# EXERCISE #7
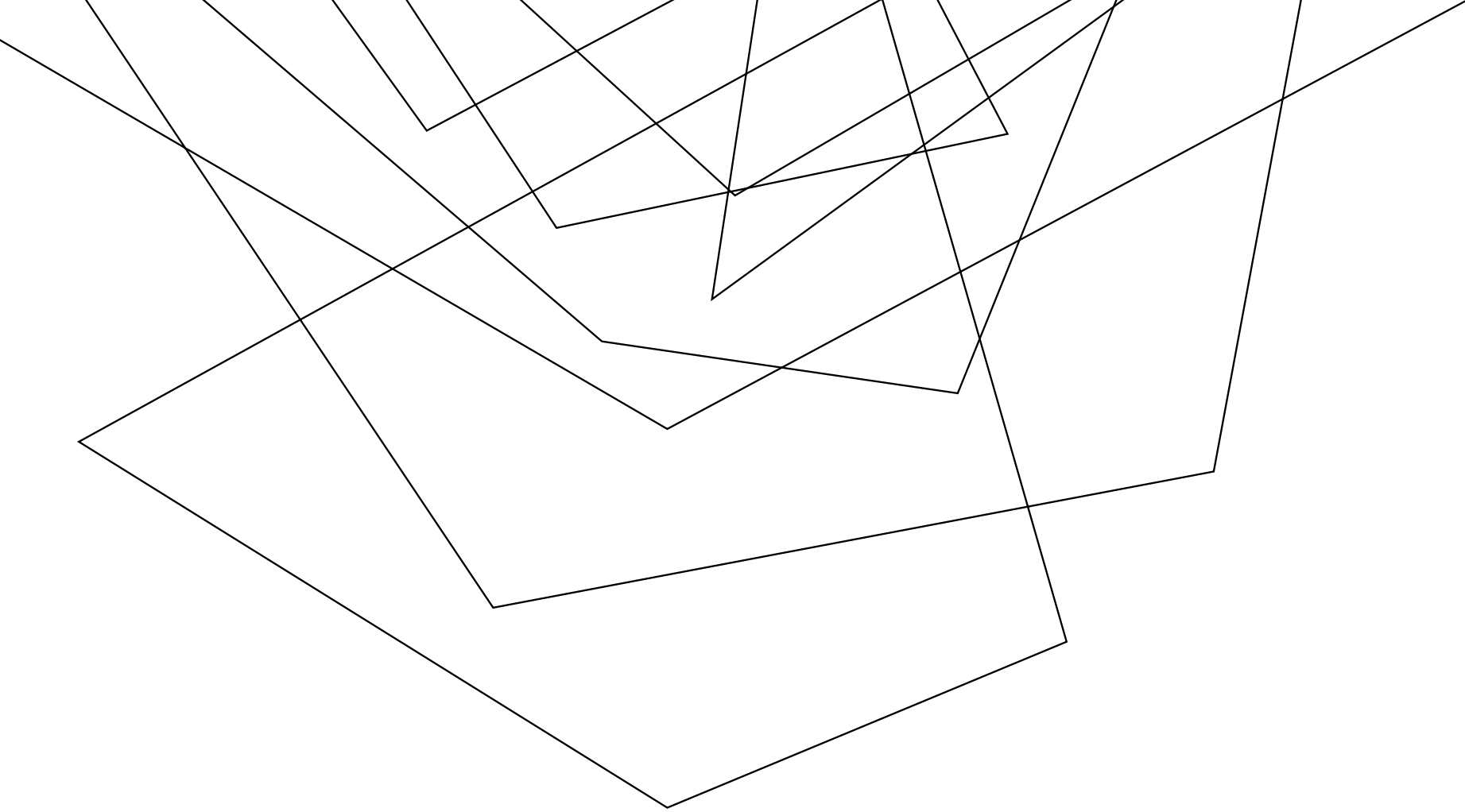
**Write your name and answer the following on a piece of paper**

Consider a simple bug-finding analysis that looks for null pointer deferences in C programs. The analysis raises an alert on any program that has ANY pointer operation, and does not raise an alert on any other program.
Is this analysis sound, complete, neither, or both? Justify your answer.

