# Foundations of Software Engineering

Taint Analysis

Miguel Velez

isr | institute for SOFTWARE RESEARCH

# Learning goals

- Define taint analysis.
- Compare the dynamic and static approaches, as well as their benefits and limitations.
- Apply the analysis to several examples
- Understand how dynamic and static analyses can be combined to overcome the limitations of each other.

institute for
SOFTWARE
RESEARCH

# DYNAMIC ANALYSIS

# Dynamic Analysis

- Learn about program's properties by executing it.
- Examine program state throughout/ after execution by gathering additional information.

# Performance Analysis

How would you learn about method execution time?

```
1. void main(a) {
2.    if(a > 0) {
3.       sleep_ms(a);
4.    else {
5.       sleep_ms(1000);
6.    }
7. }
```

institute for
SOFTWARE
RESEARCH

```
1. void main(a) {
2.    start("main");
3.    if(a > 0) {
4.        sleep_ms(a);
5.    else {
6.        sleep_ms(1000);
7.    }
8.    end("main");
9. }
```

# Benefits

# Benefits

- Analyzes the state of the program in a runtime environment.
- If the property we are looking for is found, we can be sure that it exists.
- Validate static analysis findings.

# Limitations

# Limitations

- Input dependent
- Cannot explore all paths
- Cost of tracking information
- Heisenbuggy behavior

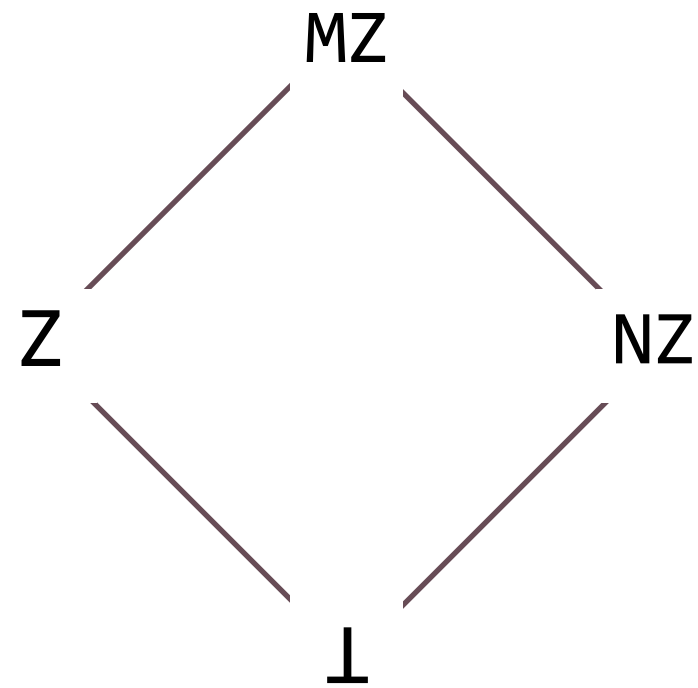# STATIC ANALYSIS

# Static Analysis

- Learn about program's properties without executing it.
- Systematic examination of an abstraction of a program

institute for
SOFTWARE
RESEARCH

# Zero Analysis

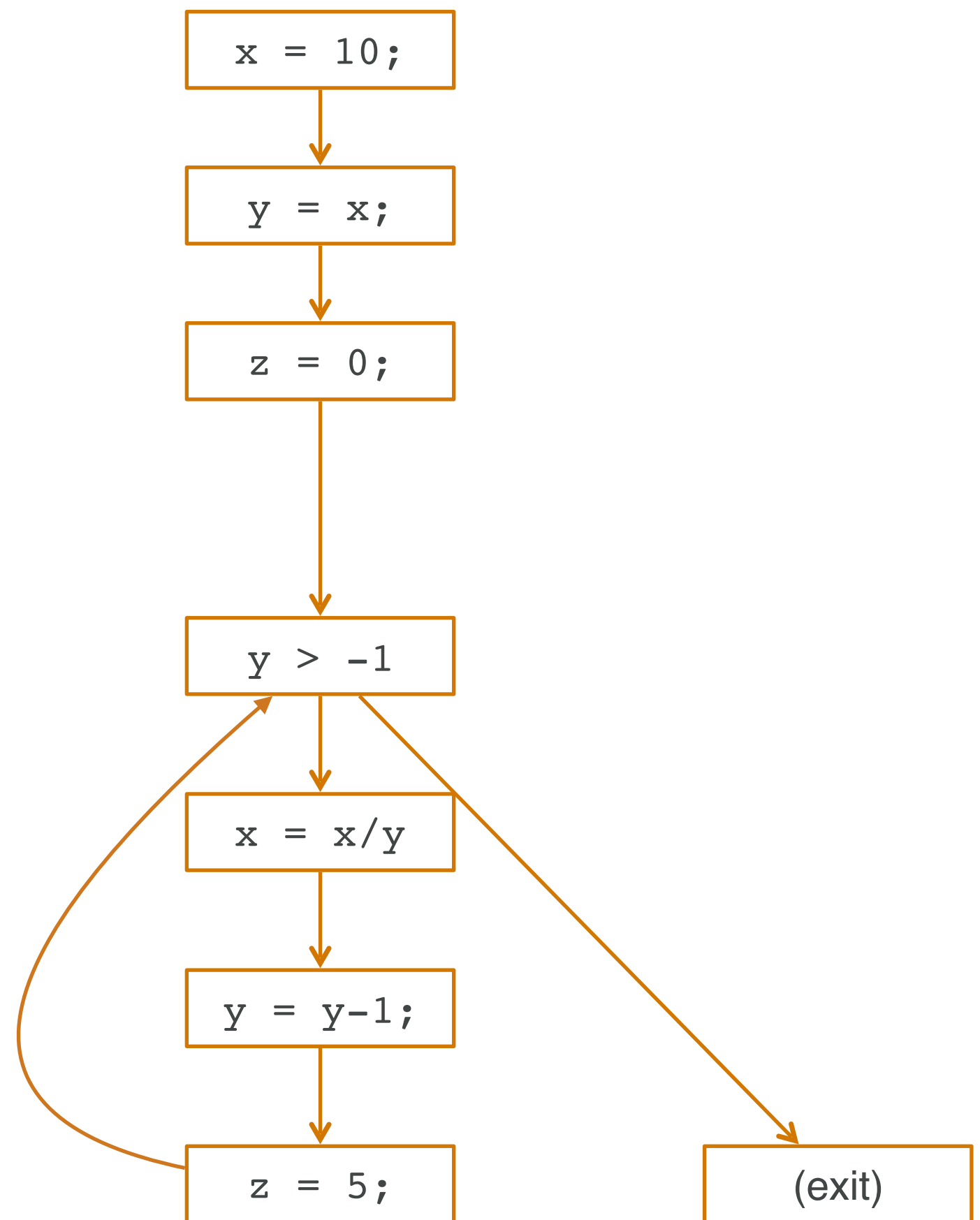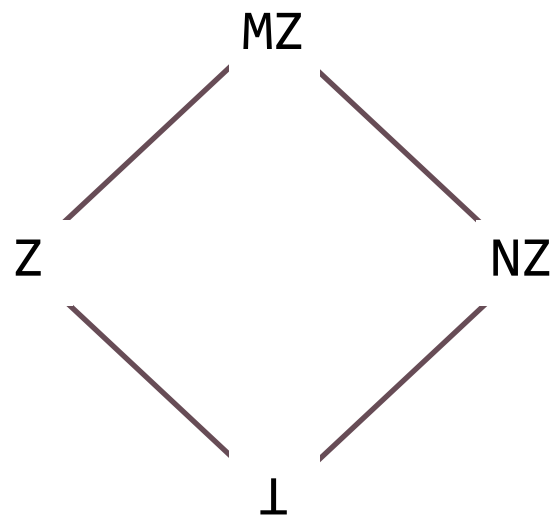How would you learn if you divide by 0?

```
1. x = 10;
2. y = x;
3. z = 0;
4. while(y > -1) {
5.    x = x/y;
6.    y = y-1;
7.    Z = 5;
8. }
```

```
1. x = 10;
2. y = x;
3. z = 0;
4. while(y > -1) {
5.    x = x/y;
6.    y = y-1;
7.    z = 5;
8. }
```

