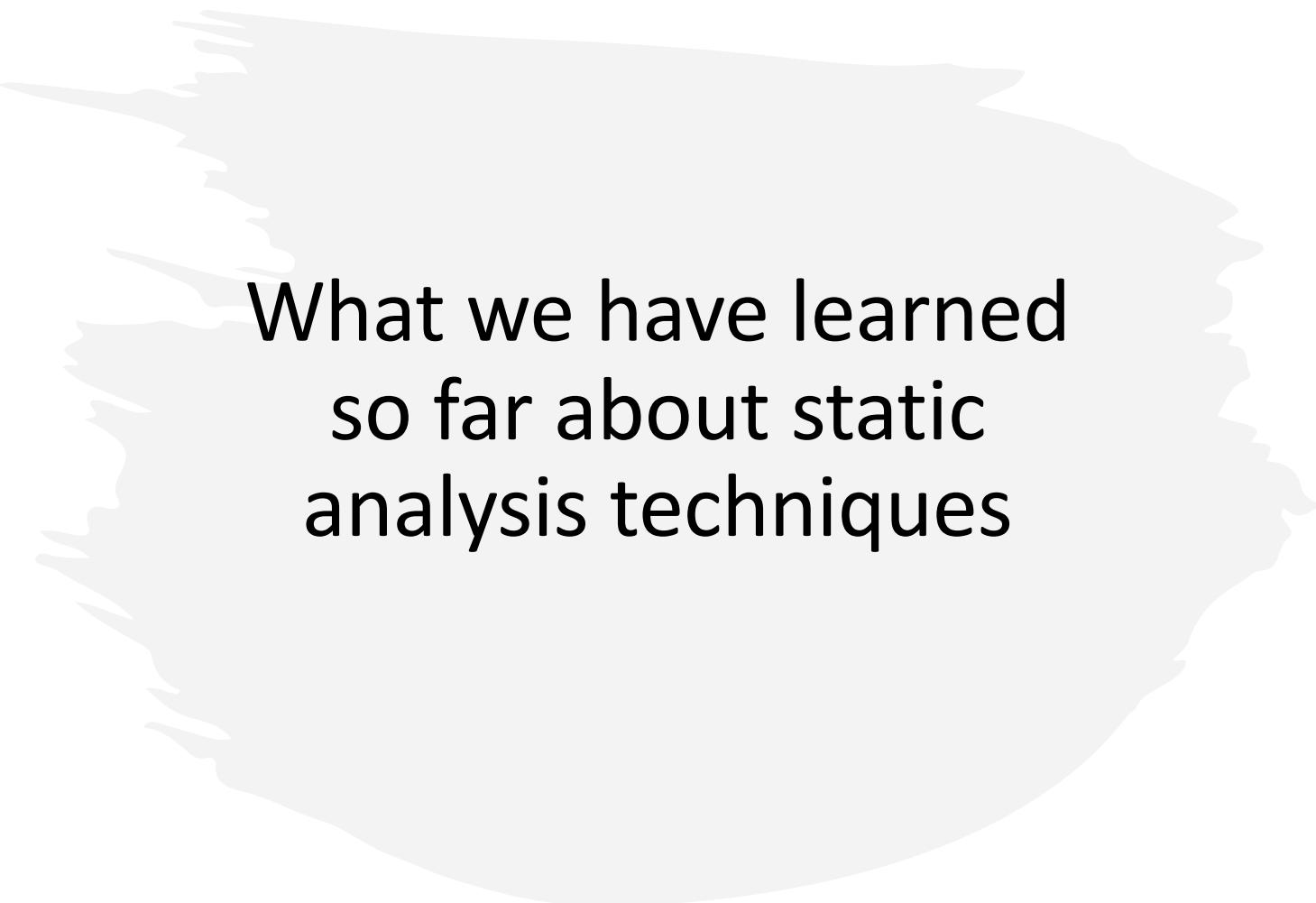


Static Taint Analysis



What we have learned
so far about static
analysis techniques

Recap

Static analysis ≠ perfect precision

Need to choose between:

- **Safety** (conservatively warn of all software vulnerabilities)
- **Certainty** (guarantee properties for compiler optimization)

Recap

Static analysis ≠ perfect precision

Need to choose between:

- **Safety** (conservatively warn of all software vulnerabilities)
- **Certainty** (guarantee properties for compiler optimization)

Reconcile two views

- **Over-approximation**: conservative, safe
- **Under-approximation**: precise

Over-approximation

May Analysis: What Might Be True

- Describes properties that **may hold on some paths**

Over-approximates possible runtime behavior

Key characteristics:

- Powerset Lattice order: \subseteq (least element = \emptyset)
- Flow Equation: use **union (U)**

Examples:

- **Live Variables Analysis**

Under- approximation

Must Analysis: What Must Be True

- Describes properties that **hold on all paths**

Under-approximates – only keeps what's guaranteed

Key characteristics:

- Powerset Lattice order: \supseteq (greatest element = universal set)
- Flow Equation: use **intersection** (\cap)

Examples:

- **Available Expressions**

Comparison

Feature	May Analysis	Must Analysis
Interpretation	May be true	Must be true
Approximation type	Over-approximation	Under-approximation
Powerset lattice order	\subseteq (least = \emptyset)	\supseteq (greatest = universal)
Combine operator	\cup (union)	\cap (intersection)
Typical goal	Vulnerability detection	Optimization

Comparison

Adopt **May** analysis when:

- We care about **possible behaviors** (e.g., potential vulnerabilities)

Adopt **Must** analysis when:

- We need **guaranteed facts** (e.g., safe optimizations)

The main difference lies in:

- **Lattice structure** (\subseteq vs \supseteq)
- **Combining constraints** (\cup vs \cap)

Question:
What structure would you
choose for static taint
analysis?

Static Taint Analysis

Interpretation	??
Approximation	??
Lattice order	??
Combine operator	??
Goal	??

Static Taint Analysis

Interpretation	May
Approximation	Over-approximation
Lattice order	$\text{untaint} \sqsubseteq \text{taint}$
Combine operator	\sqcup Least upper bound
Goal	Taint detection

From Dynamic to Static taint tracking

- Dynamic taint analysis operates at runtime: It tracks taint information as the program executes.
 - It ensures precise detection of vulnerabilities but can incur performance overhead.
 - Incomplete Coverage: Only the executed paths are analyzed.
 - If a piece of code or a dangerous path is never executed during runtime, it is never checked.
 - This leads to false negatives (i.e., vulnerabilities may go undetected).

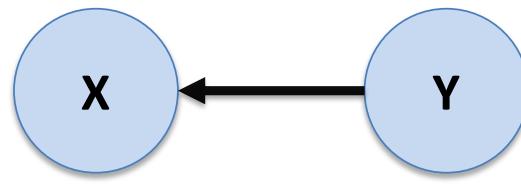
From Dynamic to Static taint tracking

- Static taint analysis ensures a broad coverage, works before execution
 - No runtime costs
 - Other issues ... later

(Static) Intuition: typing flow of assignments

X := Y

The flow from Y to X is legal whenever
the (type of) value of variable Y (τ_Y must be compatible
with the (type of) value of variable x (τ_X))



FLOW

$$\frac{\tau_Y < \tau_X}{X := Y \text{ is ok}}$$

TYPING RULE

(Static) Intuition: typing flow of assignments

taint



untaint

The degree of
taint lattice

(Static) Intuition: typing flow of assignments

taint

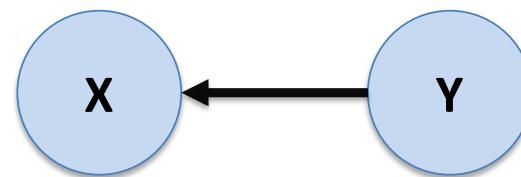


untaint

The degree of
Taint lattice

X := Y

The flow from Y to X is legal whenever
the the **degree of taint** τ_Y of variable Y is less than the
degree of taint τ_X of variable x.



$$\frac{\tau_Y < \tau_X}{X := Y \text{ is ok}}$$

FLOW

TYPING RULE

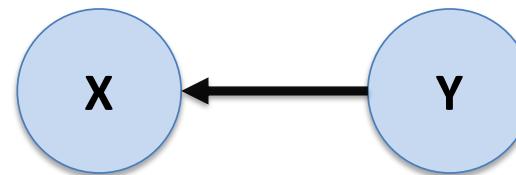
(Static) Intuition: typing flow of assignments

taint



untaint

$$\tau_y = \text{untaint}$$



$$\frac{\tau_Y < \tau_X}{X := Y \text{ is } ok}$$

TYPING RULE

The degree of
Taint lattice

(Static) Intuition: typing flow of assignments

taint

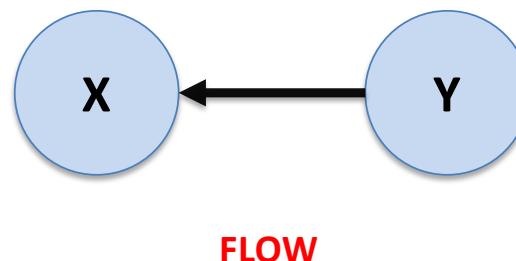


untaint

The degree of
Taint lattice

$$\tau_y = \text{taint}$$

$$\tau_x = \text{untaint}$$



$$\frac{\tau_y < \tau_x = \text{false}}{X := Y \text{ is not ok}}$$

The ingredients of
the analysis



The analysis by examples

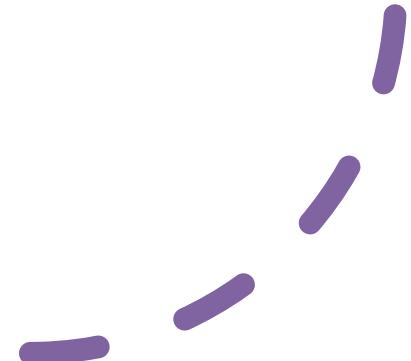
```
string name = source()  
string x = name;  
sink(x)
```

sink (untainted string s) = printf(untainted string s) { ... // trusted sink};

Source() = tainted string getsFromNetwork(...); //untrusted source

Ingredient 1: the analysis

- **May analysis:**
- Lattice order is \subseteq
- Join function is **union (\cup)**
- Taint facts: set of **tainted** variables -- we track which variables are **tainted**



Ingredient 2: the flow equation

We define for each program point

- **GEN(n)**: New taint facts **generated** at statement n
- **KILL(n)**: Taint facts **removed** or overwritten at statement n
- **IN(n)**: Taint facts **before** statement n
- **OUT(n)**: Taint facts **after** statement n
- **OUT(n) = (IN(n) – KILL(n)) ∪ GEN(n)**

1. string name = **source()**
2. string x = name;
3. **sink(x)**

GEN(1) = {name}

KILL(1) = \emptyset

IN(1) = \emptyset

OUT(1) = GEN(1) \cup (IN(1) – KILL(1)) = { name }

1. string name = **source()**
2. string x = name;
3. **sink(x)**

$$IN(1) = \emptyset$$

$$GEN(1) = \{name\}$$

$$KILL(1) = \emptyset$$

$$OUT(1) = GEN(1) \cup (IN(1) - KILL(1)) = \{ name \}$$

$$IN(2) = OUT(1)$$

if name \in OUT(1) then x \in GEN(2)

$$KILL(2) = \emptyset$$

$$OUT(2) = GEN(2) \cup (IN(2) - KILL(2)) = \{ name, x \}$$

1. string name = **source()**
2. string x = name;
3. **sink(x)**

$$IN(1) = \emptyset$$

$$GEN(1) = \{name\}$$

$$KILL(1) = \emptyset$$

$$OUT(1) = GEN(1) \cup (IN(1) - KILL(1)) = \{ name \}$$

$$IN(2) = OUT(1)$$

if name \in OUT(1) then x \in GEN(2)

$$KILL(2) = \emptyset$$

$$OUT(2) = GEN(2) \cup (IN(2) - KILL(2)) = \{ name, x \}$$

$$IN(3) = OUT2) = \{ name, x \}$$

$$GEN(3) = \emptyset$$

$$KILL(2) = \emptyset$$

Warning!!! Since x \in IN(3) a tainted value reaches sink!

