# Safety of Autonomous Systems

## Stuart Reid  PhD, FBCS
([stureid.test@gmail.com](mailto:stureid.test@gmail.com) / www.stureid.info)

© Stuart Reid  2018

# Scope of the Talk

- **Introduction to Autonomous Systems**
- **Specifying Objectives (Safely)**
- **Online vs Off-Line Machine Learning**
- **Machine Learning Challenges**
- **Black Box Testing**
- **White Box Testing**
- **The Necessity of Virtual Test Environments**
- **Conclusions**

# *Introduction to Autonomous Systems*
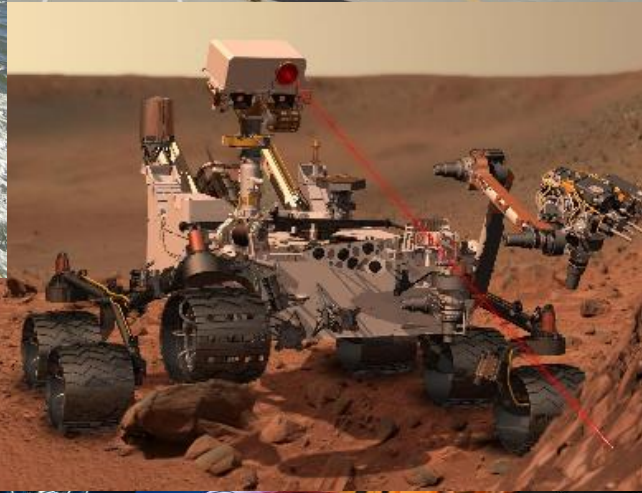
# Definition used for Autonomous System

- **Autonomy**
  - the capacity to make an informed, un-coerced decision
  - Autonomous organizations or institutions are independent or self-governing
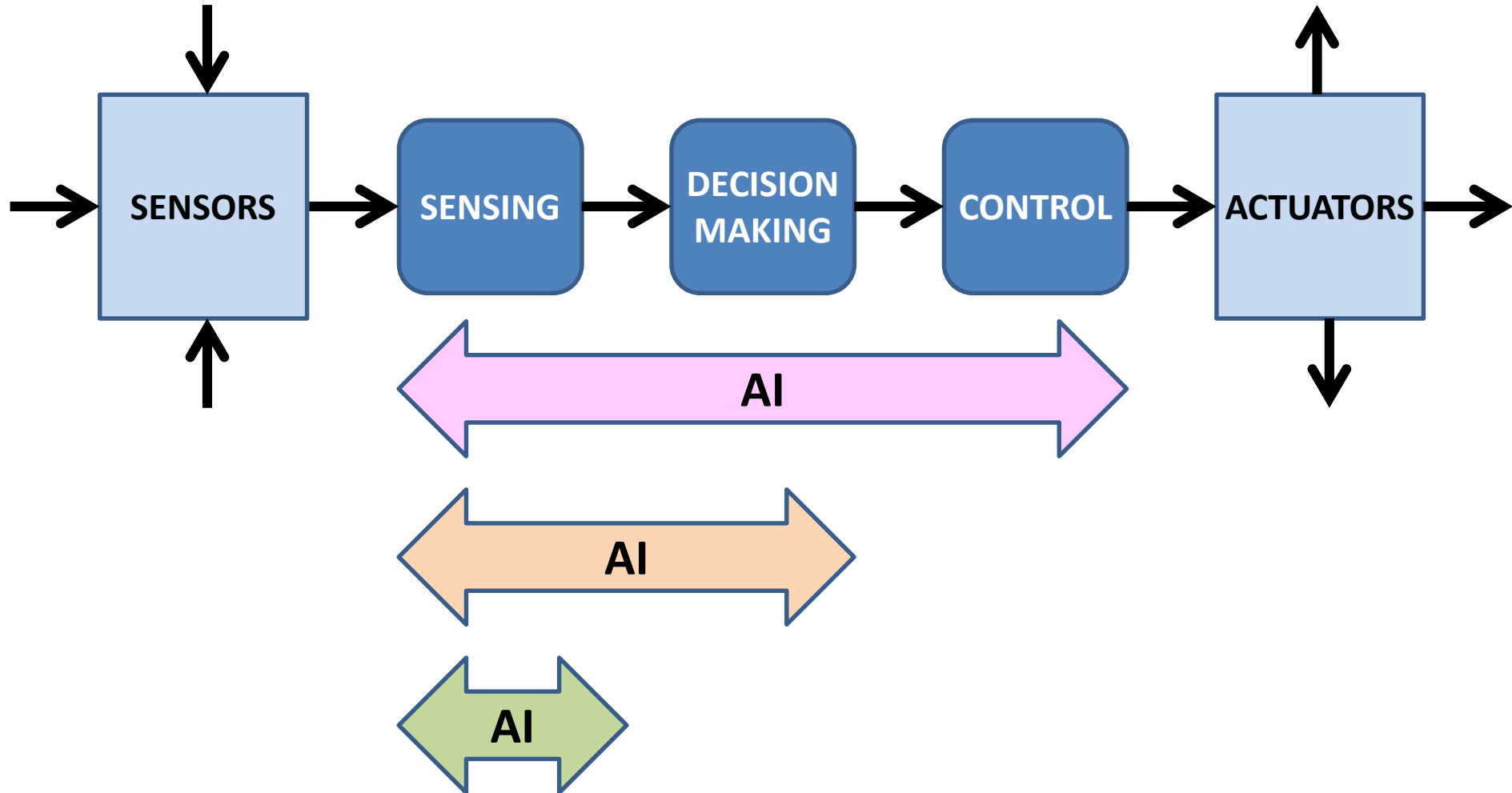  - the ability to act independently of direct human control and in unrehearsed conditions

- **Autonomous System**
  - system that changes its behaviour based on its experiences and the current situation to achieve given objectives without human control
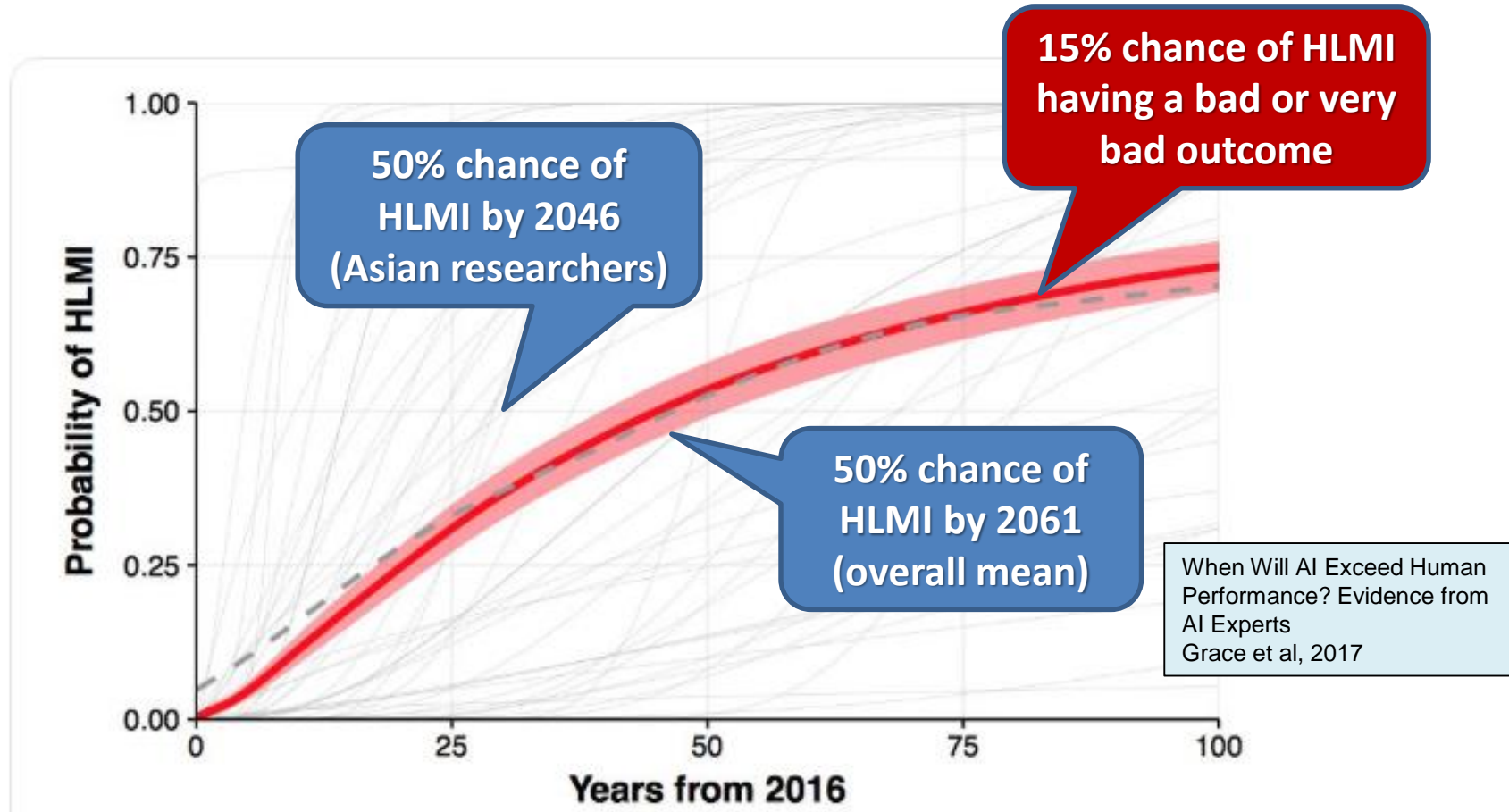
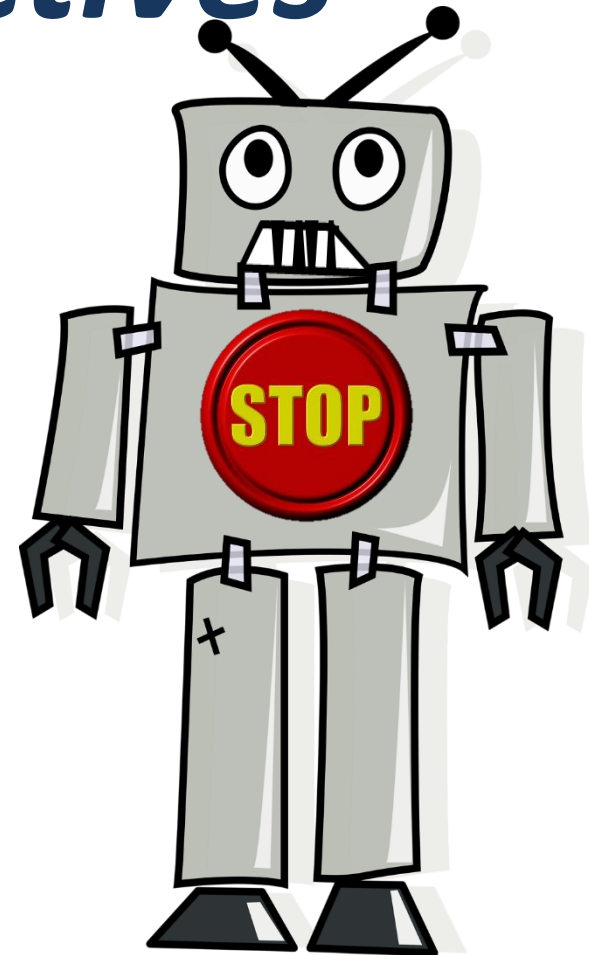# Basic Autonomous System Framework