MZ

Z             NZ
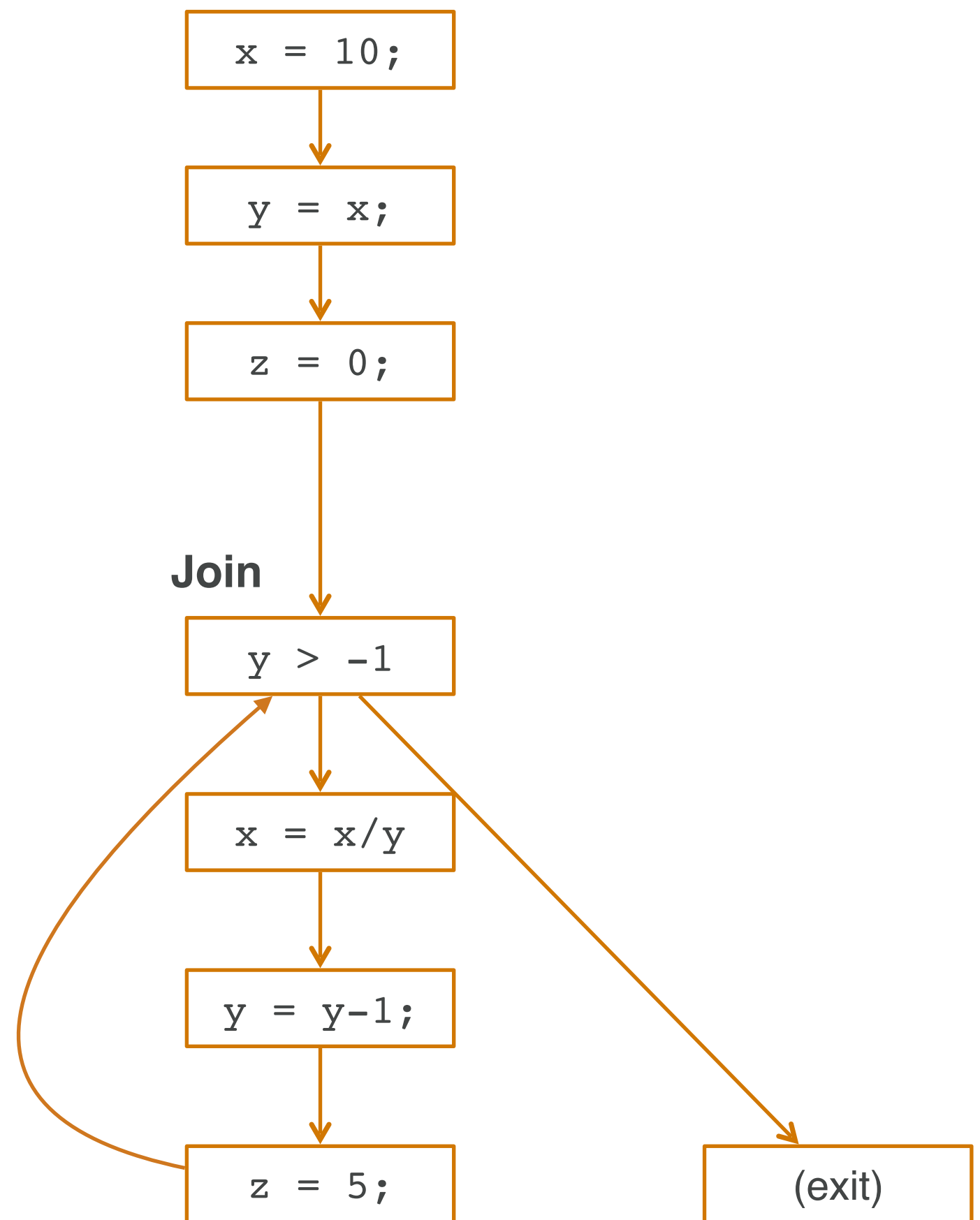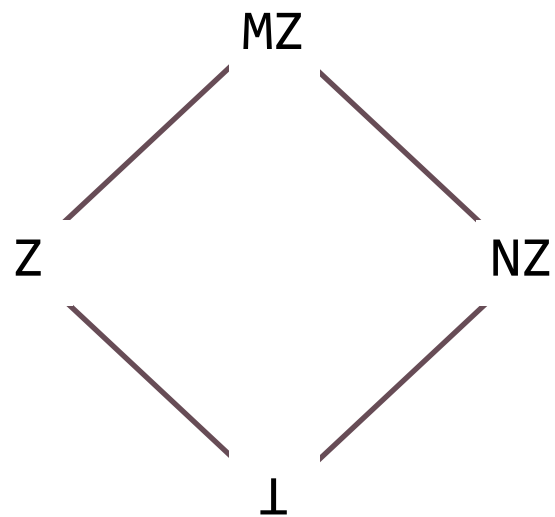
⊥

institute for
SOFTWARE
RESEARCH

```
1. x = 10;
2. y = x;
3. z = 0;
4. while(y > -1) {
5.    x = x/y;
6.    y = y-1;
7.    Z = 5;
8. }
```