# STATIC ANALYSIS

EECS 677: Software Security Evaluation

Drew Davidson

H1 - past due
H2 - coming out

**ADMINISTRIVIA
AND
ANNOUNCEMENTS**

# LAST TIME: ANALYSIS DEFINITIONS
## REVIEW: COMPUTABILITY

**Analysis Target – The system being analyzed**

- For us this will usually be a software program

**Analysis Engine – The system doing analysis**

- For us this will usually be a software program

**Analysis Goal – The phenomenon we are detecting**

- The existence of a certain (program) behavior?
- The absence of a certain (program) behavior?

# LAST TIME: ANALYSIS LIMITS
## REVIEW: COMPUTABILITY

## The limits of computability

- The Halting Problem: No decision procedure for halting
- Rice's Theorem: The Halting Problem implies no decision procedure for any reachability problem

## Analysis without decision procedures

- Approximation
- How do we approximate? Soundness / Completeness

# LAST TIME: ANALYSIS DESIGN
## REVIEW: COMPUTABILITY

**NO analysis can be both sound and complete**

**Building an analysis that is <u>either</u> sound <u>or</u> complete is trivial**

- Complete – Always report positive, no false negatives
- Sound – Always report negative, no false positives


POBODY'S NERFECT
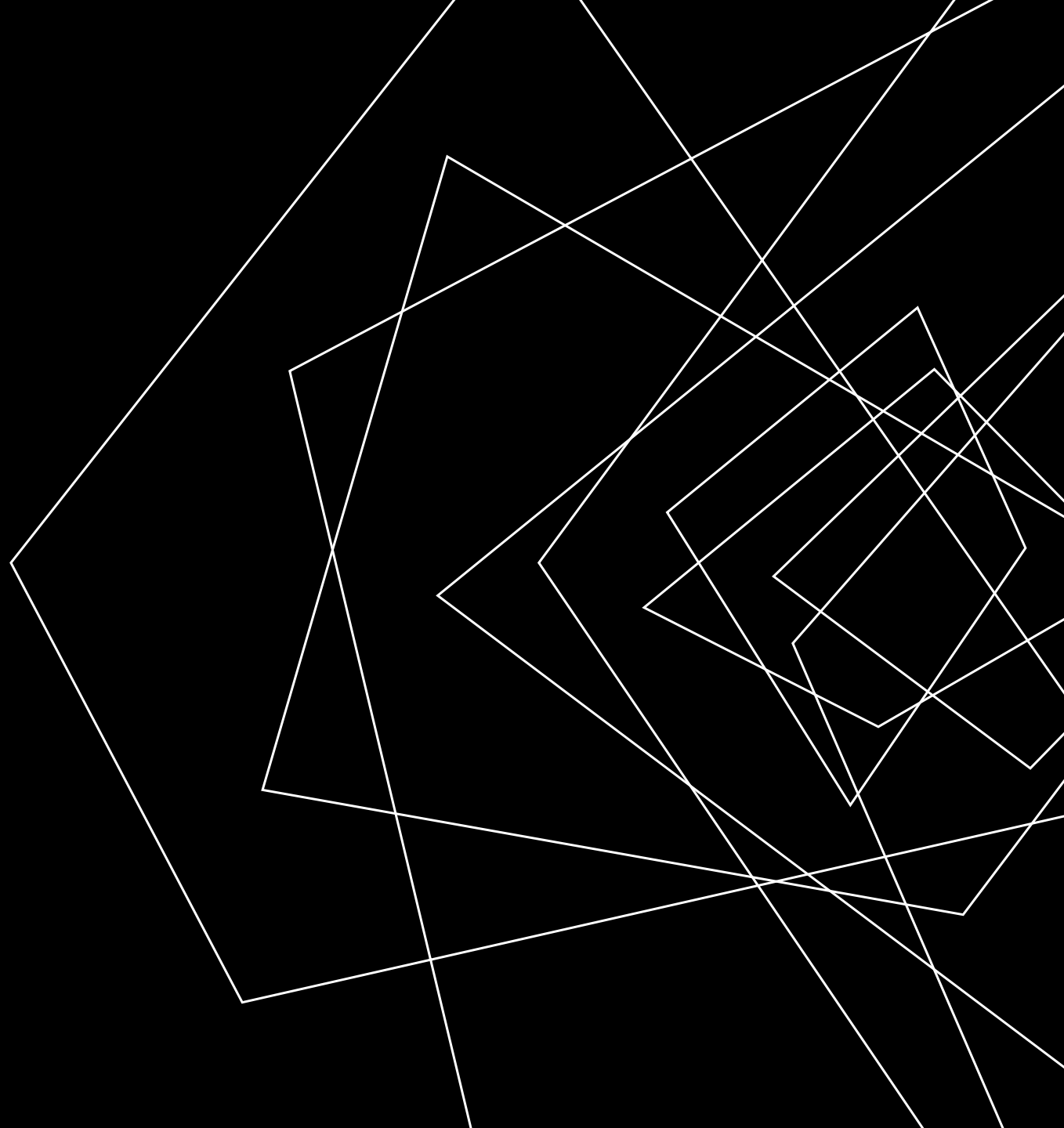
# THIS TIME: WORKING WITH CONSTRAINTS
## TODAY'S LECTURE

Given the limitations of analysis, how might we still provide useful tools for software security evaluation?



**Static Analysis, get it?**

# LECTURE OUTLINE

- Contextualizing Rice's Theorem

- Program Guarantees

- Analysis Specificity

- Dataflow analysis

# RICE'S THEOREM (NON)ASSERTIONS

## CONTEXTUALIZING RICE'S THEOREM

Excludes properties that are true or false for **every** program

Excludes properties that do not regard program's **behavior**

All non-trivial semantic properties of programs are undecidable.

Notably, the theorem ignores **syntactic** properties



*j/k: Rice's Theorem is named after Dr. Henry Gordon Rice*

# SYNTACTIC STATIC ANALYSIS
## CONTEXTUALIZING RICE'S THEOREM

Some troubling behavior of a program may be discoverable via simply observing its text

```
int main(int argc, const char * argv[]){
  const char * password = argv[1];
  int good;
  if (password == "supersecret"){
      good = 1;
  } else {
      good = 0;
  }
  authenticate(good);

}
```

# INSIGHTS
## CONTEXTUALIZING STATIC ANALYSIS

Software engineering "code smells" / stats