Key Issue: Flow Sensitivity in Static Taint Analysis

Flow-dependent (flow-sensitive) taint analysis considers the **order of statements** in a program.

The analysis tracks **how taint facts change over time** as execution flows through the control-flow graph.

GEN and KILL in Flow Sensitivity

GEN(s): What is generated (newly tainted) by basic block s

GEN(s): adds taint to variables.

- $x = \text{source}()$ then $\text{GEN} = \{ x \}$
- $x = E$ and $\exists y \in PVar(E)$ and $y \in IN(s)$ then $y \in GEN(s)$
- because x becomes tainted at this point

KILL(s): What is killed (make untainted) or overwritten by block s

KILL(s) removes taint from variables, either by sanitizing or overwriting.

- $x = 0$ then $\text{KILL}(s) = \{ x \}$
- $x = E$ and $\forall y \in PVar(E) y \notin IN(s)$ then $\text{KILL}(s) = \{x\}$
- the taint degree that was on x is now gone.

The flow equation

- At each program point, the analysis computes
$$OUT(s) = (IN(s) - KILL(s)) \cup GEN(s)$$
 - $IN(s)$: taint info **before** s
 - $OUT(s)$: taint info **after** s
- The flow equation precisely captures **flow sensitivity**:
- If a variable was tainted earlier ($IN(s)$), but gets overwritten or sanitized ($KILL(s)$), it will no longer be tainted in $OUT(s)$.
- If a taint is introduced ($GEN(s)$), it appears in $OUT(s)$ even if it was not previously present.

Why this matters?

Without GEN/KILL, we would have to assume taints never disappear. This may lead to:

- **Overtainting** (false positives)
- Poor precision

Using GEN/KILL sets enables:

- Accurate **taint propagation**
- **Precise flow-sensitive reasoning** about taint evolution

```
1: x = source(); // GEN = { x }
2: x = 0; // KILL = { x }
3: sink(x);
```

$x \in OUT(1)$



$x \notin OUT(2)$

sink(x) is sanitized!!

What about the monotone framework?

The Monotone Framework

- Analysis Lattice: $L = \textcolor{green}{untainted} \leq \textcolor{red}{tainted}$
- Analysis State: $[i] \in PVar \rightarrow L$
- The Join operation: $JOIN(i) = \bigcup [j] j \in pred(i)$

The Monotone Framework

$$\frac{i.a = E}{[i] = JOIN(i)[a = eval(JOIN(i), E)]}$$

$$eval(\sigma, a) = \sigma(a)$$

$$eval(\sigma, CstInt) = \textcolor{teal}{untainted}$$

$$eval(\sigma, input) = \textcolor{red}{tainted}$$

$$eval(\sigma, E1 \ op \ E2) = eval(\sigma, E1) \ \widehat{\op} \ eval(\sigma, E2)$$

Example: $\widehat{+}$ is the abstract additions over the abstract values given by

$$\textcolor{teal}{untainted} \widehat{+} \textcolor{teal}{untainted} = \textcolor{teal}{untainted}$$

$$\textcolor{red}{tainted} \widehat{+} v = \textcolor{red}{tainted}$$

$$v \widehat{+} \textcolor{red}{tainted} = \textcolor{red}{tainted}$$

Example

1: string name = **source()**

2: string x = name;

3: **sink(x)**

[1] = [name = **tainted**]

[2] = JOIN(2)[x = eval(JOIN(2), name)] = [1][x = eval([1], name)]

= [name = **tainted**][eval(x = [name = **tainted**](name))]

= [name = **tainted**][x = **tainted**] = [name = **tainted**, x = **tainted**]

Conditionals

```
1: string name = source();    // taint source
2: string x;
3: if (...)                  // conditional
4:   x = name;                // taint flows here
5: else x = "ciao";           // constant, untainted
6: sink(x);                  // potential leak
```

GOAL: Determine if tainted data may reach sink(x), using GEN/KILL sets

```
1: string name = source() ←
2: string x;
3: If ()
4:   x = name
5: else x = "ciao";
6: sink(x)
```

$$GEN(1) = \{name\}$$

$$KILL(1) = \emptyset$$

$$OUT(1) = GEN(1) \cup (IN(1) - KILL(1)) = \{name\}$$

```
1: string name = source()  
2: string x; ←  
3: If ()  
4:   x = name  
5: else x = "ciao";  
6: sink(x)
```

$$\begin{aligned}GEN(1) &= \{name\} \\KILL(1) &= \emptyset \\OUT(1) &= GEN(1) \cup (IN(1) - KILL(1)) = \{name\}\end{aligned}$$
$$\begin{aligned}GEN(2) &= \emptyset \\KILL(1) &= \emptyset \\OUT(2) &= IN(2) = \{name\}\end{aligned}$$

```
1: string name = source()  
2: string x;  
3: If () ←—————  
4:   x = name  
5: else x = "ciao";  
6: sink(x)
```

$GEN(1) = \{name\}$
 $KILL(1) = \emptyset$
 $OUT(1) = GEN(1) \cup (IN(1) - KILL(1)) = \{name\}$

$GEN(2) = \emptyset$
 $KILL(1) = \emptyset$
 $OUT(2) = IN(2) = \{name\}$

$GEN(3) = \emptyset$
 $KILL(3) = \emptyset$
 $OUT(3) = IN(3) = \{name\}$

```
1: string name = source()  
2: string x;  
3: If ()  
4:   x = name ← blue arrow  
5: else x = "ciao";  
6: sink(x)
```

$$\begin{aligned}GEN(1) &= \{name\} \\KILL(1) &= \emptyset \\OUT(1) &= GEN(1) \cup (IN(1) - KILL(1)) = \{name\}\end{aligned}$$
$$\begin{aligned}GEN(2) &= \emptyset \\KILL(1) &= \emptyset \\OUT(2) &= IN(2) = \{name\}\end{aligned}$$
$$\begin{aligned}GEN(3) &= \emptyset \\KILL(3) &= \emptyset \\OUT(3) &= IN(3) = \{name\}\end{aligned}$$
$$\begin{aligned}IN(4) &= \{name\} \\GEN(4) &= \{x\} \\KILL(4) &= \emptyset \\OUT(4) &= \{x\} \cup IN(4) = \{x, name\}\end{aligned}$$

```
1: string name = source()  
2: string x;  
3: If ()  
4:   x = name  
5: else x = "ciao"; ←  
6: sink(x)
```

$GEN(1) = \{name\}$	$KILL(1) = \emptyset$	$OUT(1) = GEN(1) \cup (IN(1) - KILL(1)) = \{name\}$
$GEN(2) = \emptyset$	$KILL(1) = \emptyset$	$OUT(2) = IN(2) = \{name\}$
$GEN(3) = \emptyset$	$KILL(3) = \emptyset$	$OUT(3) = IN(3) = \{name\}$
$IN(4) = \{name\}$	$GEN(4) = \{x\}$	$IN(5) = \{name\}$
$KILL(4) = \emptyset$		$GEN(5) = \emptyset$
$OUT(4) = \{x\} \cup IN(4) = \{x, name\}$		$KILL(5) = \{x\}$
		$OUT(5) = IN(5) - \{x\} = \{name\}$

```

1: string name = source()
2: string x;
3: If ()
4:   x = name
5: else x = "ciao";
6: sink(x) ←

```

$IN(1) = \{name\}$	$GEN(1) = \{name\}$
	$KILL(1) = \emptyset$
	$OUT(1) = GEN(1) \cup (IN(1) - KILL(1)) = \{name\}$
	$GEN(2) = \emptyset$
	$KILL(2) = \emptyset$
	$OUT(2) = IN(2) = \{name\}$
	$GEN(3) = \emptyset$
	$KILL(3) = \emptyset$
	$OUT(3) = IN(3) = \{name\}$
$IN(4) = \{name\}$	$IN(5) = \{name\}$
$GEN(4) = \{x\}$	$GEN(5) = \emptyset$
$KILL(4) = \emptyset$	$KILL(5) = \{x\}$
$OUT(4) = \{x\} \cup IN(4) = \{x, name\}$	$OUT(5) = IN(5) - \{x\} = \{name\}$
	$IN(6) = OUT(4) \cup OUT(5) = \{name x\}$

Warning!!! Since $x \in IN(6)$ a tainted value reaches sink!

SUMMARY

	IN	GEN	KILL	OUT
1	\emptyset	{ name }	\emptyset	{ name }
2	{ name }	\emptyset	\emptyset	{ name }
3	{ name }	\emptyset	\emptyset	{ name }
4	{ name }	{ x }	\emptyset	{ name, x }
5	{ name }	\emptyset	{ x }	{ name }
6	{ name, x }	\emptyset	\emptyset	{ name, x }

SUMMARY

	IN	GEN	KILL	OUT
1	\emptyset	{ name }	\emptyset	{ name }
2	{ name }	\emptyset	\emptyset	{ name }
3	{ name }	\emptyset	\emptyset	{ name }
4	{ name }	{ x }	\emptyset	{ name, x }
5	{ name }	\emptyset	{ x }	{ name }
6	{ name, x }	\emptyset	\emptyset	{ name, x }

Analysis Result

Variable x **may be tainted** at the sink (line 6)

Static analysis **reports a possible taint flow**

Dropping Conditional

```
1: string name = source();  
2: string x;  
3: if (...)  
4:   x = name;  
5: else x = "ciao";  
6: sink(x);
```



```
1: string name = source();  
2: string x1;  
3: string x2;  
4: x1 = name;  
5: x2 = "ciao";  
6: sink(2);
```

GOAL: Determine if tainted data may reach sink(), using GEN/KILL sets

SUMMARY

	IN	GEN	KILL	OUT
1	\emptyset	{ name }	\emptyset	{ name }
2	{ name }	\emptyset	\emptyset	{ name }
3	{ name }	\emptyset	\emptyset	{ name }
4	{ name }	{ x1 }	\emptyset	{ name, x1 }
5	{ name, x1 }	\emptyset	{ x2 }	{ name, x1 }
6	{ name, x1 }	\emptyset	\emptyset	{ name, x1 }

