# PhysCov: Physical Test Coverage for Autonomous Vehicles

**Carl Hildebrandt**, Meriel von Stein, and Sebastian Elbaum

hildebrandt.carl@virginia.edu

# Motivation

## Autonomous Systems are here

Waymo [1]

Oxa [2]

[1] https://waymo.com/
[2] https://oxa.tech/

# Motivation

## Soon they will be common place

Waymo [1]

Oxa [2]

AutoX [3]

Cruise [4]

May Mobility [5]

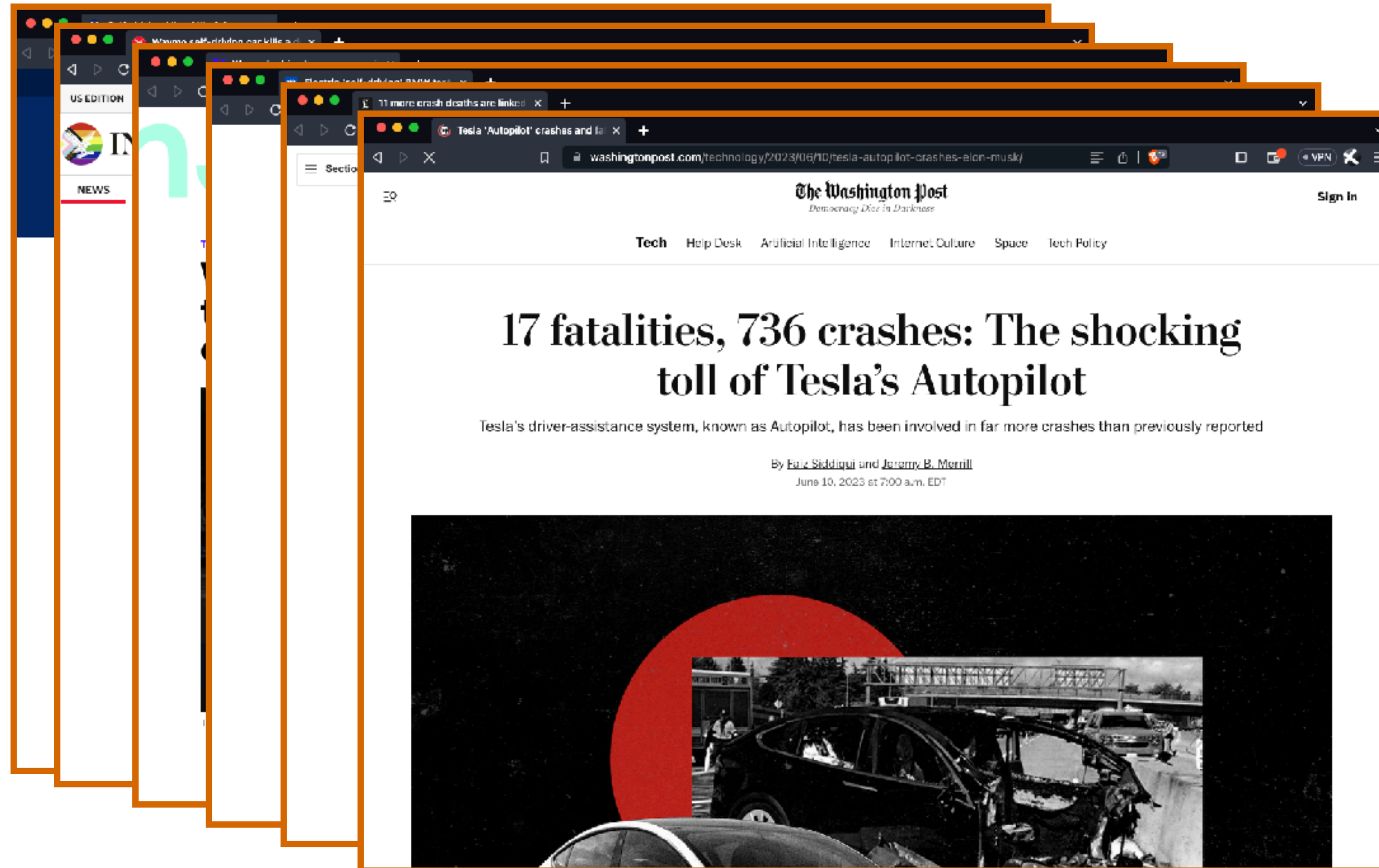Motional [6]

Pony AI [7]

Zoom [8]

Tesla [9]

• • •

[1] https://waymo.com/
[2] https://oxa.tech/
[3] https://www.autox.ai
[4] https://getcruise.com
[5] https://maymobility.com
[6] https://motional.com
[7] https://www.pony.ai
[8] https://zoox.com/
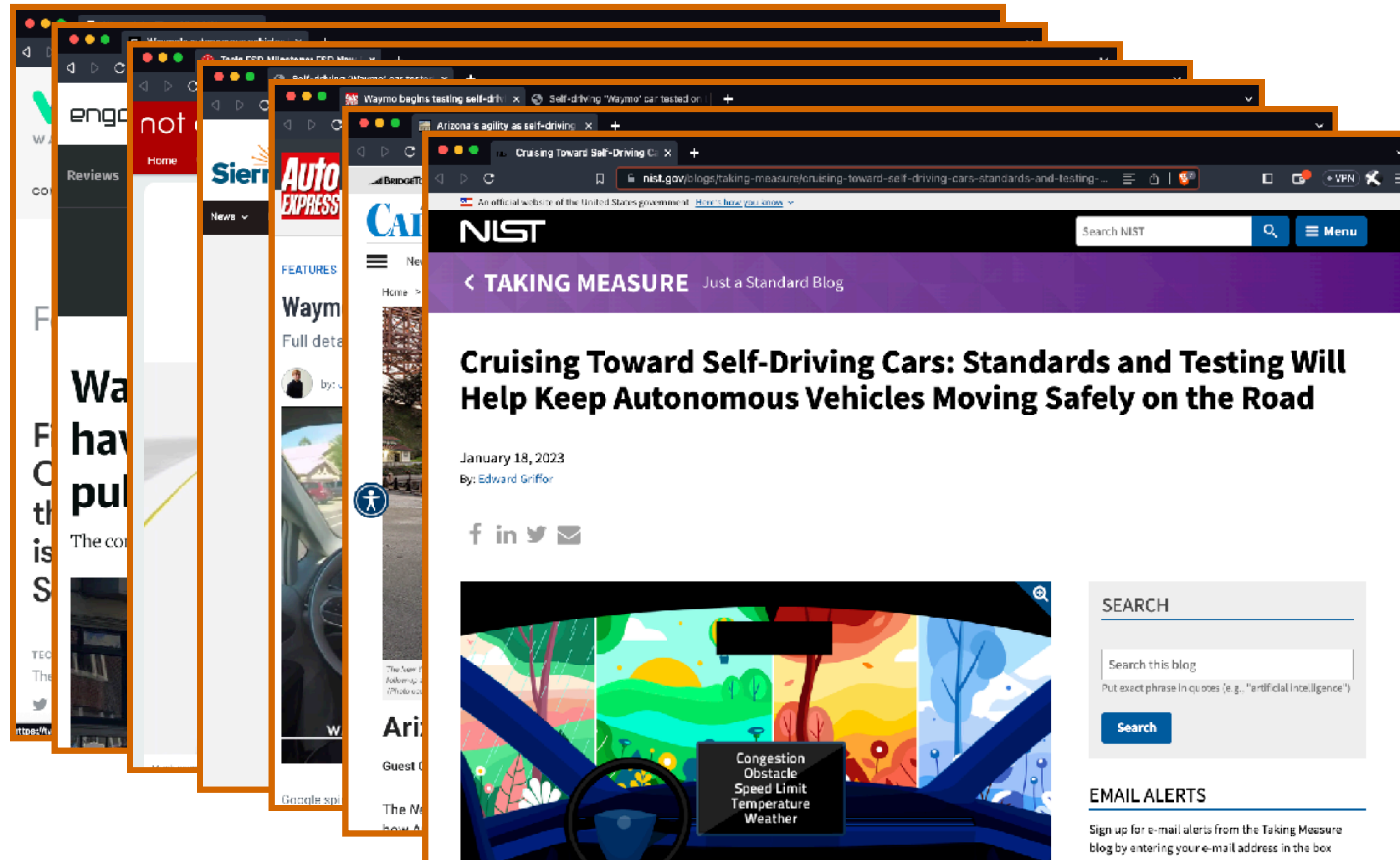[9] https://www.tesla.com/autopilot

