



Prof. Philip Koopman

Carnegie  
Mellon  
University

# Autonomous Vehicle Safety and Perception Robustness Testing

Australian System Safety Conference  
May 23, 2019



@PhilKoopman



## ■ Making safe robots

- Doer/Checker safety

## ■ Edge cases matter

- Robust perception matters

## ■ The heavy tail distribution

- Fixing stuff you see in testing isn't enough

## ■ Perception stress testing

- Finding the weaknesses in perception

## ■ UL 4600: autonomy safety standard



[General Motors]



TRIP COMPLETE !!!  
2797/2849 miles (98.2%)

## ■ Washington DC to San Diego

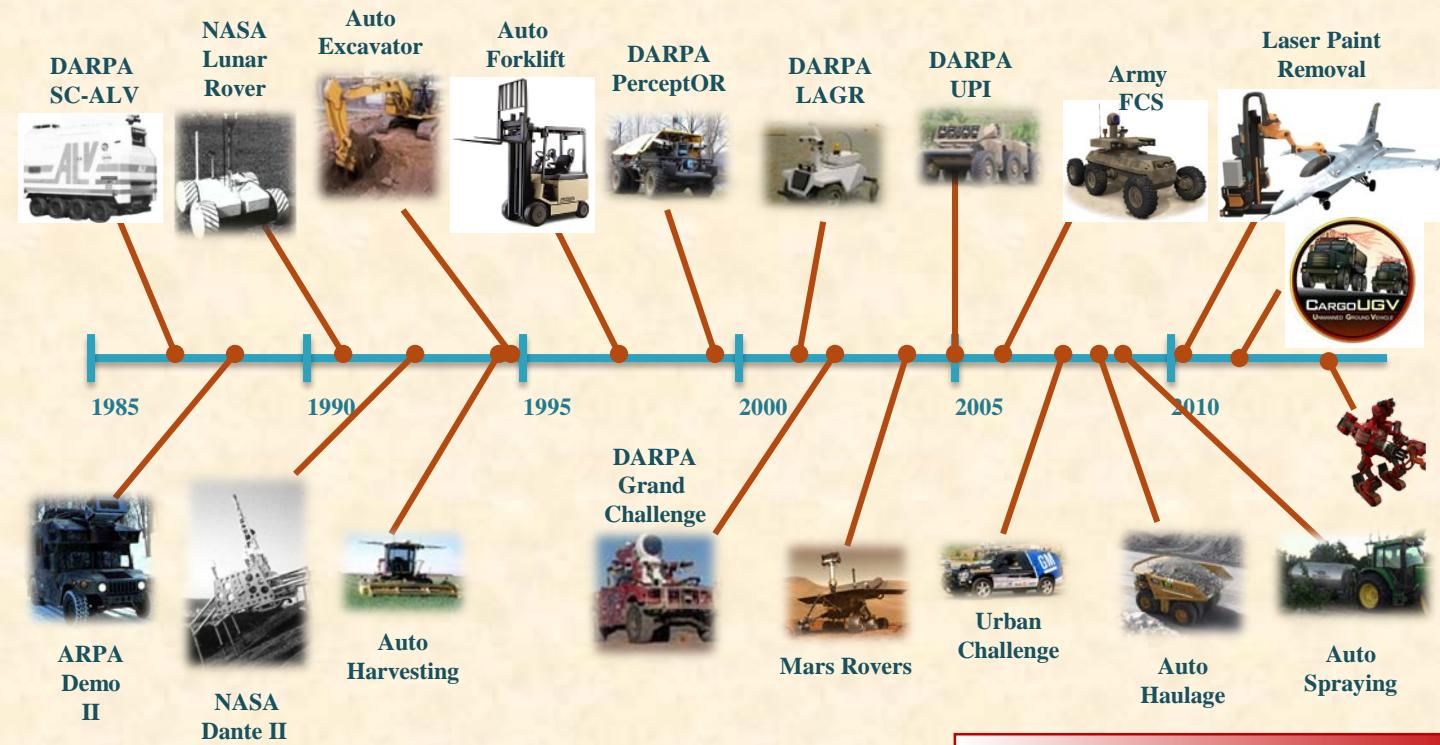
- CMU Navlab 5
- Dean Pomerleau
- Todd Jochem

[https://www.cs.cmu.edu/~tjochem/nhaa/nhaa\\_home\\_page.html](https://www.cs.cmu.edu/~tjochem/nhaa/nhaa_home_page.html)

## ■ AHS San Diego demo Aug 1997



# NREC: 30+ Years Of Cool Robots



Carnegie Mellon University Faculty, staff, students  
Off-campus Robotics Institute facility

**Software  
Safety**

# Before Autonomy Software Safety

## ■ The Big Red Button era



# Traditional Validation Vs. Machine Learning

- Use traditional software safety where you can

..BUT..