Analysis Result

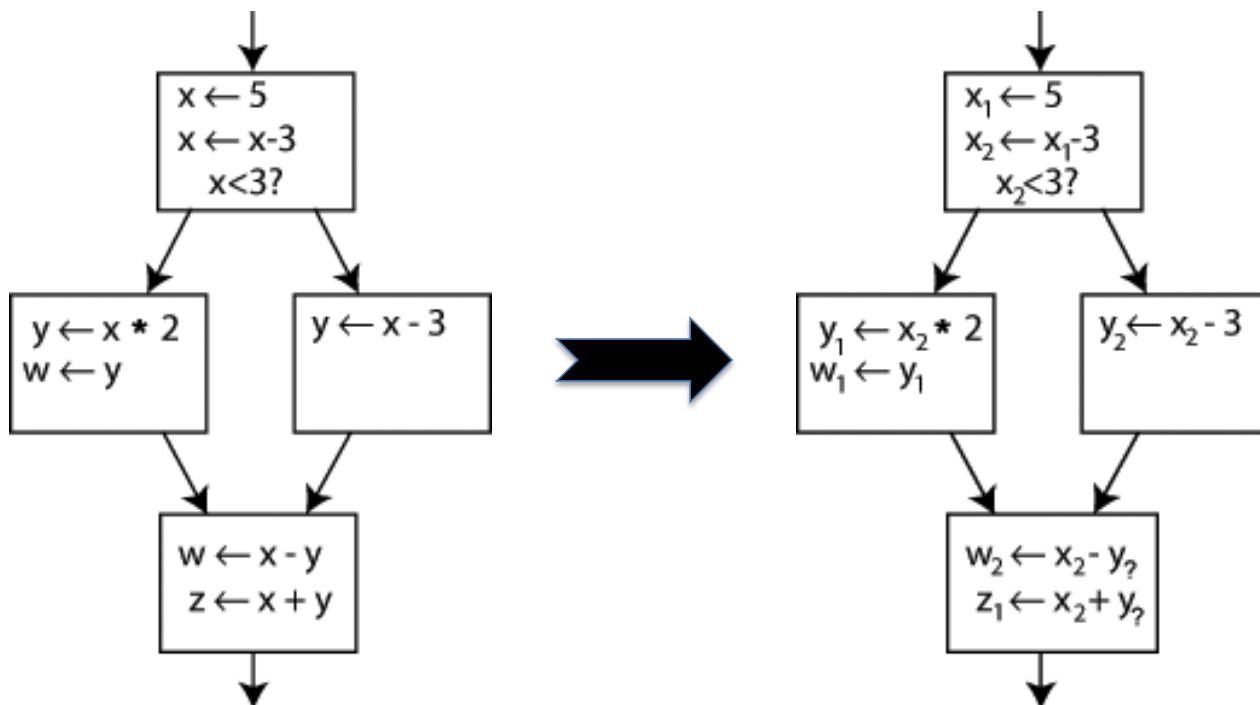
Variable x2 **is not tainted** at the sink (line 6)

Static analysis **does not report a possible taint flow**

Static Single Assignment (SSA)

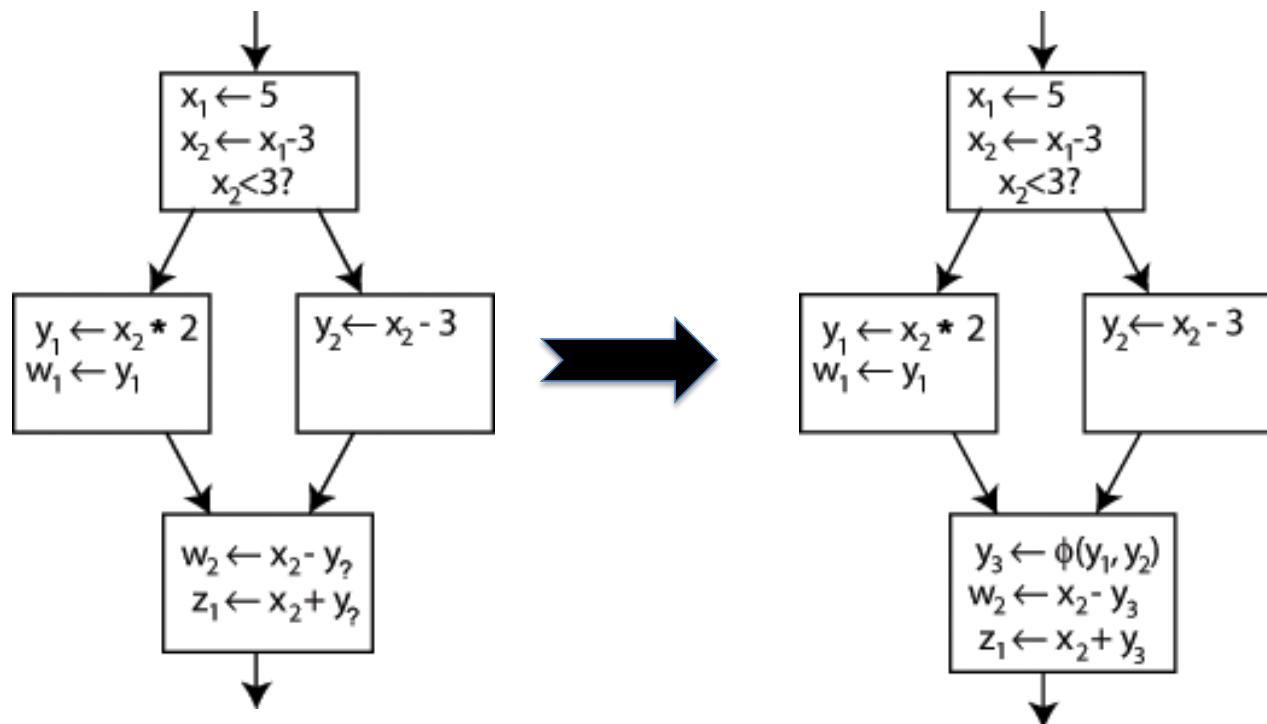
- SSA is a way of structuring the intermediate representation (IR) of programs so that **every variable is assigned exactly once** and **every variable is defined before it is used**
- Intuition: Existing variables in the original IR are split into versions, new variables typically indicated by the original name with a subscript, so that every definition gets its own version.
- This is formally equivalent to continuation-passing style (CPS) translation

SSA Example



y in the bottom block could be referring to either y_1 or y_2 ,

SSA Example



ϕ (*Phi*) function generates a new definition of y called y_3 by "choosing" either y_1 or y_2 , depending on the control flow.

SSA: discussion

- Given an arbitrary control-flow graph, it can be difficult to tell where to insert Φ functions, and for which variables.
 - this general question has an efficient solution that can be computed using a concept called *dominance frontiers*
- A compiler can implement a Φ function by inserting "move" operations at the end of every predecessor block.
 - The compiler might insert a move from y_1 to y_3 at the end of the left block and a move from y_2 to y_3 at the end of the right block.



SSA

- The LLVM Compiler Infrastructure uses SSA form
- The GNU Compiler Collection makes extensive use of SSA.
- Oracle's HotSpot Java Virtual Machine uses an SSA-based intermediate language in its JIT compiler.
- Microsoft Visual C++ compiler (2015 Update) uses SSA

Multiple conditionals

```
int printfun(untainted string) { ... };  
tainted string getsFromNetwork(...);
```

```
void f (int x) {  
    string y;  
     If (x) y = "ciao"  
    else y = getsFromNetwork()  
    if (x) printfun(y);  
}  
y = untainted
```

Multiple conditionals

```
int printfun(untainted string) { ... };  
tainted string getsFromNetwork(...);
```

```
void f (int x) {  
    a string y;  
    If (x) y = "ciao"  
     else y = getsFromNetwork()  
    if (x) printfun(y);  
}  
    y = tainted
```

Multiple conditionals

```
int printfun(untainted string) { ... };  
tainted string getsFromNetwork(...);
```

```
void f (int x) {  
    a string y;  
    If (x) y = "ciao"  
    else y = getsFromNetwork()  
    if (x) printfun(y);  
}  
    y = untainted  
    y = tainted
```



Multiple conditionals

```
int printfun(untainted string) { ... };  
tainted string getsFromNetwork(...);
```

```
void f (int x) {  
    a string y;  
    If (x) y = "ciao"  
    else y = getsFromNetwork()  
    if (x) printfun(y);  
}  
    y = untainted  
    y = tainted
```



No solution

Multiple conditionals

```
int printfun(untainted string) { ... };  
tainted string getsFromNetwork(...);
```

```
void f (int x) {  
    a string y;  
    If (x) y = "ciao"  
    else y = getsFromNetwork()  
    if (x) printfun(y)  
}
```



No solution

False Alarm: because of the conditions on the guard

Path sensitivity

- The problem is that the constraints we generates do not correspond to **feasible paths** (i.e. **feasible executions**)
- **Solution:** *We develop an analysis which considers the feasibility of paths when generating constraints*

TEST

A simple while program

```
string input = source();
char[] buffer = new char[32];
int i = 0;

while (i < input.length() && i < 32) {
    buffer[i] = input.charAt(i);
    i = i + 1;
}

string userInput = new String(buffer);
grantAccess(userInput);
}
```

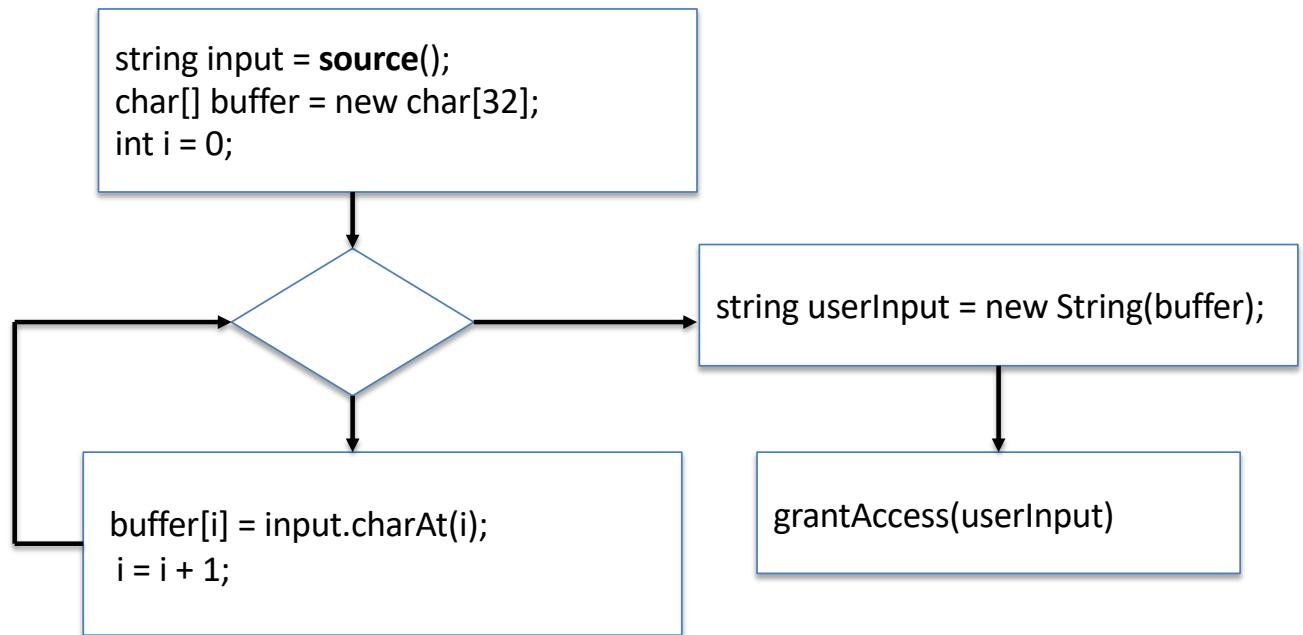
GOAL: Determine if tainted data may reach grantAccess