# Motivation

These vehicles fail, resulting in the <u>loss of life</u>

# Motivation

## They are being tested

# Motivation
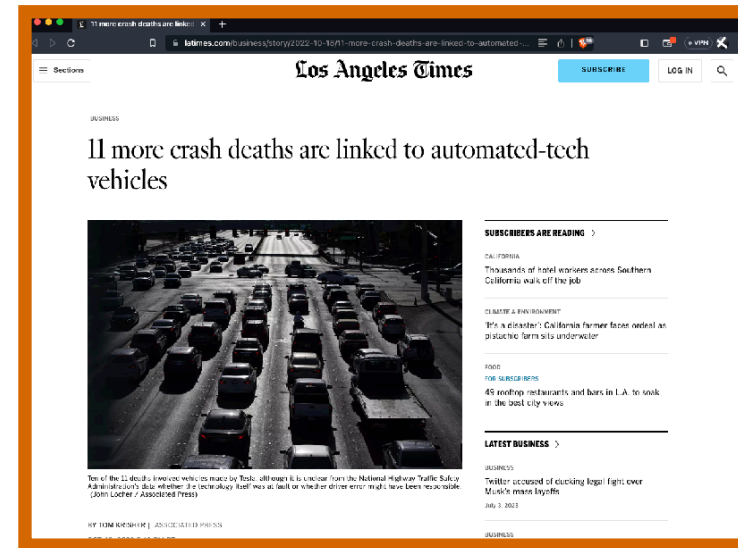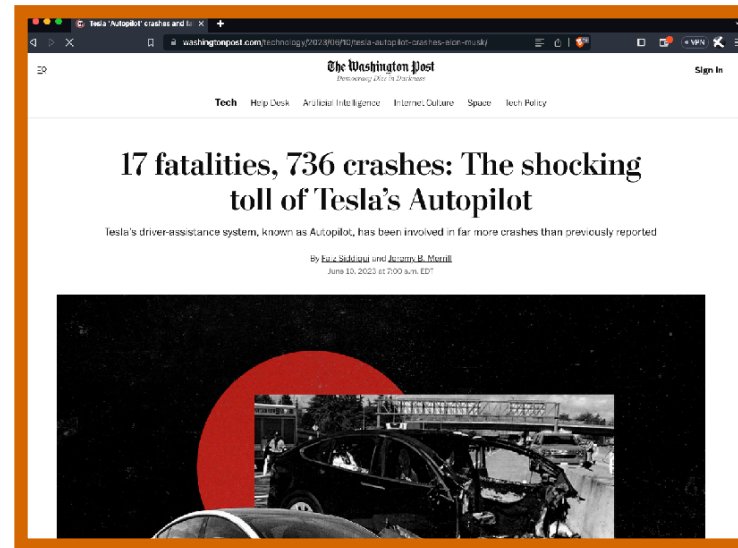
## Failures



## Testing

## Where is the disconnect?

# Disconnect

*"Was the previous test useful?"*

*"How thoroughly is the current system tested?"*

*"When is it safe to stop testing?"*

# Disconnect

*"How do we quantify an autonomous vehicle's test adequacy?"*

# Problem

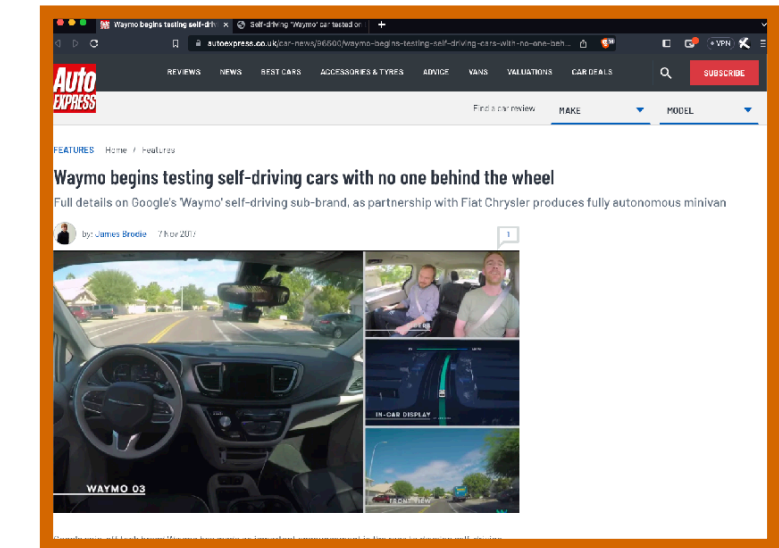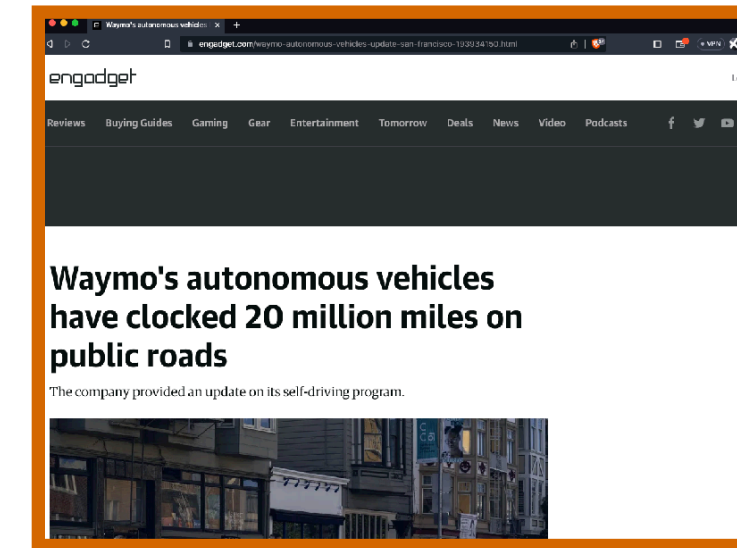## Traditional software uses test adequacy metrics

Input → System → Output

# Problem

## Traditional software uses test adequacy metrics

System

Input

Output

Software

How much of the input space have we seen?

# Problem

## Traditional software uses test adequacy metrics

System



Input → Output

Software

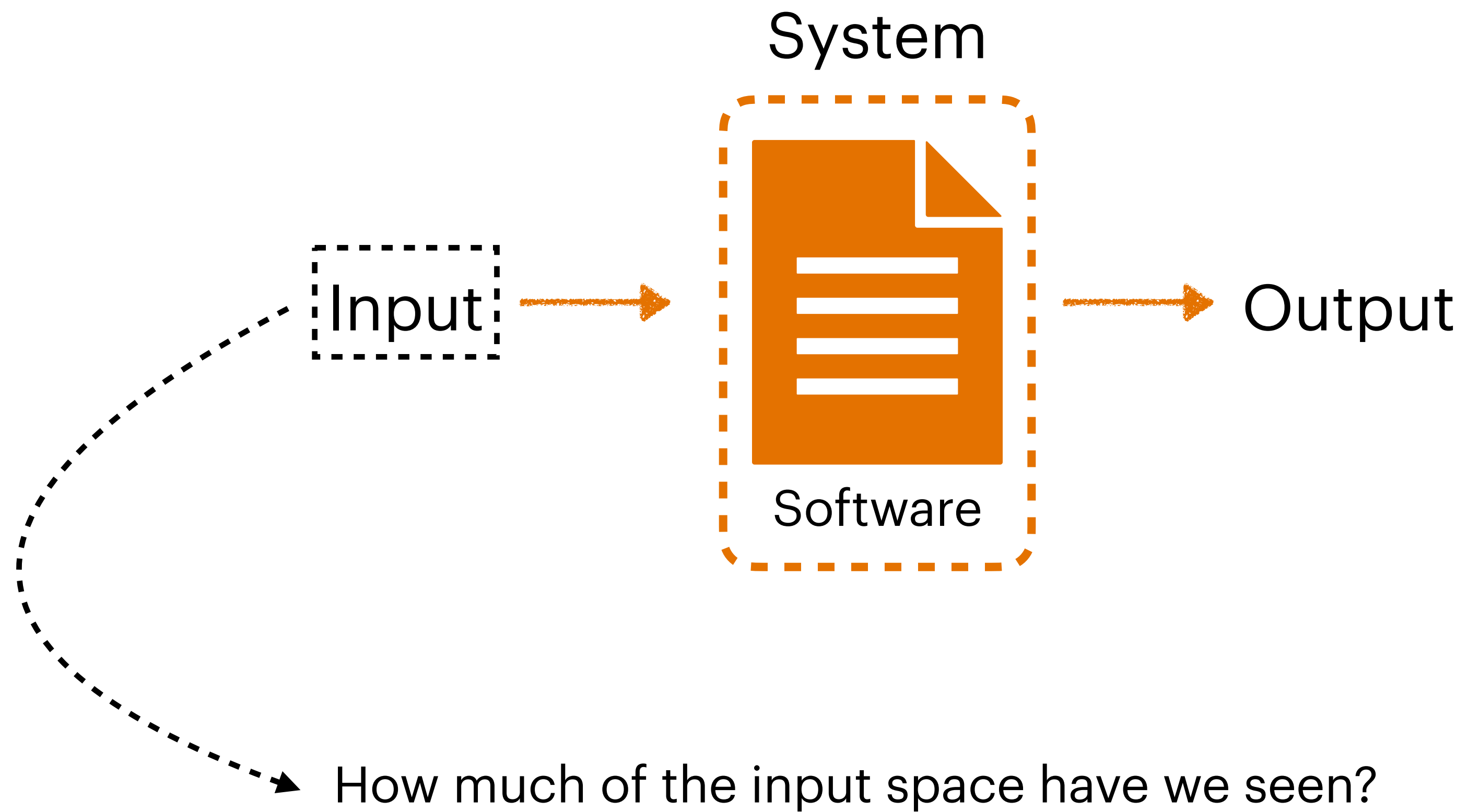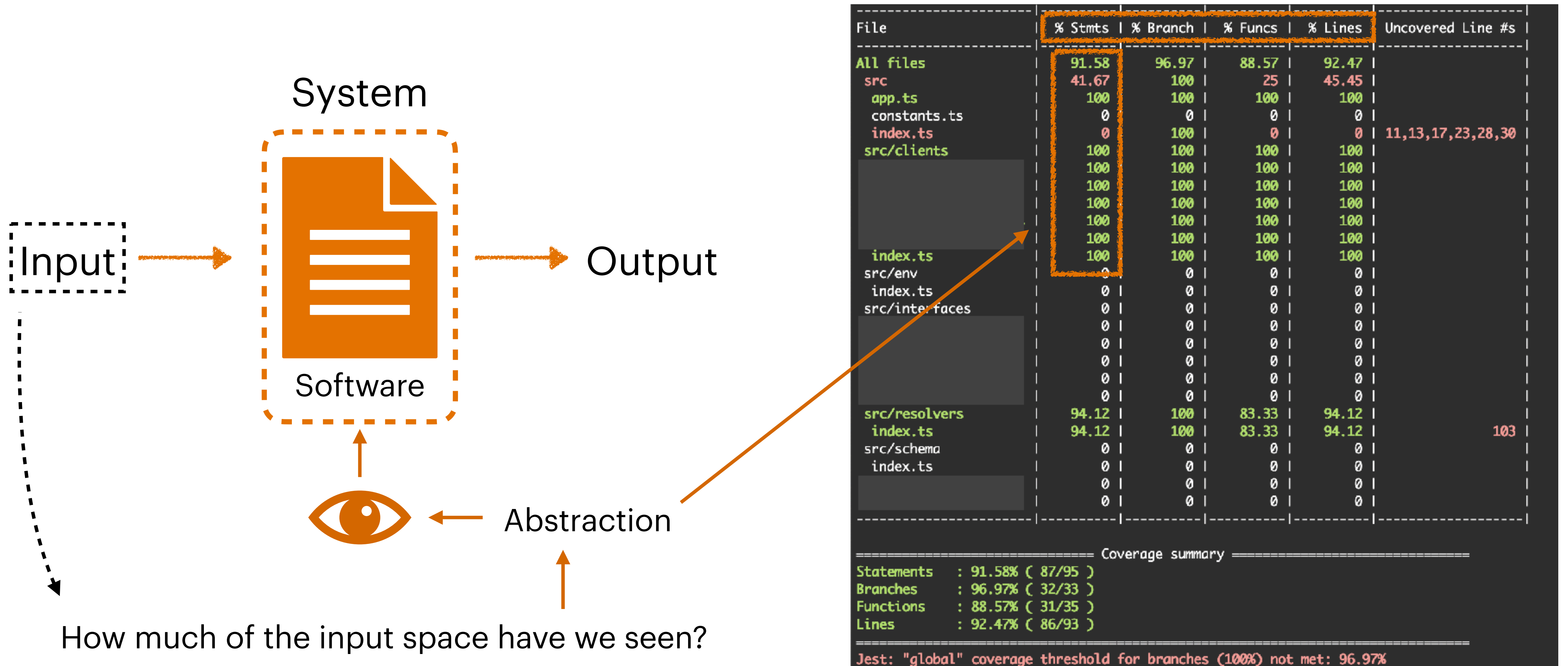👁 ← Abstraction