MZ

Z          NZ

⊥

x = 10;

y = x;

z = 0;

y > -1

x = x/y

y = y-1;

z = 5;

(exit)

institute for
SOFTWARE
RESEARCH

```
1. x = 10;
2. y = x;
3. z = 0;
4. while(y > -1) {
5.     x = x/y;
6.     y = y-1;
7.     Z = 5;
8. }
```

MZ

Z          NZ

⊥

```
x = 10;
```

```
y = x;
```

```
z = 0;
```

**Join**

```
y > -1
```

```
x = x/y
```

```
y = y-1;
```

```
z = 5;
```

(exit)
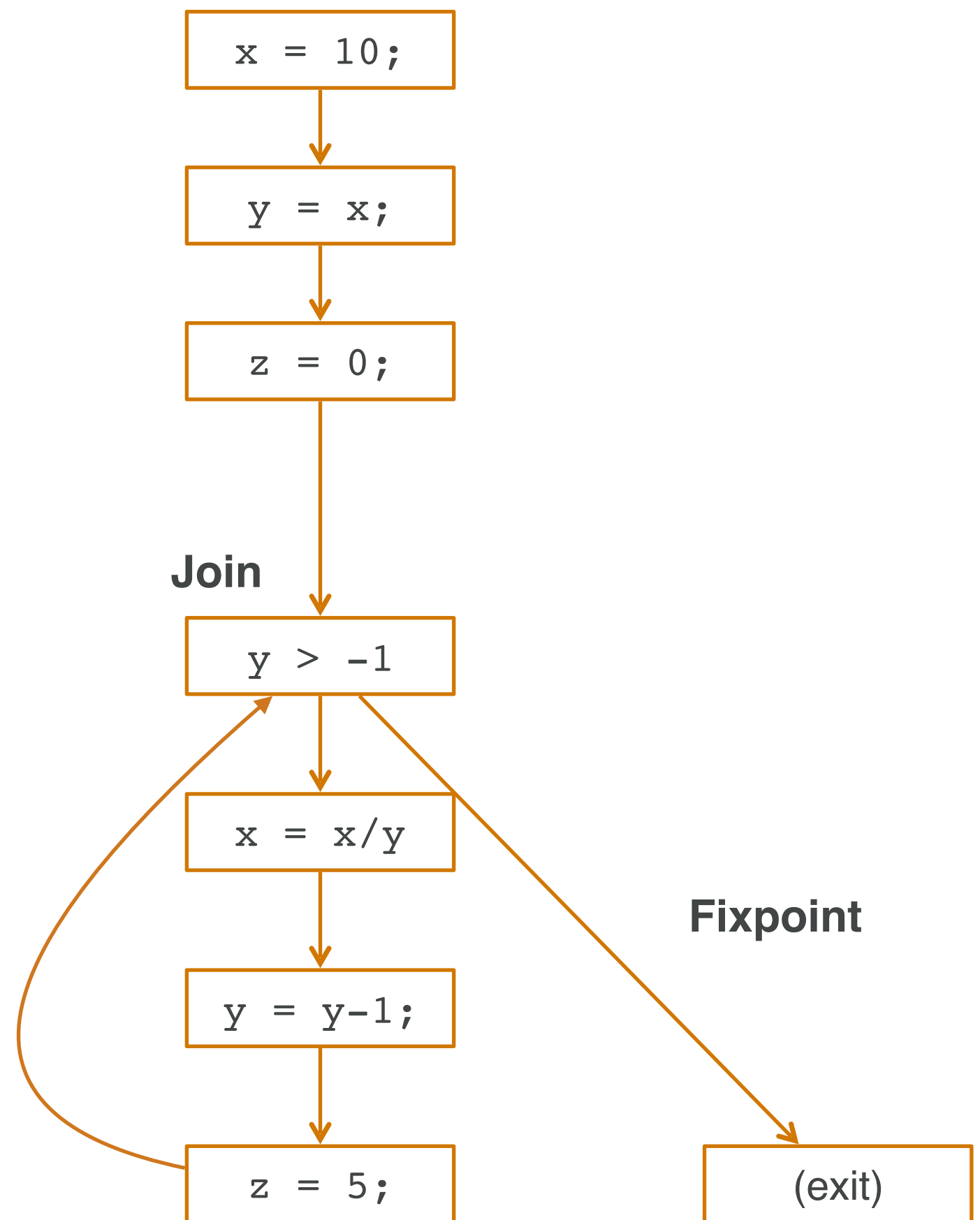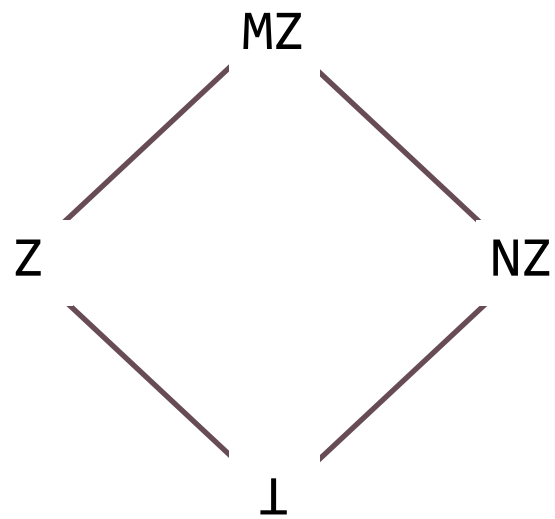
institute for
SOFTWARE
RESEARCH

```
1. x = 10;
2. y = x;
3. z = 0;
4. while(y > -1) {
5.    x = x/y;
6.    y = y-1;
7.    Z = 5;
8. }
```

MZ

Z          NZ

⊥

x = 10;

y = x;

z = 0;

**Join**

y > -1

x = x/y

**Fixpoint**

y = y-1;

z = 5;