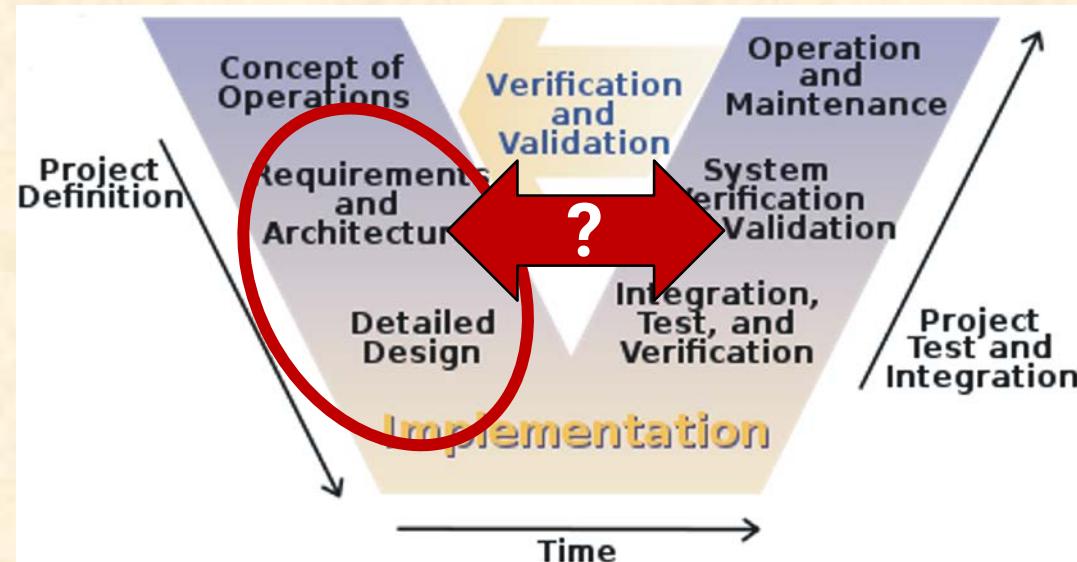
- Machine Learning (inductive training)

- **No requirements**

- Training data is difficult to validate

- **No design insight**

- Generally inscrutable; prone to gaming and brittleness

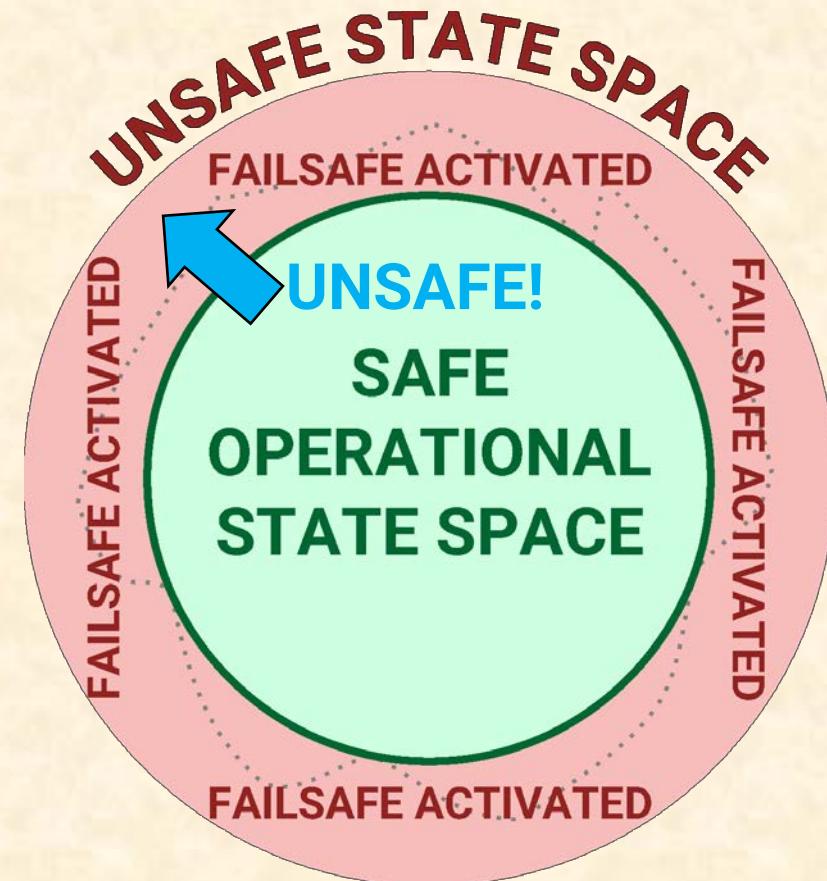


# APD (Autonomous Platform Demonstrator)



# Safety Envelope Approach to ML Deployment

- Specify unsafe regions
- Specify safe regions
  - Under-approximate to simplify
- Trigger system safety response upon transition to unsafe region