A simple while program and its CFG

```
string input = source();
char[] buffer = new char[32];
int i = 0;

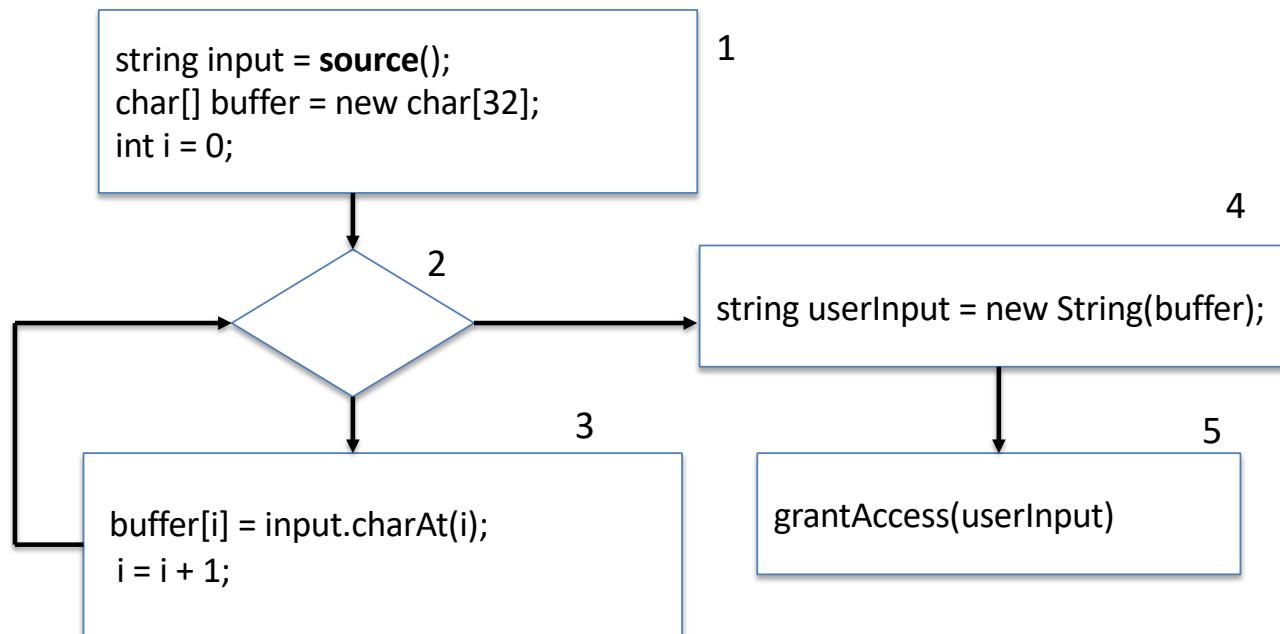
while (i < input.length() && i < 32) {
    buffer[i] = input.charAt(i);
    i = i + 1;
}

string userInput = new String(buffer);
grantAccess(userInput);
}
```

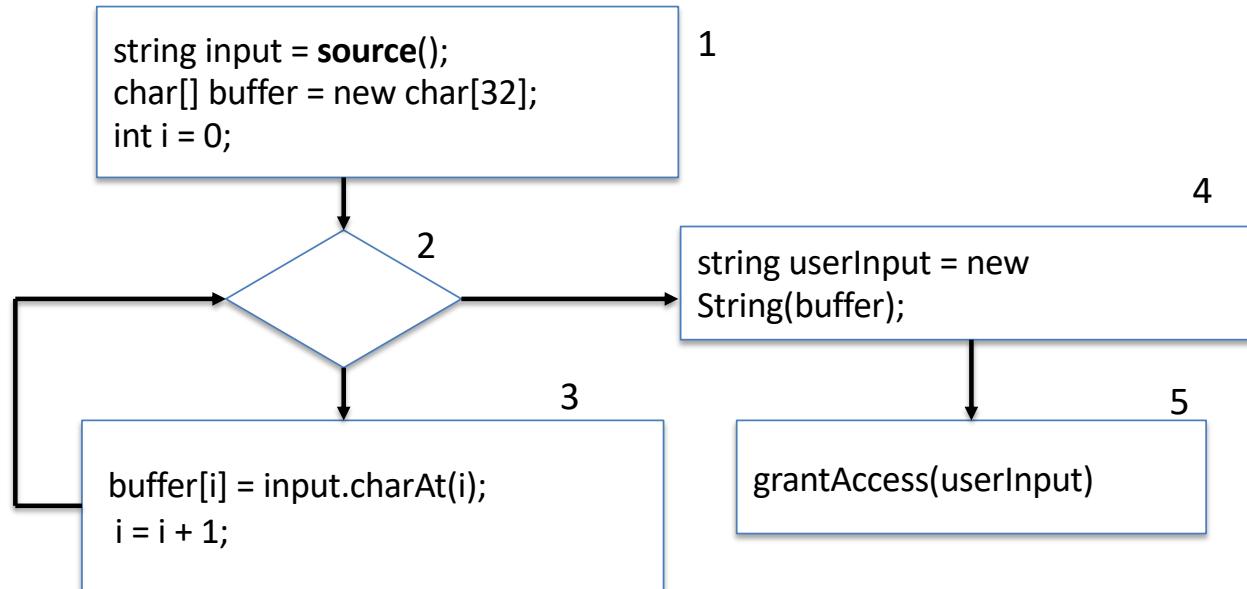


GOAL: Determine if tainted data may reach grantAccess

The CFG



GOAL: Determine if tainted data may reach `grantAccess` at block 5



BASIC BLOCK	GEN SET	KILL SET
1	{input}	\emptyset
2	\emptyset	\emptyset
3	$buffer[i] \in GEN(3)$	\emptyset
4	$\{userInput, buffer\} \subseteq GEN(4)$	\emptyset
5	\emptyset	\emptyset

The Constraints

IN	OUT
IN(1) = ??	OUT(1) =??
IN(2) = ??	OUT(2) = ??
IN(3) = ??	OUT(3) = ??
IN(4) = ??	OUT(4) = ??
IN(5) = ??	OUT(5) = ??

	GEN $\{\text{input}\}$	KILL
1		\emptyset
2	\emptyset	\emptyset
3	buffer[1] if input \subseteq IN(3)	\emptyset
4	userinput if buffer \subseteq IN(4)	\emptyset
5	\emptyset	\emptyset

$$OUT(S) = (IN(S) - KILL(S)) \cup GEN(S)$$

IN(1)	\emptyset	OUT(1)	$\{\text{input}\}$
IN(2)	$OUT(1) \cup OUT(3)$	OUT(2)	IN(2)
IN(3)	OUT(2)	OUT(3)	$IN(3) \cup \{\text{buffer}\}$ if input \subseteq IN(3)
IN(4)	OUT(2)	OUT(4)	$IN(3) \cup \{\text{userinput}\}$ if buffer \subseteq IN(4)
IN(5)	OUT(4)	OUT(5)	IN(5)

IT. 0	IN(S)	OUT(S)	IT. 1
1	\emptyset	\emptyset	1
2	\emptyset	\emptyset	2
3	\emptyset	\emptyset	3
4	\emptyset	\emptyset	4
5	\emptyset	\emptyset	5

IT. 2	IN(S)	OUT(S)	IT. 3	IN(S)	OUT(S)
1	\emptyset	$\{\text{input}\}$	1	\emptyset	$\{\text{input}\}$
2	$\{\text{input}\}$	\emptyset	2	$\{\text{input}\}$	$\{\text{input}\}$
3	\emptyset	\emptyset	3	\emptyset	\emptyset
4	\emptyset	\emptyset	4	\emptyset	\emptyset
5	\emptyset	\emptyset	5	\emptyset	\emptyset

IT. 4

	IN(s)	OUT(s)
1	\emptyset	{input}
2	{input}	{input}
3	{input}	\emptyset
4	{input}	\emptyset
5	\emptyset	\emptyset

IT. 5

	IN(s)	OUT(s)
1	\emptyset	{input}
2	{input}	{input}
3	{input}	{input, buffer}
4	{input}	{input}
5	\emptyset	\emptyset

IT. 6

	IN(s)	OUT(s)
1	\emptyset	{input}
2	{input, buffer}	{input}
3	{input}	{input, buffer}
4	{input}	{input}
5	{input}	\emptyset

IT. 7

	IN(s)	OUT(s)
1	\emptyset	{input}
2	{input, buffer}	{input, buffer}
3	{input}	{input, buffer}
4	{input}	{input}
5	{input}	{input}

IT. 8

	IN(s)	OUT(s)
1	\emptyset	{input}
2	{input, buffer}	{input, buffer}
3	{input, buffer}	{input, buffer}
4	{input, buffer}	{input}
5	{input}	{input}

IT. 9

	IN(s)	OUT(s)
1	\emptyset	{input}
2	{input, buffer}	{input, buffer}
3	{input, buffer}	{input, buffer}
4	{input, buffer}	{input, buffer, userinput}
5	{input, buffer, userinput}	{input}

IT. 9

	IN(s)	OUT(s)
1	\emptyset	{input}
2	{input, buffer}	{input, buffer}
3	{input, buffer}	{input, buffer}
4	{input, buffer}	{input, buffer, userinput}
5	{input, buffer, userinput}	{input}

IT. 10

	IN(s)	OUT(s)
1	\emptyset	{input}
2	{input, buffer}	{input, buffer}
3	{input, buffer}	{input, buffer}
4	{input, buffer}	{input, buffer, userinput}
5	{input, buffer, userinput}	{input, buffer, userinput}

- Abstraction step: if I have a buffer and each of the elements of buffer is known, then buffer is known. This is safe because now buffer, as a whole is given as input. So we say "because all the elements of buffer are known then buffer is known".

The Constraints

IN	OUT
$IN(1) = \emptyset$	$OUT(1) = ??$
$IN(2) = ??$	$OUT(2) = ??$
$IN(3) = ??$	$OUT(3) = ??$
$IN(4) = ??$	$OUT(4) = ??$
$IN(5) = ??$	$OUT(5) = ??$

The Constraints

IN	OUT
$IN(1) = \emptyset$	$OUT(1) = (IN(1) - KILL(1)) \cup GEN(1)$
$IN(2) = OUT(1) \cup OUT(3)$	$OUT(2) = (IN(2) - KILL(2)) \cup GEN(2)$
$IN(3) = OUT(2)$	$OUT(3) = (IN(3) - KILL(3)) \cup GEN(3)$
$IN(4) = OUT(2)$	$OUT(4) = (IN(4) - KILL(4)) \cup GEN(4)$
$IN(5) = OUT(4)$	$OUT(5) = (IN(5) - KILL(5)) \cup GEN(5)$

The Constraints

IN	OUT
$IN(1) = \emptyset$	$OUT(1) = (\emptyset - \emptyset) \cup \{\text{input}\}$
$IN(2) = OUT(1) \cup OUT(3)$	$OUT(2) = (IN(2) - KILL(2)) \cup GEN(2)$
$IN(3) = OUT(2)$	$OUT(3) = (IN(3) - KILL(3)) \cup GEN(3)$
$IN(4) = OUT(2)$	$OUT(4) = (IN(4) - KILL(4)) \cup GEN(4)$
$IN(5) = OUT(4)$	$OUT(5) = (IN(5) - KILL(5)) \cup GEN(5)$