(exit)

```
1. x = 10;
2. y = x;
3. z = 0;
4. while(y > -1) {
5.     x = x/y;
6.     y = y-1;
7.     Z = 5;
8. }
```
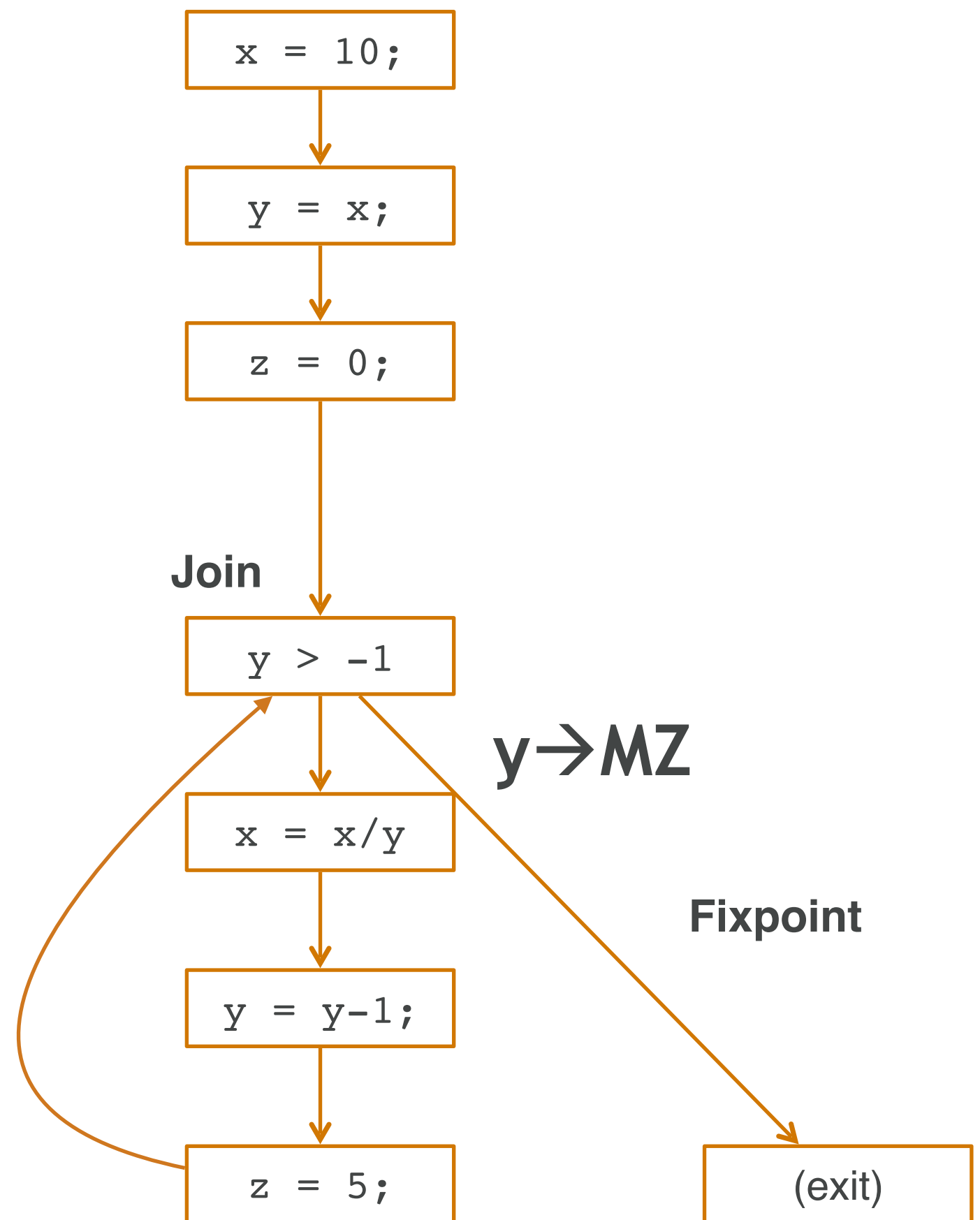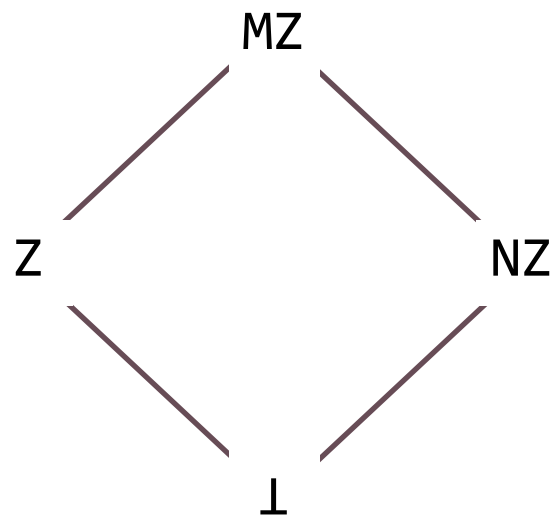
MZ

Z                    NZ

⊥

```
x = 10;
```

```
y = x;
```

```
z = 0;
```

**Join**

```
y > -1
```

**y→MZ**

```
x = x/y
```

**Fixpoint**

```
y = y-1;
```

```
z = 5;
```

(exit)

institute for
SOFTWARE
RESEARCH

# Benefits

# Benefits

- Analyzes all possible executions of the program.
- Pinpoint in code where issues occur.
- Detects issues in the early stages of development.

# Limitations

# Limitations

- Rice's Theorem: Every static analysis is necessarily incomplete or unsound or undecidable (or multiple of these).
- Difficult to track runtime properties.
- Can analyze parts of the program that are never executed.
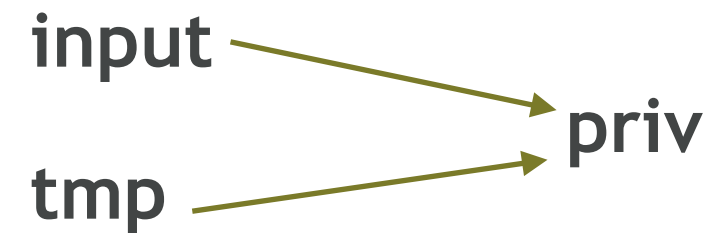
# TAINT ANALYSIS

# Taint Analysis

- Information flow analysis.
- Used in the security domain.
- Tracking how private information flows through the program and if it is leaked to public observers.

institute for SOFTWARE RESEARCH

# Example

```
1. input = get_input();
2. tmp = "select …" + input;
3. query(tmp);
4. log(tmp);
```

institute for
SOFTWARE
RESEARCH

# Example

1. `input = get_input();`
2. `tmp = "select …" + input;`
3. `query(tmp);`
4. `log(tmp);`

input

tmp

priv

Warning!

institute for
SOFTWARE
RESEARCH

# Terminology

- Sources
  - Private data of interest

- Sinks
  - Locations of interest
  - Check taints of incoming information
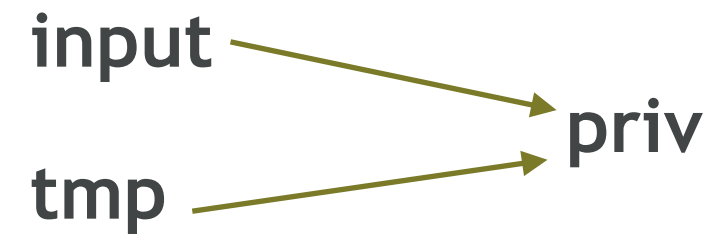  - Determines if there is a leak in the program.