Use of the forbidden / arcane constructs

Cyclomatic complexity

Long functions

# RICE'S THEOREM (NON)ASSERTIONS
## CONTEXTUALIZING RICE'S THEOREM

Excludes properties that are true or false for **every** program

Excludes properties that do not regard program's **behavior**

All non-trivial semantic properties of programs are undecidable.

Notably, the theorem ignores **syntactic** properties

No decision procedure: that's a high bar!

We can design analyses that work perfectly in some cases, but admit uncertainty in others

*j/k: Rice's Theorem is named after Dr. Henry Gordon Rice*

# LECTURE OUTLINE

- Contextualizing Rice's Theorem
- Program Guarantees
- Preciseness vs Efficiency
- Dataflow analysis

# STATIC ANALYSIS – OPPORTUNITIES
## STATIC ANALYSIS PHILOSOPHY



For security analysis, we want to lock out "bad" programs (even at the cost of locking out some "good" programs)

**Provide assurances about what a program will NEVER or ALWAYS do**

- Static analysis might report EVERY program that (possibly) has a problem
- Static analysis might certify EVERY program that (definitely) has no problem

# STATIC ANALYSIS – OPPORTUNITIES
## STATIC ANALYSIS PHILOSOPHY

**Provide assurances about what a program will NEVER or ALWAYS do**

- Static analysis might report EVERY program that (possibly) has a problem
- Static analysis might certify EVERY program that (definitely) has no problem

Target Program → Analysis

Definitely Good — Complete bug finder – no flag

¯\\_(ツ)_/¯ — Complete bug finder - flag

Definitely Bad

lump "unsure" into "bad"

# STATIC ANALYSIS – OPPORTUNITIES
## STATIC ANALYSIS PHILOSOPHY



Target Program → Analysis

Definitely Good — Complete bug finder – no flag

¯\_(ツ)_/¯

Definitely Bad

Complete bug finder - flag

**lump "unsure" into "bad"**

**Goal: minimize the uncertainty**

# LECTURE OUTLINE

- Contextualizing Rice's Theorem

- Program Guarantees

- Analysis over CFG

- Preciseness vs Efficiency

# STATIC ANALYSIS – OPPORTUNITIES
## STATIC ANALYSIS PHILOSOPHY



**Definitely Good** — Complete bug finder – no flag

¯\\_(ツ)_/¯ — Complete bug finder - flag

**Definitely Bad**

**lump "unsure" into "bad"**

**Goal: minimize the uncertainty
(comes in part from reachability)**

Key idea: Build the Control-Flow Graph,
explore routes through the graph to
(over)approximate reachability