# Architecting A Safety Envelope System

## ■ “Doer” subsystem

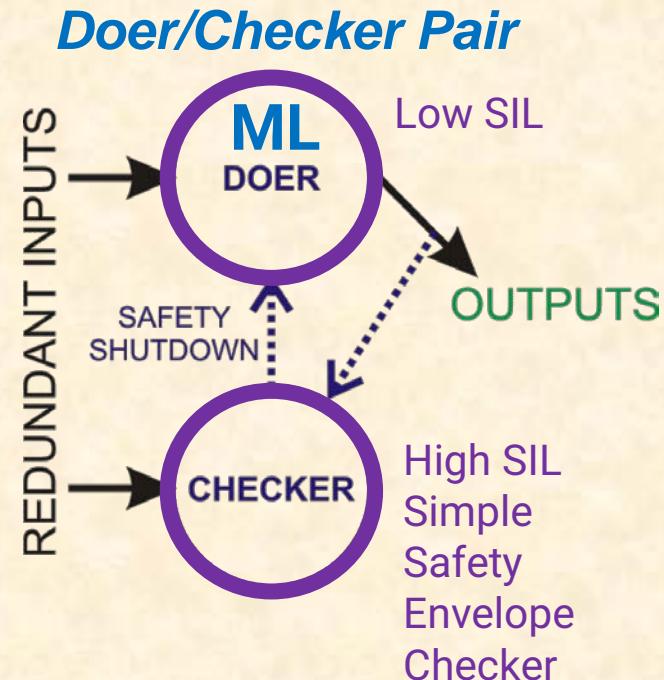
- Implements normal, untrusted functionality

## ■ “Checker” subsystem – Traditional SW

- Implements failsafes (safety functions)

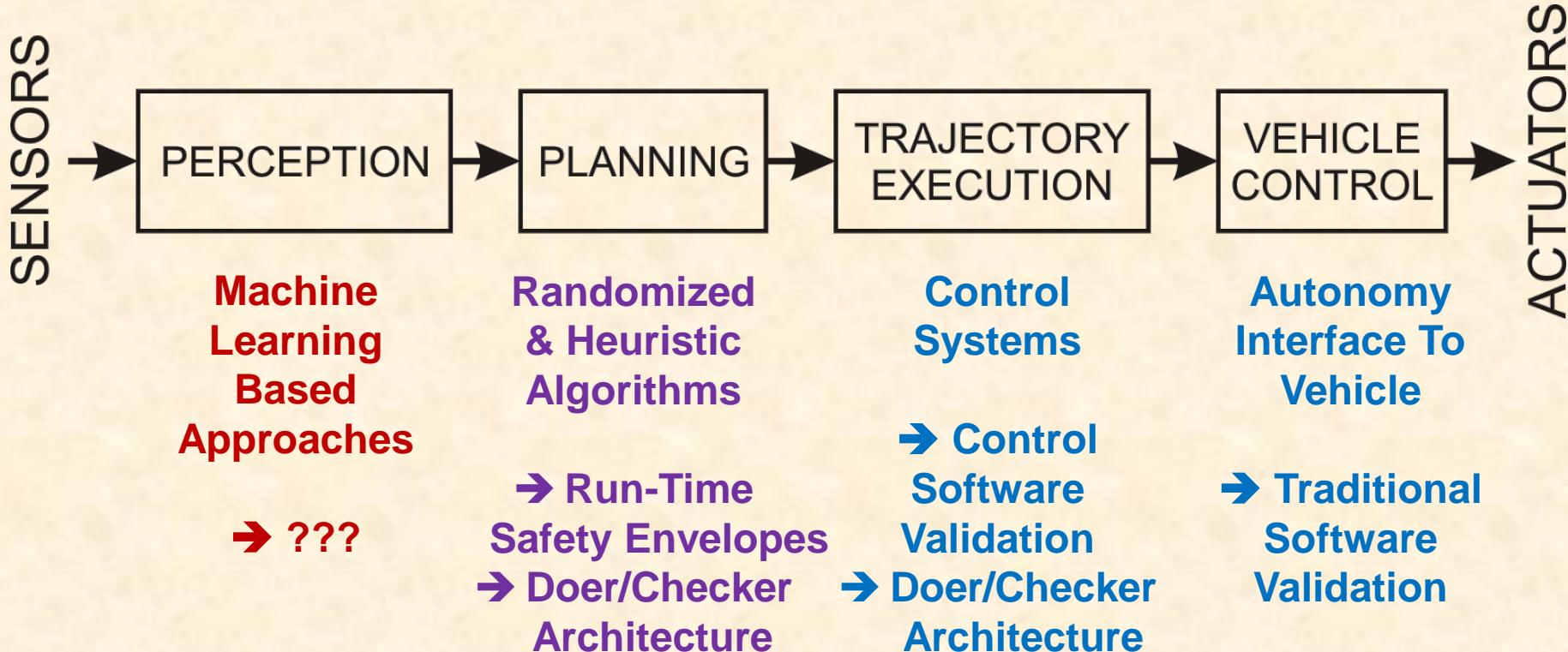
## ■ Checker entirely responsible for safety

- Doer can be at low Safety Integrity Level
- Checker must be at higher SIL



(Also known as a “safety bag” approach;  
also monitor/actuator pair)