How much of the input space have we seen?

| File | % Stmts | % Branch | % Funcs | % Lines | Uncovered Line #s |
|------|---------|----------|---------|---------|-------------------|
| All files | 91.58 | 96.97 | 88.57 | 92.47 | |
| src | 41.67 | 100 | 25 | 45.45 | |
| app.ts | 100 | 100 | 100 | 100 | |
| constants.ts | 0 | 0 | 0 | 0 | |
| index.ts | 0 | 100 | 0 | 0 | 11,13,17,23,28,30 |
| src/clients | 100 | 100 | 100 | 100 | |
| | 100 | 100 | 100 | 100 | |
| | 100 | 100 | 100 | 100 | |
| | 100 | 100 | 100 | 100 | |
| | 100 | 100 | 100 | 100 | |
| | 100 | 100 | 100 | 100 | |
| index.ts | 100 | 100 | 100 | 100 | |
| src/env | 0 | 0 | 0 | 0 | |
| index.ts | 0 | 0 | 0 | 0 | |
| src/interfaces | 0 | 0 | 0 | 0 | |
| | 0 | 0 | 0 | 0 | |
| | 0 | 0 | 0 | 0 | |
| | 0 | 0 | 0 | 0 | |
| | 0 | 0 | 0 | 0 | |
| src/resolvers | 94.12 | 100 | 83.33 | 94.12 | |
| index.ts | 94.12 | 100 | 83.33 | 94.12 | 103 |
| src/schema | 0 | 0 | 0 | 0 | |
| index.ts | 0 | 0 | 0 | 0 | |
| | 0 | 0 | 0 | 0 | |

```
========================= Coverage summary =========================
Statements   : 91.58% ( 87/95 )
Branches     : 96.97% ( 32/33 )
Functions    : 88.57% ( 31/35 )
Lines        : 92.47% ( 86/93 )
====================================================================
Jest: "global" coverage threshold for branches (100%) not met: 96.97%
```

11

# Problem

Why can't we do this with autonomous systems?

Input → System → Output

# Problem

Why can't we do this with autonomous systems?

Input → System → Output

Sensor Data ← Environment

# Problem

## Why can't we do this with autonomous systems?

Input → System → Output

Environment

Sensor Data

# Problem

Why can't we do this with autonomous systems?

Input → System → Output

Sensor Data ← Environment

# Problem

Why can't we do this with autonomous systems?



Input → System → Output

Sensor Data

Environment

All Possible Scenarios …

$-\infty$ ⟷ $+\infty$

# Problem

Why can't we do this with autonomous systems?



Input → System → Output

Sensor Data

Environment

# Problem

Why can't we do this with autonomous systems?



State

Input

Software          Hardware

Output

Environment

Sensor
Data

# Problem

Why can't we do this with autonomous systems?



State

Input

Software

Hardware

Output

Environment

Sensor
Data

# Problem

Why can't we do this with autonomous systems?

State

Input → Software → Hardware → Output

Environment

Sensor Data

# Current Solutions

Current approaches are not cognizant of the environment and state

| Coverage Metric | Account for Environments | Account for State |
|---|---|---|
| **Structural Code Coverage** | ❌ | ❌ |
| **Miles Driven / Incident per Miles** | ✅ | ❌ |
| **Requirement Coverage** | ✅ | ✅ |
| **Scenario Coverage** | ✅ | ❌ |
| **Trajectory Coverage** | ✅ | ✅ |
| **Physical Coverage** | ✔️ | ✔️ |

# Insight

**1)** The environment is highly complex and practically infinite:
Only the sensed environment, which the vehicle can reach is important to the vehicles current behavior.

**2)** The vehicles state is dependent the specific systems hardware:
Kinematic models offer a way to abstract the state for any vehicle.

# PhysCov: Approach

# PhysCov: Approach

**SUT**

**RRS Pipeline**

$\tau_k$

AV

Reduction to Reachable Set

Reduction to Sensed Reachable Set

Sensed Reachable Set Vectorization

$f$

**Physical Coverage**

**Test Suite Coverage**

RRS [13, 16, 17, 20, ...., 25, ..., 25]

Environment

Reachable Set

Sensed Reachable Set

Vectorization

Trailer

Truck

Av

Obstacle

People

Car    Car

Av

Av

13.4
15.9
17.3
19.7

25

Av

25

# PhysCov: Approach



SUT

RRS Pipeline

$\tau_k$

AV

Reduction to Reachable Set → Reduction to Sensed Reachable Set → Sensed Reachable Set Vectorization

$f$

Physical Coverage

Test Suite Coverage

[ 5.  5.  15.  15.  35.  35.  25.  35.  25.  5.]

# PhysCov: Approach

$$\alpha = \{(e_1^{sen}, s_1), \dots \}$$

[ 5.  5.  15.  15.  35.  35.  25.  35.  25.  5.]



$$PhysCov = \frac{\alpha}{\beta}$$



$$\beta = E \times S$$

Start: [0, 0, 0, ... 0]

X = Max size reachable set

End: [X, X, X, ..., X]

# PhysCov: Approach

We couldn't cover all the details and we encourage you to read the paper!

# Study

We asked three different research questions:

**RQ1)** How effective RRS at grouping equivalent environment inputs such that they cause similar behaviors?

**RQ2)** How effective is PhysCov at selecting tests that induce unique failures?

**RQ3)** Can PhysCov distinguish similar from different scenarios?

# Environments

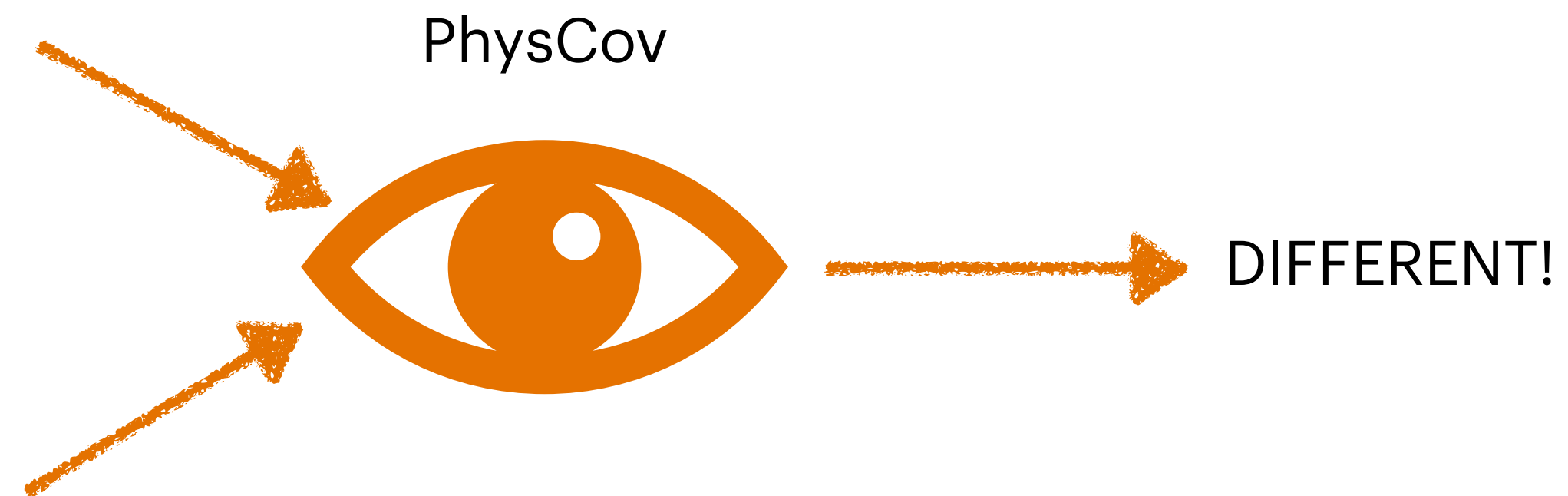HighwayEnv

1,000,000 tests

BeamNG

10,000 tests

Waymo Open Dataset

4 Hours 26 Minutes Driving

# Research Question 1

How effective RRS at grouping equivalent environment inputs such that they cause similar behaviors?