# BUILDING THE CFG
## CONTROL FLOW GRAPH ANALYSIS

### In general, an iterative process:

Pass over instructions, mark leaders/ terminators
        Might create new leaders / terminators,
        so keep doing passes until no more found
Build basic blocks by boundaries
Connect control source to control destination
Refine based on analyses

### For LLVM Bitcode, even easier

All blocks (except possibly entry) have labels
Connect control source to control destination
**One deviation from LLVM: split after a call!**
Refine based on analyses

# CFG NOTATION
## STATIC ANALYSIS PHILOSOPHY

**We'll enforce a couple of constraints on the CFG**

It must be a hammock –
One entry point, one exit block

Call sites and return points are connected via a special link edge

# VISUALIZING THE CFG: DOT
## CONTROL FLOW GRAPH ANALYSIS

FLOW-SENSITIVE ANALYSIS RELIES HEAVILY ON THE CONTROL-FLOW GRAPH CONCEPT

**It's pretty helpful to have a CFG in hand**

Good news! You know how to automatically induce the CFG structure

Gooder news! There's a format to visualize CFGs

**File graph.dot**

```
digraph name {
   nodeA […];
   nodeB […];
   nodeA -> nodeB […];
}
```

**cmdline**

```
dot -Tpdf graph.dot -o graph.pdf
```

**output**

# EXISTING CFG TOOLS

## CONTROL FLOW GRAPH ANALYSIS

Good news! You know how to automatically induce the CFG structure

Gooder news! There's a format to visualize CFGs

Goodest news! llvm can output a dot-format CFG for .ll-format code

```
opt -dot-cfg prog.ll > /dev/null

opt -passes=dot-cfg prog.ll > /dev/null
```

# OVERAPPROXIMATING REACHABILITY
## STATIC ANALYSIS PHILOSOPHY

```llvm
1 define i32 @main() {
2   ret i32 1
3 blk3:
4   %val = sdiv i32 1, 0
5   ret i32 %val
6 }
```

```llvm
 1 define i32 @main(i32 %argc) {
 2   %pred = icmp sgt i32 %argc, 2
 3   br i1 %pred, label %blkTrue, label %blkFalse
 4
 5 blkTrue:
 6   ret i32 1
 7
 8 blkFalse:
 9   ret i32 2
10
11 blkAfter:
12   %val = sdiv i32 1, 0
13   ret i32 %val
14 }
```
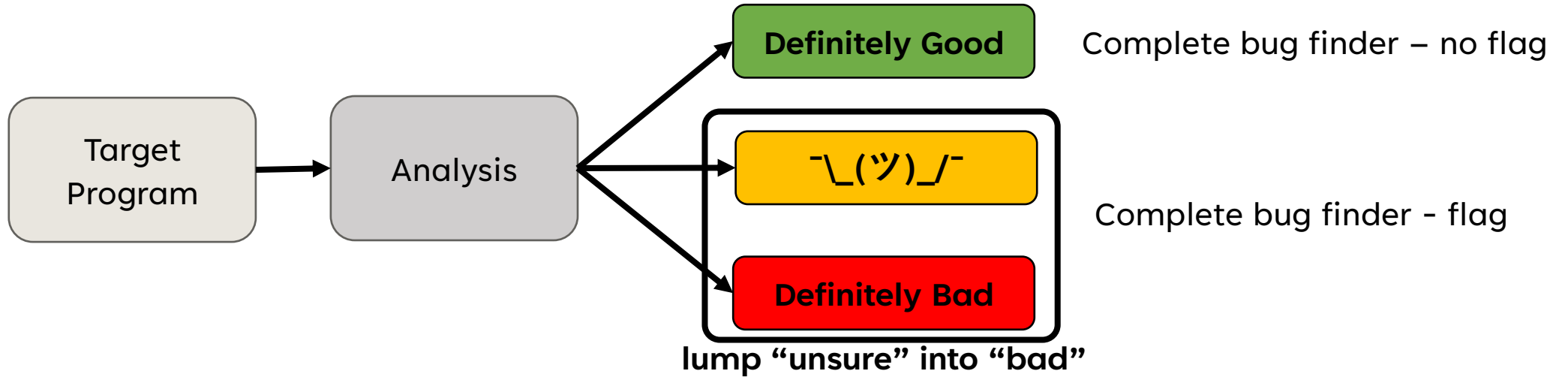
# LECTURE OUTLINE

