

Safety of Autonomous Systems

Stuart Reid PhD, FBCS

stureid.test@gmail.com / www.stureid.info

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

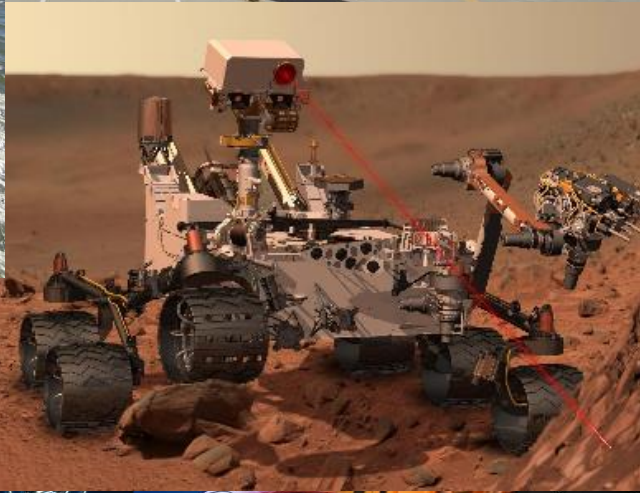
- **Autonomy**

- 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

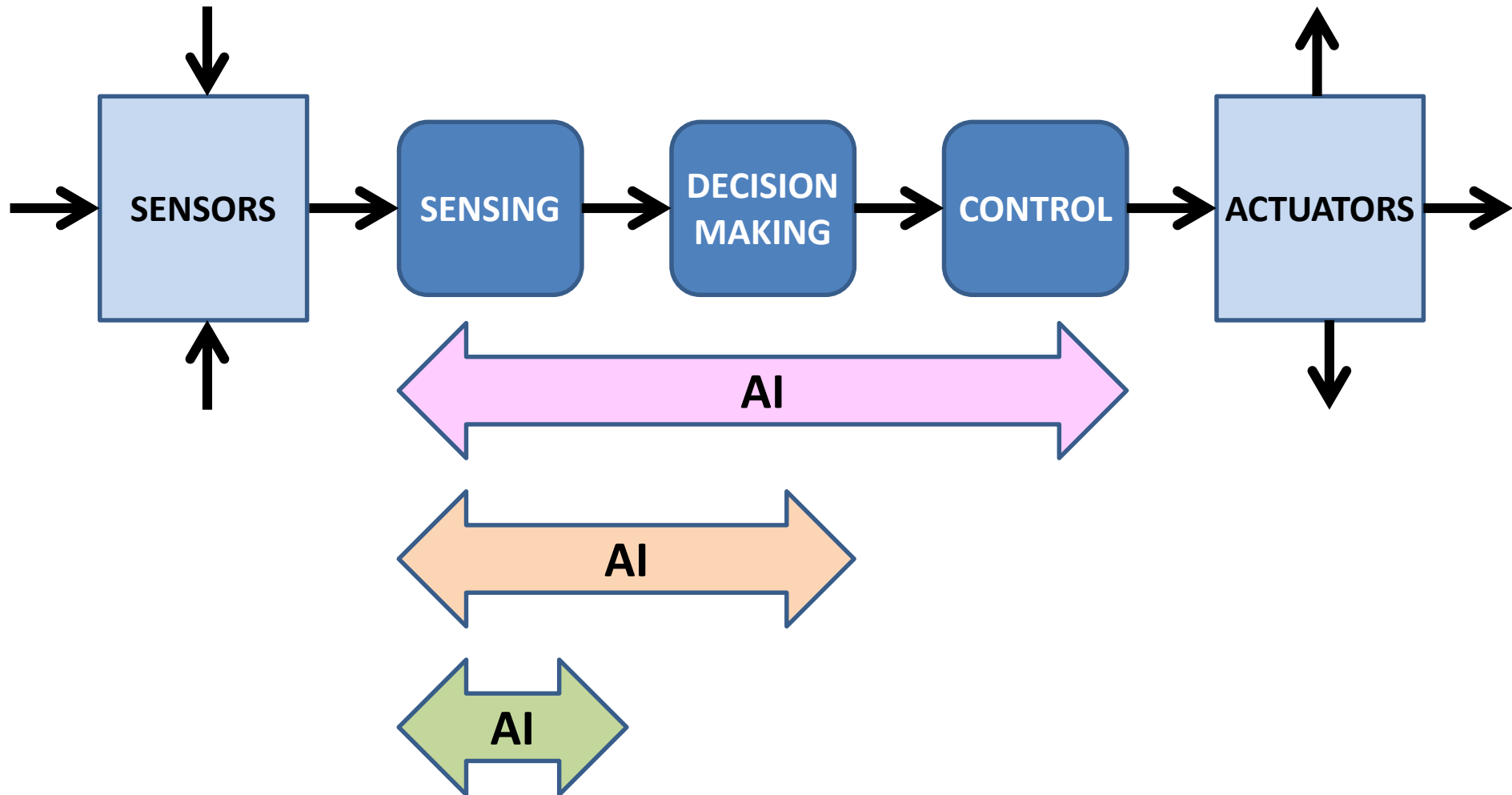
- **Autonomous System**

- system that changes its behaviour based on its experiences and the current situation to achieve given objectives without human control