PhysCov

DIFFERENT!
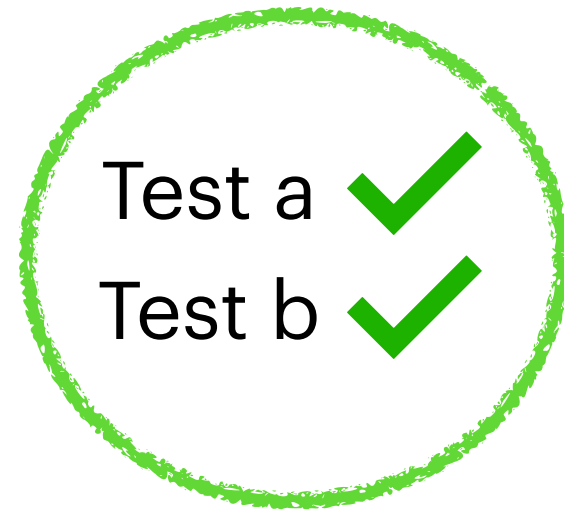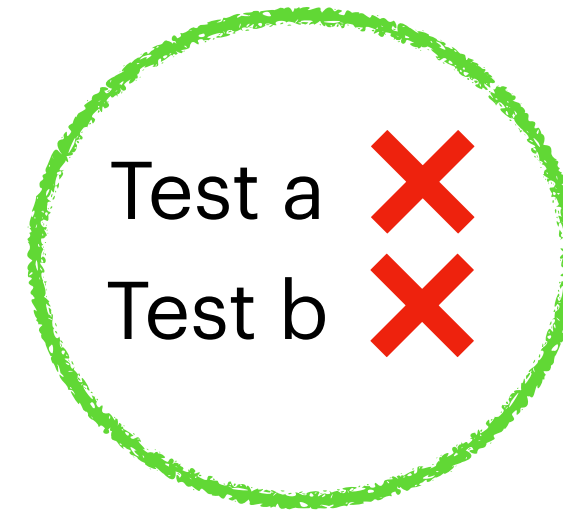
# Research Question 1

How effective RRS at grouping equivalent environment inputs such that they cause similar behaviors?

# Research Question 1

How effective RRS at grouping equivalent environment inputs such that they cause similar behaviors?
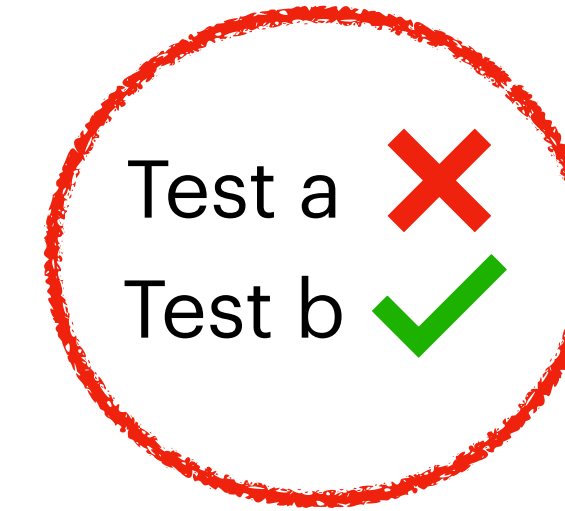
## Structural Code Coverage

- Line Coverage
- Branch Coverage
- Intraprocedural prime path coverage
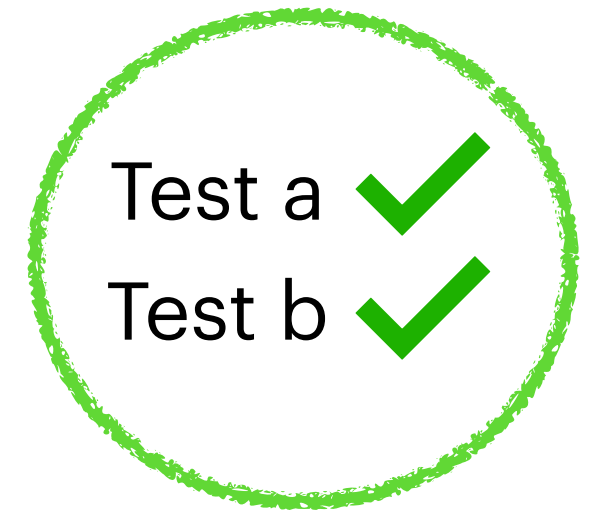- Intraprocedural path coverage
- Absolute path coverage

10,000 tests →
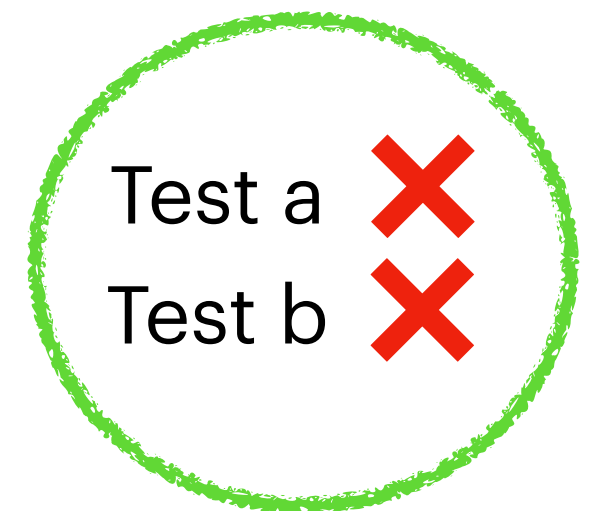
## Trajectory Coverage

- Improved to include irregular maps

## PhysCov

- $\Psi_1$ - RRS of length 1
- $\Psi_5$ - RRS of length 5
- $\Psi_{10}$ - RRS of length 10

Class 1

Test a ✓
Test b ✓

Class 2

Test a ✗
Test b ✗

Class 3

Test a ✗
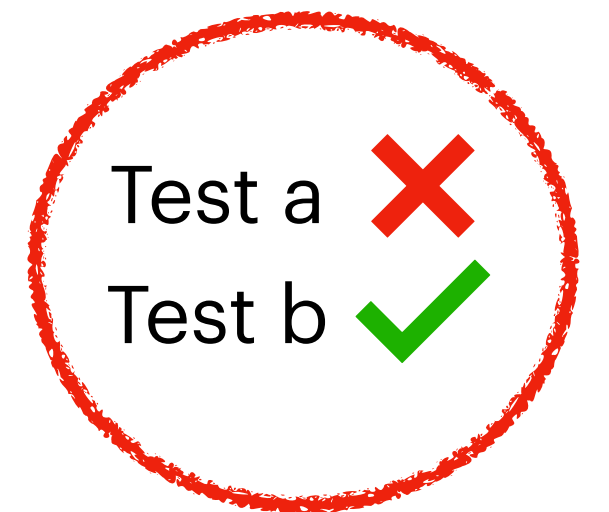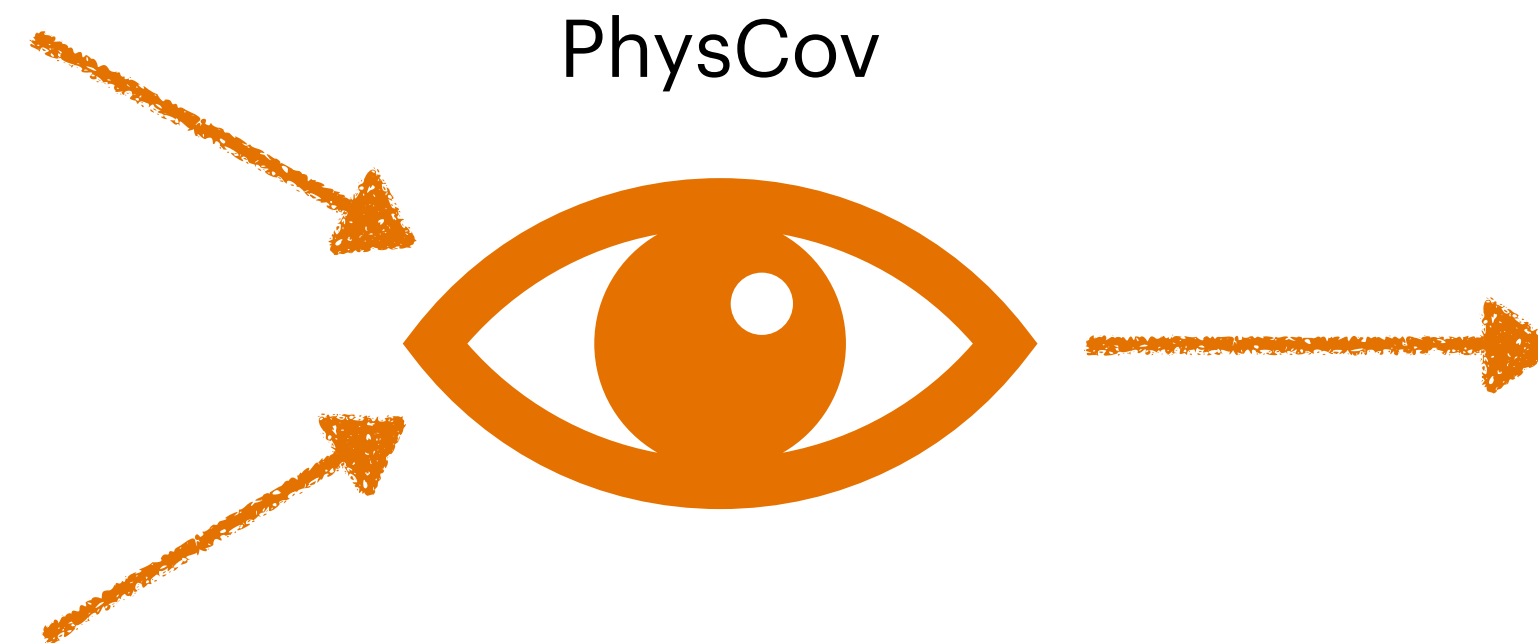Test b ✓

# Research Question 1

How effective RRS at grouping equivalent environment inputs such that they cause similar behaviors?