# High-Level Machine Intelligence
## - as predicted by published AI researchers



**15% chance of HLMI having a bad or very bad outcome**

**50% chance of HLMI by 2046 (Asian researchers)**

**50% chance of HLMI by 2061 (overall mean)**

When Will AI Exceed Human Performance? Evidence from AI Experts
Grace et al, 2017

"High-level machine intelligence" (HLMI) is achieved when unaided machines can accomplish every task better and more cheaply than human workers.

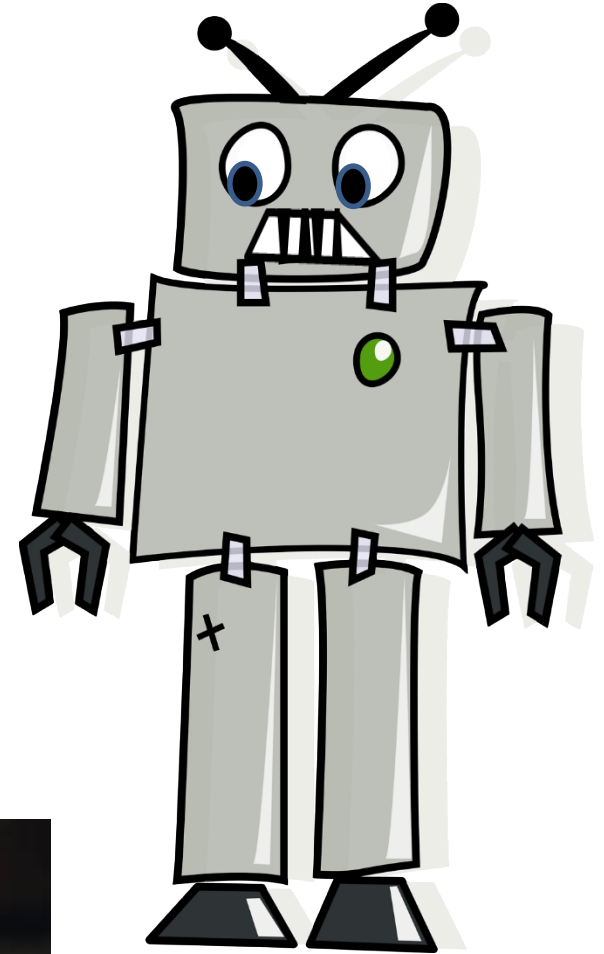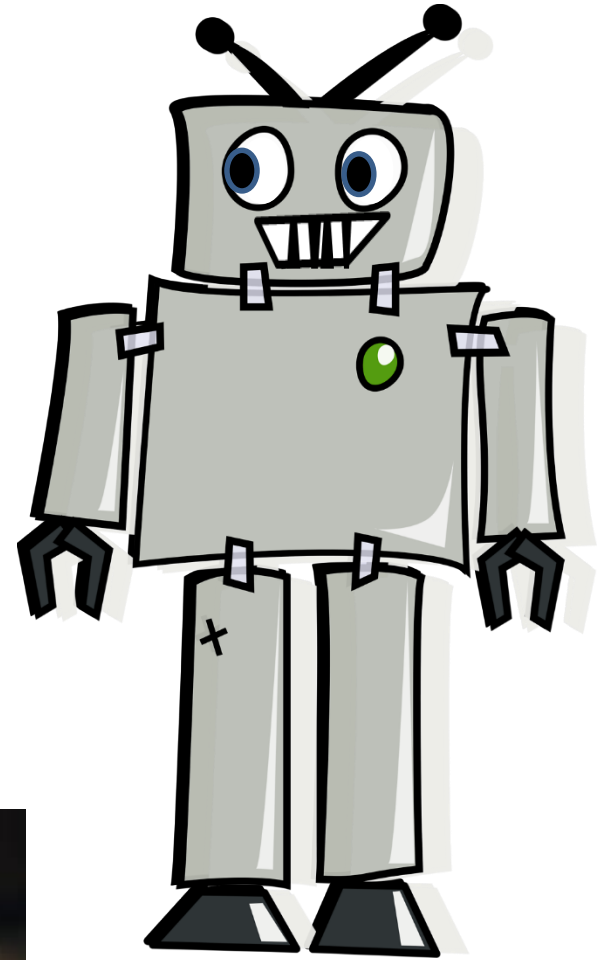# *Specifying Objectives (Safely)*

# The Midas Problem "마이더스의 손"



MIDAS' DAUGHTER TURNED TO GOLD
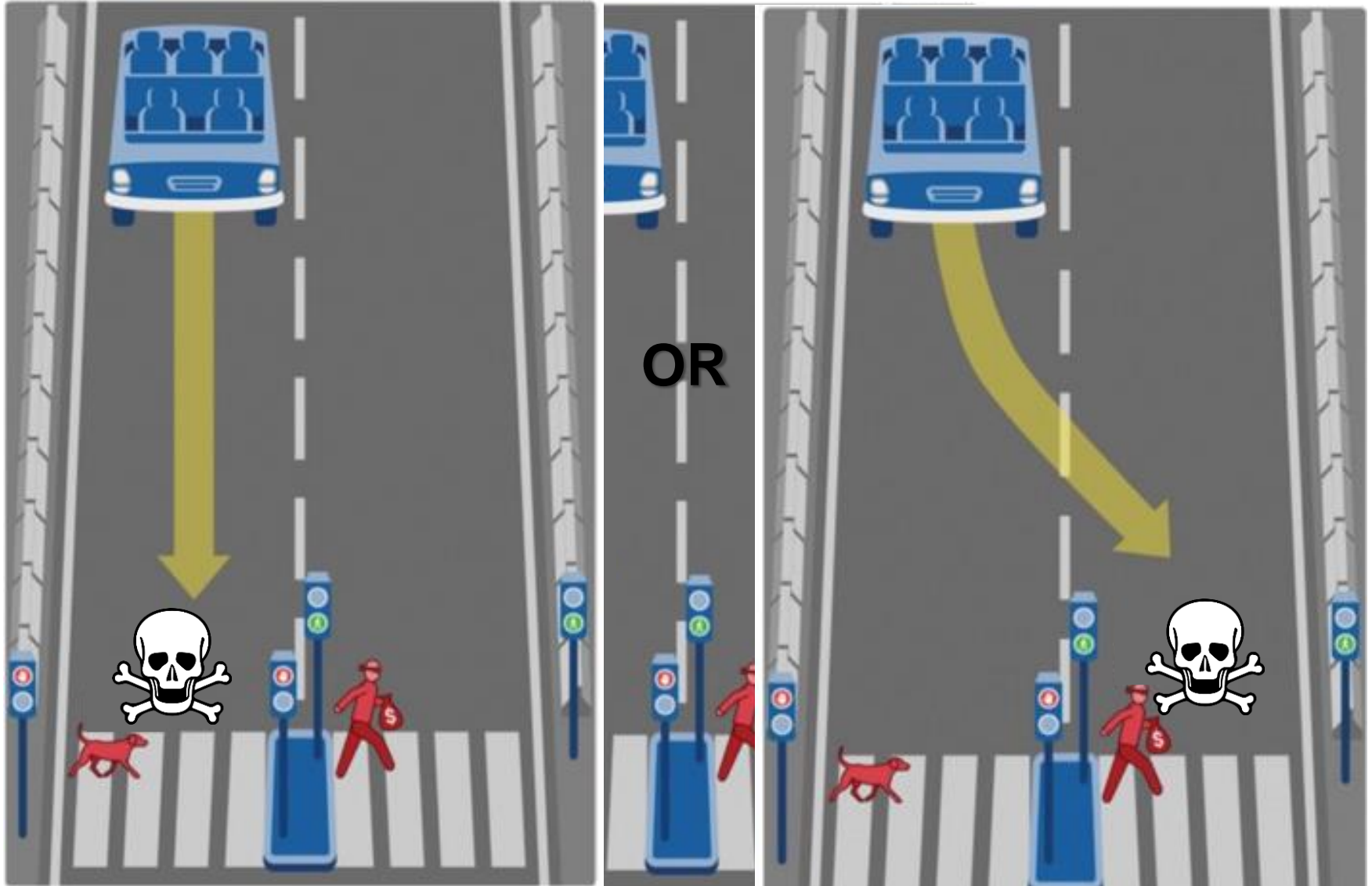
# "I'm hungry!   Make me dinner"

# "Keep the kitchen clean"

# Side-Effects, Reward Hacking and Role Models

- **Reinforcement learning involves the system being rewarded for achieving objectives**
  - must be aware of side-effects
  - however problems can arise with 'reward hacking' when the system 'hacks' the objectives
- **Instead, we can get systems to learn from human demonstrations**
  - and get feedback from humans
- **BUT**
  - make sure the humans are representative
  - recognize that human values change over time
  - humans aren't always the best role models...