Examples: Autonomous Systems



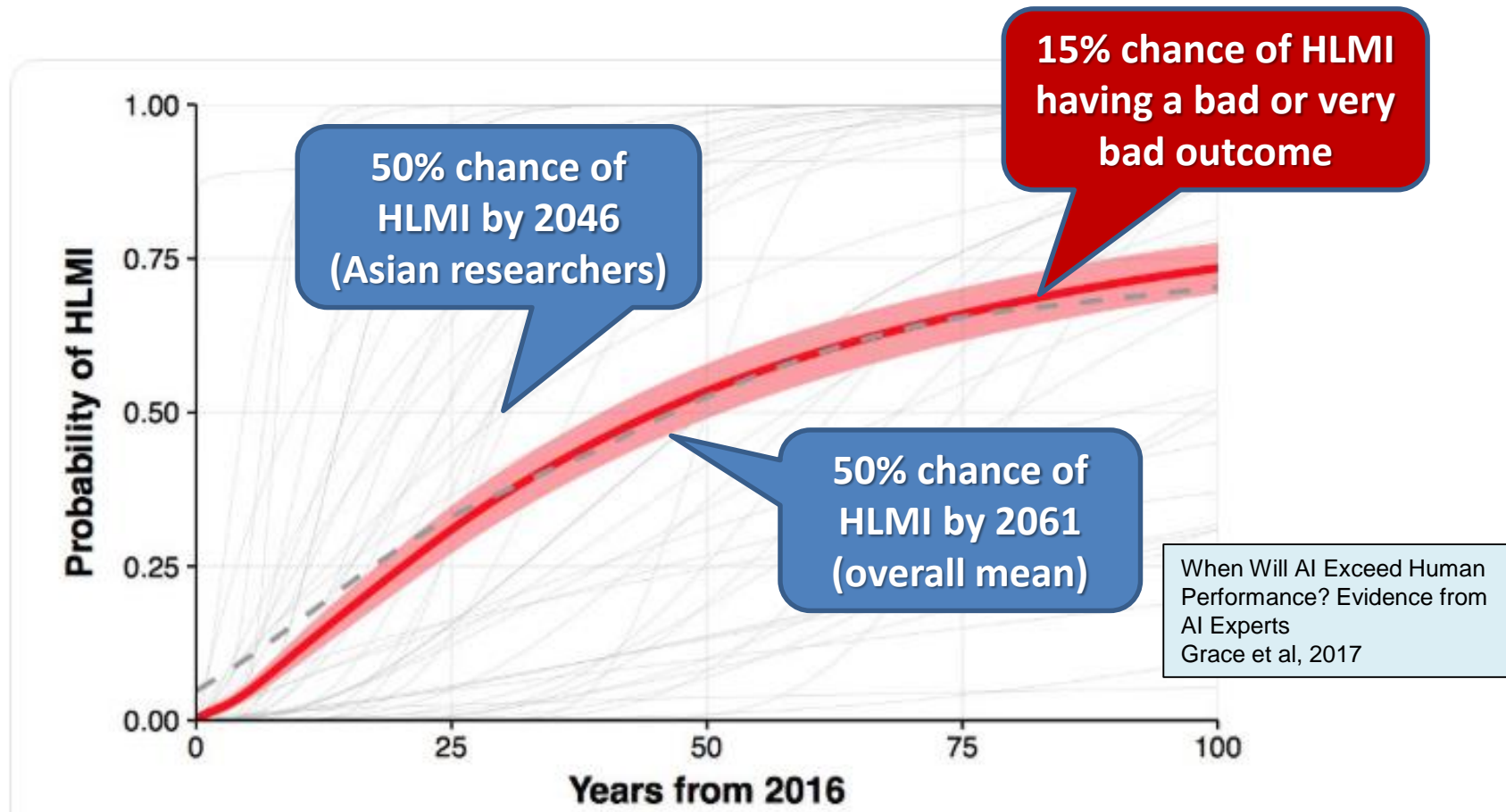
Basic Autonomous System Framework



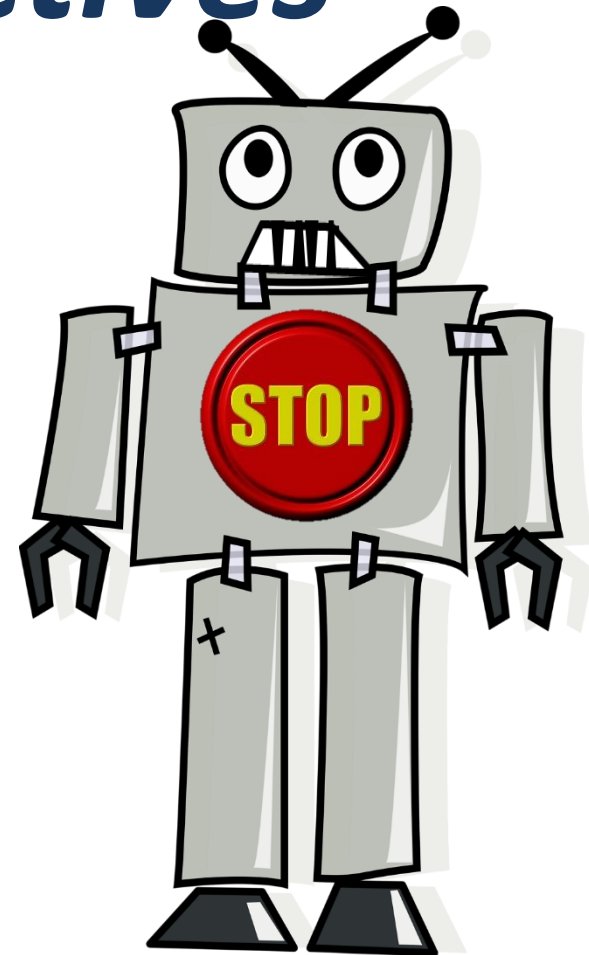
High-Level Machine Intelligence

- as predicted by published AI researchers

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



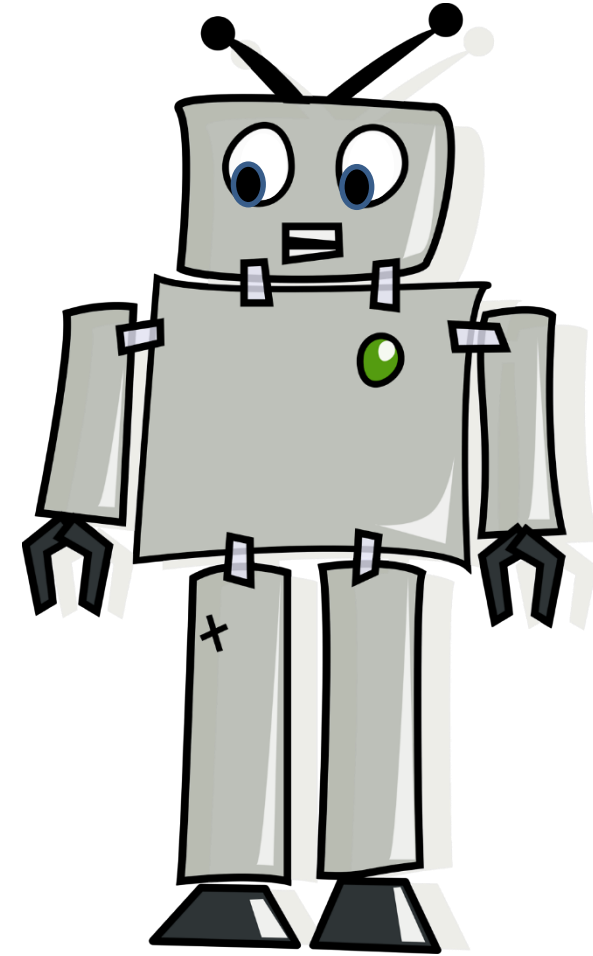
Specifying Objectives (Safely)



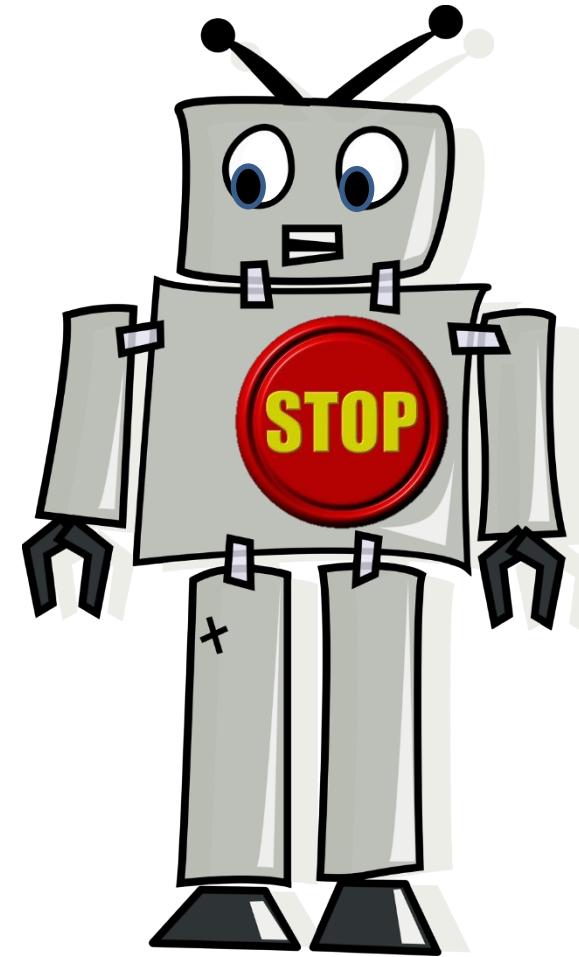
The Midas Problem “마이더스의 손”



“I’m hungry! Make me dinner”



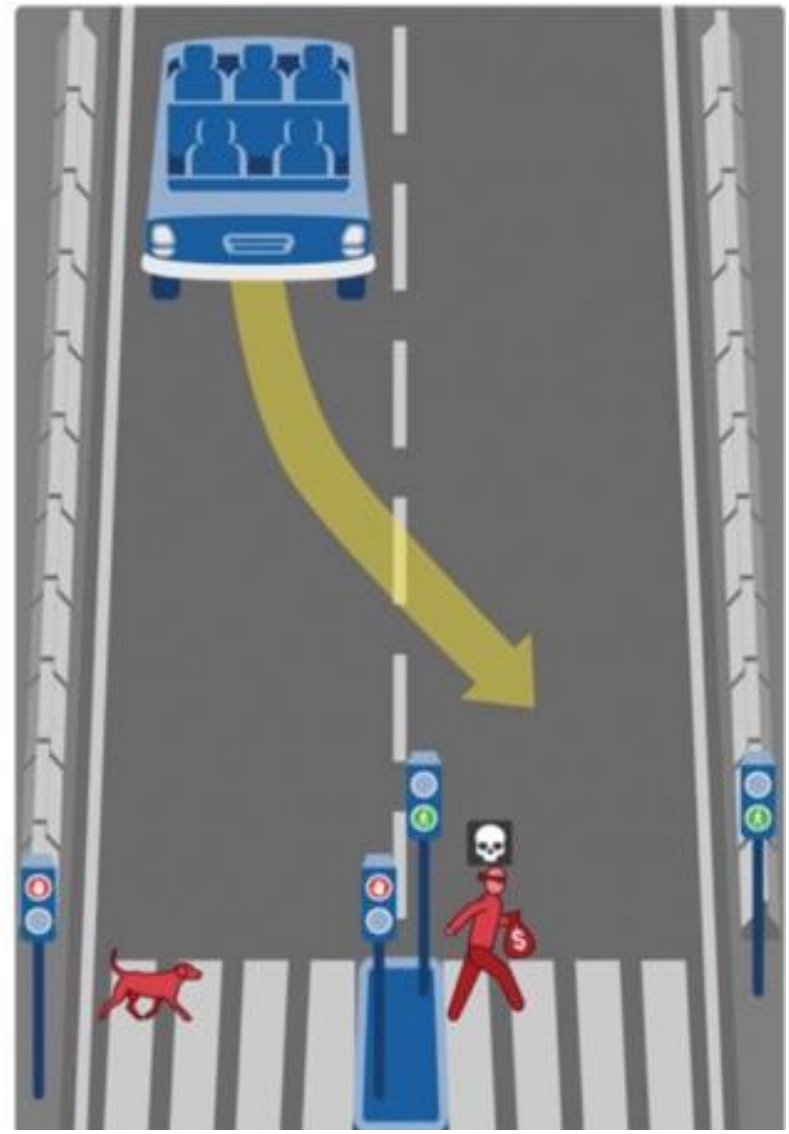
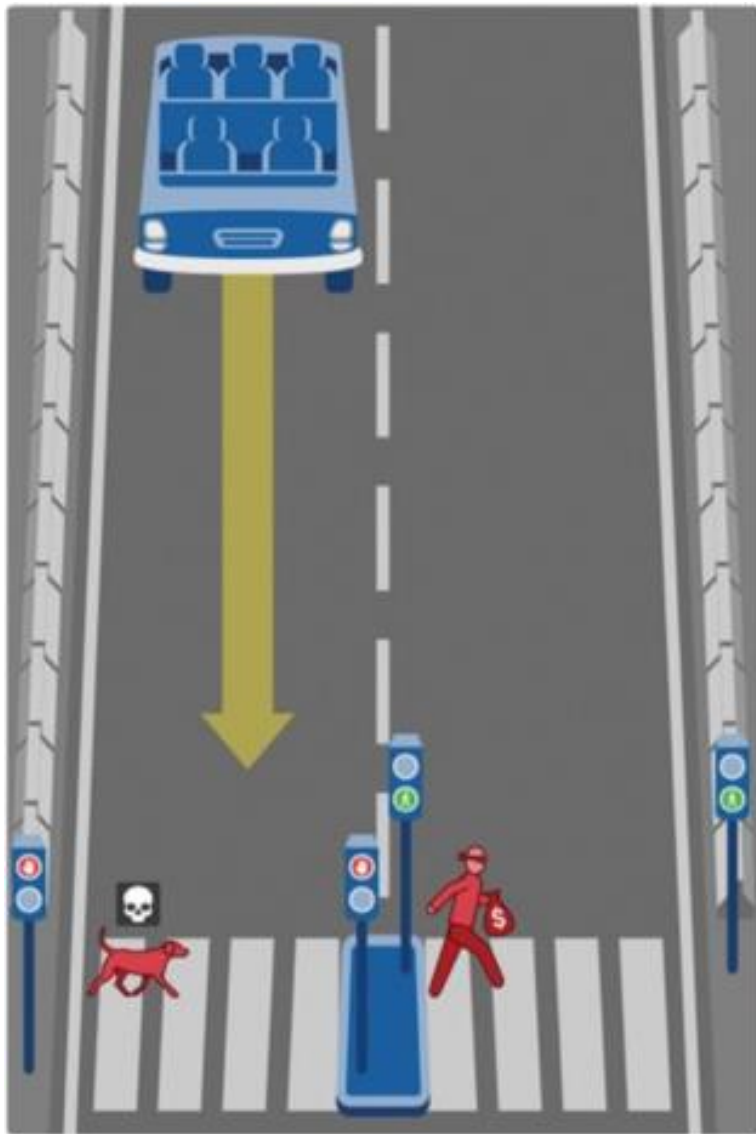
“Keep the kitchen clean”



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)