# Validating an Autonomous Vehicle Pipeline



Perception presents a uniquely difficult assurance challenge

# Brute Force AV Validation: Public Road Testing

- Good for identifying “easy” cases
  - Expensive and potentially **dangerous**



# Validation Via Brute Force Road Testing?

## ■ If 100M miles/critical mishap...

- Test 3x–10x longer than mishap rate  
→ Need 1 Billion miles of testing

## ■ That's ~25 round trips on every road in the world

- With fewer than 10 critical mishaps

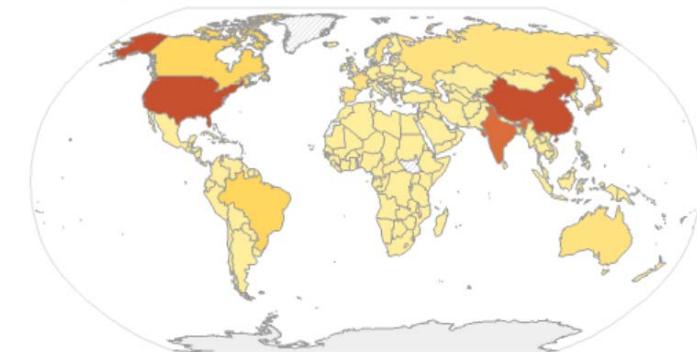
...

WolframAlpha computational knowledge engine.

miles of roads	
Summary:	
total	20.46 million mi
median	11 630 mi
highest	4.03 million mi (United States)
lowest	4.97 mi (Tuvalu)

(1994 to 2008)  
(based on 225 values; 24 unavailable)

Total road length map:



□ (no data available)  
□ 0  
□ 4 to 360 000  
□ 360 000 to 720 000  
□ 720 000 to 1.1 million  
□ 1.1 million to 1.4 million  
□ 1.4 million to 1.8 million  
□ 1.8 million to 2.1 million  
■ > 2.1 million  
(in miles)

# Closed Course Testing

## ■ Safer, but expensive

- Not scalable
- Only tests things you have thought of!



# Simulation

## ■ Highly scalable; less expensive

- Scalable; need to manage fidelity vs. cost
- Only tests things you have thought of!



# What About Edge Cases?

## ■ Gaps in training data can lead to perception failure

- Safety needs to know: “Is that a person?”
- Machine learning provides: “Is that thing like the people in my training data?”



<http://bit.ly/2ln4rzj>

PREDICTED CONCEPT	PROBABILITY
bird	0.997
no person	0.990
one	0.975
feather	0.970
nature	0.963
poultry	0.954
outdoors	0.936
color	0.910
animal	0.908

<https://www.clarifai.com/demo>

## ■ Edge Case are surprises

- You won't see these in training or testing
- Edge cases are the stuff you didn't think of!

# Need An Edge Case “Zoo”

- Novel objects (missing from zoo) are triggering events



# Why Edge Cases Matter

■ Where will you be after 1 Billion miles of validation testing?

■ Assume 1 Million miles between unsafe “surprises”

- Example #1:

**100 “surprises” @ 100M miles / surprise**

- All surprises seen about 10 times during testing
- With luck, all bugs are fixed

- Example #2:

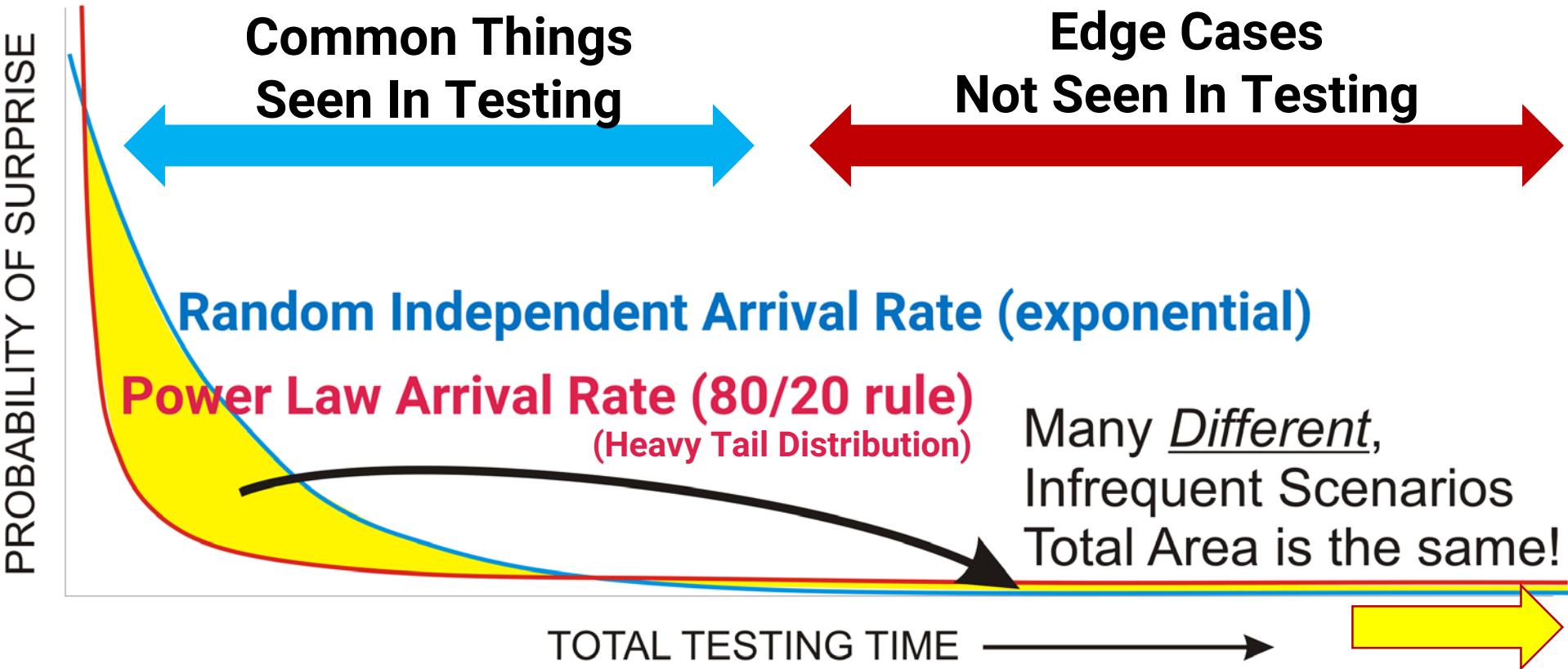
**100,000 “surprises” @ 100B miles / surprise**