# MIT's Moral Machine (moralmachine.mit.edu)

# Who should die?

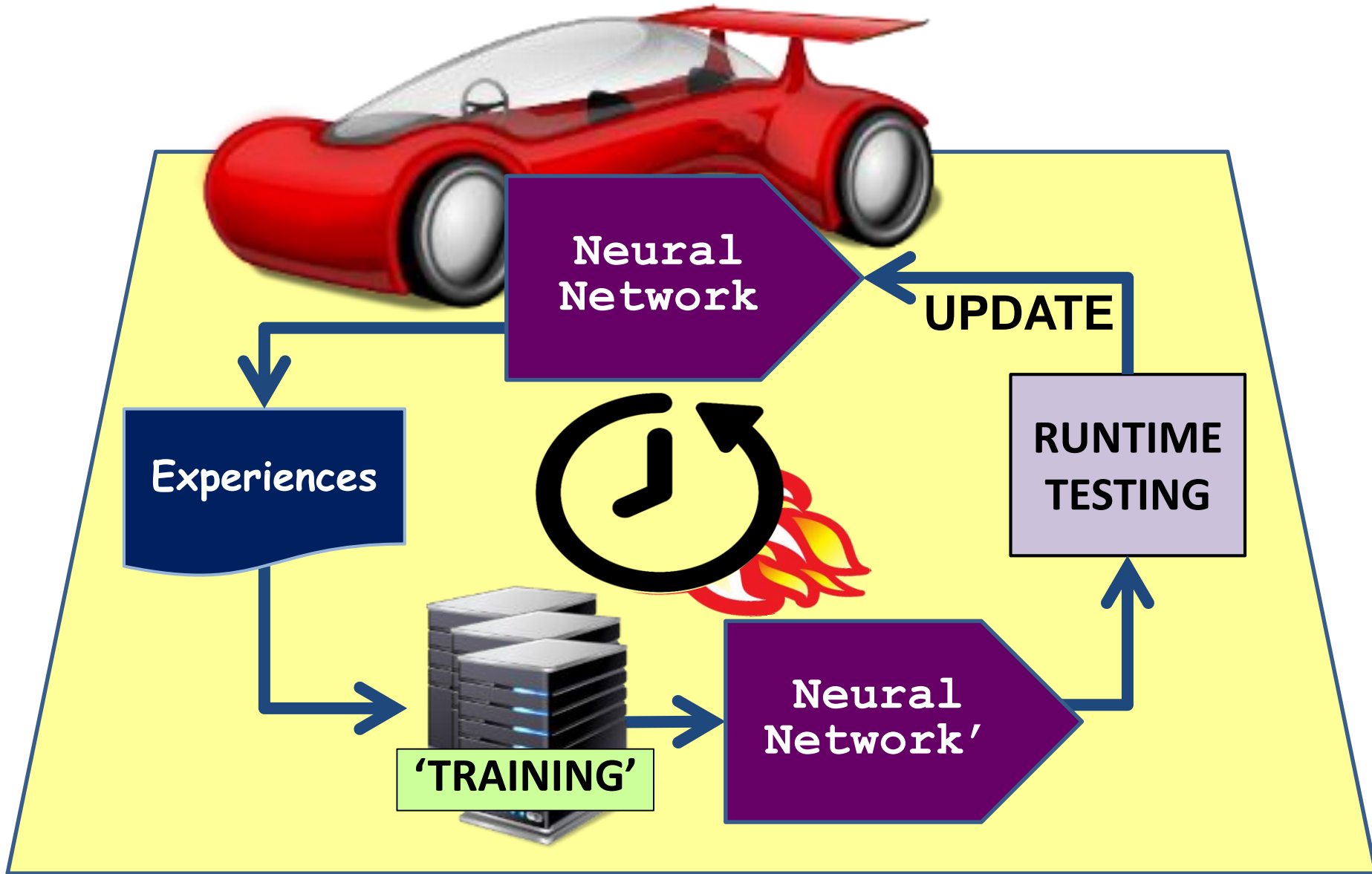| 0 | 0 |
|---|---|
| Dog | Criminal |

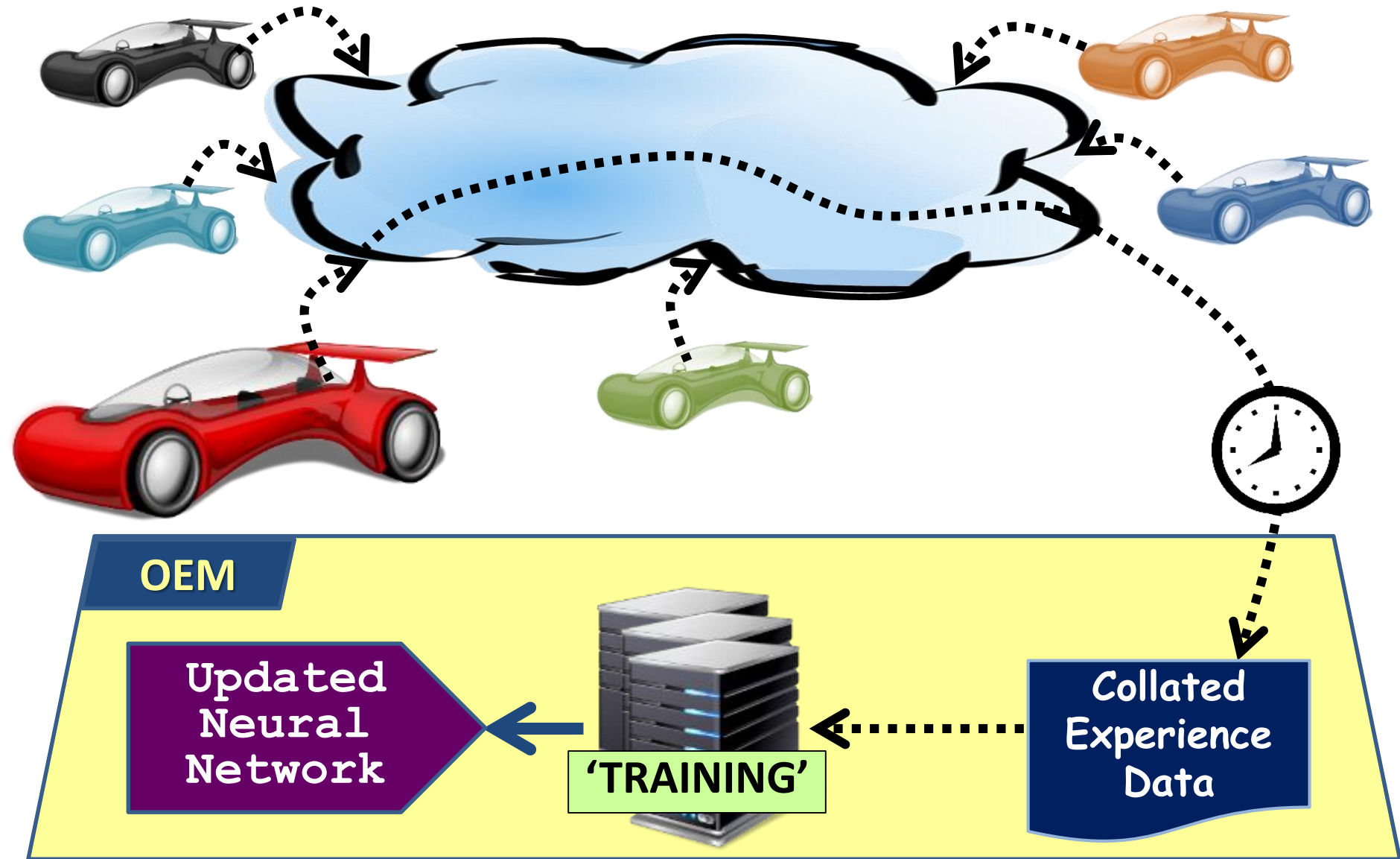Slide is not active **Activate**

0

# Better than Humans?

# Online vs Off-Line Machine Learning

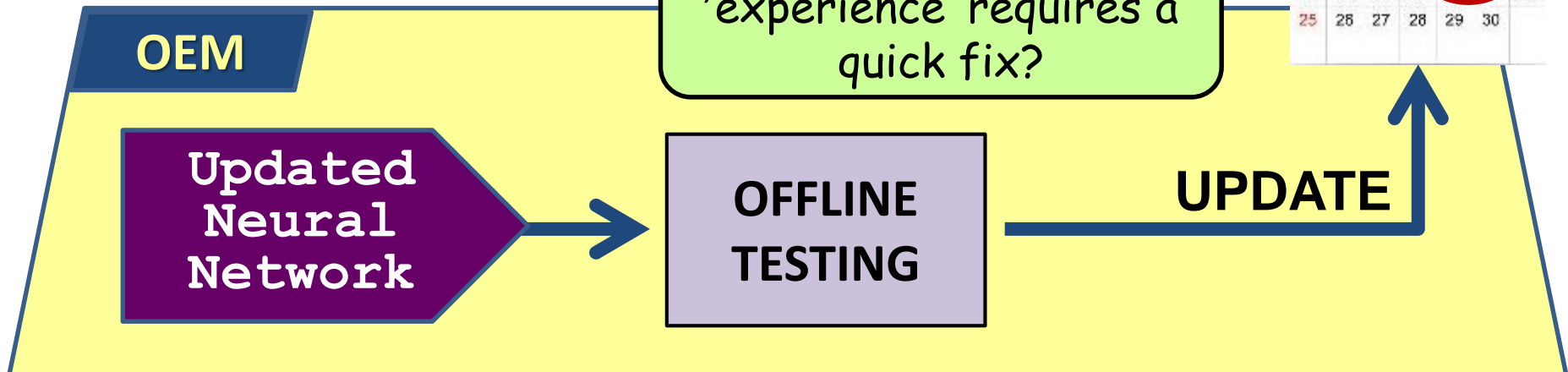# Continuous Online Learning

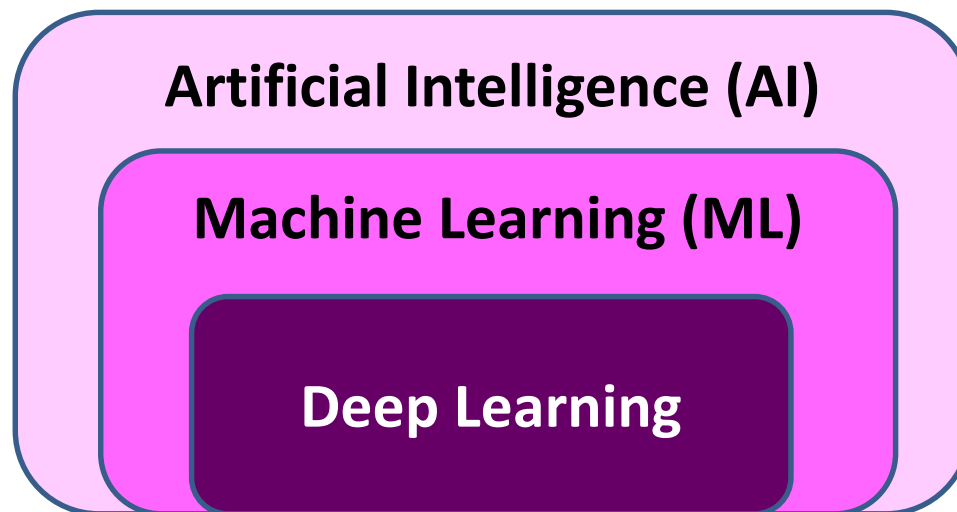# Off-Line Learning – from Day-to-Day Use

# Performance Updates - Over-The-Air

# *Machine Learning Challenges*

# Deep Learning Systems

Artificial Intelligence (AI)