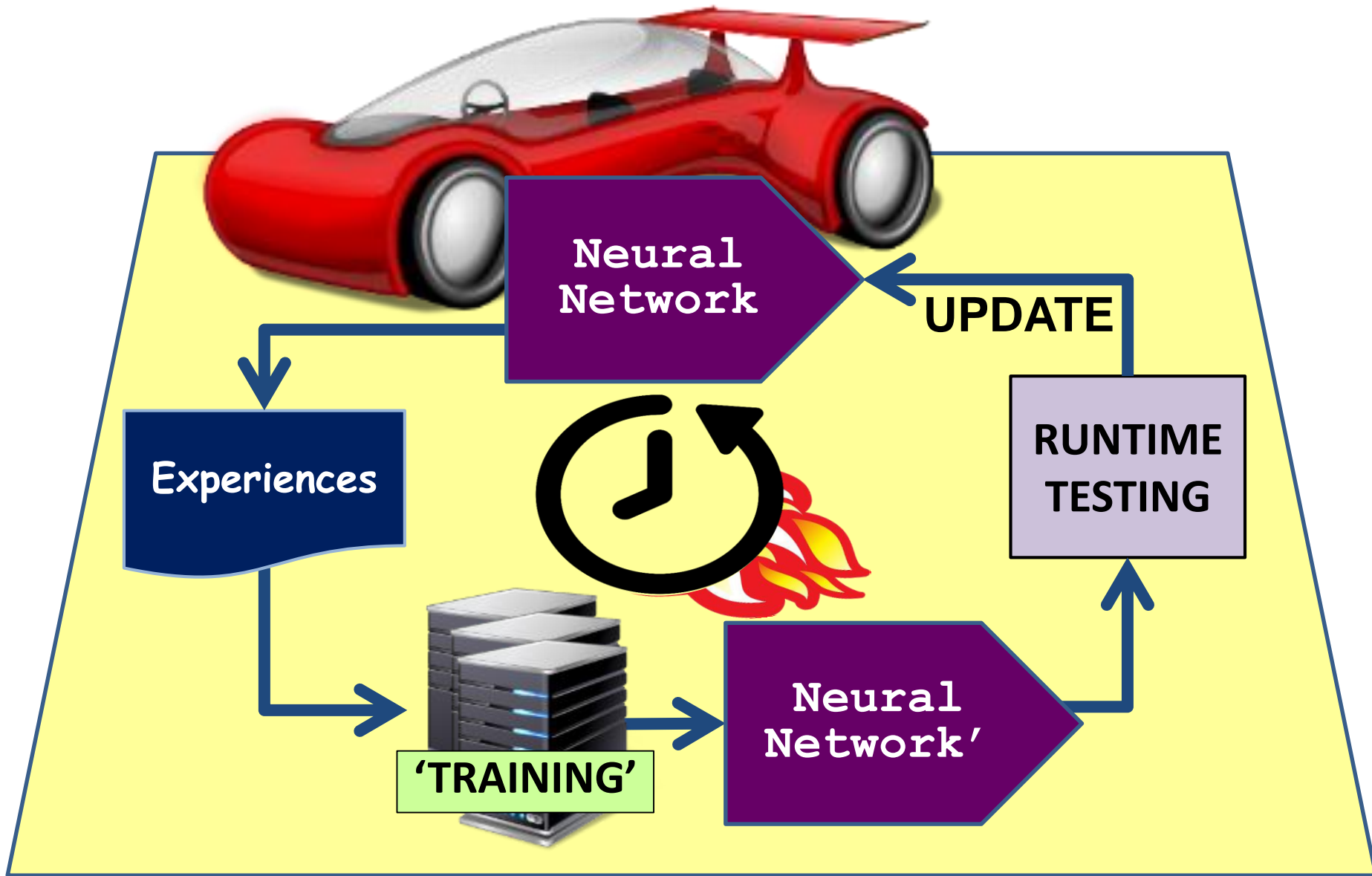


Better than Humans?