The Constraints

IN	OUT
$IN(1) = \emptyset$	$OUT(1) = (\emptyset - \emptyset) \cup \{input\}$
$IN(2) = \{input\} \cup OUT(3) = \{input\}$	$OUT(2) = \{input\}$
$IN(3) = \{input\}$	$OUT(3) = (IN(3) - KILL(3)) \cup GEN(3)$
$IN(4) = OUT(2)$	$OUT(4) = (IN(4) - KILL(4)) \cup GEN(4)$
$IN(5) = OUT(4)$	$OUT(5) = (IN(5) - KILL(5)) \cup GEN(5)$

The Solution

IN	OUT
IN(1) = \emptyset	OUT(1) = {input}
IN(2) = {input, buffer}	OUT(2) = {input, buffer}
IN(3) = {input, buffer}	OUT(3) = {input, buffer}
IN(4) {input, buffer}	OUT(4) = {input, buffer, userInput}
IN(5) = {input, buffer, userInput}	OUT(5) = {input, buffer, userInput}

The Solution

IN	OUT
IN(1) = \emptyset	OUT(1) = {input}
IN(2) = {input, buffer}	OUT(2) = {input, buffer}
IN(3) = {input, buffer}	OUT(3) = {input, buffer}
IN(4) {input, buffer}	OUT(4) = {input, buffer, userInput}
IN(5) = {input, buffer, userInput}	OUT(5) = {input, buffer, userInput}

The Constraints

IN	OUT
$IN(1) = \emptyset$	$OUT(1) = \{$
$IN(2) = OUT(1) \cup OUT(3)$	$OUT(2) = (IN(2) - KILL(2)) \cup GEN(2)$
$IN(3) = OUT(2)$	$OUT(3) = (IN(3) - KILL(3)) \cup GEN(3)$
$IN(4) = OUT(2)$	$OUT(4) = (IN(4) - KILL(4)) \cup GEN(4)$
$IN(5) = OUT(4)$	$OUT(5) = (IN(5) - KILL(5)) \cup GEN(5)$



WARNING

The algorithm

```
↑ Killn generated by gen  
Out(s) := Gen(s) for all basic block s  
W := {all basic block}  
repeat {  
    take s from W↑ block from worklist  
    In(s) :=  $\cup$  s'  $\in$  pred(s) Out(s')  
    temp := Gen(s)  $\cup$  (In(s) - Kill(s))  
    if (temp != Out(s)) { → Need to update  
        Out(s) := temp  
        W := W  $\cup$  succ(s)}  
} until W = empty → if not, fixed point not reached  
↑ In(s) = Out(s)  
Mr is exactly worklist  
algorithm adapted
```

Path sensitivity

- The problem is that the constraints we generates do not correspond to **feasible paths** (i.e. **feasible executions**)
- **Solution:** *We develop an analysis which considers the feasibility of paths when generating constraints*

Path sensitivity: a running example

```
1: string name = source(); // taint source
2: string x;
3: if (...) {
4:   x = name;
5: } else {
6:   x = "ciao";
7: }
8: sink(x);           // potential taint sink
```

We analyze each path through the program separately:

Path P1 (True branch):

1 → 2 → 3 → 4 → 7 → 8

First flow of execution

Path sensitivity

```
1: string name = source(); // taint source
2: string x;
3: if (...) {
4:   x = name;
5: } else {
6:   x = "ciao";
7: }
8: sink(x);           // potential taint sink
```

We analyze each path through the program separately:

Path P2 (False branch):

1 → 2 → 3 → 5 → 6 → 7 → 8

Second, possible flow of execution

Path sensitivity

```
1: string name = source(); // taint source
2: string x;
3: if (...) {
4:   x = name;
5: } else {
6:   x = "ciao";
7: }
8: sink(x);           // potential taint sink
```

We analyze each path through the program separately:

Path P1 (True branch):

1 → 2 → 3 → 4 → 7 → 8

Path P2 (False branch):

1 → 2 → 3 → 5 → 6 → 7 → 8

We maintain separate **IN/OUT sets per block per path.**

```

1: string name = source(); // taint source
2: string x;
3: if (...) {
4:   x = name;
5: } else {
6:   x = "ciao";
7: }
8: sink(x);           // potential taint sink

```

GEN/KILL SETS

Line	Statement	GEN	KILL
1	name = source();	{ name}	Ø
2	string x;	Ø	Ø
4	x = name;	{ x } if name tainted	Ø
6	x = "ciao";	Ø	{ x }
8	sink(x);	Ø	Ø

Step by step analysis

```
1: string name = source(); // taint source
2: string x;
3: if (...) {
4:   x = name;
5: } else {
6:   x = "ciao";
7: }
8: sink(x);           // potential taint sink
```

Lines 1–2 (both paths)

Line 1:

$$\text{IN}(1) = \emptyset$$

$$\text{OUT}(1) = \text{GEN}(1) \cup (\text{IN}(1) - \text{KILL}(1)) = \{ \text{name} \}$$

Line 2:

$$\text{IN}(2) = \text{OUT}(1) = \{ \text{name} \}$$

$$\text{OUT}(2) = \{ \text{name} \}$$

Step by step analysis

```
1: string name = source(); // taint source
2: string x;
3: if (...) {
4:   x = name;
5: } else {
6:   x = "ciao";
7: }
8: sink(x);           // potential taint sink
```

True branch (path P1)

↗ we do not generate new paths (no C assumption)

$$\text{IN}(4) = \text{OUT}(3) = \text{OUT}(2) = \{ \text{name} \}$$

$$\text{Since } \text{name} \in \text{IN}(4), \text{ then } \text{GEN}(4) = \{ x \}$$

$$\text{KILL}(4) = \emptyset$$

$$\text{OUT}(4) = \{ \text{name}, x \}$$

$$\text{IN}(P1:8) = \text{OUT}(4) = \{ \text{name}, x \}$$

Since $x \in \text{IN}(P1:8)$!!!warning tainted value reaches sink

Path P1: x is tainted at the sink

Step by step analysis

```
1: string name = source(); // taint source
2: string x;
3: if (...) {
4:   x = name;
5: } else {
6:   x = "ciao";
7: }
8: sink(x);           // potential taint sink
```

False branch (path P2)

IN(6) = OUT(3) = OUT(2) = { name }

Since "ciao" is a constant, variable x is overwritten with untainted value

GEN(6) = \emptyset , KILL(6) = { x }

OUT(6) = IN(6) – KILL(6) = { name } (no taint value on variable x)

IN(P2:8) = OUT(6) = { name }

Since $x \notin \text{IN}(P2:8)$ then untainted value reaches sink

Path P2: x is untainted at the sink

SUMMARY of the function
We have analyzed

```
1: string name = source(); // taint source
2: string x;
3: if (...) {
4:   x = name;
5: } else {
6:   x = "ciao";
7: }
8: sink(x);           // potential taint sink
```

Path	vulnerability at sink?	Explanation
Path P1	Yes	Tainted name copied to x
Path P2	No	x overwritten with untainted value

Evaluation

- Extending the analysis with **path sensitivity**, allows us to **distinguish** that **sink(x)** is only unsafe on some paths.
- The (standard) conservative (safe) analyzer would still **report a potential taint violation** at **sink(x)**.
- But a **path-sensitive analyzer** could:
 - **Mark the violation with path conditions** (e.g., *only if true branch is taken*)
 - **Use this precision to reduce false positives.**

Analysis is definitely more precise than what we did before. But what is the strategy you have to define to perform this analysis? You need to keep track of feasibility and conditions for reaching a certain path. We need a precise strategy.

Path sensitivity

A **path-sensitive analysis** tracks data flow facts depending on the path taken

- Path often represented by which branches of conditionals taken

A **path-sensitive analysis** can reason more accurately about correlated conditionals

The issue: how can we make a path sensitive analysis?

- a dataflow analysis where we track facts for each possible path



Handling Path Sensitivity

Approach: The analysis exploits **path-sensitive predicates** in order to associate taint facts with **symbolic conditions** representing the execution paths.

Incorporate symbolic conditions that represent the path you are taking to the taint facts that we are gathering

Symbolic Condition: Path predicates

- Use **path predicates** to:
 - Distinguish between different execution paths
 - Track **conditional taint facts** of variables
 - Evaluate whether **taint values** may reach the **sink** under specific conditions

Path condition: conditions to be satisfied in order to reach that point.

Question now is: is sink reached by tainted value? Yes/No depending on the conditions for path traversal so we will say "sink is reached by tainted value if..."

Conditional taint facts

- We associate each taint fact with a **path predicate** φ (a symbolic condition).
- A path predicate φ represents **under what condition** a taint fact holds.
- We perform a **forward symbolic propagation** of pairs (φ, x)
meaning "variable x is tainted under condition φ "

Taint facts: set of tainted variables. To take into account paths we define φ , path predicate representing conditions under which a certain fact about knowability of a variable holds.

Our running example

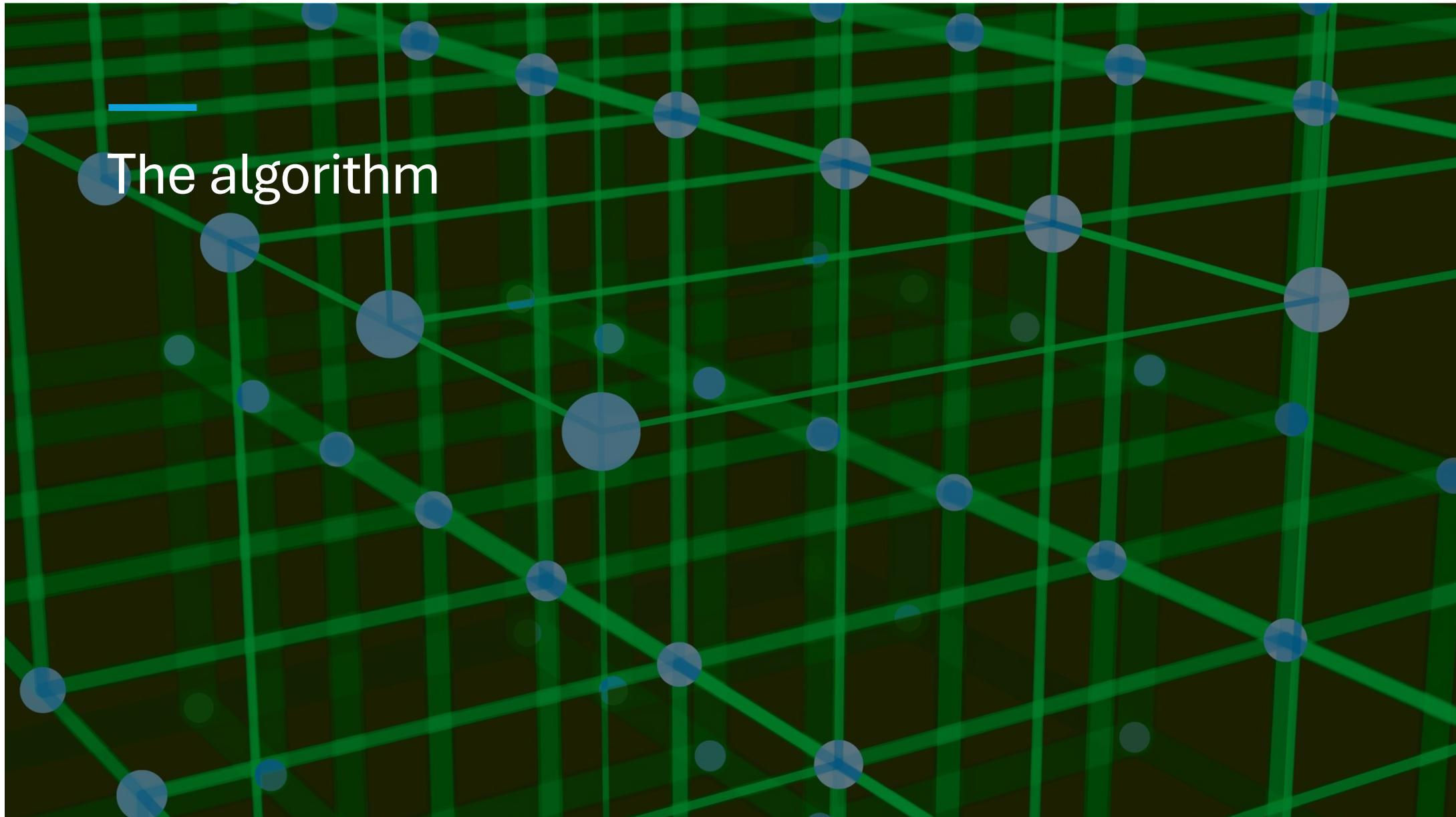
```
1: string name = source(); // taint source
2: string x;
3: if (...) {
4:   x = name;
5: } else {
6:   x = "ciao";
7: }
8: sink(x);           // potential taint sink
```

Path Conditions:

φ_1 condition for then branch

φ_2 condition for else branch

We know in this case $\varphi_2 = !\varphi_1$.