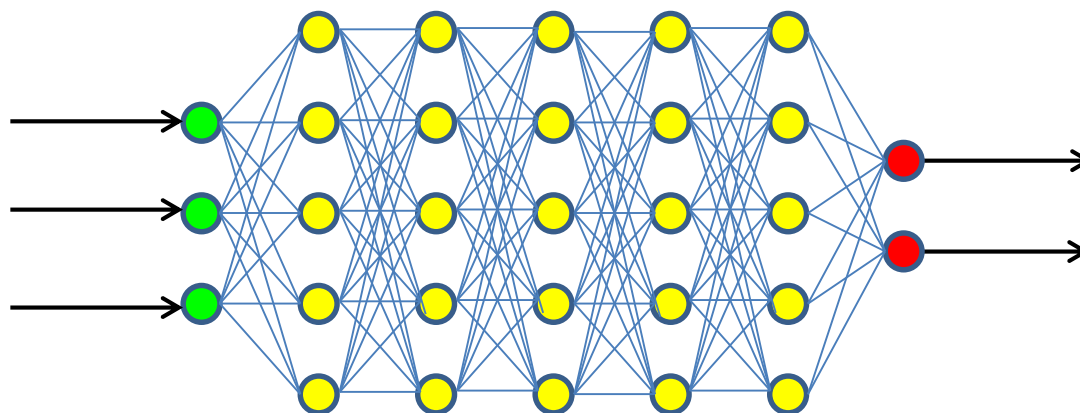
Machine Learning (ML)

Deep Learning

Deep Neural Network

# Example of Machine Learning

# Supervised Machine Learning

# Mis-Classification

# Checking the Training Set

Need to look for biased training data, overfitting, underfitting, extraneous data, outliers, etc.

Model & Parameter Selection

Training Set

Network Training

Test Set

Neural Network

Accuracy

# Misunderstanding – Data Bias
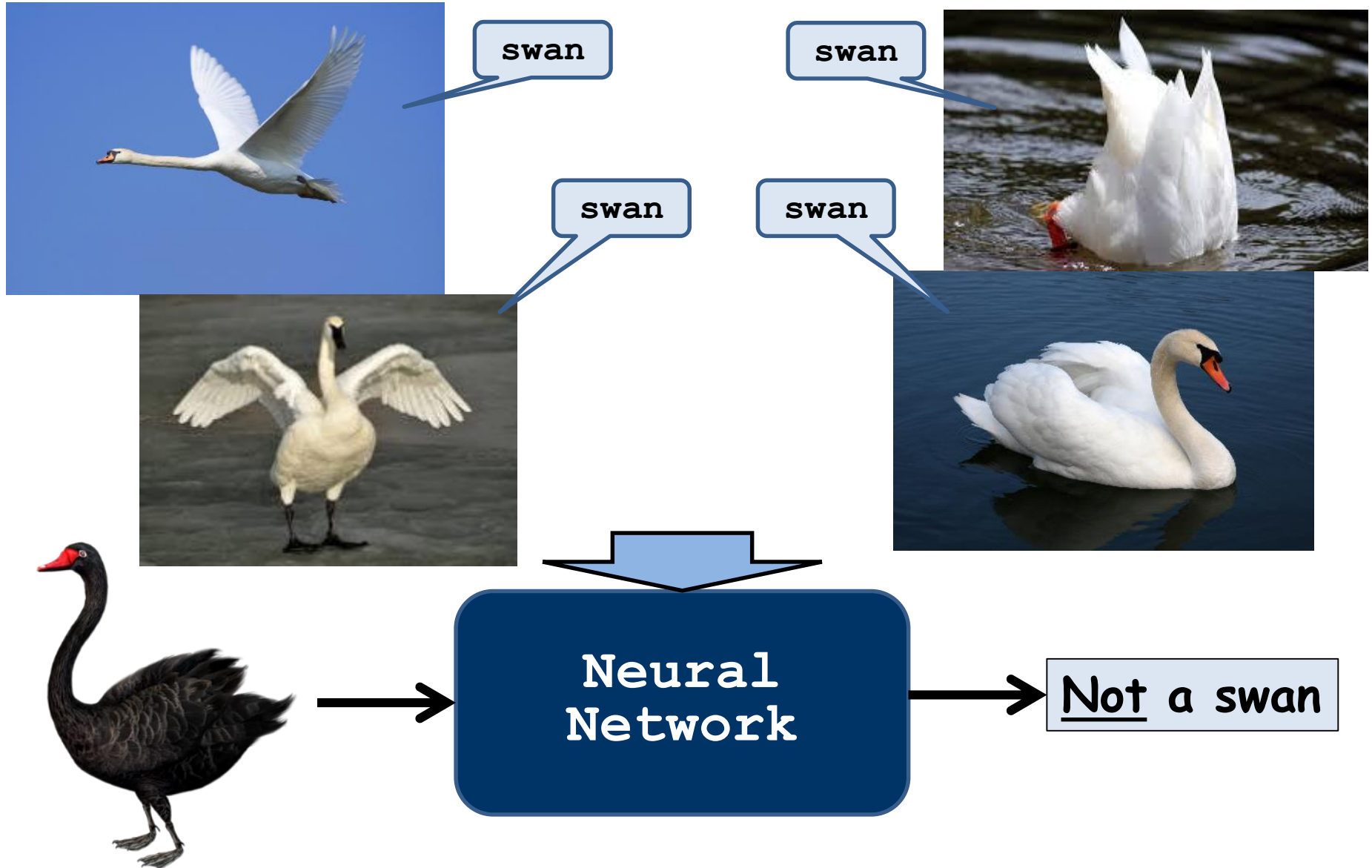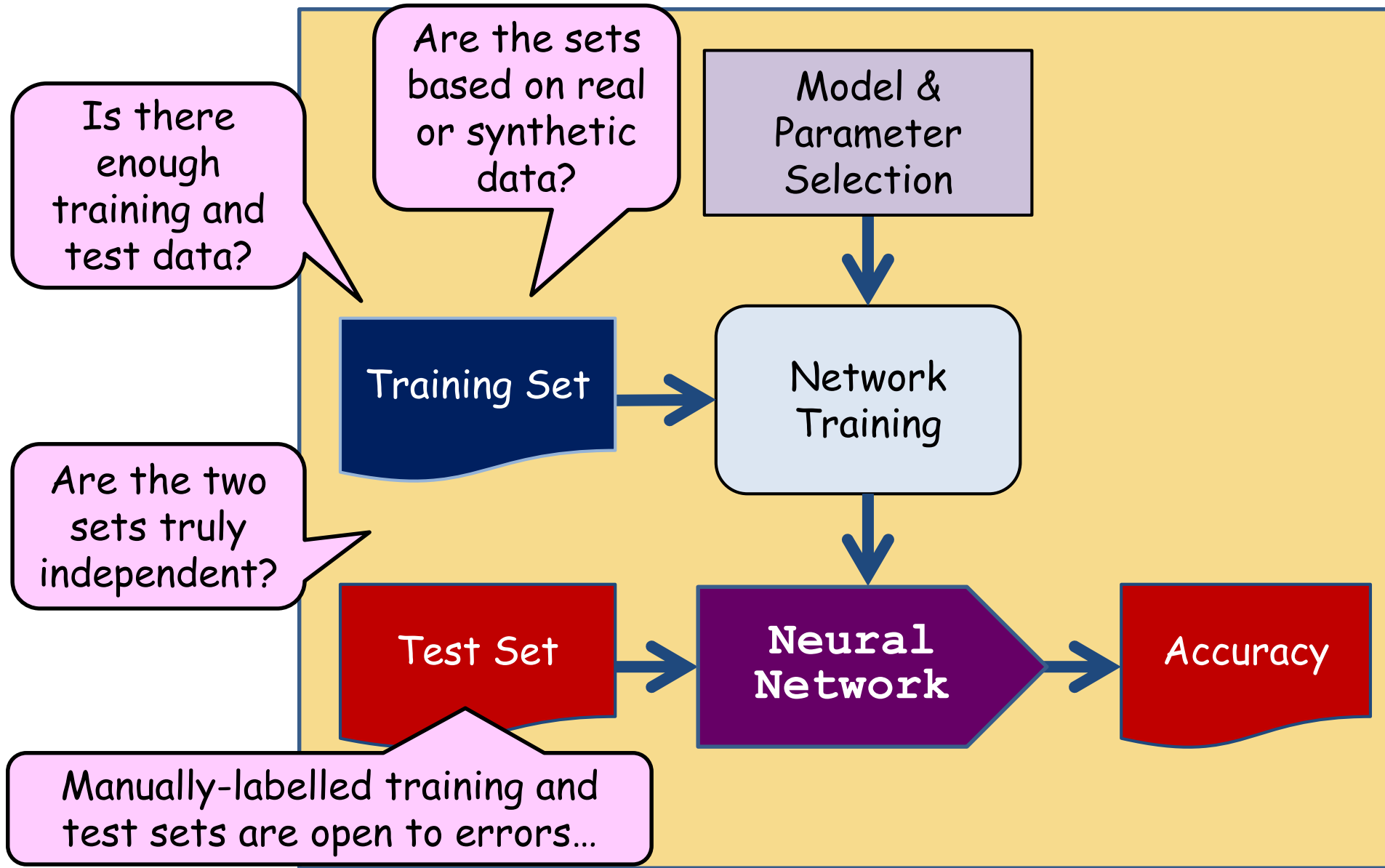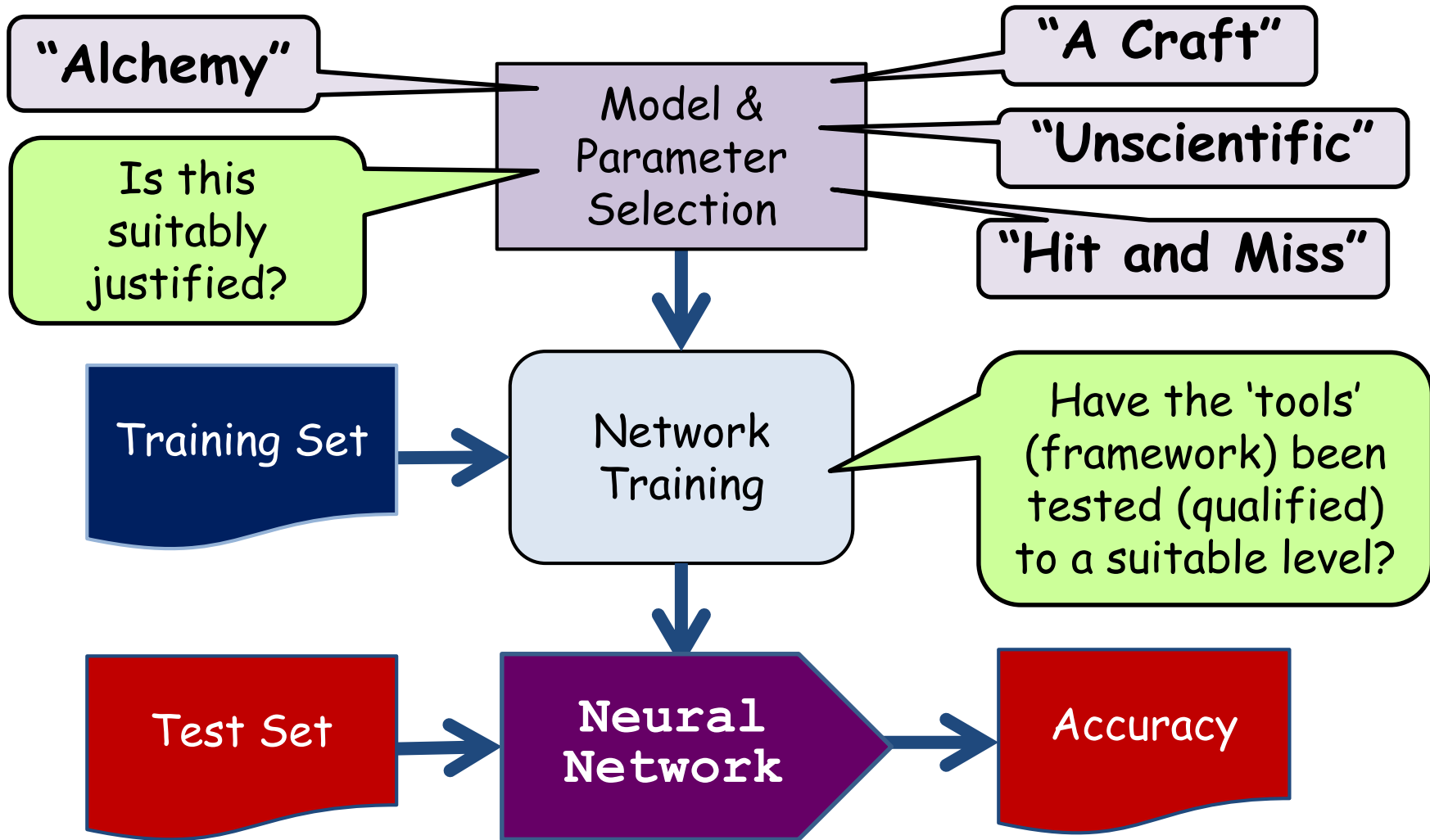
# Incomplete Training Set