- Contextualizing Rice's Theorem

- Program Guarantees

- Analysis over CFG

- Preciseness vs Efficiency

# INFEASIBLE PATHS

### STATIC ANALYSIS: DATAFLOW

**Is this program buggy?**

Just check every path(?!?!?!)

```c
int f(int arg) {
    int b = 1;
    int c = 0;
    if (arg > 0){
        b = 0;
        c = 1;
    }
    if (b && c) {
        arg = 2 / 0;
    }
    return arg;
}
```

```llvm
1 define i32 @f(i32 %arg) {
2 entry:
3    %b0 = add i32 0, 1
4    %c0 = add i32 0, 0
5    %argPos = icmp sgt i32 %arg, 0
6    br i1 %argPos, label %block1, label %block2
7
8 block1:
9    %b1 = add i32 0, 0
10   %c1 = add i32 0, 1
11   br label %block2
12
13 block2:
14   %bJoin = phi i32 [%b0, %entry], [%b1, %block1]
15   %cJoin = phi i32 [%b0, %entry], [%b1, %block1]
16   %prod = mul i32 %bJoin, %cJoin
17   %bothNonZero = icmp ne i32 %prod, 0
18   br i1 %bothNonZero, label %block3, label %block4
19
20 block3:
21   %argDiv = sdiv i32 2, 0
22   br label %block4
23
24 block4:
25   %argJoin = phi i32 [%arg, %entry], [%argDiv, %block4]
26   ret i32 %arg
27 }
```
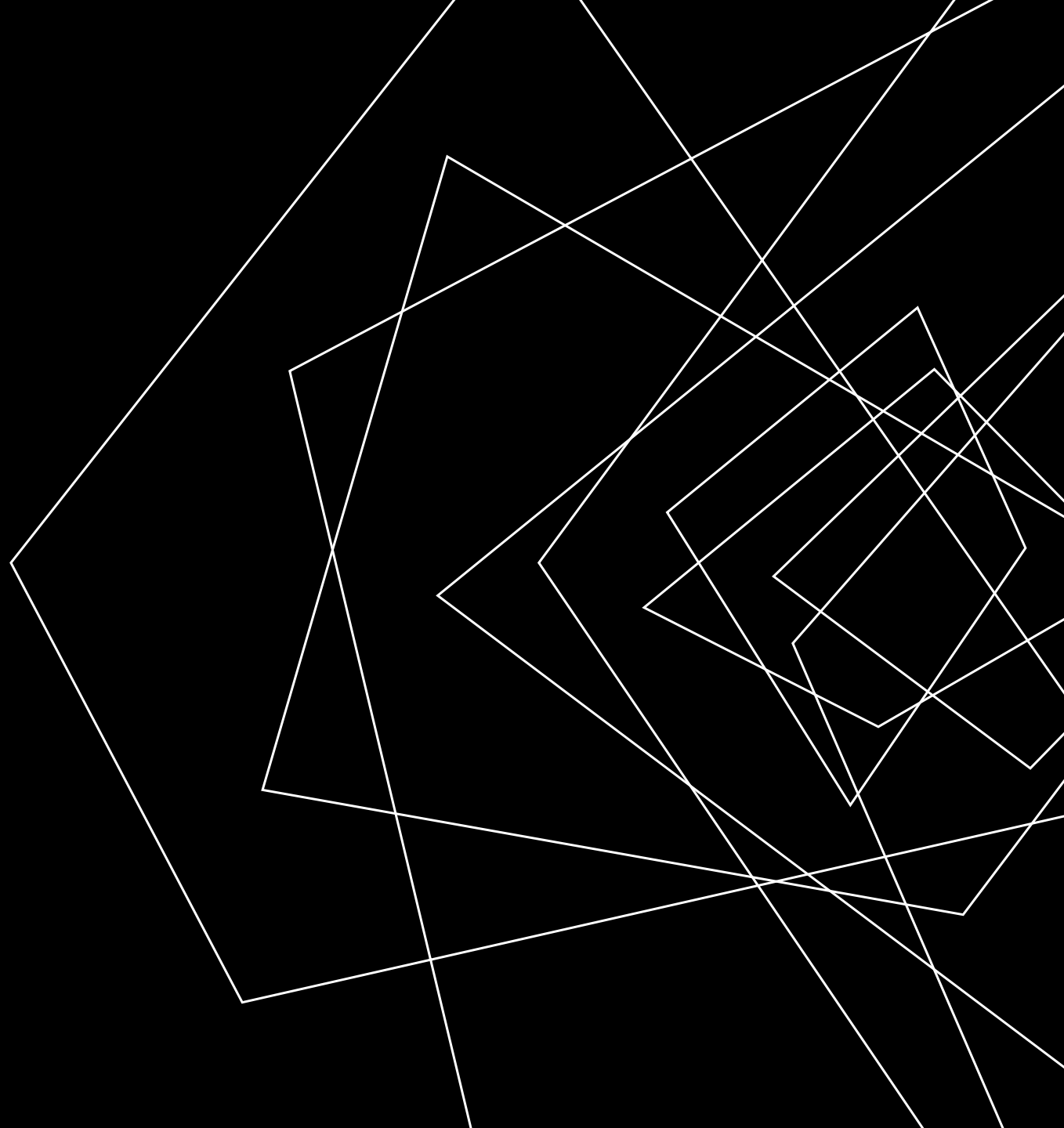
# WRASSLIN' WITH STATE SPACE
## PRECISENESS VS EFFICIENCY