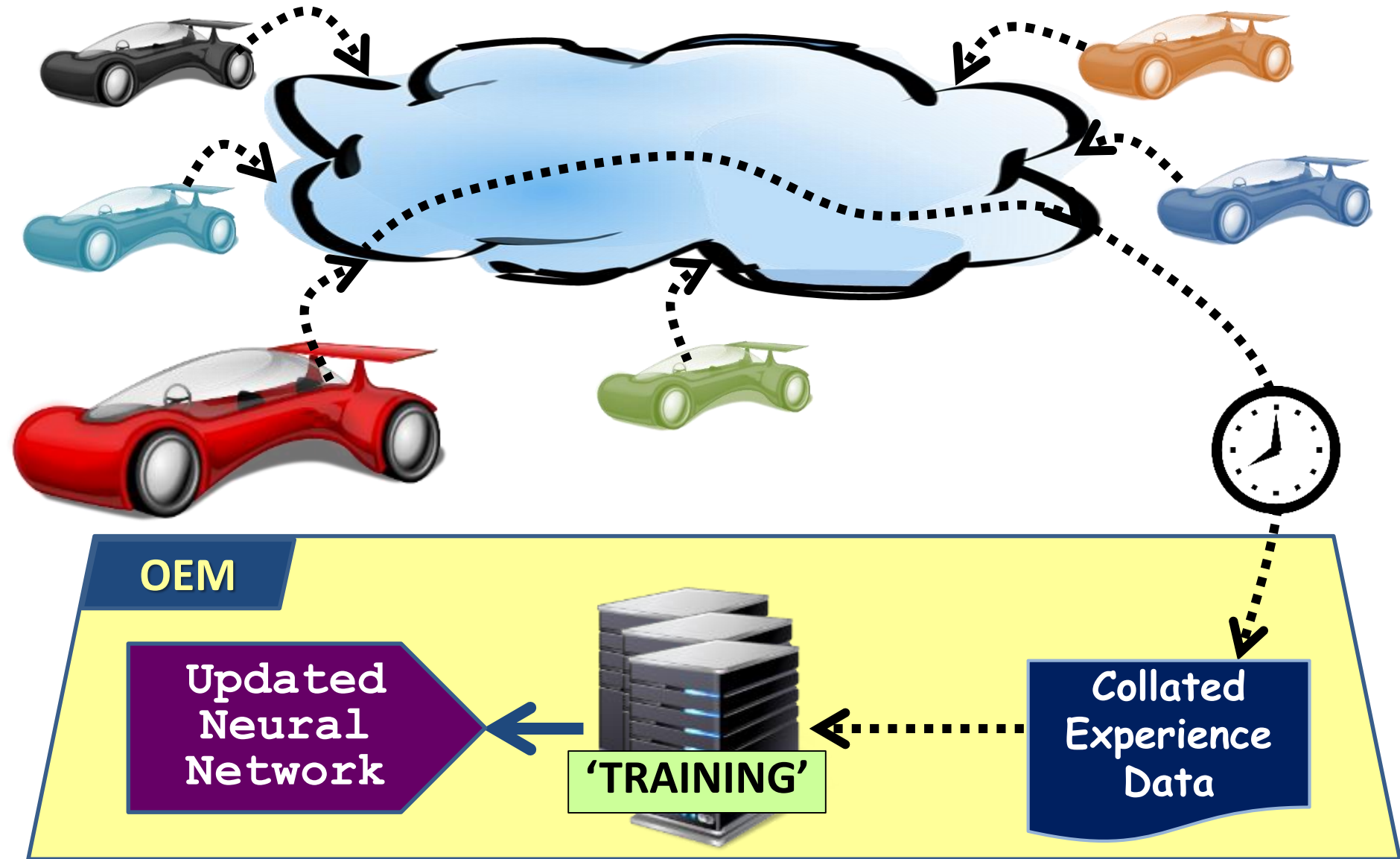


Online vs Off-Line Machine Learning

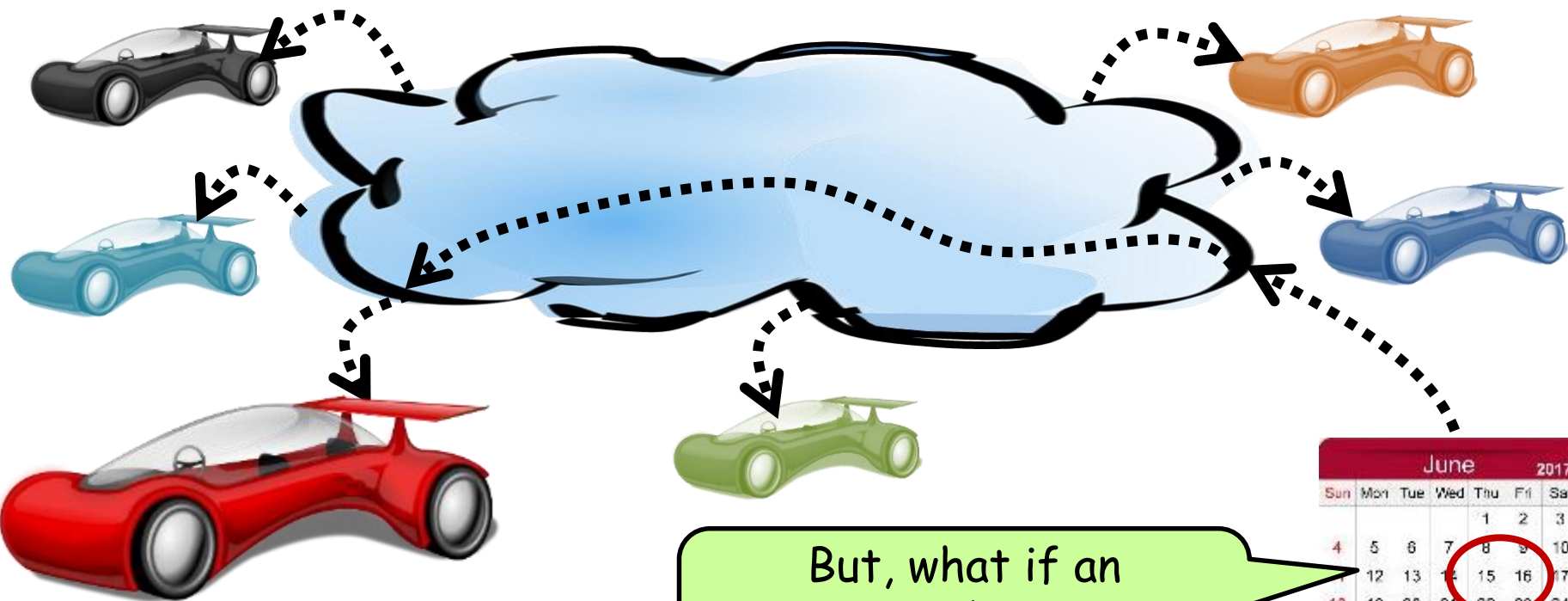
Continuous Online Learning



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

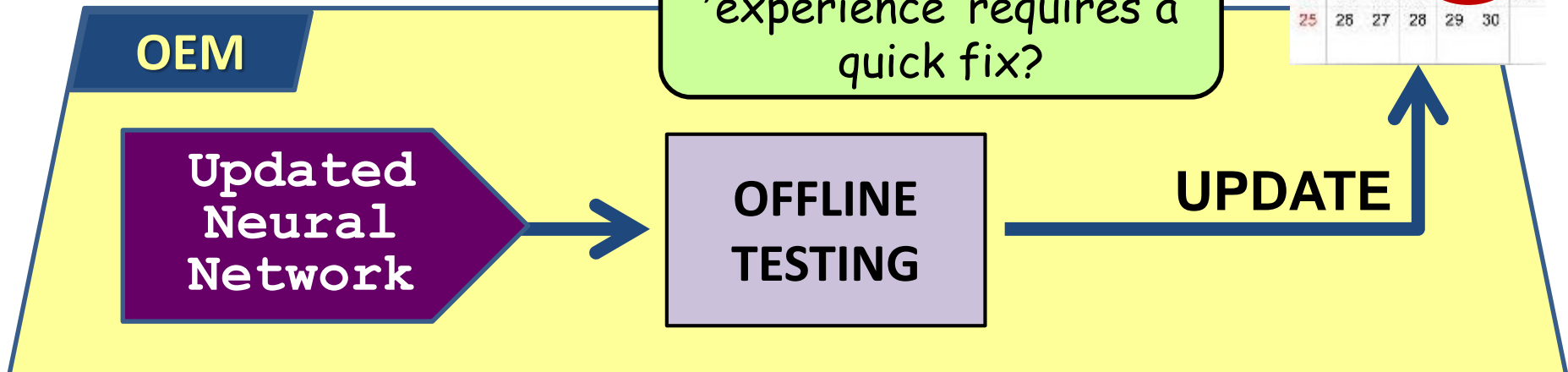


Performance Updates - Over-The-Air



June 2017						
Sun	Mon	Tue	Wed	Thu	Fri	Sat
				1	2	3
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	

But, what if an 'experience' requires a quick fix?