- Only 1% of surprises seen during 1B mile testing
- Bug fixes give no real improvement (1.01M miles / surprise)



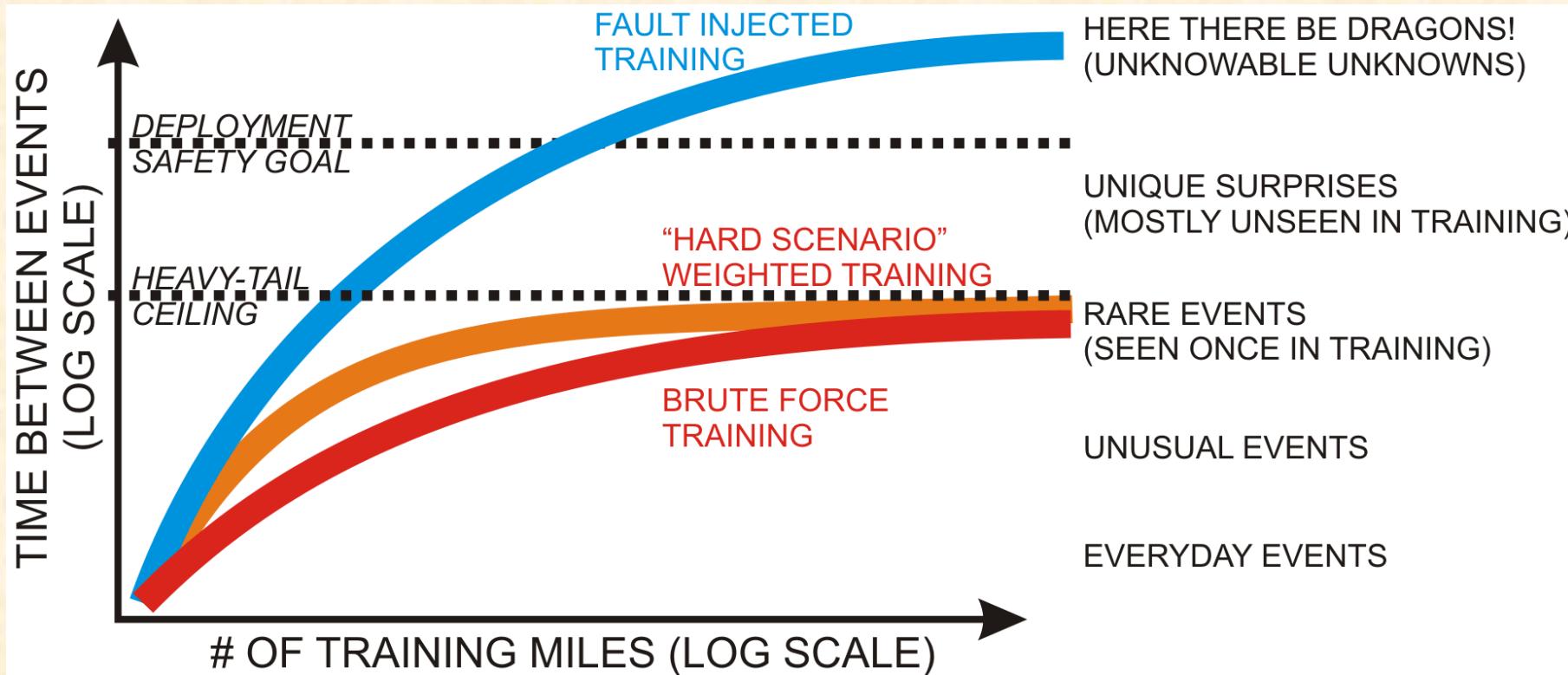
<https://goo.gl/3dzguf>

# Real World: Heavy Tail Distribution(?)



# The Heavy Tail Testing Ceiling

- Need to find “Triggering Events” to inject into sims/testing



# Edge Cases Pt. 1: Triggering Event Zoo

## ■ Need to collect surprises

- Novel objects
- Novel operational conditions



<https://goo.gl/Ni9HhU>

## ■ Issue:

**novel for person  $\neq$  novel for  
Machine Learning**

- ML can have “edges” in unexpected places
- ML might train on features that seem irrelevant to people