**State space:** the set of all possible configurations of the target model

**Naïve state space representation:** Consider each path separately

# PATH-SENSITIVE DATAFLOW ANALYSIS
## STATIC ANALYSIS: DATAFLOW

**How many paths are there in this program?**

```
 1 define i32 @main(i32 %argc) {
 2 entry:
 3   br label %loop
 4
 5 loop:
 6   %argcJoin = phi i32 [ %argc, %entry ], [ %argcInc, %loop ]
 7   %argcInc = add i32 %argcJoin, -1
 8   %mychar = call i32 (...) @getchar()
 9   %is97 = icmp eq i32 %mychar, 97
10   br i1 %is97, label %after, label %loop
11
12 after:
13   %res = sdiv i32 2, %argcJoin
14   ret i32 %res
15 }
16
17 declare i32 @getchar(...)
```

```
extern int getchar();

int main(int argc){
        int c;
        do {
                argc--;
                c = getchar();
        } while (c != 'a');
        return 2 / argc;
}
```

entry:
br label %loop

**false**

Loop:
%argcJoin = phi i32 [%argc, %entry] , [%argcInc, %loop]
%argcInc = add i32 %argcJoin, -1
%mychar = call i32 (...) @getchar();
%is97 = icmp eq i32 %mychar, 97
Br i1 %is97, label %after, label %loop

**true**

after:
%res = sdiv i32 2, %argcJoin
%argcJoin
ret i32 %res

# WRASSLIN' WITH STATE SPACE
## PRECISENESS VS EFFICIENCY

**State space:** the set of all possible configurations of the analysis target
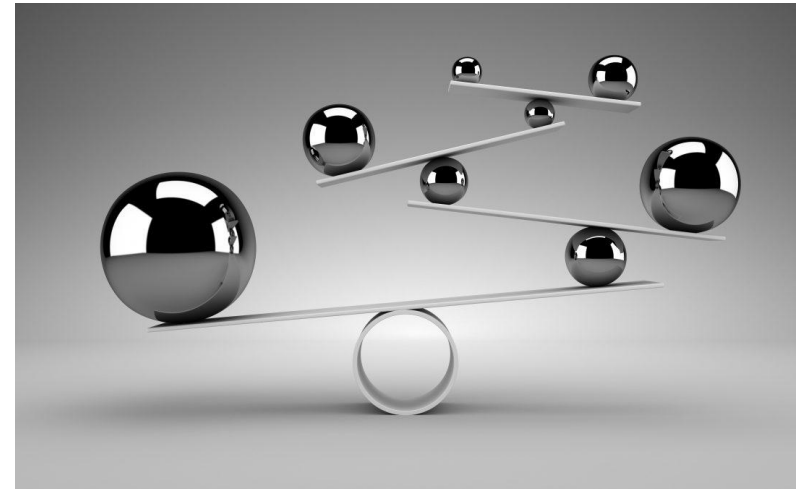
**Naïve state space representation:** enumerate all configurations of a program
- For $n$ bits of memory: $2^n$ states