OEM

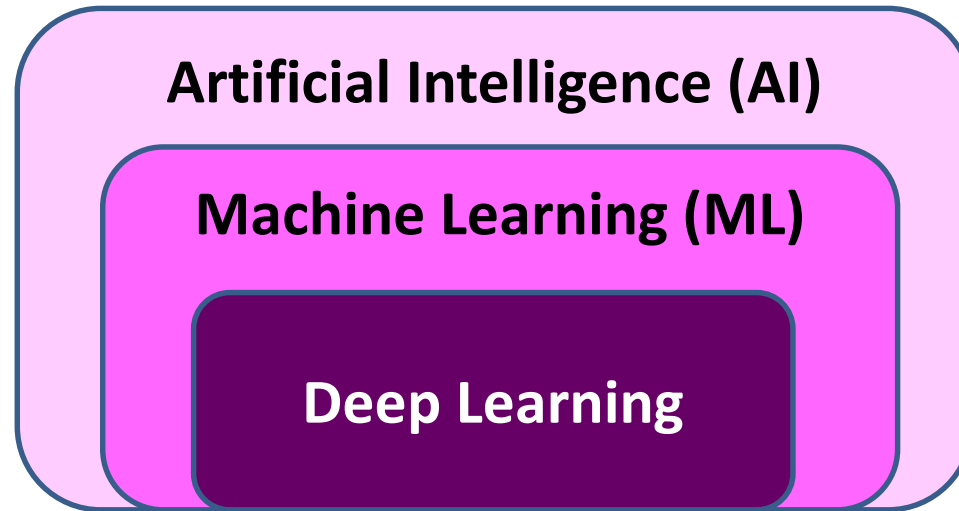
Updated Neural Network

OFFLINE TESTING

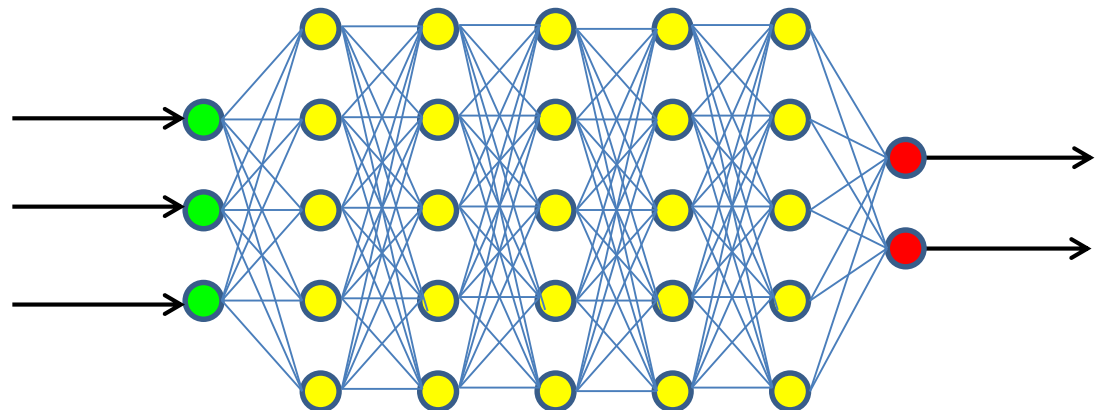
UPDATE

Machine Learning Challenges

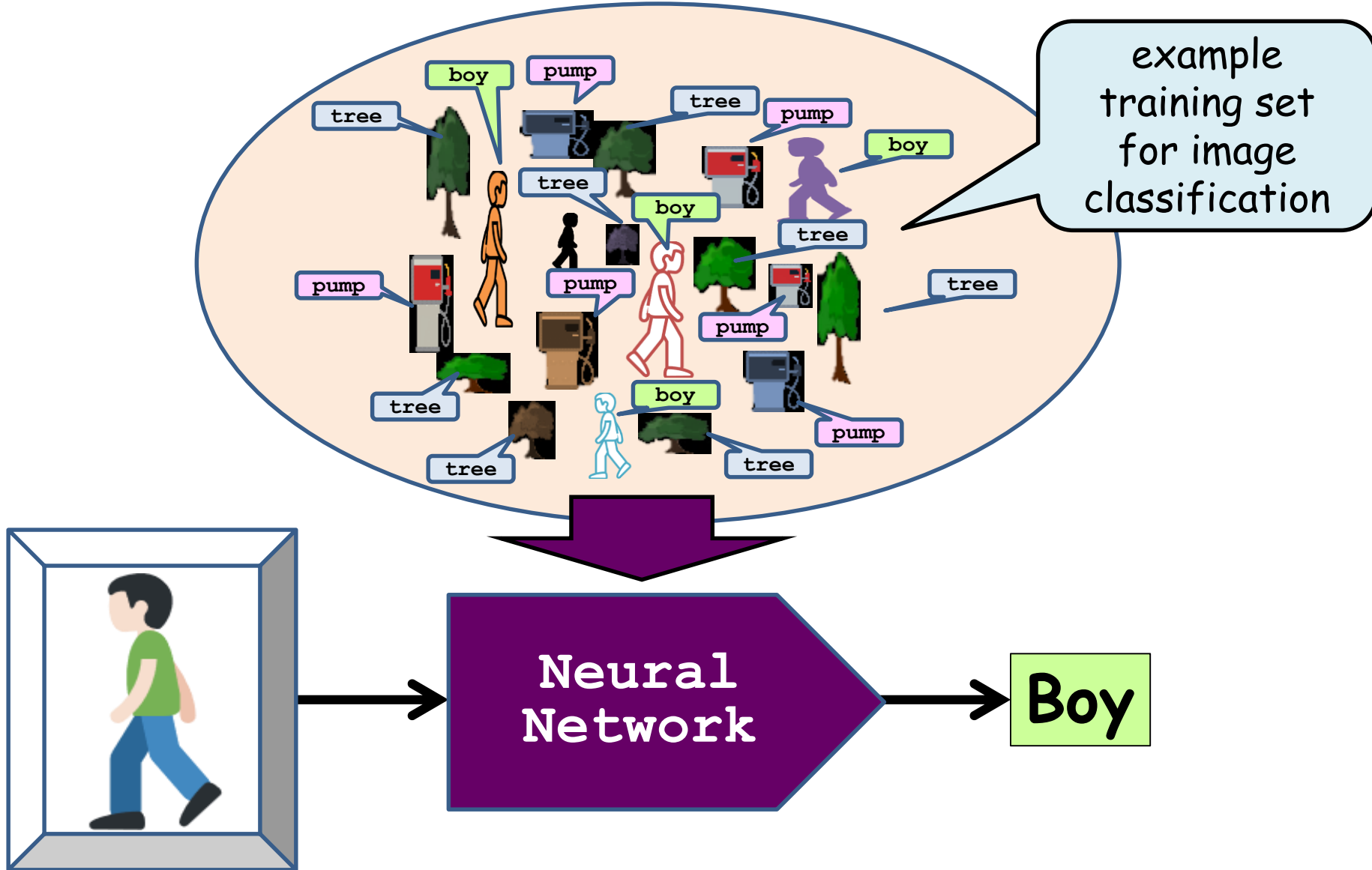
Deep Learning Systems



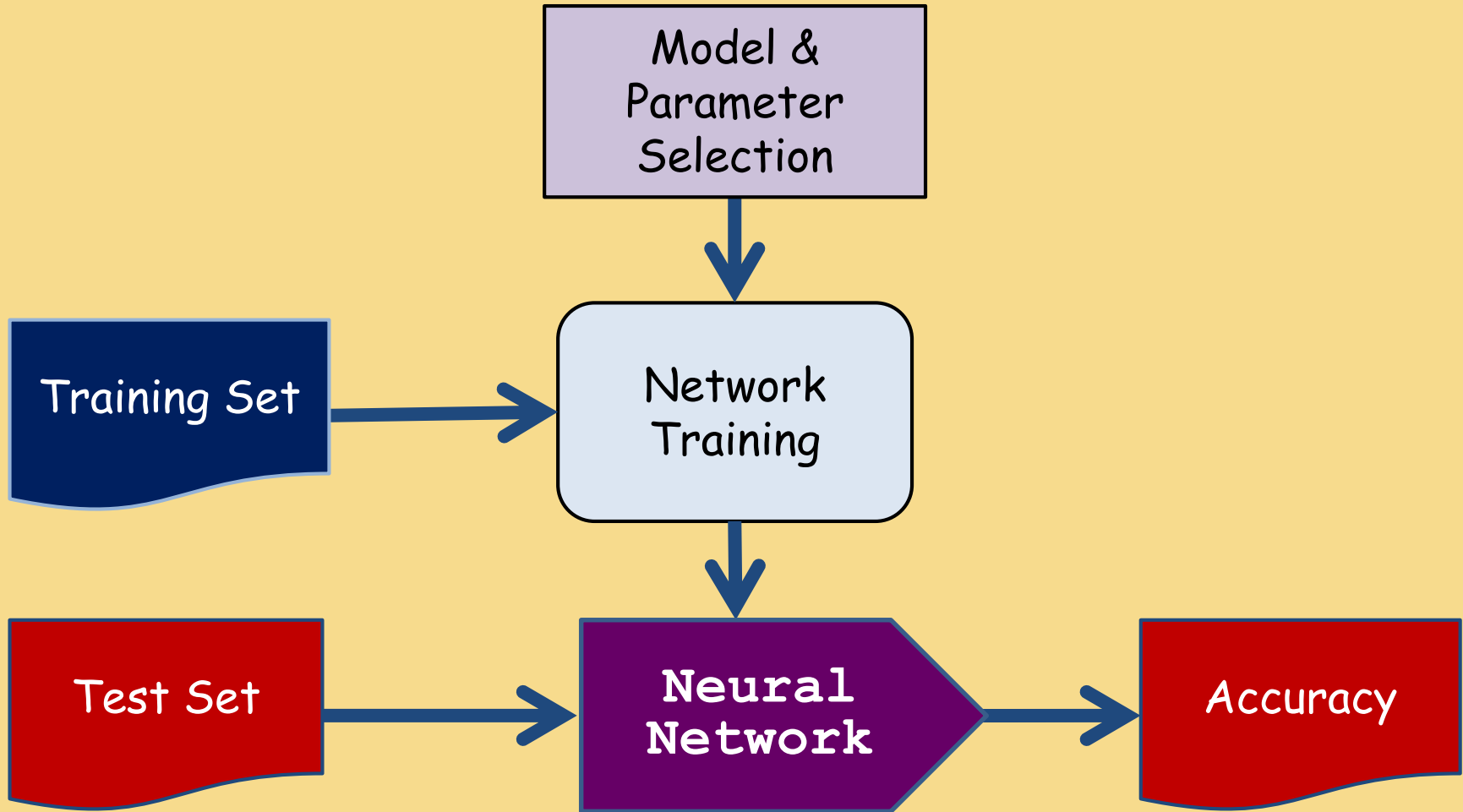
Deep
Neural
Network



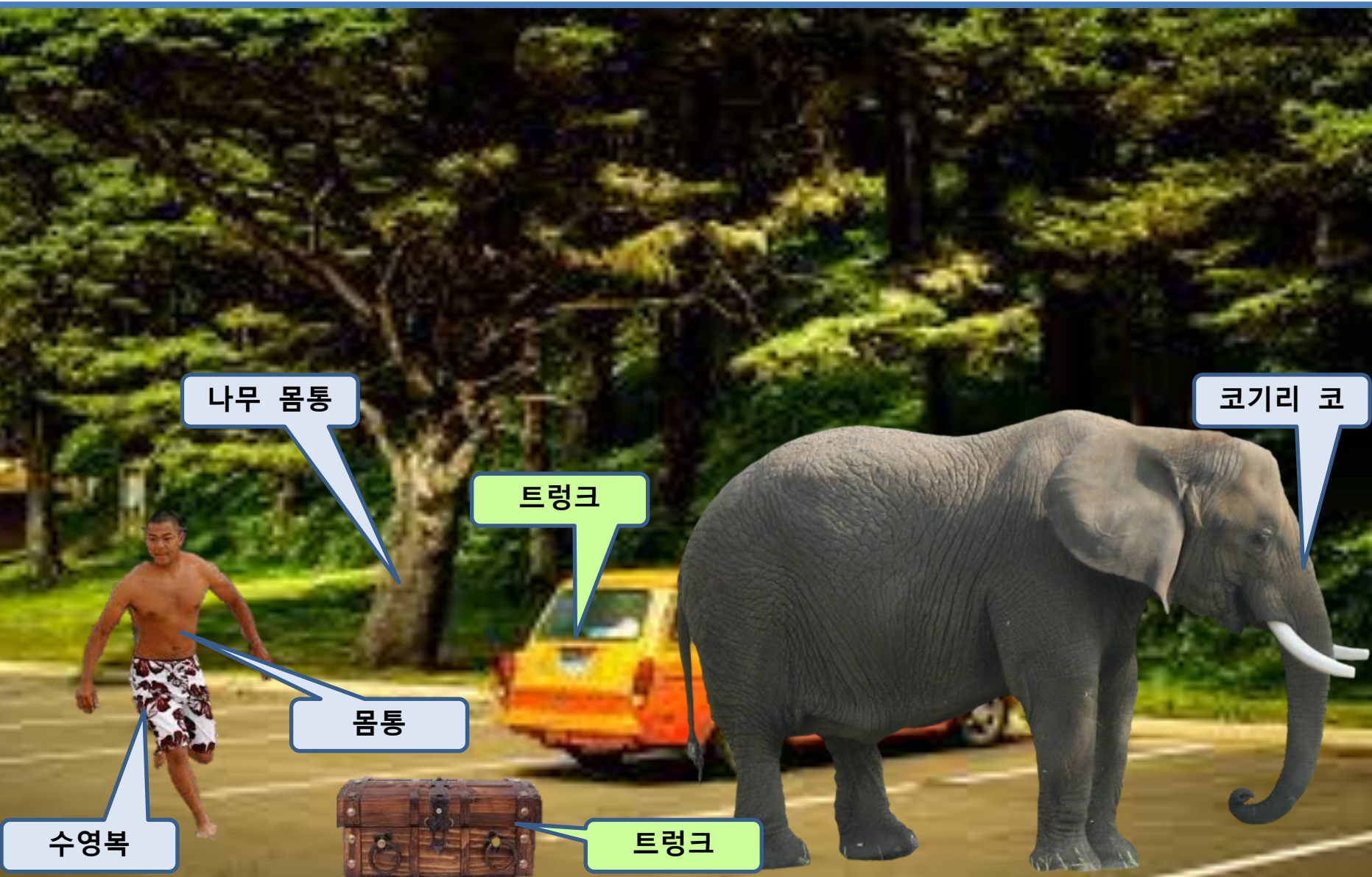
Example of Machine Learning



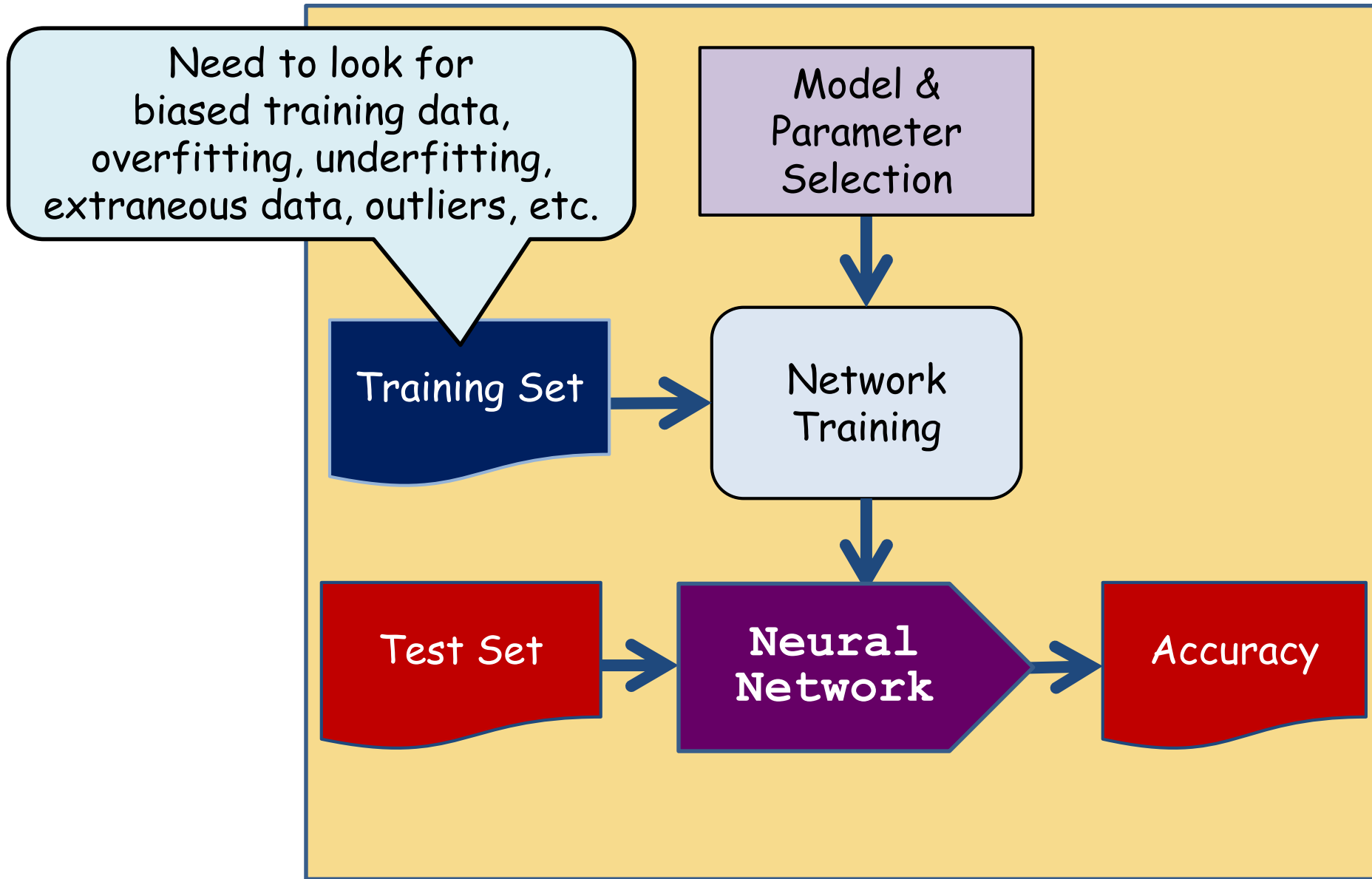
Supervised Machine Learning



잘못된 분류 - 한글?