# Example

```
1. input = get_input();
2. tmp = "select …" + input;
3. query(tmp);
4. log(tmp);
```

# Example

```
1. input = Source();
2. tmp = "select …" + input;
3. Sink(tmp);
4. log(tmp);
```

# Example

1. input = **Source();**
2. tmp = "select …" + input;
3. **Sink(**tmp**);**
4. log(tmp);

input

tmp

priv

Warning!

isr institute for SOFTWARE RESEARCH

# Example

1. input = **Source()**;
2. tmp = "select …" + input;
3. tmp = **encode(**tmp**)**
4. **Sink(**tmp**)**;
5. log(tmp);

input ⟶ priv

tmp ⟶ …

OK

# DYNAMIC TAINT ANALYSIS

isr institute for SOFTWARE RESEARCH

# Dynamic Taint Analysis

- Track what are the taints that are influencing the values of the program.

# Example

```
1. x = get_input();
2. y = 1;
3. z = x;
4. w = y + z;
5. print(w);
```

# Example

```
1. x = Source(0);
2. y = 1;
3. z = x;
4. w = y + z;
5. Sink(w);
```

# Example

1. x = Source(0);
2. y = 1;
3. z = x;
4. w = y + z;
5. Sink(w);

$x \longrightarrow 0 \rightarrow T$

# Example

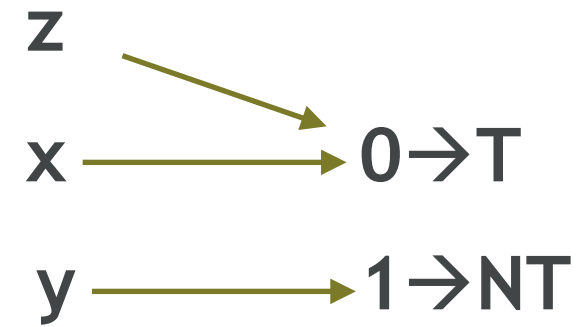1. `x = Source(0);`
2. `y = 1;`
3. `z = x;`
4. `w = y + z;`
5. `Sink(w);`
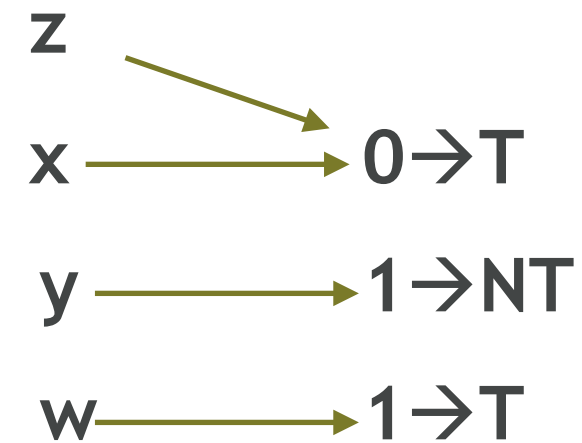
x ⟶ **0→T**

y ⟶ **1→NT**

# Example

1. x = Source(0);
2. y = 1;
3. z = x;
4. w = y + z;
5. Sink(w);

z

x ⟶ 0→T

y ⟶ 1→NT

# Example

1. x = Source(0);
2. y = 1;
3. z = x;
4. w = y + z;
5. Sink(w);

z

x ⟶ **0→T**

y ⟶ **1→NT**

w ⟶ **1→T**

# Example

1. x = Source(0);
2. y = 1;
3. z = x;
4. w = y + z;
5. Sink(w);

z

x ⟶ 0→T

y ⟶ 1→NT

w ⟶ 1→T

42

# Example

1. `x = Source(0);`
2. `y = 1;`
3. `z = x;`
4. `w = y + z;`
5. `Sink(w);`

z

x ⟶ 0→T

y ⟶ 1→NT

w ⟶ 1→T

Leak in the program!

# Is there a leak? Why? Why not?

```
1.    x = Source(0);
2.    y = x;
3.    if(y == 0) {
4.      z = 2
5.    }
6.    else {
7.      z = 1
8.    }
9.    Sink(z);
```

# Implicit Flows

- Tainted data affects the value of another variable indirectly.
- Needed for sound analysis.

# Implicit Flows

```
1.    x = Source(0);
2.    y = x;              Explicit information flow
3.    if(y == 0) {
4.       z = 2
5.    }                   Implicit information flow
6.    else {
7.       z = 1
8.    }
9.    Sink(z);
```

# Implicit Flows

```
1.    x = Source(0);
2.    y = x;
3.    if(y == 0) {
4.      z = 2
5.    }
6.    else {
7.      z = 1
8.    }
9.    Sink(z);
```

x ⟶ 0→T
y ⟋

# Implicit Flows

```
1.    x = Source(0);
2.    y = x;
3.    if(y == 0) {
4.       z = 2
5.    }
6.    else {
7.       z = 1
8.    }
9.    Sink(z);
```

x ⟶ 0→T
y
z ⟶ 2→T

# Implicit Flows
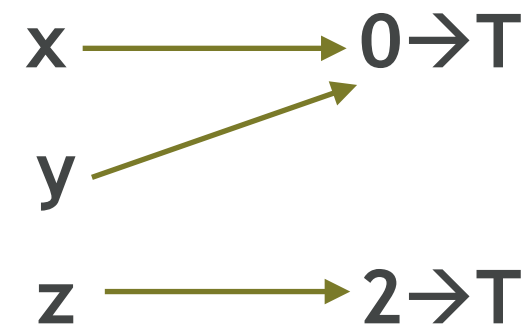
```
1.    x = Source(0);
2.    y = x;
3.    if(y == 0) {
4.      z = 2
5.    }
6.    else {
7.      z = 1
8.    }
9.    Sink(z);
```

x ⟶ 0→T

y

z ⟶ 2→T

Leak in the program!

isr institute for SOFTWARE RESEARCH

# Try it yourself

```
1.    x = Source(1);
2.    y = 0;
3.    while(x > 0) {
4.        y = y + 1;
5.        x = x - 1;
6.    }
7.    z = y;
8.    Sink(y);
9.    Sink(z);
```

# Try it yourself

```
1.    x = Source(1);
2.    y = 0;
3.    while(x > 0) {
4.        y = y + 1;
5.        x = x - 1;
6.    }
7.    z = y;
8.  Sink(y);
9.  Sink(z);
```

x $\longrightarrow$ 0→T

y $\longrightarrow$ 1→T

z

0→NT

Leaks in the program!

isr institute for SOFTWARE RESEARCH

# Limits of Dynamic Analysis

- Results are input dependent.
- Implicit flows needed for sound analysis, but difficult to track*.


- *Stayed tuned for the end of lecture.