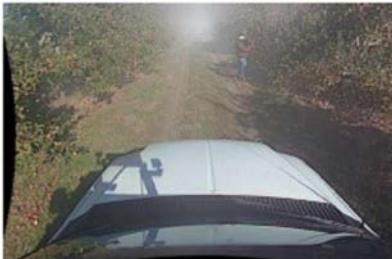


# Edge Cases Part 2: Brittleness

## ■ Sensor data corruption experiments



$u_f = 1\text{m}$ ,  $\kappa = 2$   
Defocus



$u_V = 97.8\text{m}$   
Haze

Contextual Mutators

Defocus & haze are  
a significant issue

Exploring the response of a DNN to environmental perturbations from “Robustness Testing for Perception Systems,” RIOT Project, NREC, DIST-A.

## Synthetic Equipment Faults



Correct detection



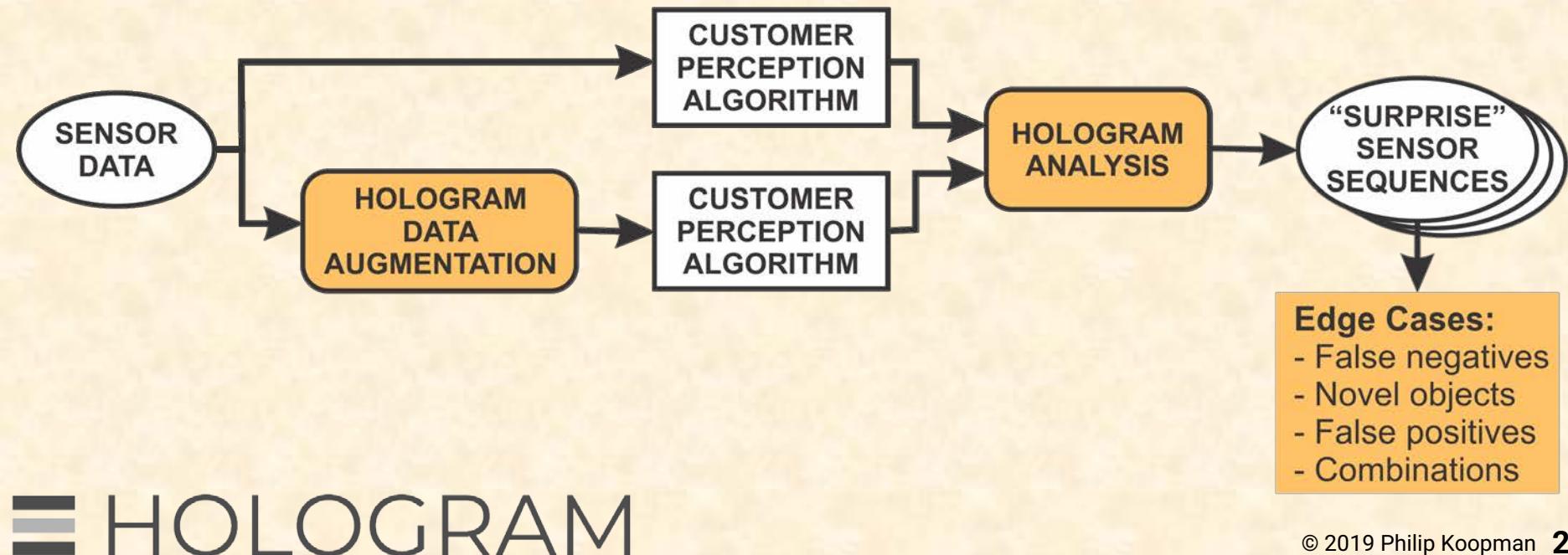
False negative

Gaussian blur

Gaussian Blur &  
Gaussian Noise cause  
similar failures

# Hologram Detects Edge Cases

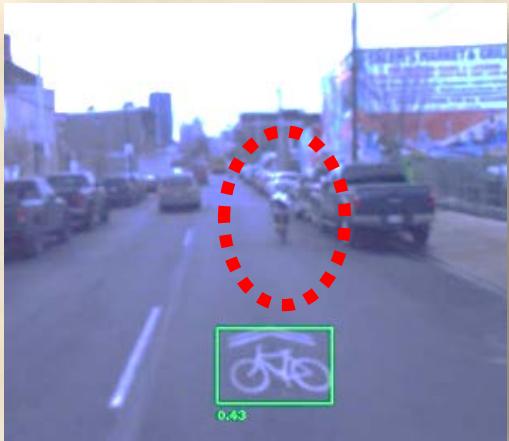
- **Brittle perception behavior indicates Edge Cases**
  - Can uncover false negatives and detect novel objects