Checking the Training Set



Misunderstanding – Data Bias



tank



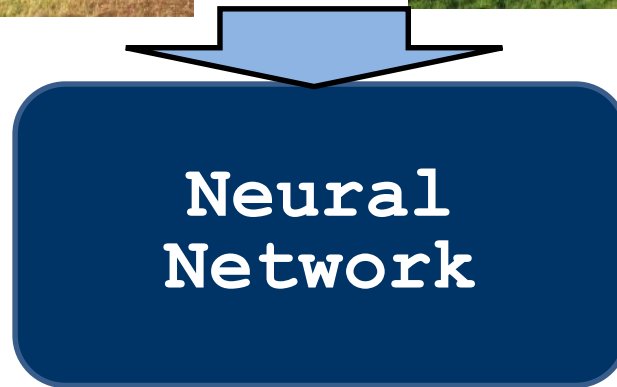
tank



tank

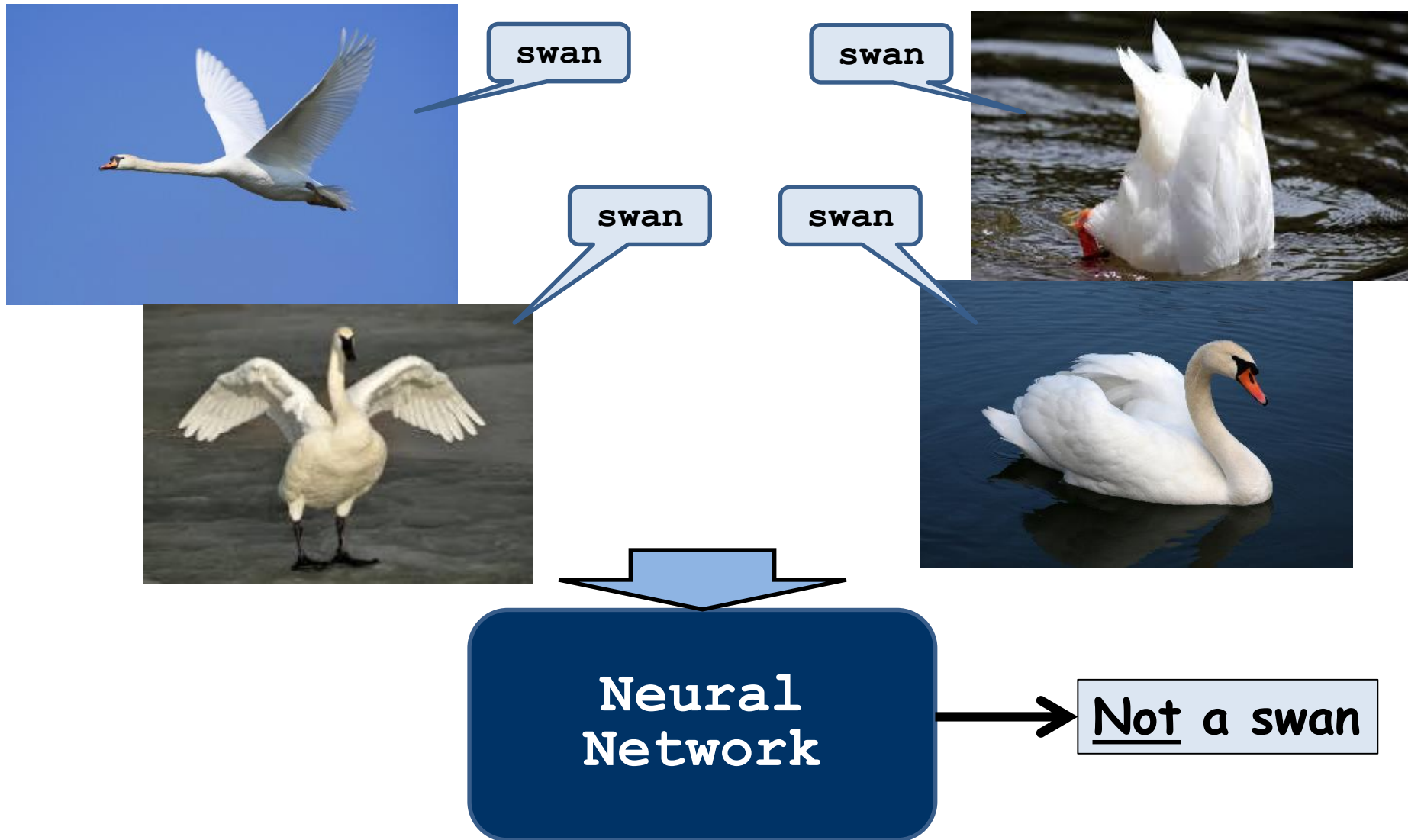


tank

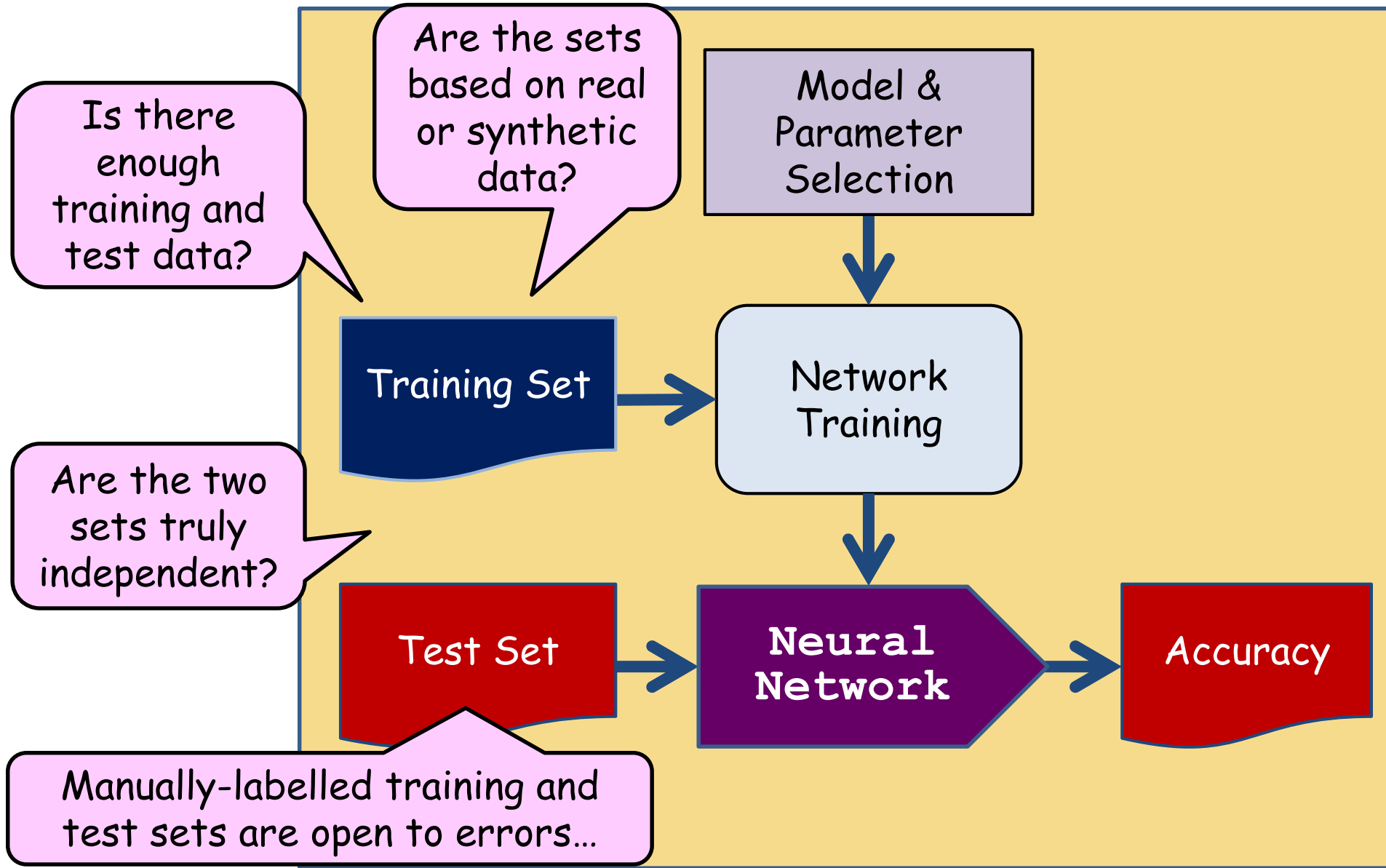


→ Not a tank

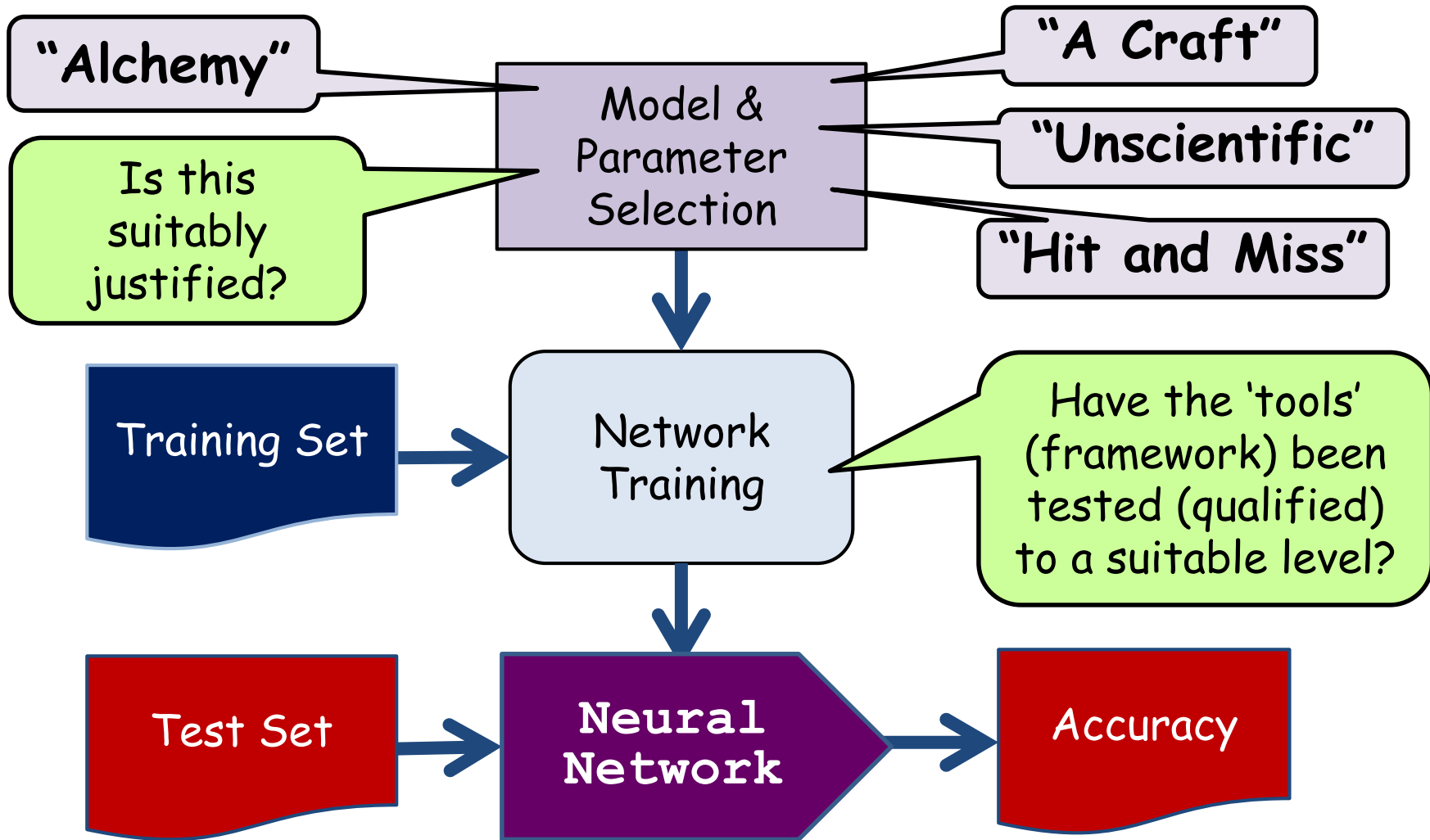
Incomplete Training Set