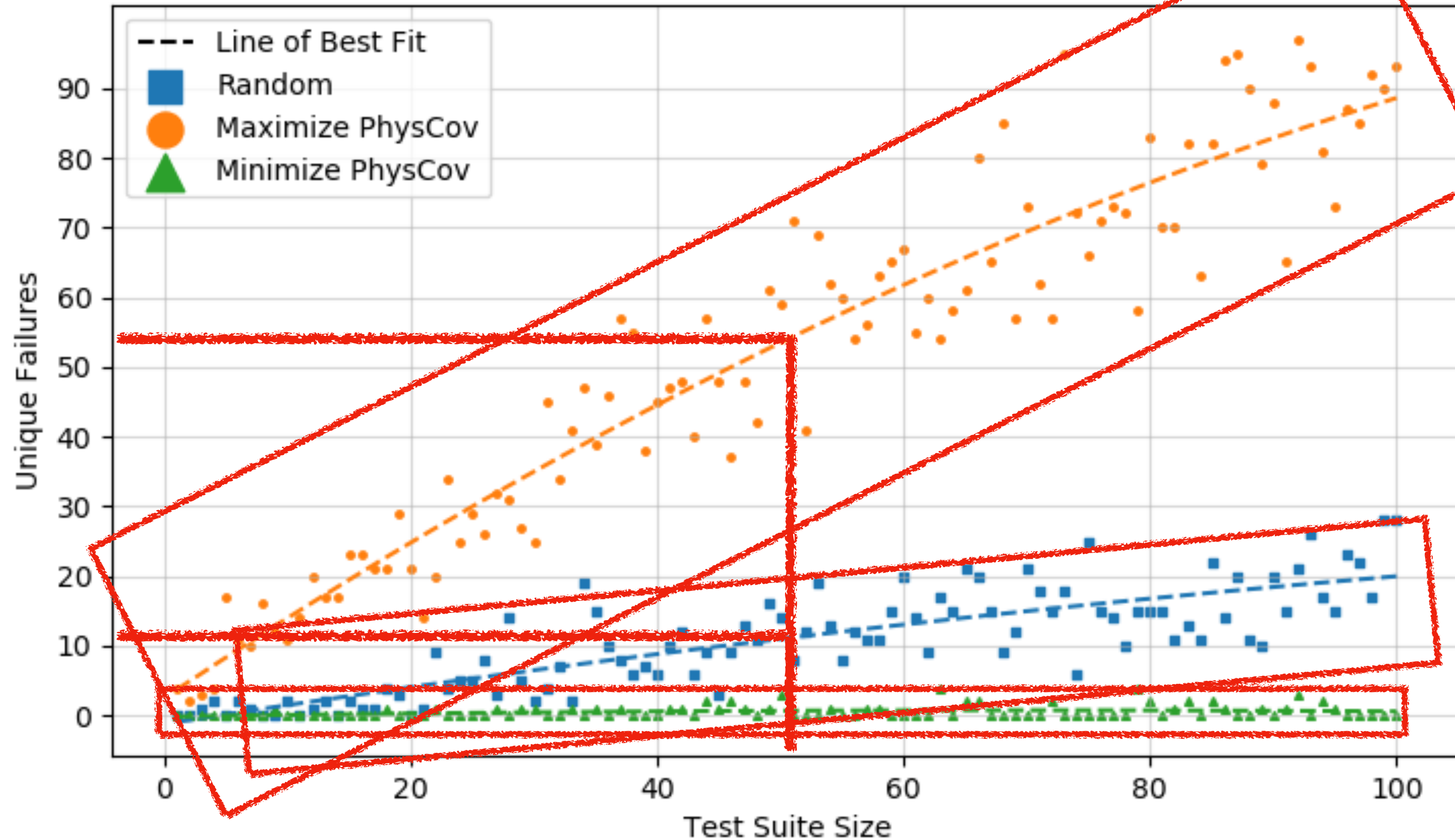| Coverage Metric | Equivalent Classes | Percentage Inconsitency |
|---|---|---|
| Line | 151 | 65% |
| Branch | 146 | 58% |
| Intraprocedural Prime Path Coverage | 421 | 75% |
| Intraprocedural Path Coverage | 10000 | —— |
| Absolute Path Coverage | 10000 | —— |
| Trajectory Coverage | 10000 | —— |
| Physical Coverage: $\Psi_1$ | 682 | 57% |
| Physical Coverage: $\Psi_5$ | 1594 | 40% |
| Physical Coverage: $\Psi_{10}$ | 3628 | 32% |

# Research Question 2

How effective is PhysCov at selecting tests that induce unique failures?



PhysCov

# Research Question 2

How effective is PhysCov at selecting tests that induce unique failures?
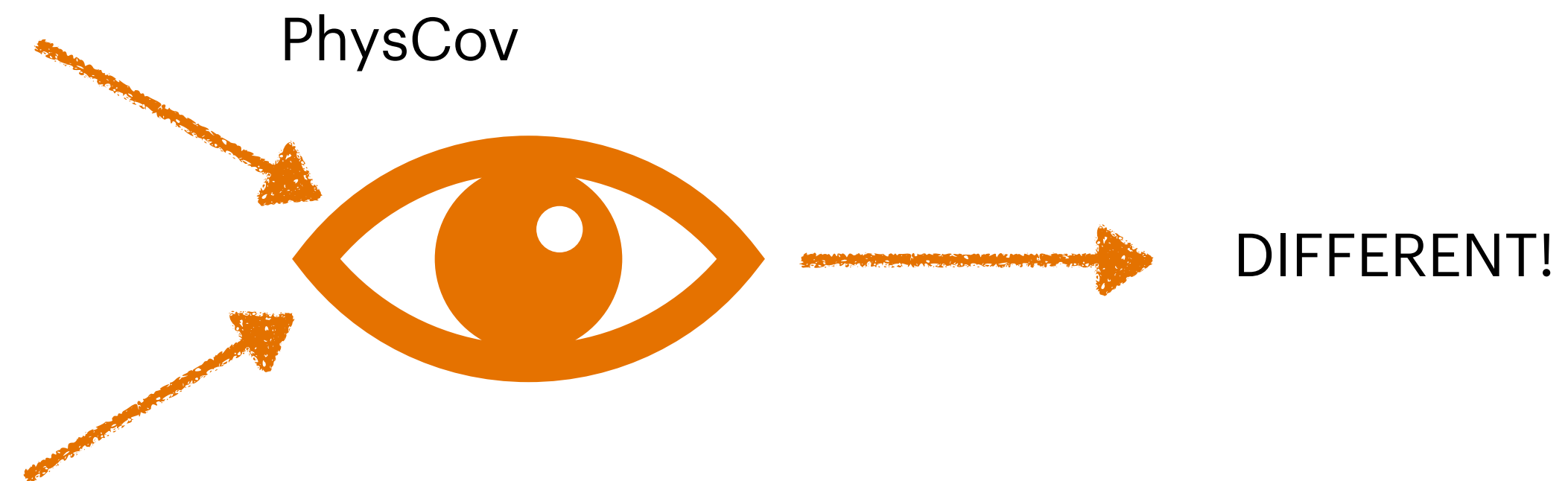
# Research Question 3