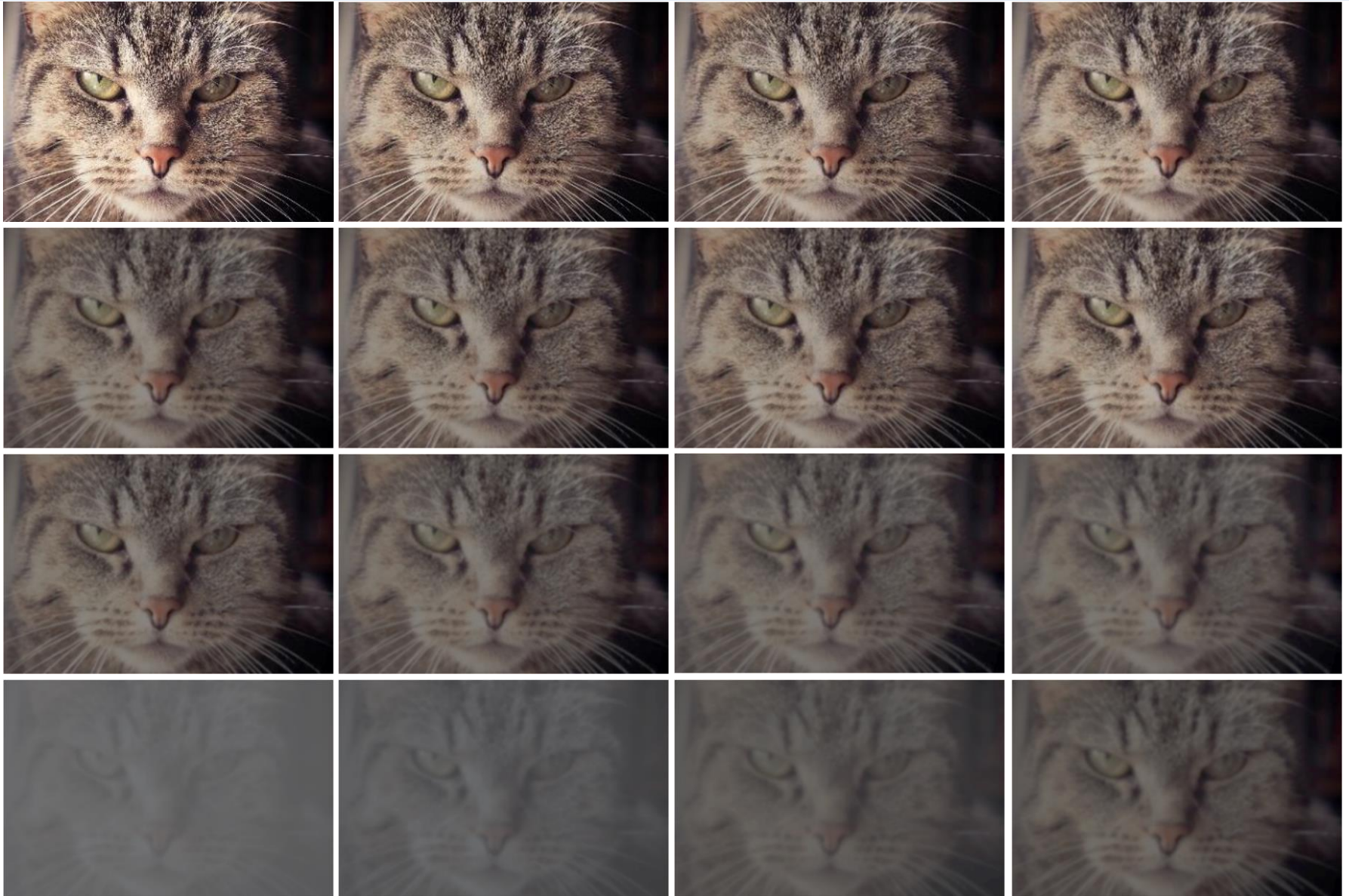# Checking the Training & Test Sets

# Checking the Training

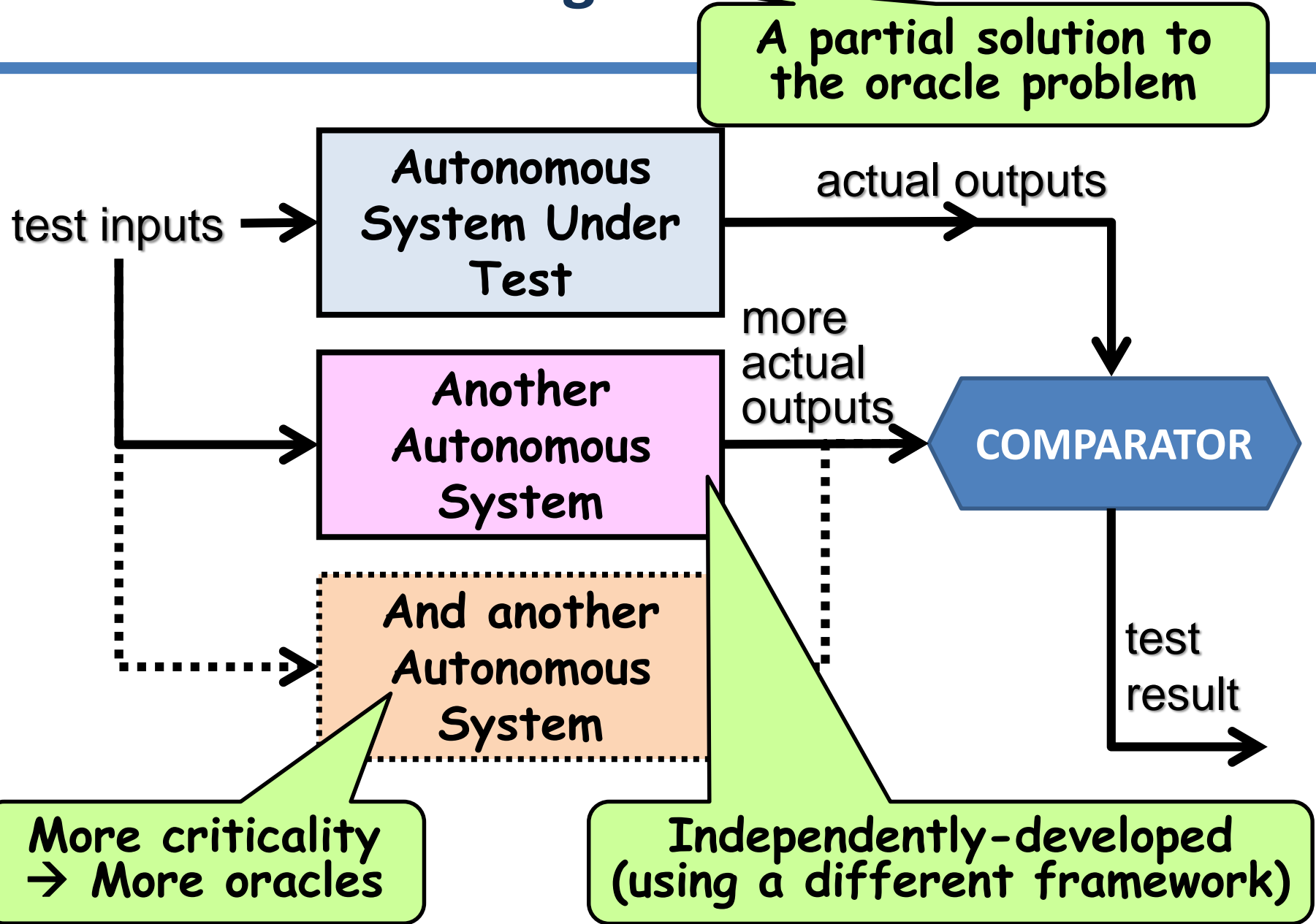# *Black Box Testing of Autonomous Systems*

# Test Challenges of Autonomous Systems

- **Expected Results (Test Oracle)**
  - if we struggle to set the objectives, then determining expected results will be equally difficult

- **Probabilistic Systems and Non-Determinism**
  - the probabilistic nature means that predicting expected results is difficult
    - we need many more tests to be statistically confident
  - non-determinism causes real problems for regression testing

- **Complexity**