Checking the Training & Test Sets



Checking the Training

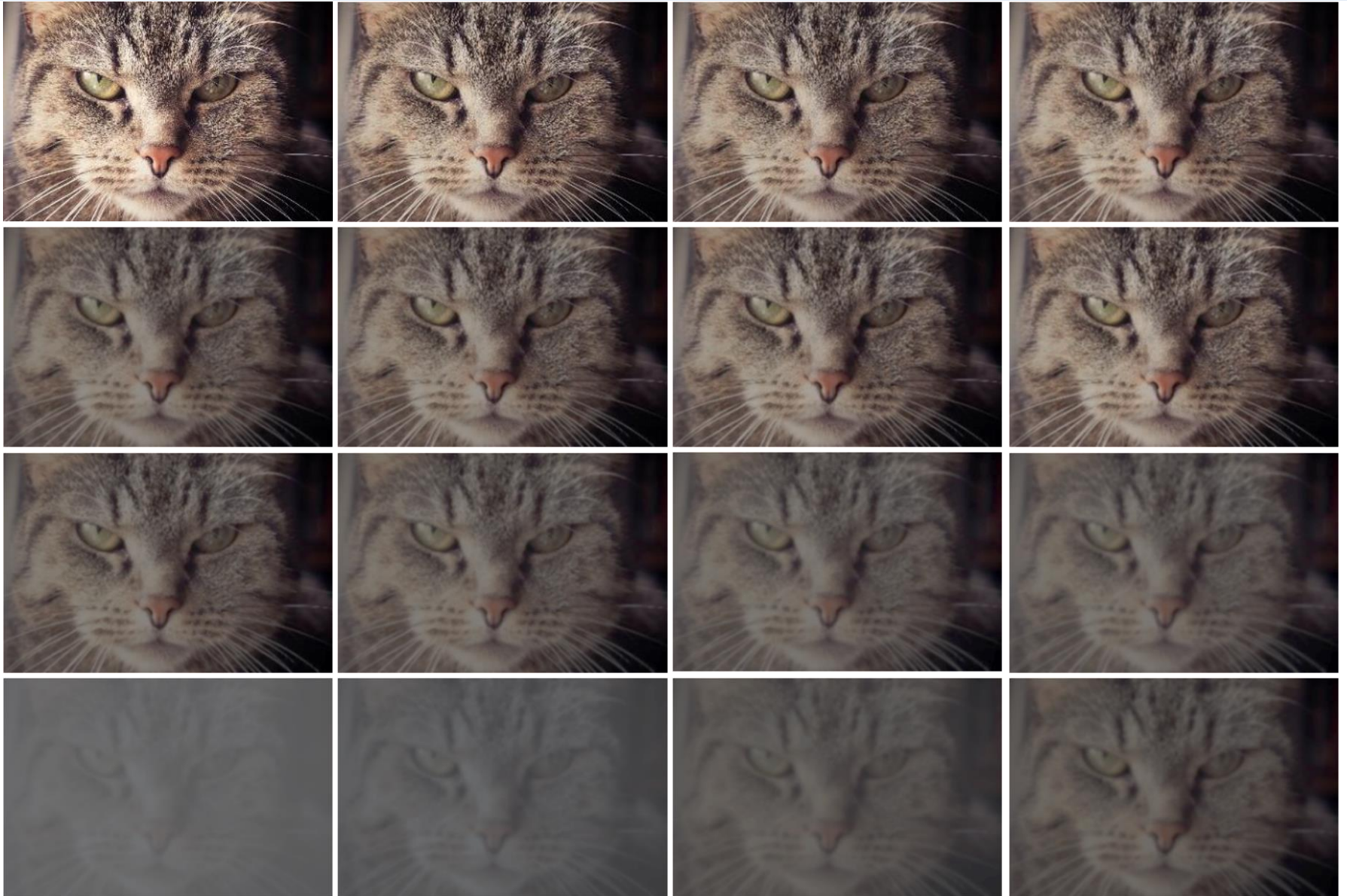


Black Box Testing of Autonomous Systems

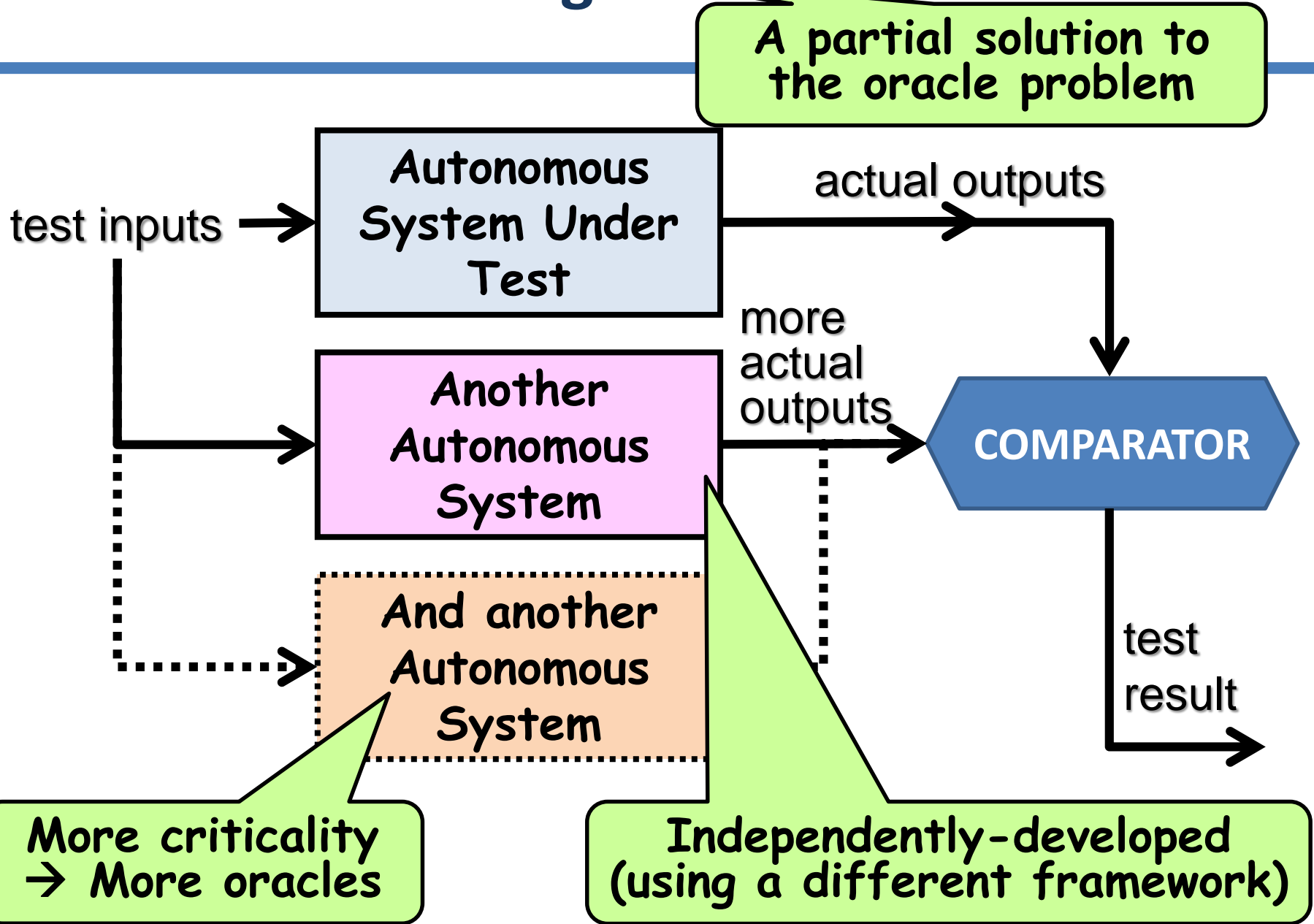
Test Challenges of Autonomous Systems

- **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...