The algorithm

Design ideas

- The algorithm keeps **separate taint facts for each path** by attaching them to **path predicates**
- This allows us to express properties of the form:
 - Variable x is tainted if φ holds
 - Variable x is untainted if φ holds
- At a sink, the algorithm checks:
 - Does any taint fact of the form (φ, x) hold?
- If yes then a taint value may reach sink under that path condition

We look out if taint values reach sinks under certain paths

The (ideal) data structures

TaintFacts: a set of pairs (φ, x) , where:

φ is a Boolean expression (path condition)

x is a program variable set of program variables

IR is a group of different analyses for each path.

PathCond

symbolic conditions accumulated during analysis traversal

- I have 2 data structures: 1 collecting Taint facts, a set of pairs (path condition, program variable)
 X is tainted if φ is true.
- Path conditions: conditions to be true to follow a path (that taints x)
conditions that identify paths to which TaintFacts refer to

Some auxiliary routines

ROUTINE: GEN/KILL TOGETHER Takes known facts, statement and path condition φ

gen_kill(s, φ , facts) = <**gen(s, φ , facts)**, **kill(s, φ , facts)**>:

Computes which taint facts are added or removed at node s, under condition φ

→ Two modes of CFG

[φ is given because it is needed to check facts]

edgeCondition(n1, n2):

Returns the condition required to take the conditional branch
(e.g., φ_1, φ_2)

Restrict(facts, φ): Takes a set of taint facts and a path condition, and remove from the set, all
Keeps only facts that are consistent with the current path predicate

The taint facts that do not satisfy φ .

Used to establish if taints we have are compatible
with our current mode reachable with path condition φ .

GO, GO, GO!



$x = \text{source}()$

(true, $\{x\}$)

PATH 1, 2, (3, 4)

if $x > 0$ then

$x = 5$

$(x > 0 \wedge \text{true}, \emptyset) = (x > 0, \emptyset)$ PATH 1

else

$x = x + 1$

$(x \leq 0, \{x\})$

↗ unfeasible path PATH 2

if $x > 0$ then

$y = 0$

TWO PATHS

$(x > 0 \wedge x > 0, \emptyset)$ $\underbrace{(x > 0 \wedge x \leq 0, \{x\})}_{\text{False}}$ PATH 1, 3)

else

$y = x$

$\underbrace{(x \leq 0 \wedge x > 0, \emptyset)}_{\text{false} = \text{unfeasible path}}$

$(x \leq 0 \wedge x \leq 0, \{x, y\})$ PATH 2, 4)

$\text{sum } K(y)$

$T = \{ \text{PATH 1 } (x > 0, \emptyset), \text{PATH 2 } (x \leq 0, \{x, y\}) \}$ RESTRICT($T, x > 0$) = $\{(x > 0, \emptyset)\}$

T contains two facts associated to paths.

We stick with one path and bring with us the conditions to follow that path and memorize wanted variables. We do this for all paths. That's it.

The (pseudo) code

Input:

Program P with control-flow graph CFG G

Output: Taint facts (with path conditions) at each program point

Initialize:

For each node n in G:

TaintIn[n] = \emptyset

TaintOut[n] = \emptyset

Worklist = [(entry_node, true)] // pair of node and path predicate

↓ we start with entry node, and path condition is true because of course we follow always that path.

The code

while Worklist is not empty:

(n, φ) = Worklist.pop() *We take a node*

// Step 1: Compute TaintIn[n] under condition φ

in_facts = \bigcup [Restrict(TaintOut[p], $\varphi \wedge \text{edgeCondition}(p,n)$)]

for p in pred(n)]

① *②*

*↳ We take info about incoming taint facts, but
we filter with current path conditions +
conditions to get to n.*

// Step 2: Apply GEN/KILL for this node under φ

gen_set, kill_set = gen_kill(n, φ , in_facts)
out_facts = (in_facts - kill_set) \cup gen_set

*Note: Union of predecessors: merging of info we have with different paths
Look at example*

The code

```
while Worklist is not empty:  
    (n,φ) = Worklist.pop()  
  
    // Step 1: Compute TaintIn[n] under condition φ  
  
    in_facts = U [ Restrict(TaintOut[p], φ ∧ edgeCondition(p,n))  
                  for p in pred(n) ]  
  
    // Step 2: Apply GEN/KILL for this node under φ   
  
    gen_set, kill_set = gen_kill(n, φ, in_facts)  
    out_facts = (in_facts - kill_set) ∪ gen_set
```

JUST DON'T THINK TOO
MUCH ABOUT IT

The code

// Step 3: Update and propagate

```
if out_facts ≠ TaintOut[n]: → if out has been updated, we update  
TaintOut[n] = TaintOut[n] ∪ out_facts (put the ones not already there)  
for s in succ(n): Take successors from the working list  
 $\varphi' = \varphi \wedge \text{edgeCondition}(n, s)$  → and take conditions to get there  
Worklist.push((s,  $\varphi'$ ))
```

Return TaintOut[n] for each node n

Supported by fact that we do
Union here!

We consider a specific path
and for different paths
we should add different
elements in Worklist?

You should
look at one path at
a time since you have one path condition

Our running example

```
1: string name = source(); // taint source
2: string x;
3: if (...) {
4:   x = name;
5: } else {
6:   x = "ciao";
7: }
8: sink(x);           // potential taint sink
```

Path Conditions:

φ_1 condition for then branch

φ_2 condition for else branch

```
TaintIn[8] = {
  ( $\varphi_1$ , {name, x}),    // x is tainted in true branch
  ( $\varphi_2$ , {name})      // x is safe in false branch
}
```

An example

```
1: string name = source();      // taint source
2: string x;
3: string log = "";
4: if (name.length() > 0) {      //  $\varphi_1$ 
5:   x = name;
6: } else {                      //  $\varphi_2 = \neg\varphi_1$ 
7:   x = "default";
8: }
9: int i = 0;
10: while (i < x.length()) {     //  $\varphi_3$ . (loop condition)
11:   log = log + x.charAt(i);   // potential use of tainted x
12:   i = i + 1;
13: }
14: sink(log);                // sink
```



Step-by-Step Analysis with Path Predicates



Meta-notation

The notation

taint(x)

is used to indicate that variable x belongs to the set of tainted variables

The notation

!taint(x)

is used to indicate that variable x does not belong to the set of tainted variables (i.e. it has to be removed from that set)

Line 1: name = source();

Always executed (path predicate always true)

(true, taint(name))

```
1: string name = source();      // taint source
2: string x;
3: string log = "";
4: if (name.length() > 0) {      //  $\varphi_1$ 
5:   x = name;
6: } else {                      //  $\varphi_2 = \neg\varphi_1$ 
7:   x = "default";
8: }
9: int i = 0;
10: while (i < x.length()) {     //  $\varphi_3$ . (loop condition)
11:   log = log + x.charAt(i);   // potential use of tainted x
12:   i = i + 1;
13: }
14: sink(log);                // sink
```

Line 1: name = source();

*Always executed (path predicate always true)
(true, taint(name))*

Line 4–8: Conditional Assignment to x

True Branch ($\varphi_1 = \text{name.length()} > 0$)

Line 5: x = name

Taint from name propagates to x

$(\varphi_1, \text{taint}(x))$

False Branch ($\varphi_2 = \neg\varphi_1$)

Line 7: x = "default"

Overwrites x with untainted value !taint(x)

Removes variable x from taint facts

```
1: string name = source();      // taint source
2: string x;
3: string log = "";
4: if (name.length() > 0) {      //  $\varphi_1$ 
5:   x = name;
6: } else {                      //  $\varphi_2 = \neg\varphi_1$ 
7:   x = "default";
8: }
9: int i = 0;
10: while (i < x.length()) {     //  $\varphi_3$ . (loop condition)
11:   log = log + x.charAt(i);  // potential use of tainted x
12:   i = i + 1;
13: }
14: sink(log);                // sink
```

Line 1: name = source();

*Always executed (path predicate always true)
(true, taint(name))*

Line 4–8: Conditional Assignment to x

True Branch ($\varphi_1 = \text{name.length()} > 0$)

Line 5: x = name

Taint from name propagates to x

(φ_1 , taint(x))

False Branch ($\varphi_2 = \neg\varphi_1$)

Line 7: x = "default"

Overwrites x with untainted value **!taint(x)**

Removes variable x from taint facts

Lines 9–13: Loop Copies x into log

We treat the loop as symbolic propagation. We reason symbolically

Assume:

- On every iteration, one character from x is appended to log
- We abstract the loop effect as: "log becomes tainted if x is tainted"

```
1: string name = source();      // taint source
2: string x;
3: string log = "";
4: if (name.length() > 0) {     //  $\varphi_1$ 
5:   x = name;
6: } else {                     //  $\varphi_2 = \neg\varphi_1$ 
7:   x = "default";
8: }
9: int i = 0;
10: while (i < x.length()) {    //  $\varphi_3$ . (loop condition)
11:   log = log + x.charAt(i);  // potential use of tainted x
12:   i = i + 1;
13: }
14: sink(log);                // sink
```

At each iteration we have 1 character. Assumption:
log becomes tainted as soon as we encounter 1 char
tainted.