  - autonomous systems are difficult to understand - and to test
  - interacting autonomous systems may cause 'special' failures
  - many sensors can create many tests...

# Back-to-Back Testing

**A partial solution to the oracle problem**

test inputs →

**Autonomous System Under Test** → actual outputs

**Another Autonomous System** → more actual outputs

**And another Autonomous System**

**COMPARATOR** → test result

**More criticality → More oracles**

**Independently-developed (using a different framework)**
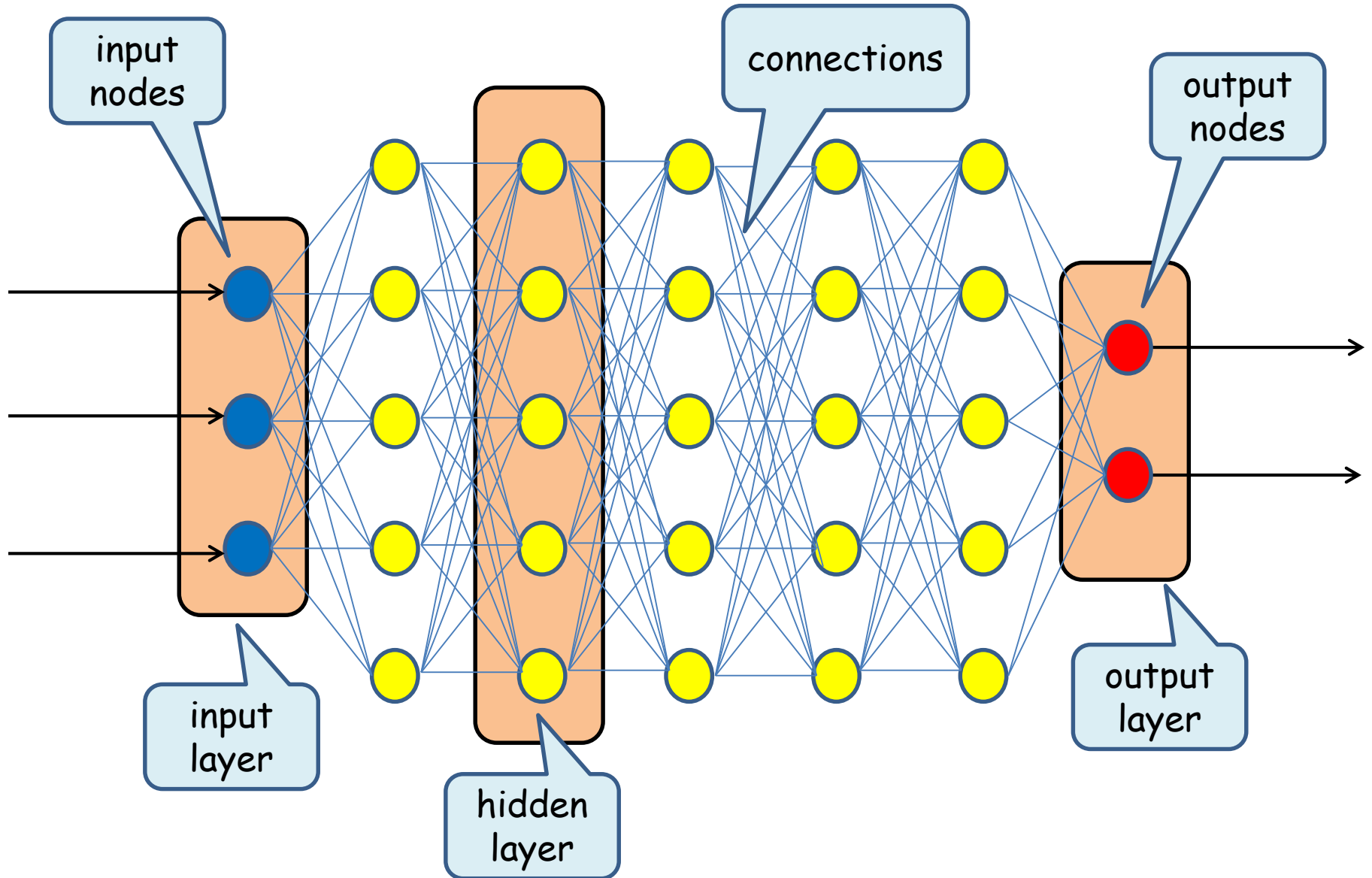
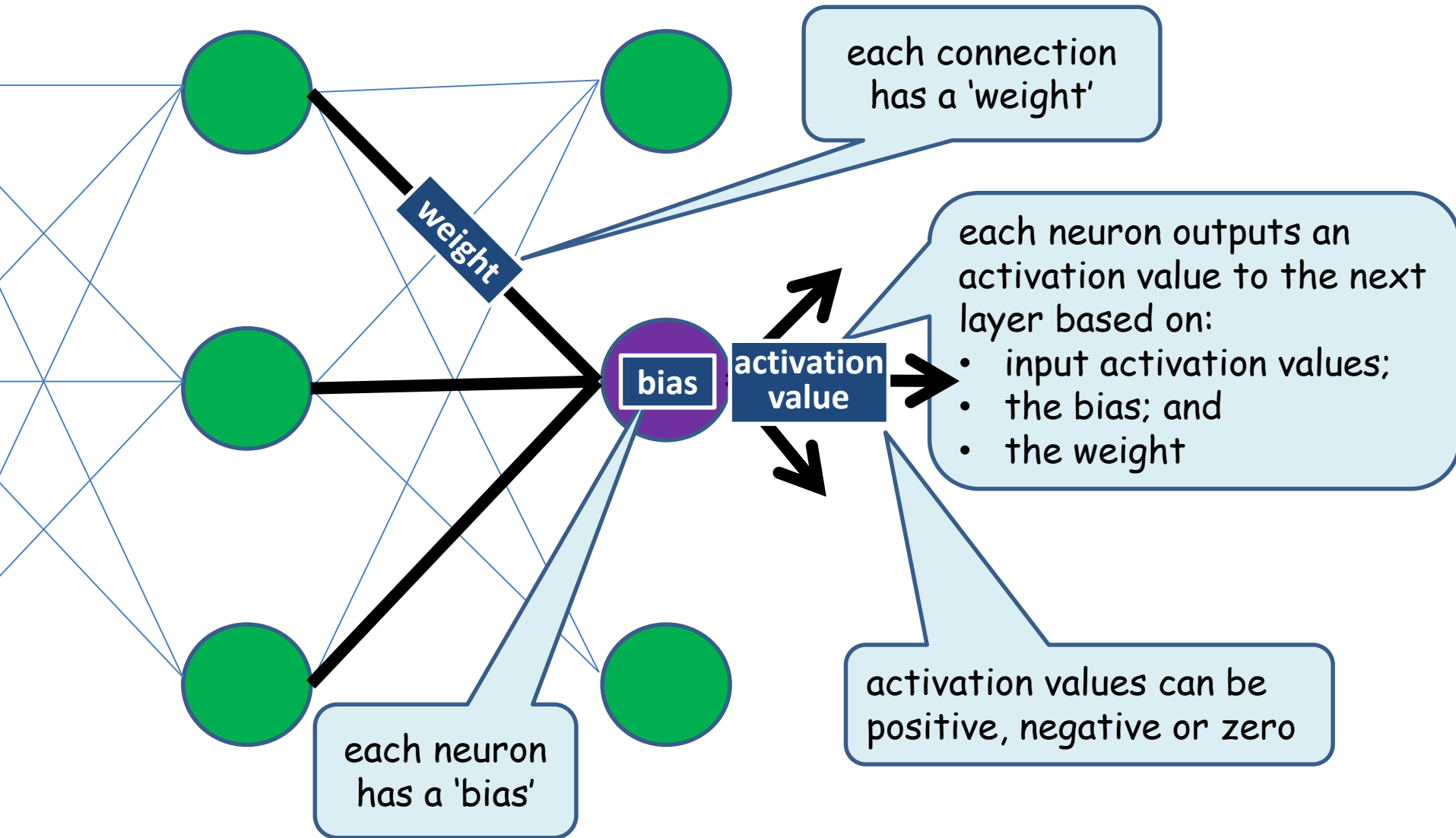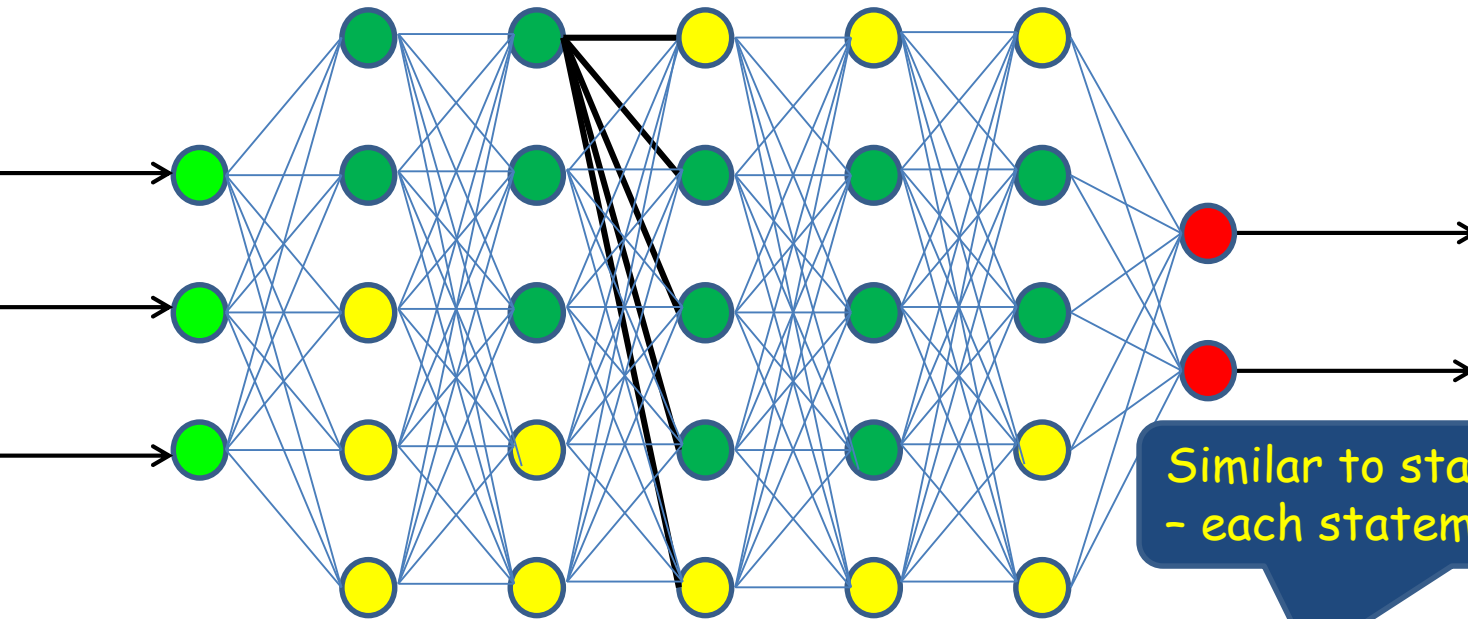# *White Box Testing of Autonomous Systems*