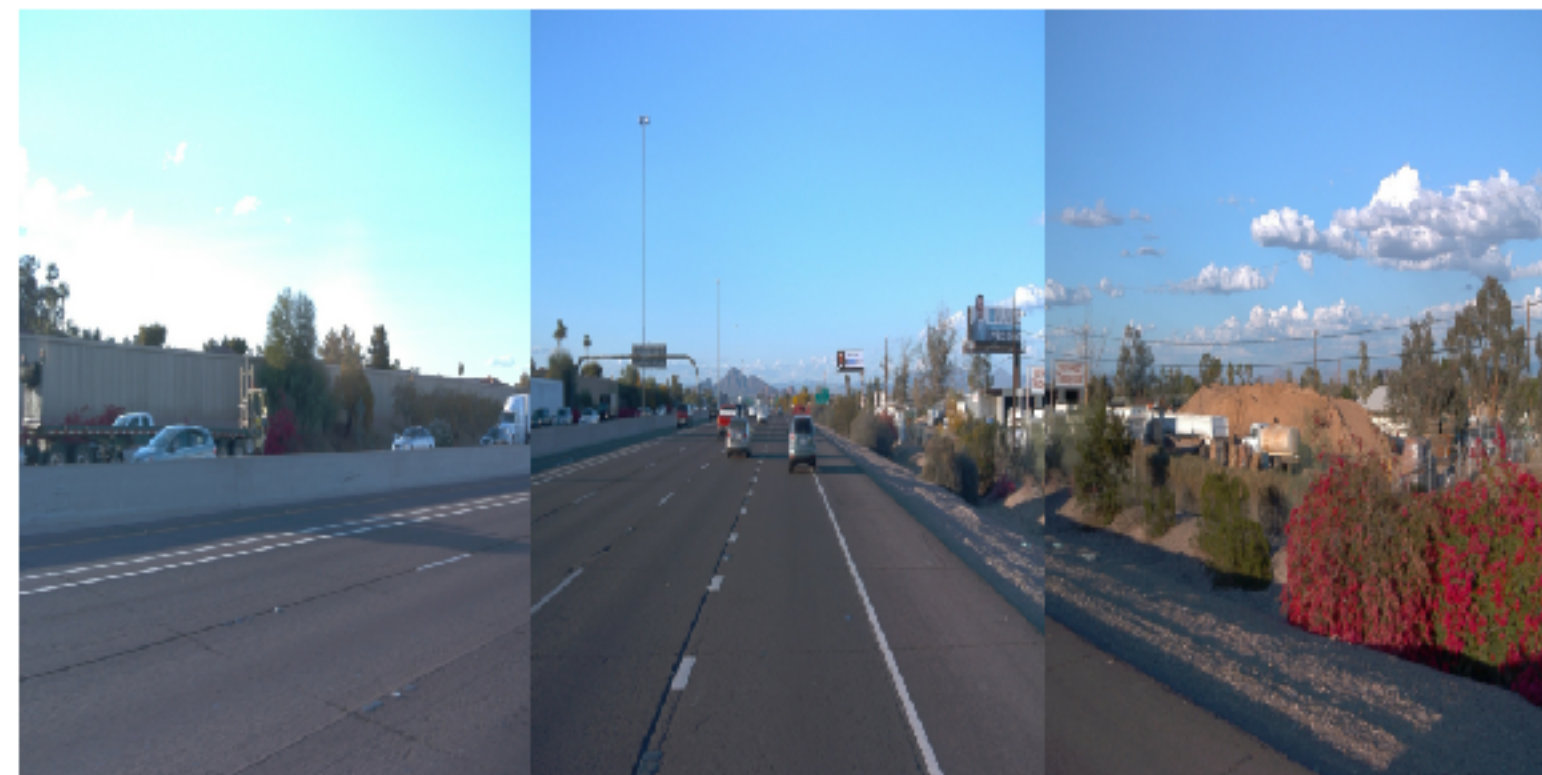Can PhysCov distinguish similar from different scenarios?
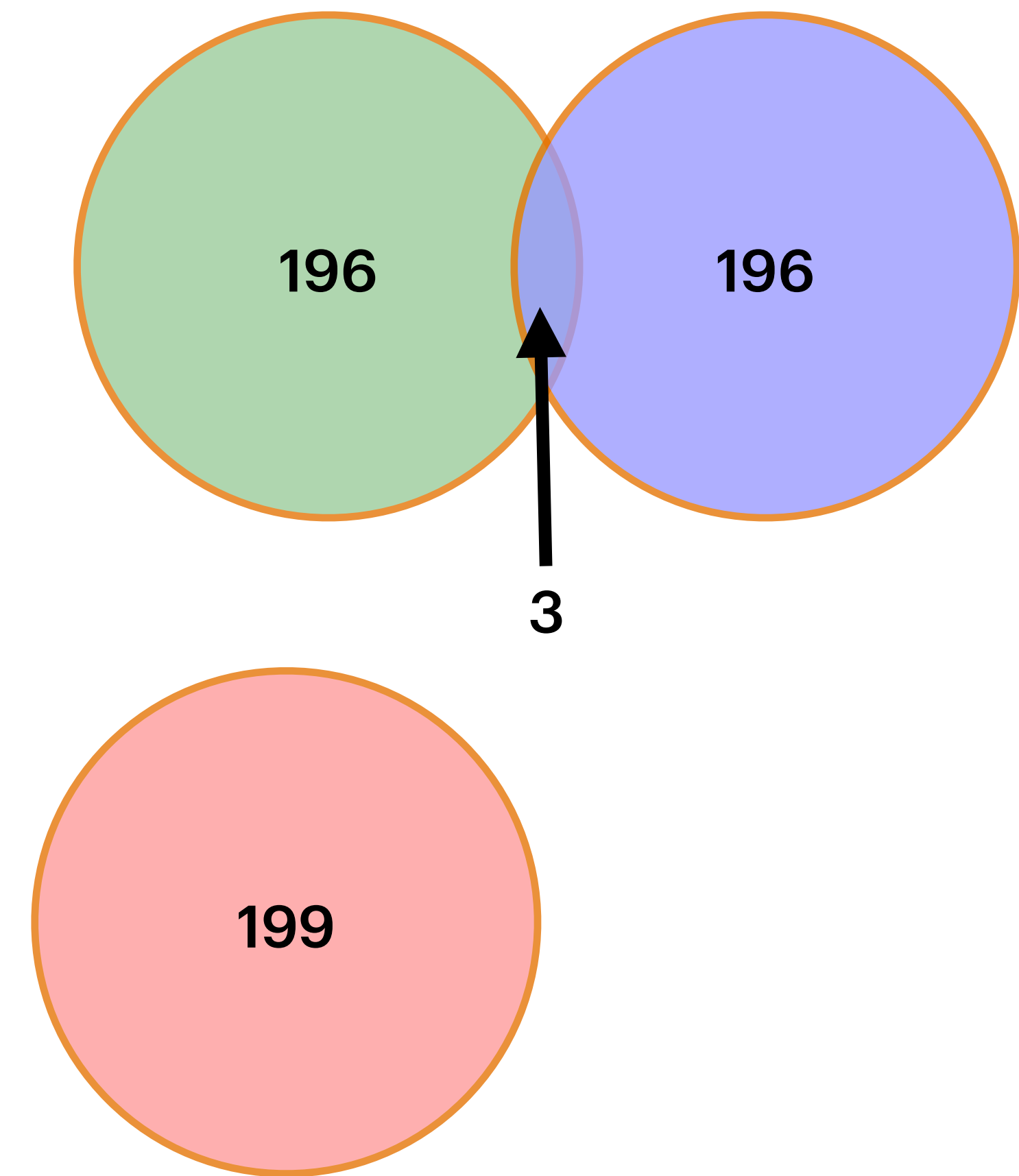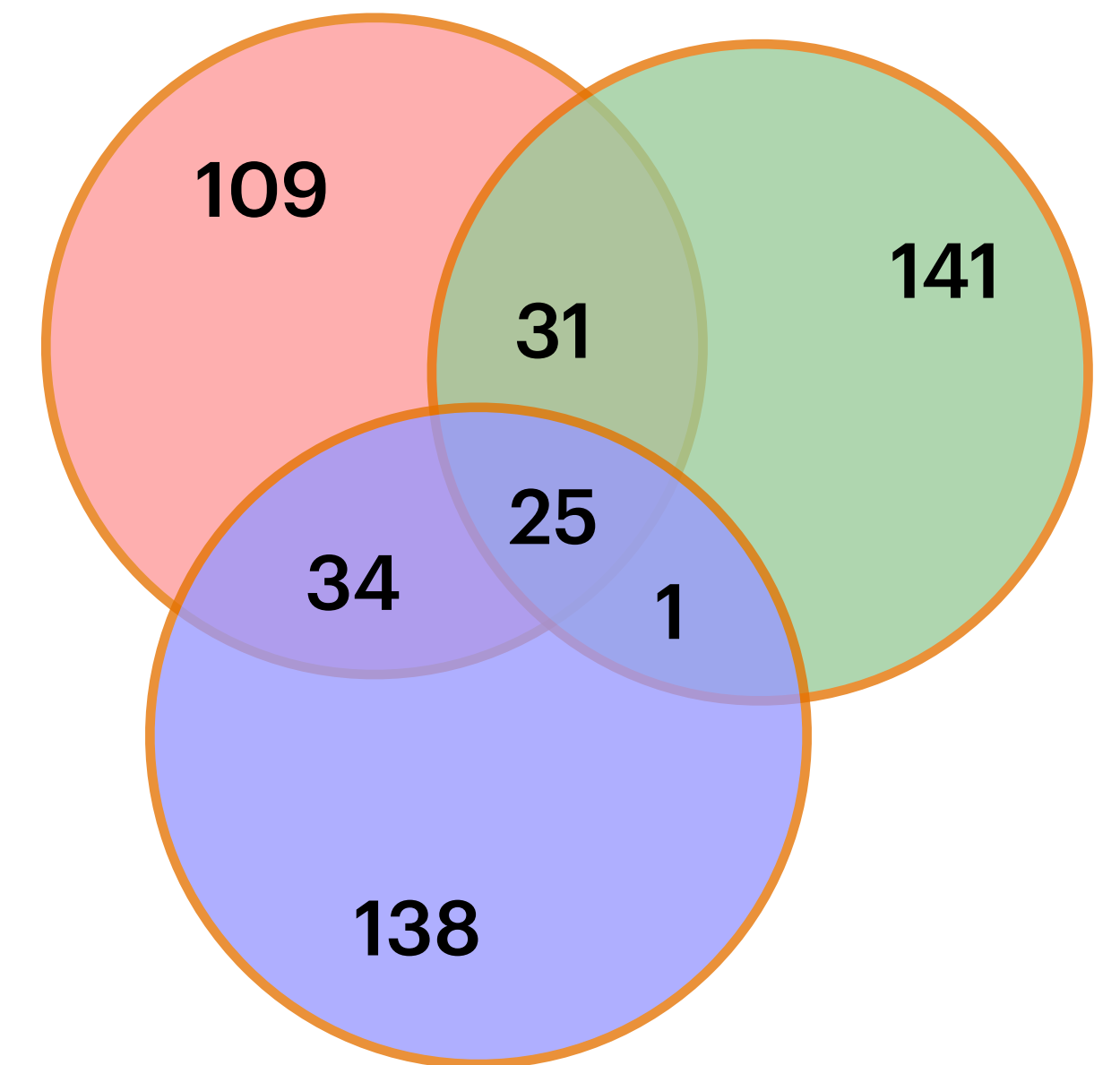


PhysCov

DIFFERENT!

# Research Question 3

Can PhysCov distinguish similar from different scenarios?

# Research Question 3

Can PhysCov distinguish similar from different scenarios?