Line 1: name = source();

*Always executed (path predicate always true)
(true, taint(name))*

Line 4–8: Conditional Assignment to x

True Branch ($\varphi_1 = \text{name.length()} > 0$)

Line 5: x = name

Taint from name propagates to x

$(\varphi_1, \text{taint}(x))$

False Branch ($\varphi_2 = \neg \varphi_1$)

Line 7: x = "default"

Overwrites x with untainted value $!\text{taint}(x)$

Removes variable x from taint facts

Lines 9–13: Loop Copies x into log

From φ_1 branch:

$(\varphi_1 \wedge \varphi_3, \text{taint}(x)) \Rightarrow (\varphi_1 \wedge \varphi_3, \text{taint}(\log))$

From φ_2 branch:

$(\varphi_2 \wedge \varphi_3, !\text{taint}(x)) \Rightarrow !\text{taint}(\log))$

```
1: string name = source();      // taint source
2: string x;
3: string log = "";
4: if (name.length() > 0) {     //  $\varphi_1$ 
5:   x = name;
6: } else {                     //  $\varphi_2 = \neg \varphi_1$ 
7:   x = "default";
8: }
9: int i = 0;
10: while (i < x.length()) {    //  $\varphi_3$ . (loop condition)
11:   log = log + x.charAt(i);  // potential use of tainted x
12:   i = i + 1;
13: }
14: sink(log);                // sink
```

$$\varphi_2 = \neg \varphi_1$$

$$\varphi_3 = i < x.\text{length}()$$

$$T = \{(\text{true}, \{\text{name}\})\}$$

$$T = \{(\varphi_1, \{x, \text{name}\}), (\varphi_2, \{\text{name}\})\}$$

Line 1: name = source();

Always executed (path predicate always true)
(true, taint(name))

Line 4–8: Conditional Assignment to x

True Branch ($\varphi_1 = \text{name.length()} > 0$)

Line 5: x = name

Taint from name propagates to x

($\varphi_1, \text{taint}(x)$)

False Branch ($\varphi_2 = \neg \varphi_1$)

Line 7: x = "default"

Overwrites x with untainted value **!taint(x)**

Removes variable x from taint facts

Lines 9–13: Loop Copies x into log

From φ_1 branch:

$(\varphi_1 \wedge \varphi_3, \text{taint}(x)) \Rightarrow (\varphi_1 \wedge \varphi_3, \text{taint}(\log))$

From φ_2 branch:

$(\varphi_2 \wedge \varphi_3, !\text{taint}(x)) \Rightarrow (\varphi_2 \wedge \varphi_3, !\text{taint}(\log))$

Line 14: sink(log);

From φ_1 : **($\varphi_1 \wedge \varphi_3, \text{taint}(\log)$) sink receive tainted data**

From φ_2 : **($\varphi_2 \wedge \varphi_3, !\text{taint}(\log)$) sink is safe**

1: string name = **source()**; // taint source

2: string x;

3: string log = "";

4: if (name.length() > 0) { // φ_1

5: x = name;

6: } else { // $\varphi_2 = \neg \varphi_1$

7: x = "default";

8: }

9: int i = 0;

10: while (i < x.length()) { // φ_3 . (loop condition)

11: log = log + x.charAt(i); // potential use of tainted x

12: i = i + 1;

13: }

14: **sink(log);** // sink

$$\varphi_2 = \neg \varphi_1$$

$$\varphi_3 = i < x.\text{length}()$$

You can also consider $\varphi_4 = \neg \varphi_3$,
so while may be executed.

Line	Path Predicate	GEN	KILL
1	true	{ name }	\emptyset
2–3	true	\emptyset	\emptyset
5	$\varphi_1 = (\text{name.length()} > 0)$	{ x }	\emptyset
7	$\varphi_2 = \neg\varphi_1$	\emptyset	{ x }
10–13	φ_3 (loop condition)	if $x \in \text{IN} \Rightarrow \{ \log \}$	\emptyset
14	—	\emptyset	\emptyset

SUMMARY

Path Condition	Line 14	Explanation
$\varphi_1 \wedge \varphi_3$	Tainted	$x = \text{name}$, loop propagates
$\neg\varphi_1 \wedge \varphi_3$	Untainted	$x = \text{"default"}$, log remains clean



Remark

The **path-sensitive analysis** shows that `sink(log)` receives **tainted data only if**:

- `name.length() > 0` (φ_1)
- and the loop executes (φ_3)

On other paths, the sink is safe

The analysis

- Avoids **false positives**
- Allows one to generate warnings with path conditions,

SUMMARY (AGAIN)

Path Condition	Line 14	Explanation
$\varphi_1 \wedge \varphi_3$	Tainted	$x = \text{name}$, loop propagates
$\neg\varphi_1 \wedge \varphi_3$	Untainted	$x = \text{"default"}$, log remains clean

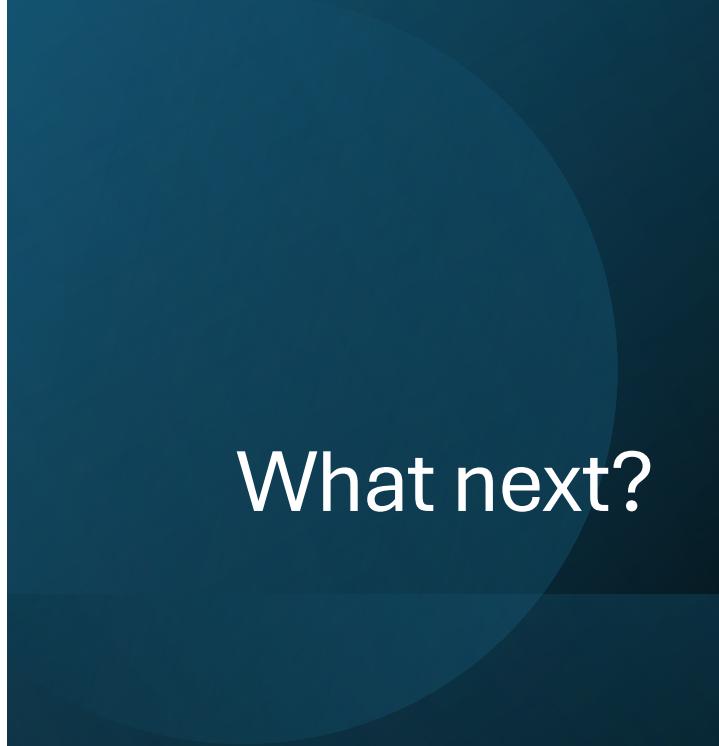
Symbolic Taint Condition

The analysis generates a **symbolic condition** under which taint reaches the sink:

$$\text{Taint at sink} \Leftrightarrow \varphi_T = \varphi_1 \wedge \varphi_3 \text{ holds}$$

$$\varphi_T = (\text{name.length(}) > 0) \wedge (x.\text{length(}) > 0)$$

A boolean φ_T that corresponds to some conditions related to path traversal.
if φ_T holds,
we have taint reaching sink.
You have to provide solution
you need routine solving constraint



What next?

Exploit Static taint analysis with **path predicates to determine the symbolic taint predicate**

Use a suitable solver to **determine if any path condition leading to taint(log) is satisfiable**

... but what is a suitable solver?



SMT Solver

SMT Solver: Theorem prover. SW which solves logical formulas. Their job is: give an input as a formula and solver tries to give a result.

- **SMT (Satisfiability Modulo Theories)** solvers are software toolkits that takes **logical formulas** as input and determines whether these formulas are:
 - **SAT** (satisfiable)
 - at least one assignment makes the formula true (for boolean variables)
 - **UNSAT** (unsatisfiable)
 - no assignment satisfies the formula
 - **UNKNOWN**
 - can't determine (e.g., due to undecidable theory or timeout)

SAT: Boolean Satisfiability Problem

What is the SAT Problem?

Input: A Boolean formula

Question: Is there an assignment of true/false values to variables that makes the formula true?

The problem proven NP-complete (Cook, 1971)

$$(x \vee \neg y) \wedge (\neg x \vee z)$$

SAT if $x = \text{true}$, $y = \text{false}$, $z = \text{true}$ Satisfied w/ ---

- No polynomial solution up to now

SAT is hard to solve!!

- The search space is exponential: 2^n possible assignments for n variables
- SAT is NP-complete
 - No known polynomial-time algorithm for all inputs
- Many real-world problems reduce to SAT:
 - Circuit verification
 - Program analysis static analysis
 - Scheduling, planning
- Even small changes to input can change the solution space drastically

Techniques to address SAT

How Modern SAT Solvers Behave

- DPLL algorithm and its modern variants
 - Heuristics for variable selection
 - Preprocessing and simplification To make several cases more tractable
 - Solvers like MiniSAT, Z3, and CryptoMiniSat make SAT feasible in many cases
- Improvements done, modern solvers effiwer and work well.

However ... DLLP first



The **DPLL algorithm** (Davis–Putnam–Logemann–Loveland) is the foundation of most modern **SAT solvers**.



DPLL is a **complete backtracking-based search algorithm** for solving the **SAT problem**. It determines whether a **Boolean formula in CNF (conjunctive normal form)** is **satisfiable**.

you try to solve, if you find conflicts you go back



DLLP improves brute-force search by introducing **smart pruning and inference techniques**.

** Main operators are AND and OR (produces same)*

In CNF, a formula is made up of **ANDs of ORs**. Think of it like this: you've got a bunch of clauses, and each clause is a group of statements connected by ORs, and then all the clauses are tied together with ANDs.

For example:

$$(A \vee \neg B) \wedge (B \vee C \vee \neg D) \wedge (\neg A \vee D)$$

How DLLP Works

Unit Propagation

- If a clause has **only one literal unassigned** (a *unit clause*), assign it in a way that satisfies the clause.
 - Clause: (x) the set $x=true$

Pure Literal Elimination (optional)

- If a literal appears with **only one polarity** (always positive or always negative) in all clauses, assign it to satisfy all such clauses.
 - If x appears only as x (never as $\neg x$), set $x = \text{true}$.

Variable Assignment (Decision Step)

- Pick an unassigned variable and try assigning true or false.

Backtracking

- If the current assignment leads to a **conflict** (unsatisfiable clause), **backtrack** and try the other assignment.
- If both fail, backtrack further.

1. Unit Propagation

This is like saying: "Hey, there's a clause that has only one variable left unassigned. So we've got no choice—assign it in a way that satisfies the clause."

Example:

- Clause: (x)
- Only literal is x , so to make this clause true, we **must** set $x = \text{true}$.

It's a kind of forced move. You're not making a decision here—you're following the only path that avoids contradiction.

2. Pure Literal Elimination (optional step)

Now you look at all the literals. If a variable only shows up **positively** (like x, x, x) or only **negatively** ($\neg x, \neg x, \neg x$), you can safely assign it in a way that satisfies every clause it's in.

Example:

- If x appears only as x (never $\neg x$), set $x = \text{true}$ and all those clauses become true.
- If x appears only as $\neg x$, set $x = \text{false}$.

This step simplifies things without risk, but not all SAT solvers do it.

3. Variable Assignment (Decision Step)

Okay, now you make a real choice. You pick an unassigned variable and just try setting it to true or false. You're kind of guessing, but you'll use logic and backtracking to correct any wrong turns later.

4. Backtracking

If your guess (plus the consequences it triggered) ends up creating a conflict—like a clause that's now impossible to satisfy—you backtrack.

- First, try the **opposite** value of the variable you chose.
- If that also leads to failure, you **backtrack even further**, undoing earlier decisions, until you find a path that works—or until you've tried everything and proven it's unsatisfiable.

We can do step 1 (setting unassigned variable to true to make clause true) and step 2 (assigning all x to true to make all clauses true) because clauses are just ORs of variables, so one true makes all true, correct?

Exactly, Giovanni—nailed it.

Clauses in CNF are **ORs of literals**, so as long as **at least one** literal in the clause is true, the whole clause is satisfied. That's the key idea behind both **unit propagation** and **pure literal elimination**.

The CNF formula

Let's use this one (in CNF form, of course):

$$(\neg A \vee B) \wedge (\neg B \vee C) \wedge (\neg C \vee D) \wedge (\neg D) \wedge (A \vee \neg B \vee C)$$

This one has 5 clauses and variables A, B, C, D.