# Deep Neural Net

# Activation Values

# 'Neuron' Coverage

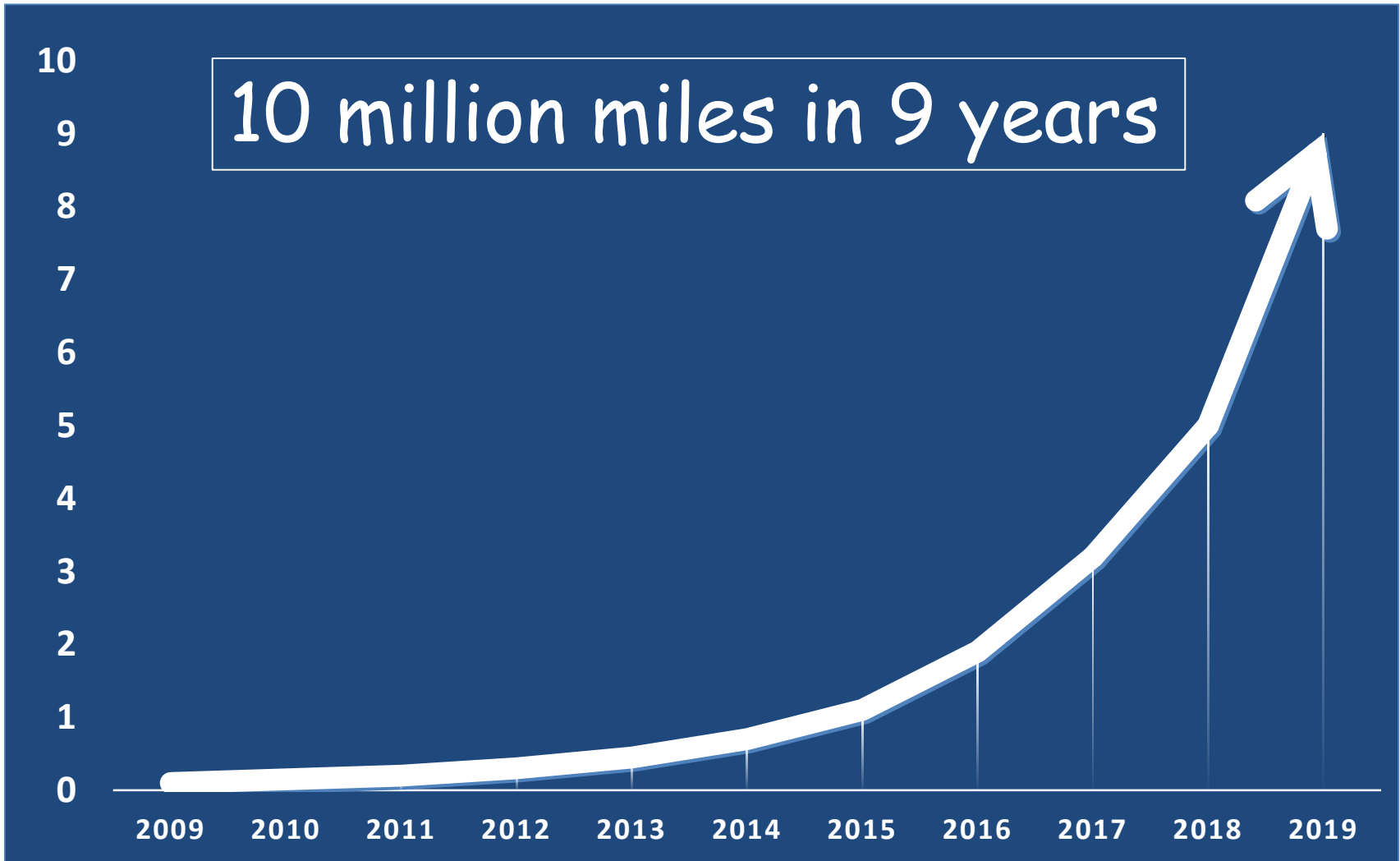

activation value is above zero

Similar to statement coverage – each statement is exercised

Full 'neuron' coverage shows that every neuron is 'activated' (value above zero) at least once (but very basic coverage – easy to achieve with a few tests)
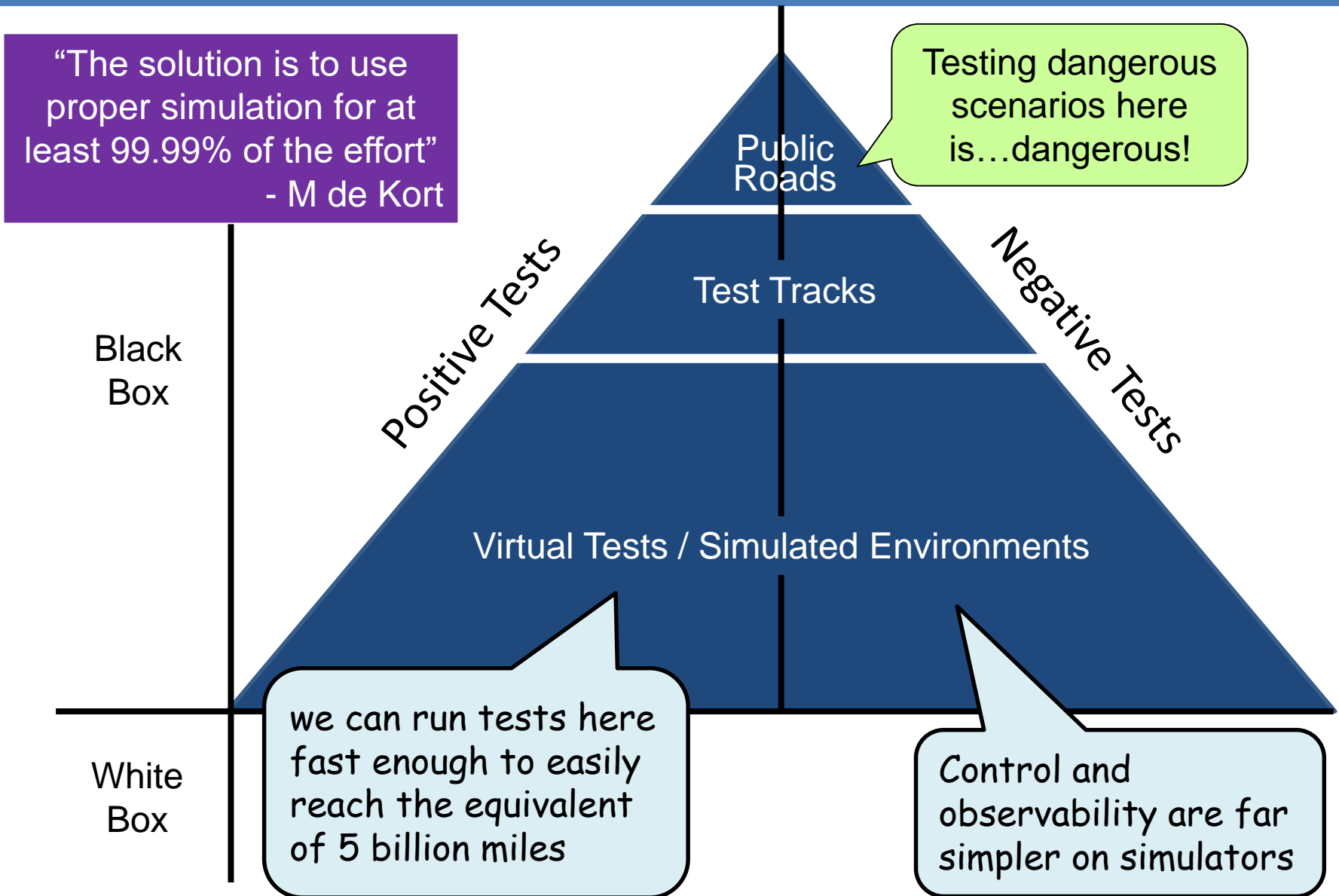
# *The Necessity of Virtual Test Environments*

# Waymo On-Road Test Miles (millions)



10 million miles in 9 years
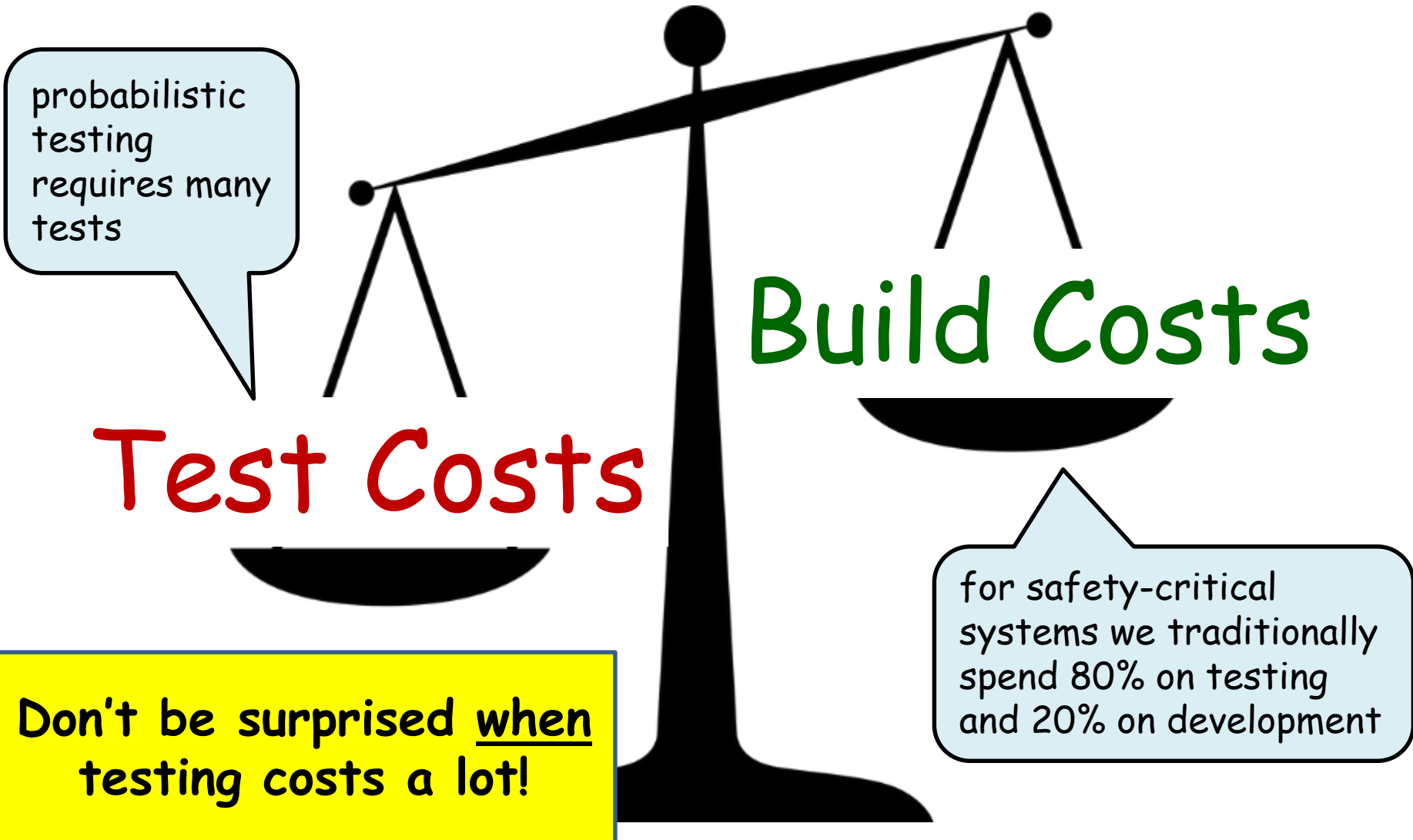
500x what has gone before = 5 Billion Miles