# Research Question 3

Can PhysCov distinguish similar from different scenarios?

PhysCov

DIFFERENT!

# Research Question 3

Can PhysCov distinguish similar from different scenarios?

# Research Question 3

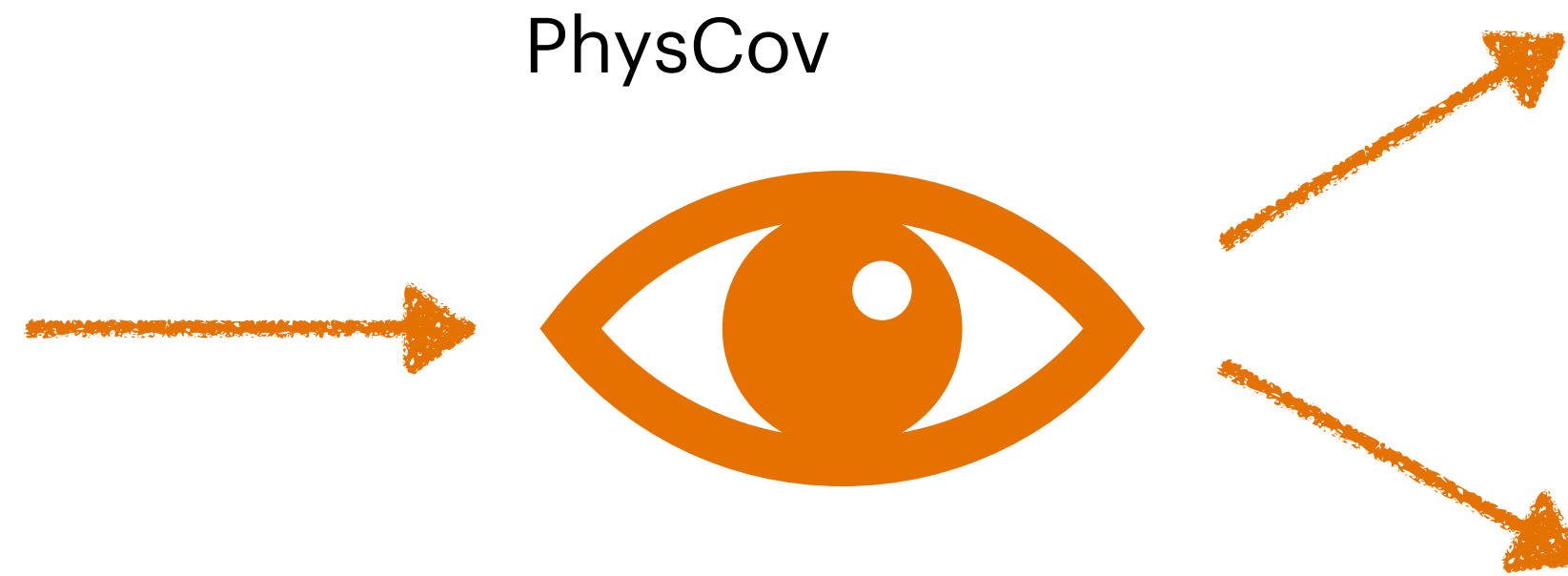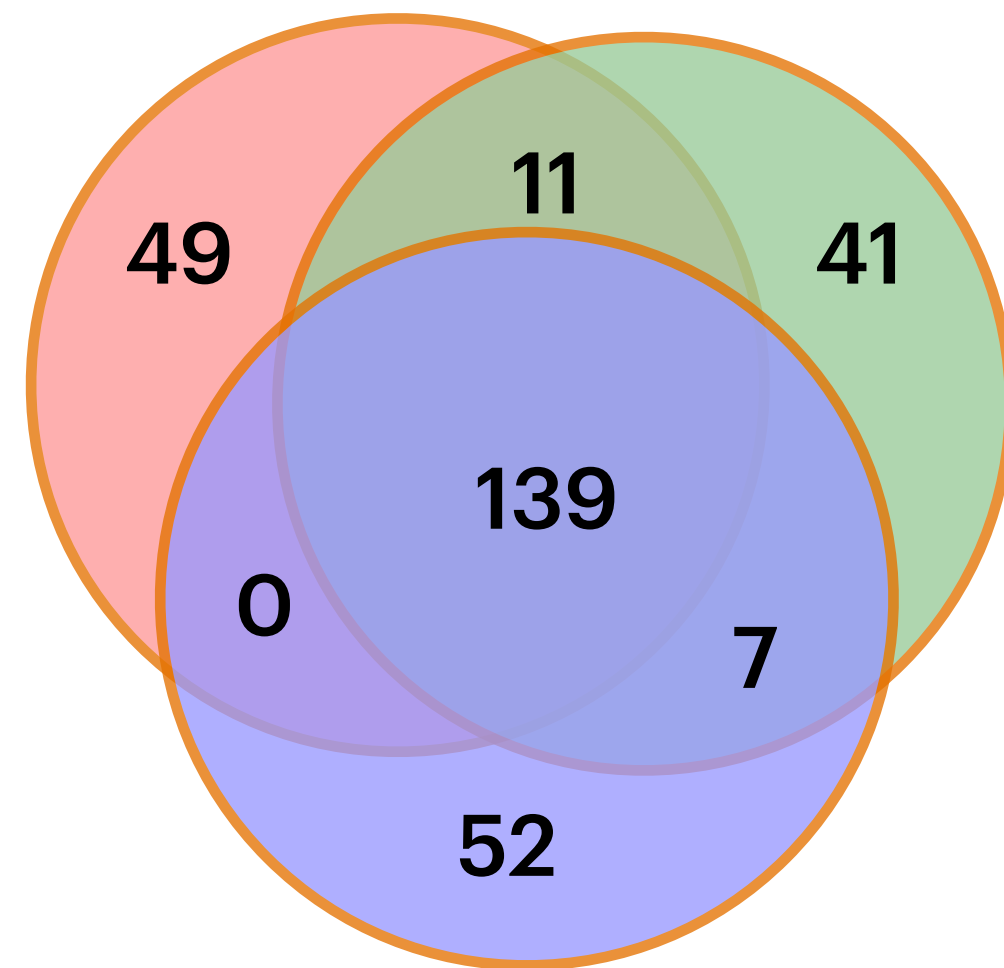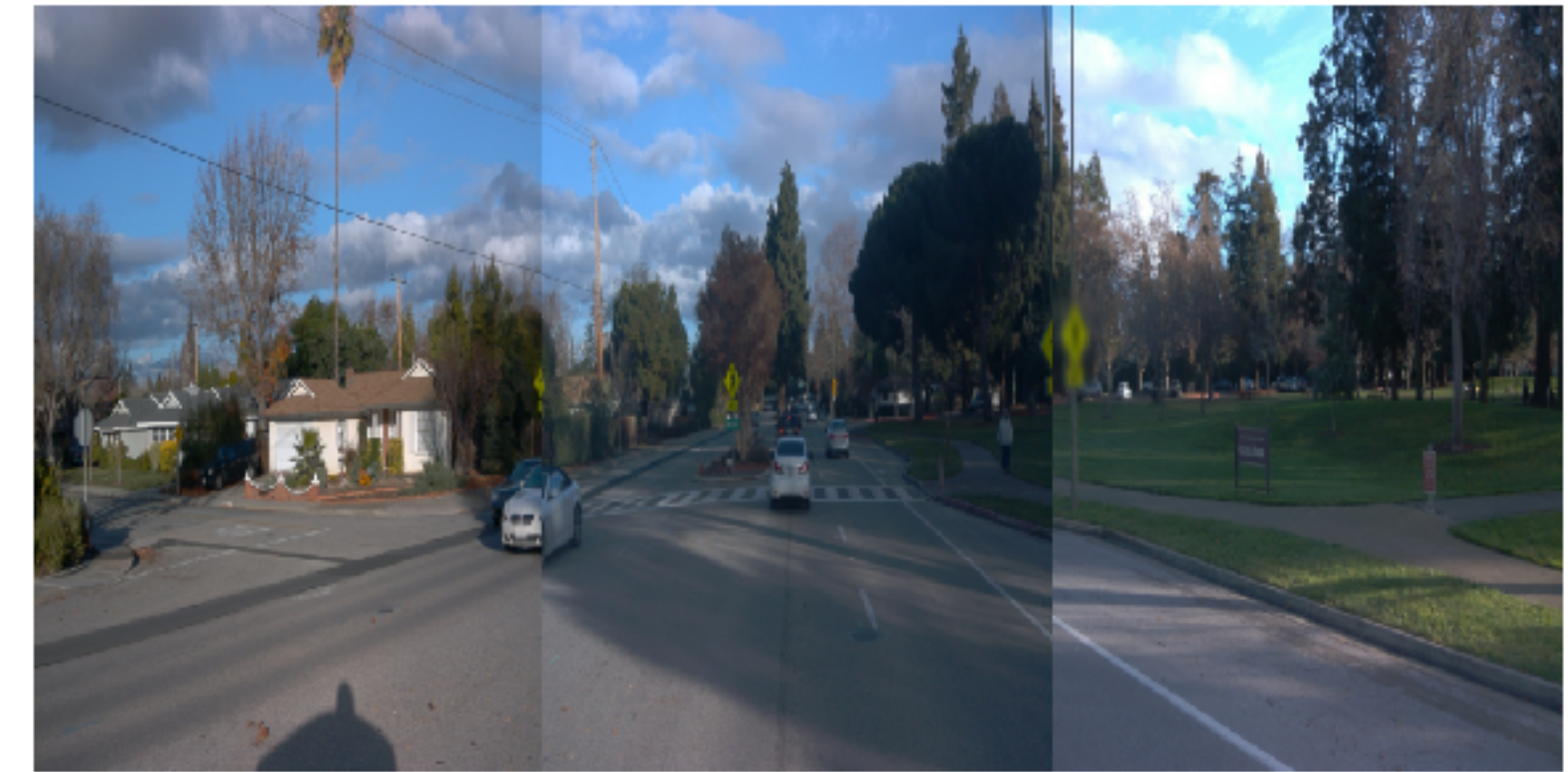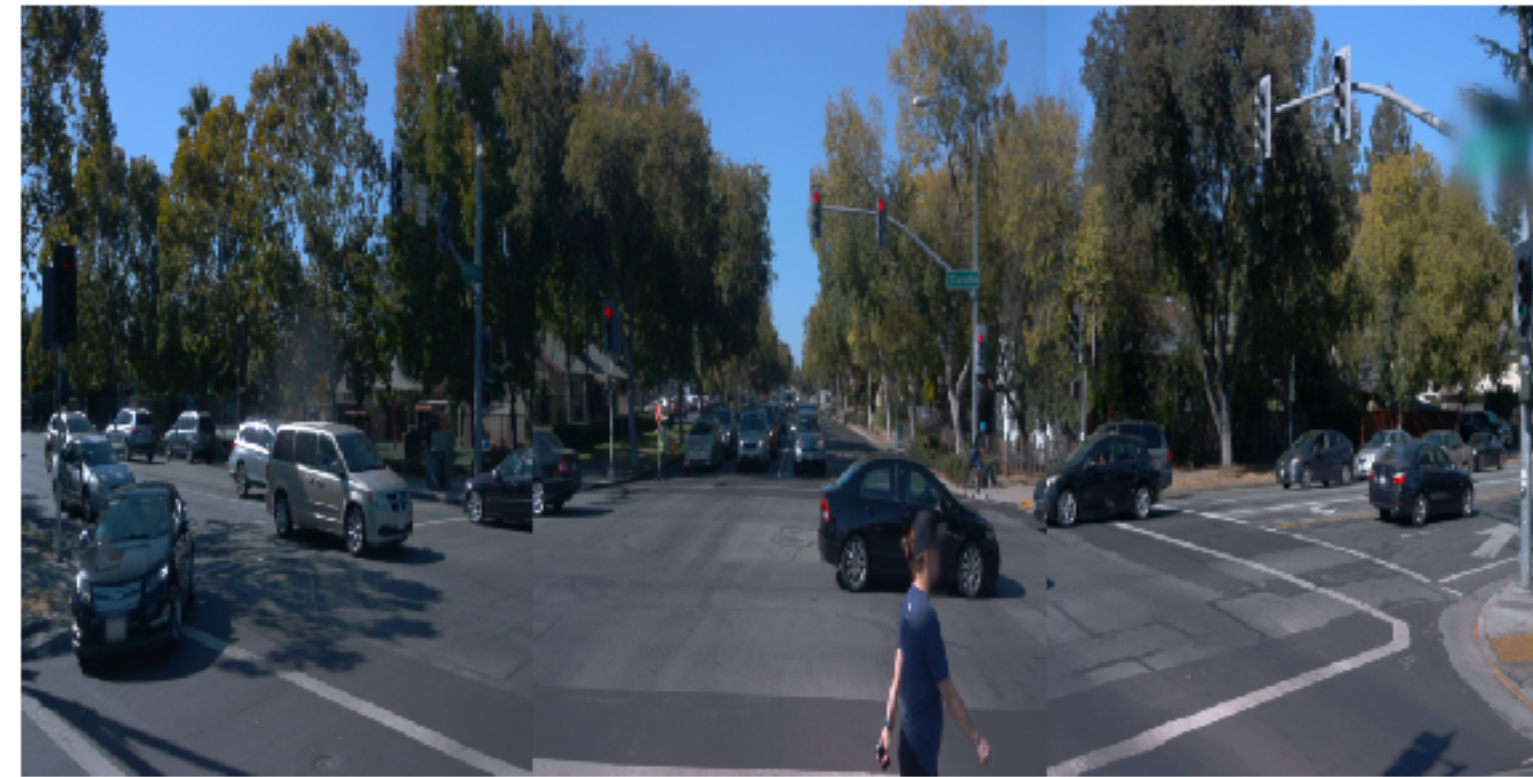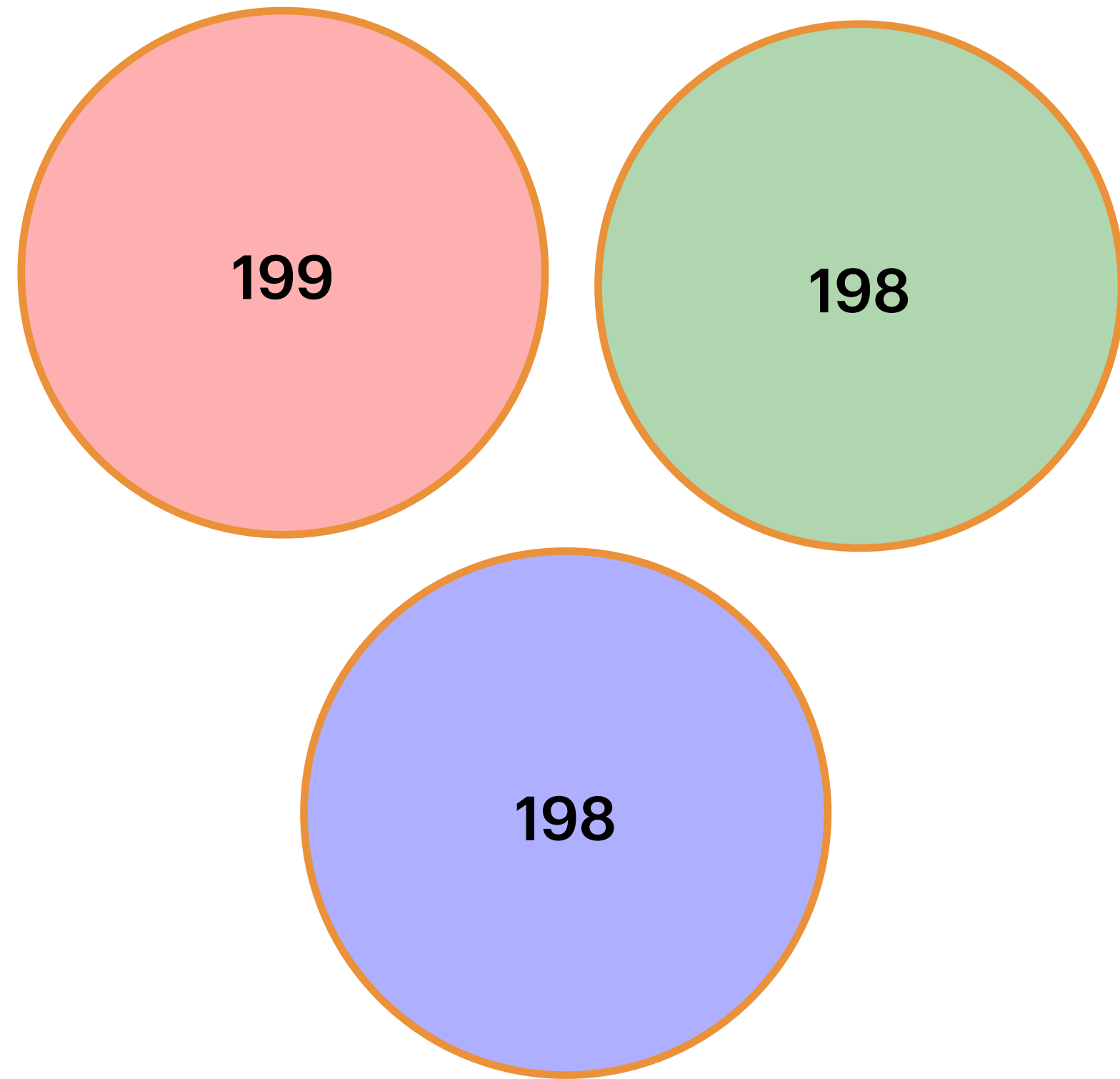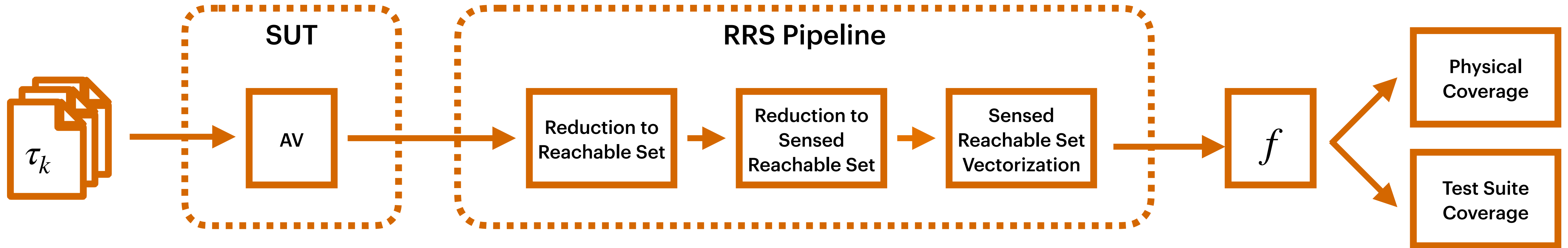Can PhysCov distinguish similar from different scenarios?

# Conclusion



| Coverage Metric | Equivalent Classes | Percentage |
|---|---|---|
| Line | 151 | 65% |
| Branch | 146 | 58% |
| Intraprocedural Prime Path Coverage | 421 | 75% |
| Intraprocedural Path Coverage | 10000 | —— |
| Absolute Path Coverage | 10000 | —— |
| Trajectory Coverage | 10000 | —— |
| Physical Coverage: | 682 | 57% |
| Physical Coverage: | 1594 | 40% |
| Physical Coverage: | 3628 | 32% |