# STATIC TAINT ANALYSIS

# Static Taint Analysis

- Track, at each instruction, what are the taints that are influencing the variables of the program.

# Example

```
1. x = Source(i);
2. y = 1;
3. z = x;
4. w = y + z;
5. Sink(w);
```

# Example

```
1. x = Source(i);        x→T
2. y = 1;                x→T
3. z = x;                x→T, z→T
4. w = y + z;            x→T, z→T, w→T
5. Sink(w);              x→T, z→T, w→T
```

# Example

1. `x = Source(i);`        x→T
2. `y = 1;`                x→T
3. `z = x;`                x→T, z→T
4. `w = y + z;`            x→T, z→T, w→T
5. `Sink(w);`             x→T, z→T, **w→T**

## Leak in the program!

# Implicit Flows

```
1.   x = Source(i);
2.   y = x;
3.   if(y == 0) {
4.     z = 0
5.   }
6.   else {
7.     z = 1
8.   }
9.   Sink(z);
```

```
1.   x = Source(i, "A");
2.   y = x;
3.   if(y == 0) {
4.     z = 0
5.   }
6.   else {
7.     z = 1
8.   }
9.   Sink(z);
```
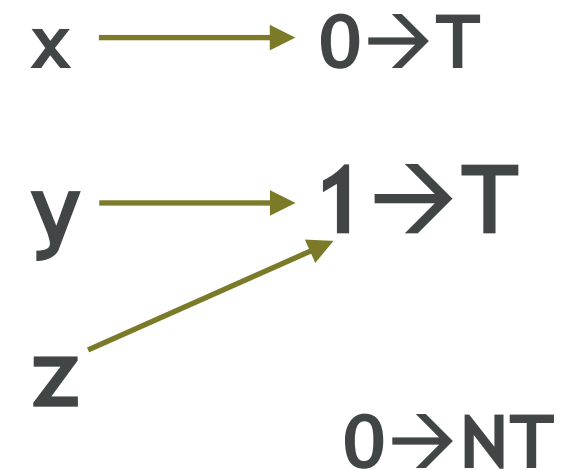
```
1.    x = Source(i);
2.    y = x;
3.    if(y == 0) {
4.       z = 0
5.    }
6.    else {
7.       z = 1
8.    }
9.    Sink(z);
```

MT

T                    NT

⊥

1: x = Source(i);

2: y = x;

3: y == 0

4: z = 0          7: z = 1

9: Sink(z);

(exit)

# Kildall's Worklist Algorithm

```
for Instruction i in program
    input[i] = ⊥
input[firstInstruction] = initialDataflowInformation
worklist = { firstInstruction }

while worklist is not empty
    take an instruction i off the worklist
    output = flow(i, input[i])
    for Instruction j in succs(i)
        if output ⋢ input[j]
            input[j] = input[j] ⊔ output
            add j to worklist
```

institute for SOFTWARE RESEARCH

```
1: x = Source(i);
        |
        v
2: y = x;
        |
        v
3: y == 0
       / \
      v   v
4: z = 0   7: z = 1
      \   /
       v v
   9: Sink(z);
        |
        v
     (exit)
```

| | Input | | |
|------|------|------|------|
| Stmt | x | y | z |
| 1 | ⊥ | ⊥ | ⊥ |
| 2 | ⊥ | ⊥ | ⊥ |
| 3 | ⊥ | ⊥ | ⊥ |
| 4 | ⊥ | ⊥ | ⊥ |
| 7 | ⊥ | ⊥ | ⊥ |
| 9 | ⊥ | ⊥ | ⊥ |

| Stmt | Worklist | x | y | z |
|------|----------|---|---|---|
| | | | | |

```
institute for
SOFTWARE
RESEARCH
```

```
1: x = Source(i);

2: y = x;

3: y == 0

4: z = 0          7: z = 1

9: Sink(z);

(exit)
```

| | Input | | |
|---|---|---|---|
| Stmt | x | y | z |
| 1 | NT | NT | NT |
| 2 | T | NT | NT |
| 3 | T | T | NT |
| 4 | T | T | NT |
| 7 | T | T | NT |
| 9 | T | T | T |

| Stmt | Worklist | x | y | z |
|---|---|---|---|---|
| 1 | 2 | T | NT | NT |
| 2 | 3 | T | T | NT |
| 3 | 4,7 | T | T | NT |
| 4 | 7,9 | T | T | T |
| 7 | 9 | T | T | T |
| 9 | | T | T | T |

institute for
SOFTWARE
RESEARCH

```
1: x = Source(i);

2: y = x;

3: y == 0

4: z = 0        7: z = 1

9: Sink(z);

(exit)
```

| Input | | | |
|-------|-----|-----|-----|
| Stmt | x | y | z |
| 1 | NT | NT | NT |
| 2 | T | NT | NT |
| 3 | T | T | NT |
| 4 | T | T | NT |
| 7 | T | T | NT |
| 9 | T | T | T |

| Stmt | Worklist | x | y | z |
|------|----------|---|---|---|
| 1 | 2 | T | NT | NT |
| 2 | 3 | T | T | NT |
| 3 | 4,7 | T | T | NT |
| 4 | 7,9 | T | T | T |
| 7 | 9 | T | T | T |
| 9 | | T | T | T |

Leak in the program!

institute for SOFTWARE RESEARCH

```
1: x = Source(i);
```

```
2: y = x;
```

```
3: y == 0
```

```
4: z = 0
```

```
7: z = 1
```

```
9: Sink(z);
```

```
(exit)
```

| | Input | | |
|------|------|------|------|
| Stmt | x | y | z |
| 1 | MT | MT | MT |
| 2 | T | MT | MT |
| 3 | T | T | MT |
| 4 | T | T | MT |
| 7 | T | T | MT |
| 9 | T | T | T |

| Stmt | Worklist | x | y | z |
|------|----------|---|---|---|
| 1 | 2 | T | MT | MT |
| 2 | 3 | T | T | MT |
| 3 | 4,7 | T | T | MT |
| 4 | 7,9 | T | T | T |
| 7 | 9 | T | T | T |
| 9 | | T | T | T |