# Autonomous Cars – Test Environments
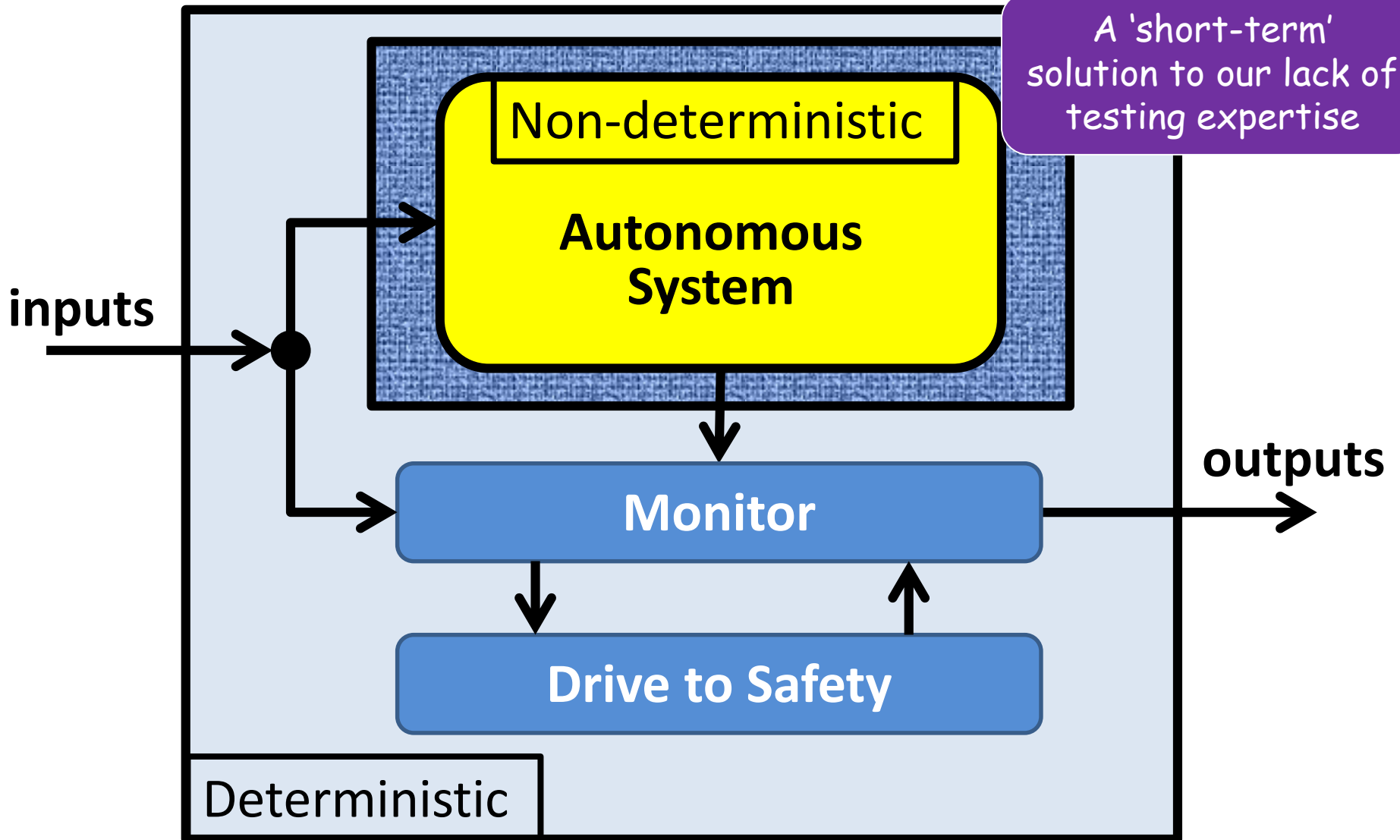
# *Conclusions*

# Conclusions – Safety of Autonomous Systems

- **For the 'simple' case of off-line systems we need:**
  - both black and white box testing
  - new test approaches and measures (with evidence)
  - more tests to assure these probabilistic systems
  - the support of sophisticated virtual test environments
- **For the learning on-line systems we need:**
  - to understand the new dangers these systems bring
- **Until we reach maturity, we should use a safety net…**

# Safety Shell Architecture

# Thank you for listening

## Any Questions?