# Context-Dependent Perception Failures

## ■ Perception failures are often context-dependent

- False positives and false negatives are both a problem



False positive on lane marking  
False negative real bicyclist



False negative when  
person next to light pole

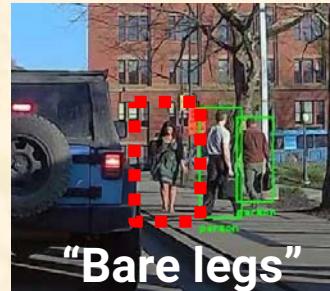


False negative when  
in front of dark vehicle

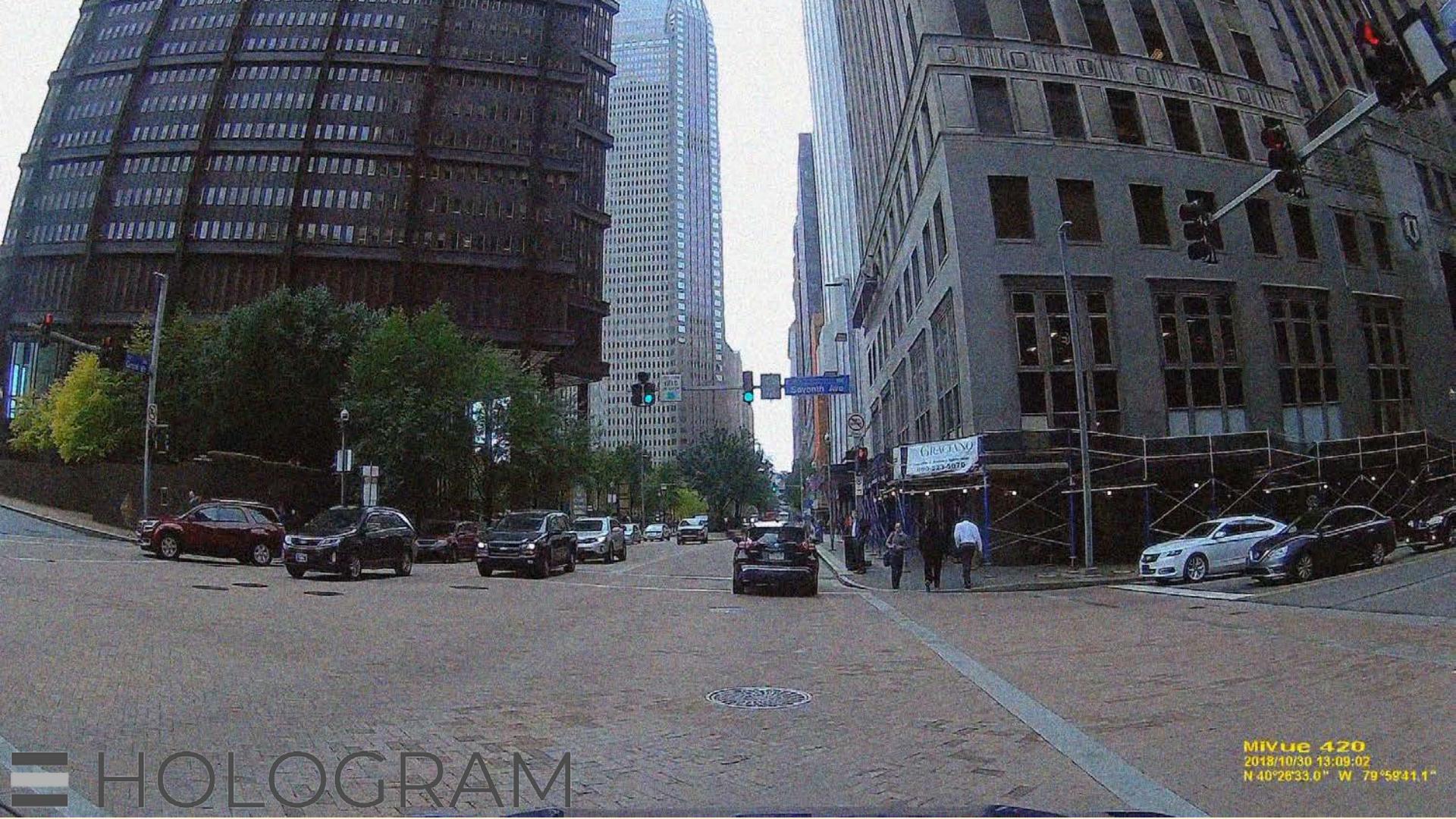
Will this pass a “vision test” for bicyclists?

# Example Triggering Events via Hologram

## ■ Mask-R CNN: examples of systemic problems we found



Notes: These are baseline, un-augmented images // Your mileage may vary.