Let's break it down clause by clause for clarity:

1. $(\neg A \vee B)$
2. $(\neg B \vee C)$
3. $(\neg C \vee D)$
4. $(\neg D)$
5. $(A \vee \neg B \vee C)$

Step 1: Unit Propagation

Check if any clause is a unit clause (only one literal):

Clause 4 is: $(\neg D) \rightarrow$ That's a unit clause \rightarrow We must assign $D = \text{false}$.

Now that $D = \text{false}$, we propagate it.

Clause 3: $(\neg C \vee D)$
 \rightarrow becomes $(\neg C \vee \text{false}) \rightarrow$ simplifies to $(\neg C)$
 Now $(\neg C)$ is a unit clause \rightarrow so $C = \text{false}$

Now that $C = \text{false}$, propagate it:

Clause 2: $(\neg B \vee C)$
 \rightarrow becomes $(\neg B \vee \text{false}) \rightarrow$ simplifies to $(\neg B)$
 \rightarrow Unit clause $\rightarrow B = \text{false}$

Now propagate $B = \text{false}$:

Clause 1: $(\neg A \vee B)$
 \rightarrow becomes $(\neg A \vee \text{false}) \rightarrow$ simplifies to $(\neg A)$
 \rightarrow Unit clause $\rightarrow A = \text{false}$

Now all variables are assigned:

- $A = \text{false}$
- $B = \text{false}$
- $C = \text{false}$
- $D = \text{false}$

Let's now check all clauses:

1. $(\neg A \vee B) \rightarrow (\text{true} \vee \text{false}) \rightarrow \text{true}$
2. $(\neg B \vee C) \rightarrow (\text{true} \vee \text{false}) \rightarrow \text{true}$
3. $(\neg C \vee D) \rightarrow (\text{true} \vee \text{false}) \rightarrow \text{true}$
4. $(\neg D) \rightarrow (\text{true}) \rightarrow \text{true}$
5. $(A \vee \neg B \vee C) \rightarrow (\text{false} \vee \text{true} \vee \text{false}) \rightarrow \text{true}$

All clauses are satisfied!

We didn't even need pure literal elimination or decision steps or backtracking—**unit propagation alone** solved the whole thing for us. That's kinda cool, actually.

Want to try another one that does need a decision and backtracking?

Can A, B etc be other CNFs or are they usually variables?

Nice question—so, Giovanni, in **CNF, A, B, C, etc. are typically variables** (or more precisely, **literals**, meaning either a variable or its negation, like A or $\neg A$).

Each clause in CNF is made of **ORs of literals**, and then the full formula is the **AND** of those clauses. That's the standard setup for SAT solvers and most logical systems using CNF.

No, they cannot

Z3 SAT SOLVER



- **Z3** is an **SMT** solver developed by Microsoft Research
 - <https://github.com/Z3Prover/z3>
- **Z3 Input:** A set of **declarative constraints**, often expressed in logic over various domains:
 - Integers, reals, booleans, bitvectors, arrays, strings, etc.
Can take a lot of different constraints in terms of
- Converts the problem into a combination of **Boolean SAT solving + theory solvers** (for strings, integers, arrays, etc.)
- Uses **efficient heuristics and decision procedures** to prune infeasible choices

Z3 in Action

Suppose you want to check if a tainted value might reach a sink **only if a certain condition is true**.

Taint Fact = $(x > 0 \&\& x < 10) \text{ AND } (x * 2 == 15)$ { $x:\text{int}$ }

Z3>

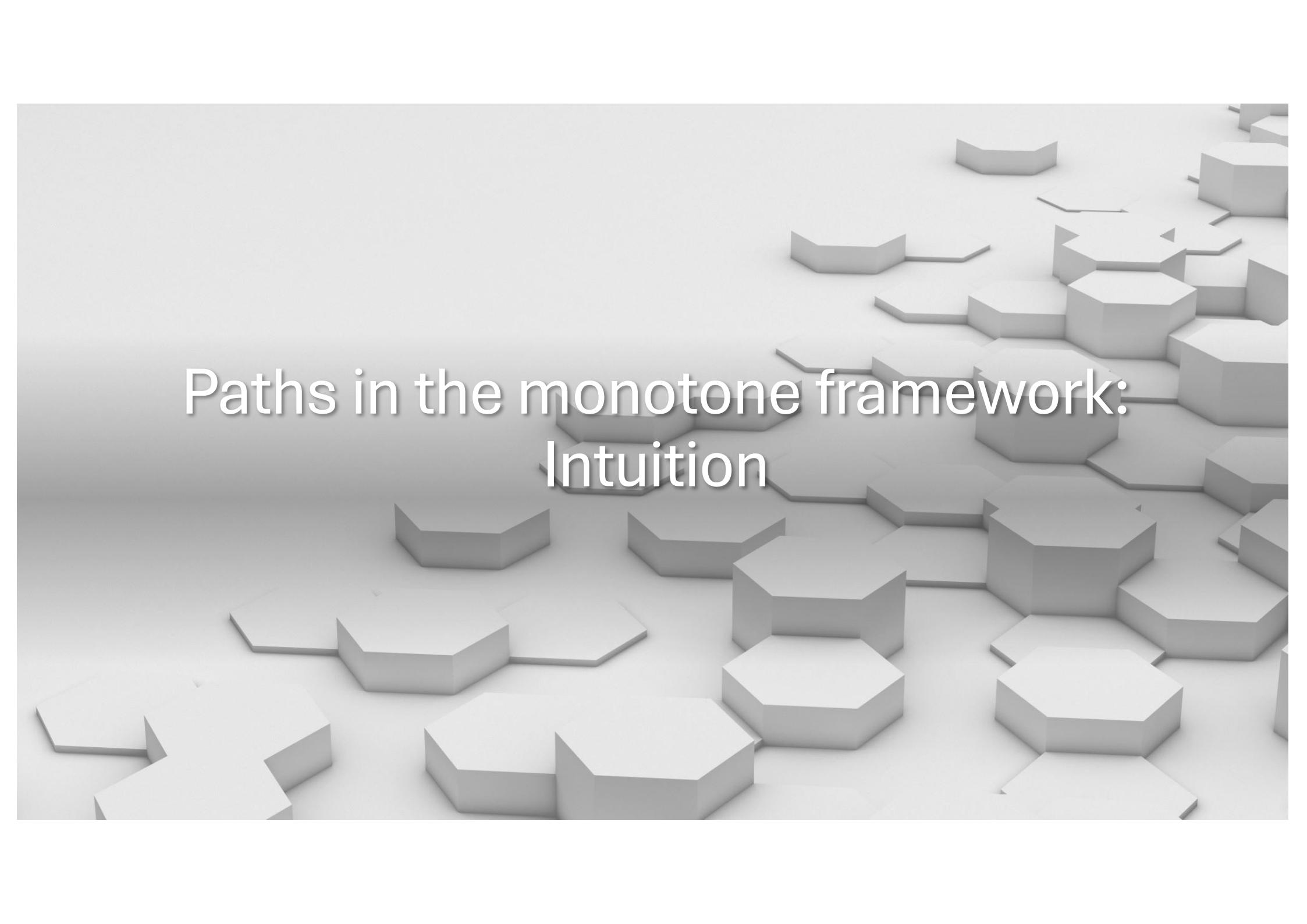
```
(declare-const x Int)
  (assert (> x 0))
  (assert (< x 10))
  (assert (= (* x 2) 15))
(check-sat)
```

Z3 >

unsat *(S expressions)*

The sink is unreachable under any input (unfeasible path)

The taint does not reach the sink



Paths in the monotone framework: Intuition

The steps

We outline **how abstract states evolve** across **branches** in a **path-sensitive taint analysis**: **path predicates** and transfer functions are exploited to trace how taint information propagates and merges across control flow branches.

Three Steps:

- 1. The Lattice Model Setup:** Formal definitions
- 2. Branch Evolution:** State splits, transfer, and joins
- 3. Worked Example:** Formal trace of a branching program

The Model

Variables: $PVar = \{x, y, z, \dots\}$ is the set of program variables.

Abstract Taint Domain

A taint state $t = PVar \rightarrow \{untaint, taint\}$ (where $untaint \leq taint$)

Path Condition Each abstract state is annotated with a **path predicate** φ (a Boolean formula)

Path-sensitive state $s = (\varphi, t)$

The overall abstract state $S \in AS$ is a **set of such path states**:

$$S = \{s_i = (\varphi_1, t_i \text{ for } i \in I)\}$$

ABSTRACT STATES

$AS = PATH \rightarrow PVar \rightarrow \{untaint \leq taint\}$

TRANSFER FUNCTION

For each command (aka basic block) **cmd** we introduce the transfer function

$$f_{cmd} : AS \rightarrow AS$$

$$f_{cmd}(S) = \bigcup \{f_{cmd}((\varphi, t) \mid (\varphi, t) \in S)\}$$

Conditional

```
if (cond) {  
    cmd1;  
} else {  
    cmd2;  
}
```

$\forall (\varphi, t) \in S$ where S is the input abstract state

Splitting paths

Then Branch: $(\varphi \wedge cond, t)$

Else Branch: $(\varphi \wedge \neg cond, t)$



Conditional Merge Point

$S' = \{(\varphi \wedge cond, t_1), (\varphi \wedge \neg cond, t_2)\}$

Transfer Functions

$$f_{cmd1}(\varphi \wedge cond, t) = (\varphi \wedge cond, t_1)$$

$$f_{cmd2}(\varphi \wedge \neg cond, t) = (\varphi \wedge \neg cond, t_2)$$

Conditional

```
if (cond) {  
    cmd1;  
} else {  
    cmd2;  
}
```

$\forall (\varphi, t) \in S$ where S is the input abstract state

Splitting paths

Then Branch: $(\varphi \wedge cond, t)$

Else Branch: $(\varphi \wedge \neg cond, t)$



Conditional Merge Point

$S' = \{(\varphi \wedge cond, t_1), (\varphi \wedge \neg cond, t_2)\}$

Transfer Functions

$$f_{cmd1}(\varphi \wedge cond, t) = (\varphi \wedge cond, t_1)$$

$$f_{cmd2}(\varphi \wedge \neg cond, t) = (\varphi \wedge \neg cond, t_2)$$

We do not join
 t_1, t_2
preserving path sensitivity

```
1: x = input();           // Tainted source
2: if (x > 0) {
3:     y = sanitize(x); // y = untaint
4: } else {
5:     y = x;           // y = taint
6: }
7: print(y);            // Sink
```

```
1: x = input();           // Tainted source
2: if (x > 0) {
3:     y = sanitize(x); // y = untaint
4: } else {
5:     y = x;           // y = taint
6: }
7: print(y);            // Sink
```

$S_7 = \{$
 $(x > 0, \{x = \text{taint}, y = \text{untaint}\}),$
 $(\neg(x > 0), \{x = \text{taint}, y = \text{taint}\})$
 $\}$

The analysis issues a **conditional warning** only on the path where $x \leq 0$.

Line	Path Predicate	Taint Status
1	true	{x = taint}
2	x > 0	{x = taint}
3	x > 0	{ x = taint, y = untaint}
4	$\neg(x > 0)$	{x = taint}
5	$\neg(x > 0)$	{ x = taint, y = taint}
6	—	Path merge
7	Both predicates	{x = taint, y = untaint or taint}

What's next?

An analysis that models only a single function at a time is **intra-procedural**

- An analysis that takes multiple functions into account is **inter-procedural**
- An inter-procedural analysis that considers the calling context when analyzing the target of a function call is **context-sensitive**