**Leak in the program!** 65

institute for
SOFTWARE
RESEARCH

# Try it yourself

```
1.    x = Source(i);
2.    y = x;
3.    if(y == 0) {
4.      z = 0
5.    }
6.    else {
7.      z = 1
8.    }
9.    Sink(z);
```

```
1: x = Source(i);

2: y = 0;

3: x > 0

4: y = y + 1

5: x = x – 1

7: z = y;

8: Sink(y);

9: Sink(z);

(exit)
```

| Stmt | Input | | |
|---|---|---|---|
| | x | y | z |
| 1 | NT | NT | NT |
| 2 | T | NT | NT |
| 3 | T | MT | NT |
| 4 | T | MT | NT |
| 5 | T | MT | NT |
| 7 | T | MT | NT |
| 8 | T | MT | MT |
| 9 | T | MT | MT |

| Stmt | Worklist | x | y | z |
|---|---|---|---|---|
| 1 | 2 | T | NT | NT |
| 2 | 3 | T | NT | NT |
| 3 | 4,7 | T | NT | NT |
| 4 | 5,7 | T | T | NT |
| 5 | 3,7 | T | T | NT |
| 3 | 4,7 | T | MT | NT |
| 4 | 5,7 | T | MT | NT |
| 5 | 7 | T | MT | NT |
| 7 | 8 | T | MT | MT |
| 8 | 9 | T | **MT** | MT |
| 9 | | T | MT | **MT** |

# Possible leak in the program!

institute for SOFTWARE RESEARCH

# Limits of Static Analysis

- Do not know what values might cause the leak.
- Overtainting

# Overtainting anti-patterns

```
1. x = Source(args[0]);
2. Object o = foo();
3. v = o.equals(x);
```

# Overtainting anti-patterns

```
1. x = Source(args[0]);
2. Object o = foo();
3. v = o.equals(x);
```

All implementation of
equals analyzed!

# Overtainting anti-patterns

```
1.    x = Source(args[0]);
2.    if(Math.max(1, x) == 0) {
3.      Sink(x);
4.    }
```

# Overtainting anti-patterns

```
1.    x = Source(args[0]);
2.    if(Math.max(1, x) == 0) {
3.      Sink(x);   x→T
4.    }
```

institute for
SOFTWARE
RESEARCH

# Overtainting anti-patterns

```
1.    i = foo();
2.    j = i + 1;
3.    a[i] = Source();
4.    a[j] = 0;
5.    Sink(a);
6.    Sink(a[i]);
7.    Sink(a[j]);
```
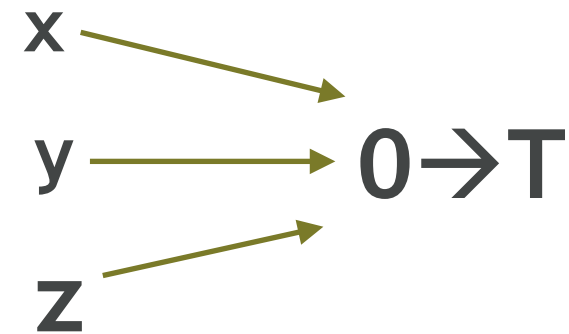
institute for
SOFTWARE
RESEARCH

# Overtainting anti-patterns

```
1.    i = foo();
2.    j = i + 1;
3.    a[i] = Source();
4.    a[j] = 0;
5.    Sink(a);
6.    Sink(a[i]);
7.    Sink(a[j]);
```

Taints the whole array

a→T

a[i]→T

a[j]→T

institute for SOFTWARE RESEARCH

# COMBINING DYNAMIC AND STATIC ANALYSIS

# Implicit Flows in Dynamic Analysis

```
1.    x = Source(0);
2.    y = x;
3.    if(y == 0) {
4.      z = 2
5.    }
6.    else {
7.      z = 1
8.    }
9.    Sink(z);
```

x

y  ⟶  **0→T**

**z**

## Leak in the program!

institute for
SOFTWARE
RESEARCH

# Implicit Flows in Dynamic Analysis

```
1.    x = Source(3);
2.    y = x;
3.    if(y == 0) {
4.      z = 2
5.    }
6.    else {
7.      z = 1
8.    }
9.    Sink(z);
```
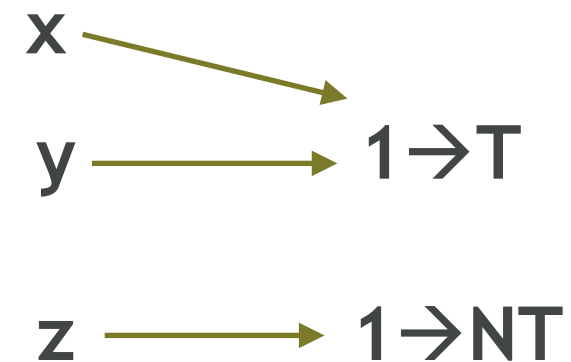
x
y ⟶ **1→T**
**z**

## Leak in the program!

institute for
SOFTWARE
RESEARCH

# Is there a leak? Why? Why not?

```
1.    x = Source(3);
2.    y = x;
3.    z = 1;
4.    if(y == 0) {
5.      z = 2
6.    }
7.    Sink(z);
```

# Is there a leak? Why? Why not?

```
1.    x = Source(3);
2.    y = x;
3.    z = 1;
4.    if(y == 0) {
5.       z = 2
6.    }
7.    Sink(z);
```

x

y $\longrightarrow$ 1$\to$T

z $\longrightarrow$ 1$\to$NT

## No leak in the program!

# Different result for Semantically the same Program?

```
1.    x = Source(3);
2.    y = x;
3.    if(y == 0) {
4.       z = 3
5.    }
6.    else {
7.       z = 1
8.    }
9.    Sink(z);
```
Leak!

```
1.    x = Source(3);
2.    y = x;
3.    z = 1;
4.    if(y == 0) {
5.       z = 2
6.    }
7.    Sink(z);
```

No Leak!

# Fundamental Issue

- In dynamic taint analysis, some implicit flows are hard to track

- If the code is not executed, we do not track its information.

# How would you solve this issue?

```
1.    x = Source(3);
2.    y = x;
3.    if(y == 0) {
4.      z = 2
5.    }
6.    else {
7.      z = 1
8.    }
9.    Sink(z);
```