HOLOGRAM

MIVUE 420  
2018/10/30 13:09:02  
N 40°26'33.0" W 79°59'41.1"



HOLOGRAM

MIVue 420  
2018/10/30 13:18:43  
N 40°26'19.7" W 80°00'8.62"



HOLOGRAM

MIVUE 420  
2018/11/05 12:50:47  
N 40°26'43.5" W 79°56'55.9"

## ■ Drivers do more than just drive

- Occupant behavior, passenger safety
- Detecting and managing equipment faults



## ■ Operational limitations & situations

- System exits Operational Design Domain
- Vehicle fire or catastrophic failure
- Post-crash passenger evacuation



## ■ Interacting with non-drivers

- Pedestrians, passengers
- Police, emergency responders

## ■ Handling updates

- Fully recertify after every weekly update?
- Security in general



## ■ Vehicle maintenance

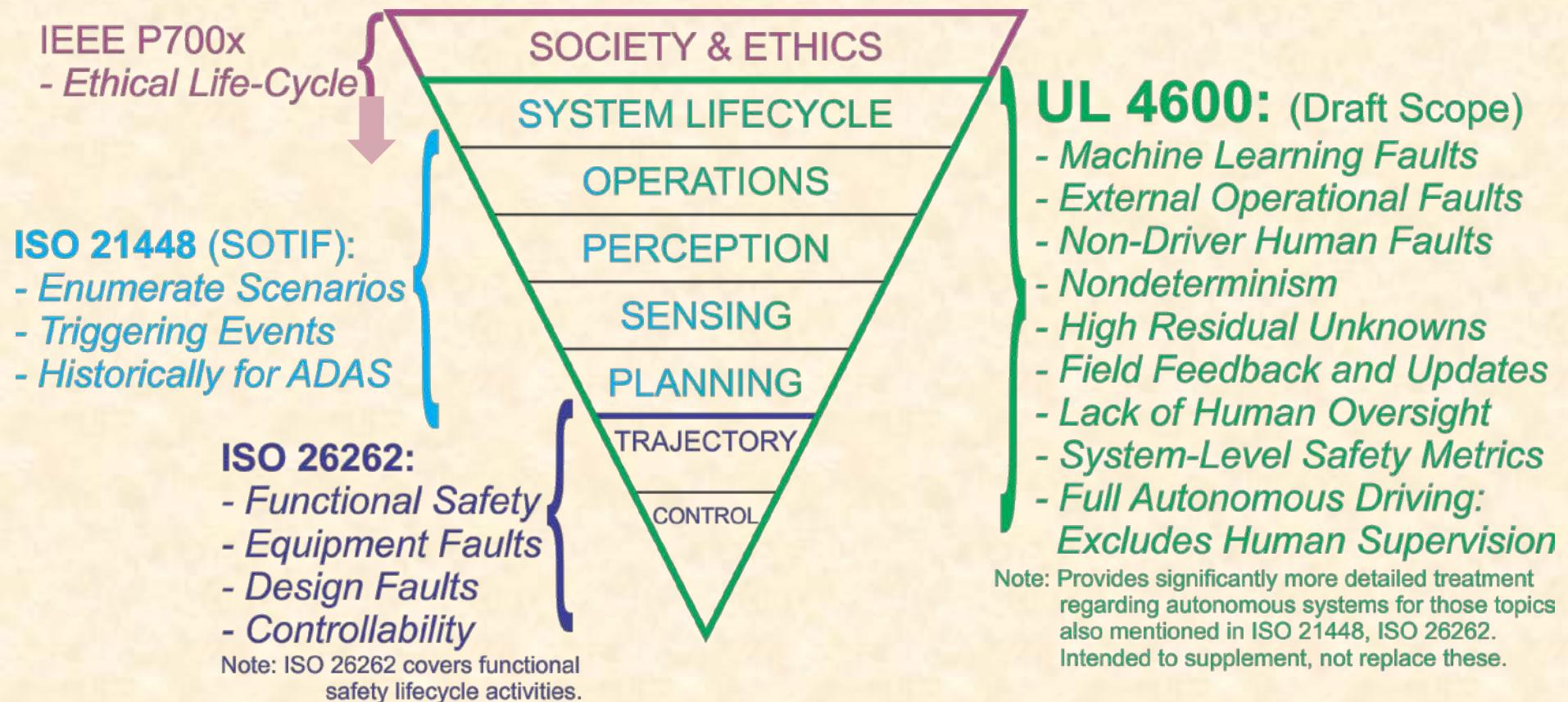
- Pre-flight checks, cleaning
- Corrective maintenance



*Is windshield cleaning fluid life critical?*

- Quality fade
- Supply chain faults

# Safety Standard Landscape



# Ways To Improve AV Safety

## ■ More safety transparency

- Independent safety assessments
- Industry collaboration on safety

## ■ Minimum performance standards

- “Driver test” is necessary – but not sufficient
  - How do you measure maturity?

## ■ Autonomy software safety standards

- ISO 26262/21448 + UL 4600 + IEEE P700x
- Dealing with uncertainty and brittleness



<http://bit.ly/2MTbT8F> (sign modified)



EDGE CASE RESEARCH

MAKING AUTONOMY SAFER  
info@ecr.guru