  *$n$ branches    $2^n$ paths*

**Practical state space representation:** Summarize sets of configurations of a program

# THE POWER OF STATIC ANALYSIS
## ANALYSIS SPECIFICITY

**The power of static analysis:**
You can see beyond the code that is executed in an individual trace

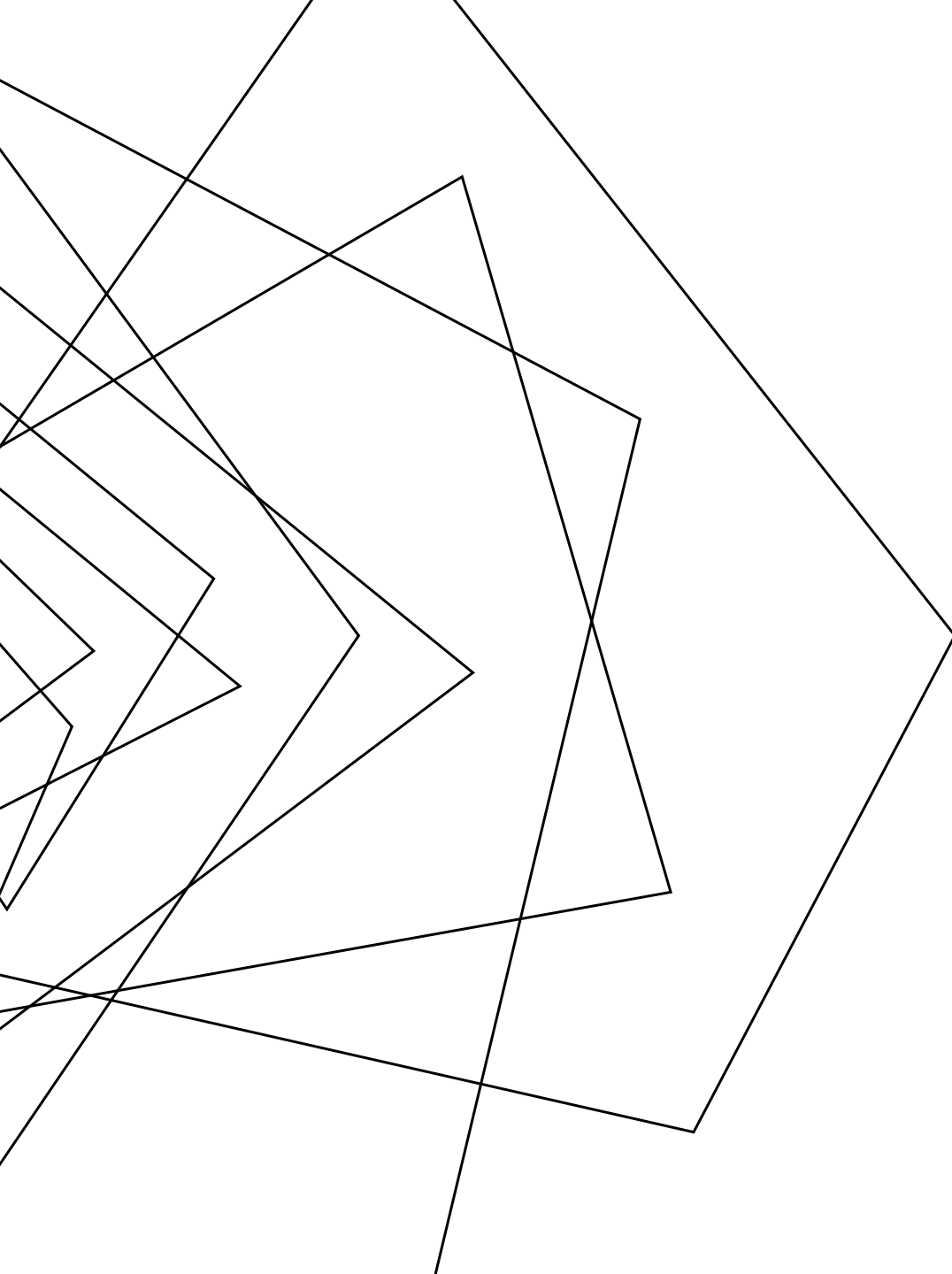**The responsibility of static analysis:**
You need to consider the conditions/circumstances/context in which code *would be* executed while keeping state space small

# OVERAPPROXIMATING REACHABILITY
## STATIC ANALYSIS PHILOSOPHY

# NEXT TIME

FLOW-SENSITIVE ANALYSIS