Leak!

```
1.    x = Source(3);
2.    y = x;
3.    z = 1;
4.    if(y == 0) {
5.      z = 2
6.    }
7.    Sink(z);
```
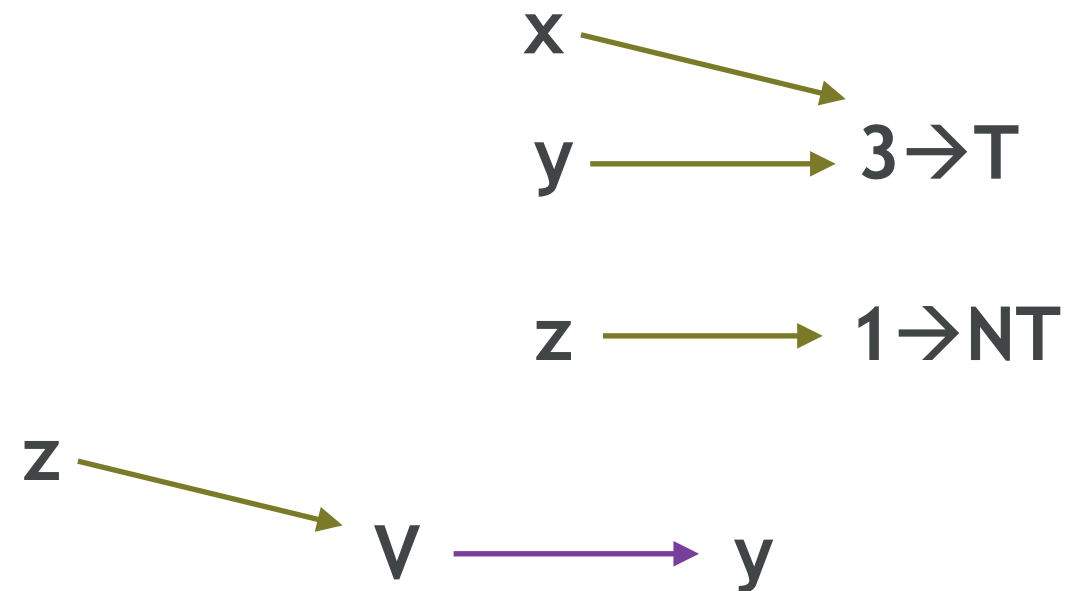
No Leak!

institute for
SOFTWARE
RESEARCH

# Branch-not-taken Analysis

```
1.    x = Source(i);
2.    y = x;
3.    z = 1;

5.    if(y == 0) {
6.      z = 2
7.    }
8.    Sink(z);
```
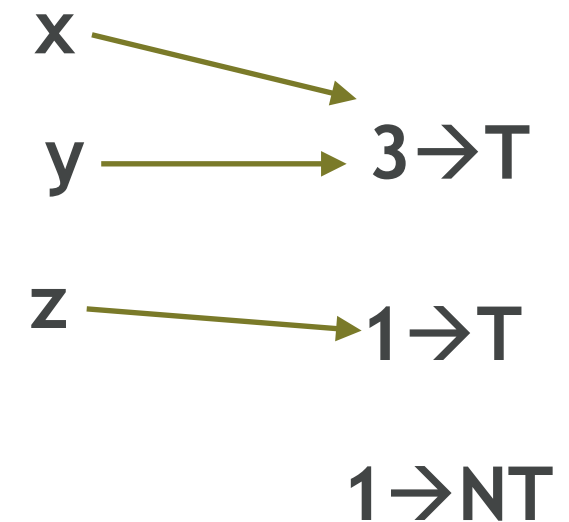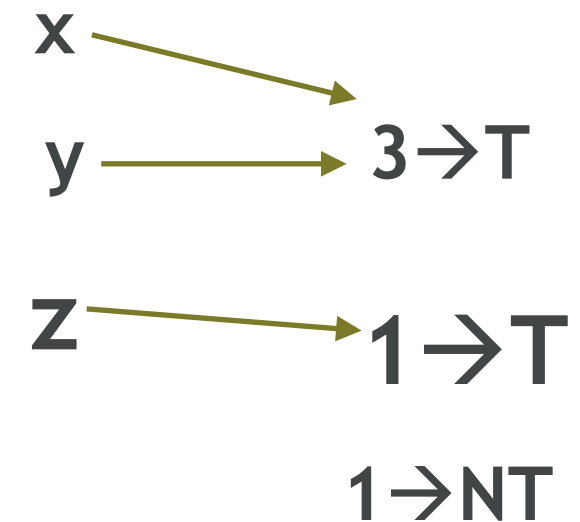
z ⟶ v ⟶ y

institute for
SOFTWARE
RESEARCH

# Branch-not-taken Analysis

```
1.    x = Source(3);
2.    y = x;
3.    z = 1;

5.    if(y == 0) {
6.       z = 2
7.    }
8.    Sink(z);
```

x

y ⟶ 3→T

z ⟶ 1→NT

z ⟶ V ⟶ y

institute for
SOFTWARE
RESEARCH

# Branch-not-taken Analysis

1.    `x = Source(3);`
2.    `y = x;`
3.    `z = 1;`

5.    `if(y == 0) {`
6.      `z = 2`
7.    `}`
8.    `Sink(z);`

x

y   → **3→T**

z  → **1→T**

**1→NT**

institute for
SOFTWARE
RESEARCH

# Branch-not-taken Analysis

```
1.    x = Source(3);
2.    y = x;
3.    z = 1;


5.    if(y == 0) {
6.       z = 2
7.    }
8.    Sink(z);
```

x

y ⟶ 3→T

z

1→T

1→NT

Leak in the program!

institute for
SOFTWARE
RESEARCH

# Is there a leak? Why? Why not?

```
1.  x = Source(3);
2.  y = x;
3.  z = 1;
4.  w = 1;
5.  if(y == 0) {
6.     z = 2
7.     if(x == 0) {
8.        w = 0;
9.     }
10. }
11. Sink(w);
```

# Limits of Branch-not-taken Analysis

```
1.  x = Source(3);
2.  y = x;
3.  z = 1;
4.  w = 1;

6.  if(y == 0) {
7.      z = 2

9.      if(x == 0) {
10.         w = 0;
11.     }
12.  }
13.  Sink(w);
```

z ⟶ V ⟶ y

w ⟶ V ⟶ x

institute for
SOFTWARE
RESEARCH

# INTERPROCEDURAL ANALYSIS

institute for
SOFTWARE
RESEARCH

# Interprocedural Analysis

```
1.  main() {
2.     x = Source(1);
3.     y = 1;
4.     z = foo(x);
5.     Sink(z);
6.     z = foo(y);
7.     Sink(z);
8.  }
```

```
1.  foo(x) {
2.     y = x * 2;
3.     return x;
4.  }
```

# Interprocedural Analysis

```
1.  main() {                1.  foo(x) {
2.    x = Source(1);        2.    y = x * 2;
3.    y = 1;                3.    return x;
4.    z = foo(x);           4. }
5.    Sink(z);
6.    z = foo(y);
7.    Sink(z);
8. }
```

Information with
context T

# Interprocedural Analysis

```
1.  main() {                    1.  foo(x) {
2.    x = Source(1);            2.    y = x * 2;
3.    y = 1;                    3.    return x;
4.    z = foo(x);              4.  }
5.    Sink(z);
6.    z = foo(y);
7.    Sink(z);
8. }
```

Information with
context T

Information with
context NT

# Summary

- Taint analysis is an information flow analysis to detect if private data is leaked in the program.

- Compare benefits and limitations of dynamic and static approaches.

- Can be combined to overcome the limitations of the other.