Example - Sensor Degradation Testing

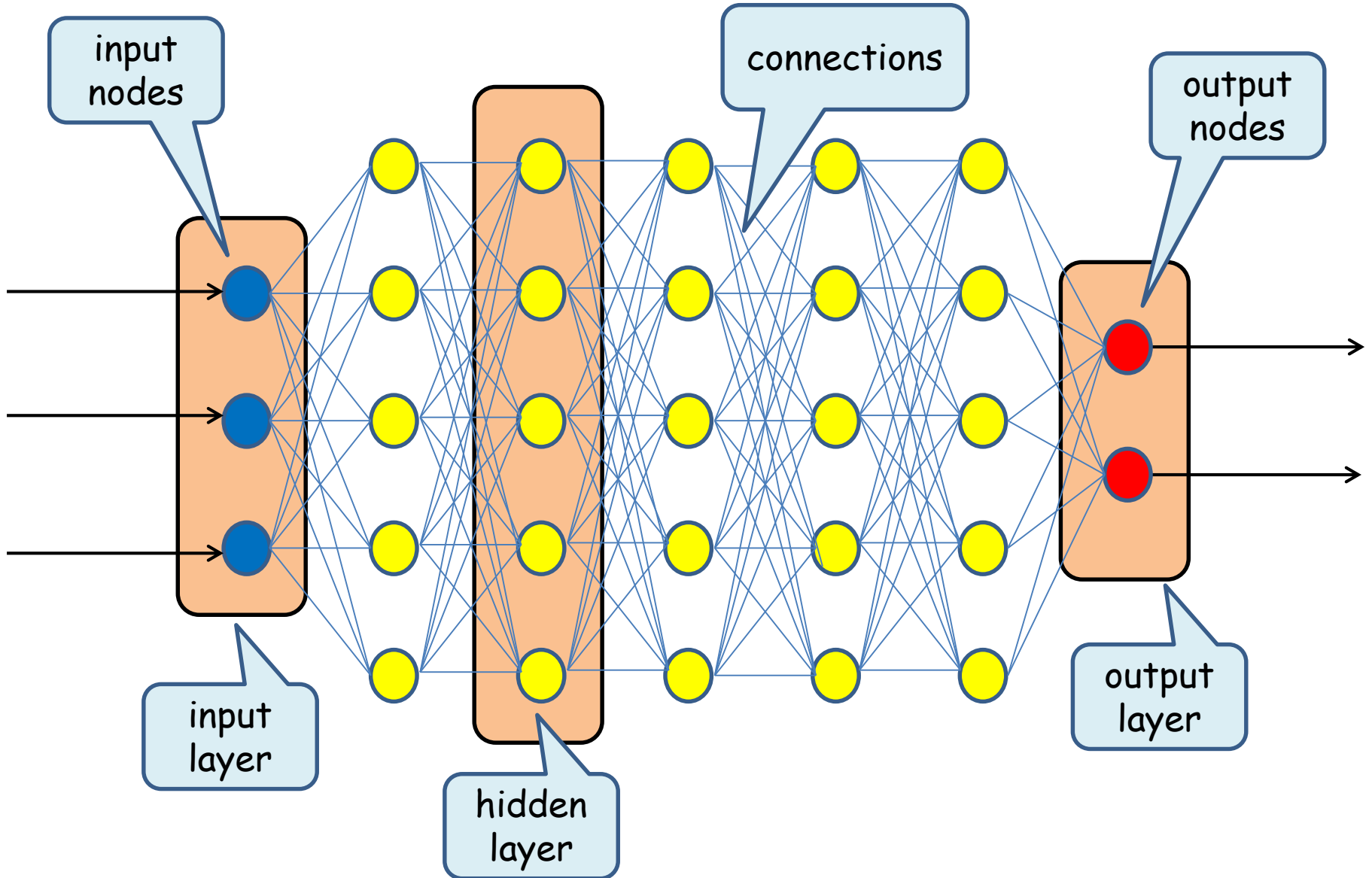


Back-to-Back Testing

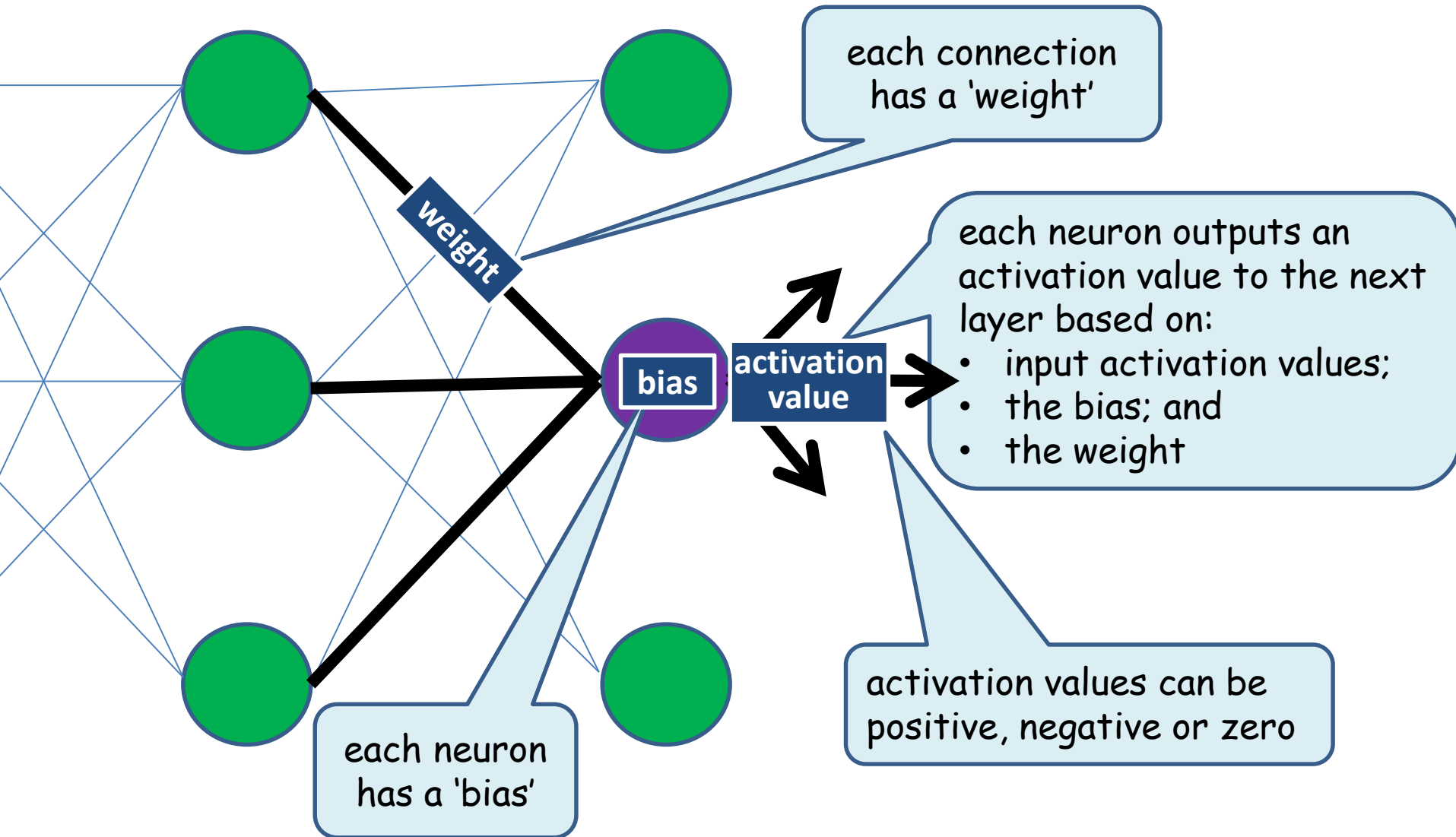


White Box Testing of Autonomous Systems

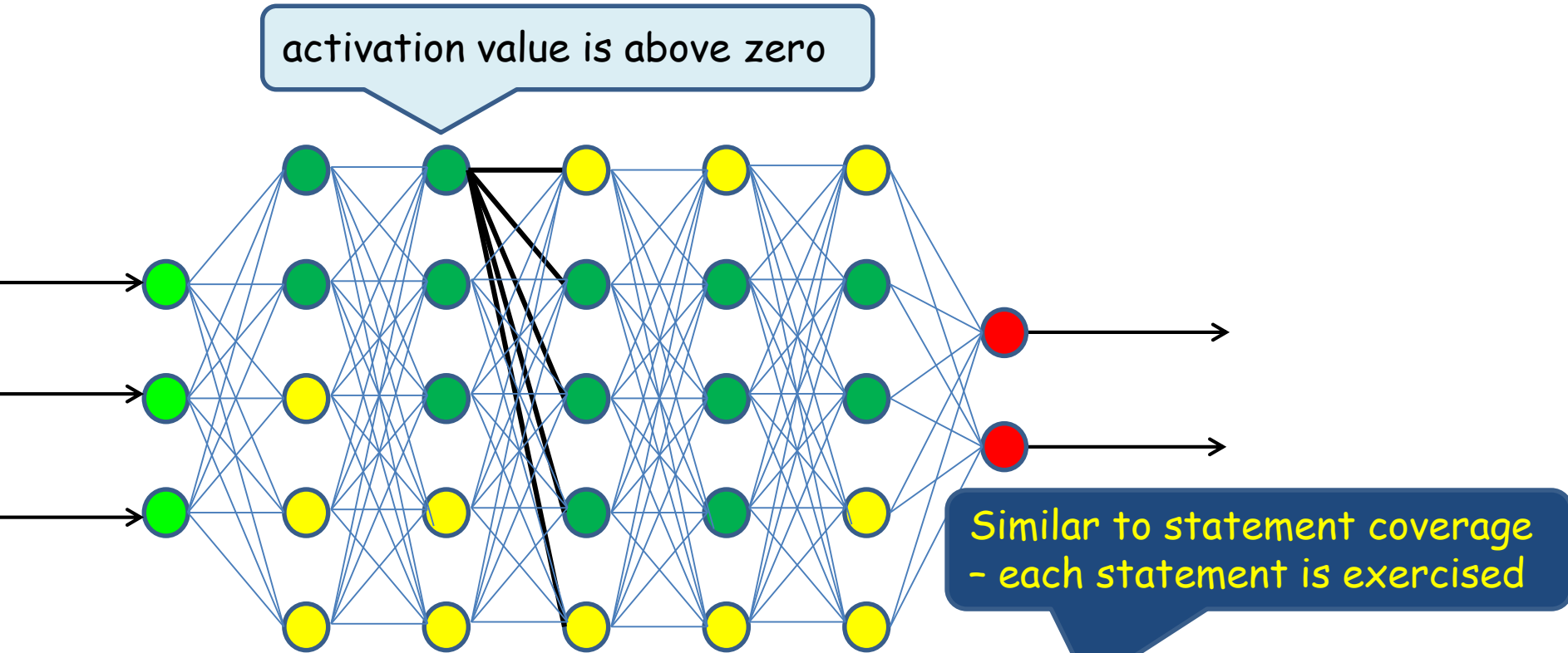
Deep Neural Net



Activation Values



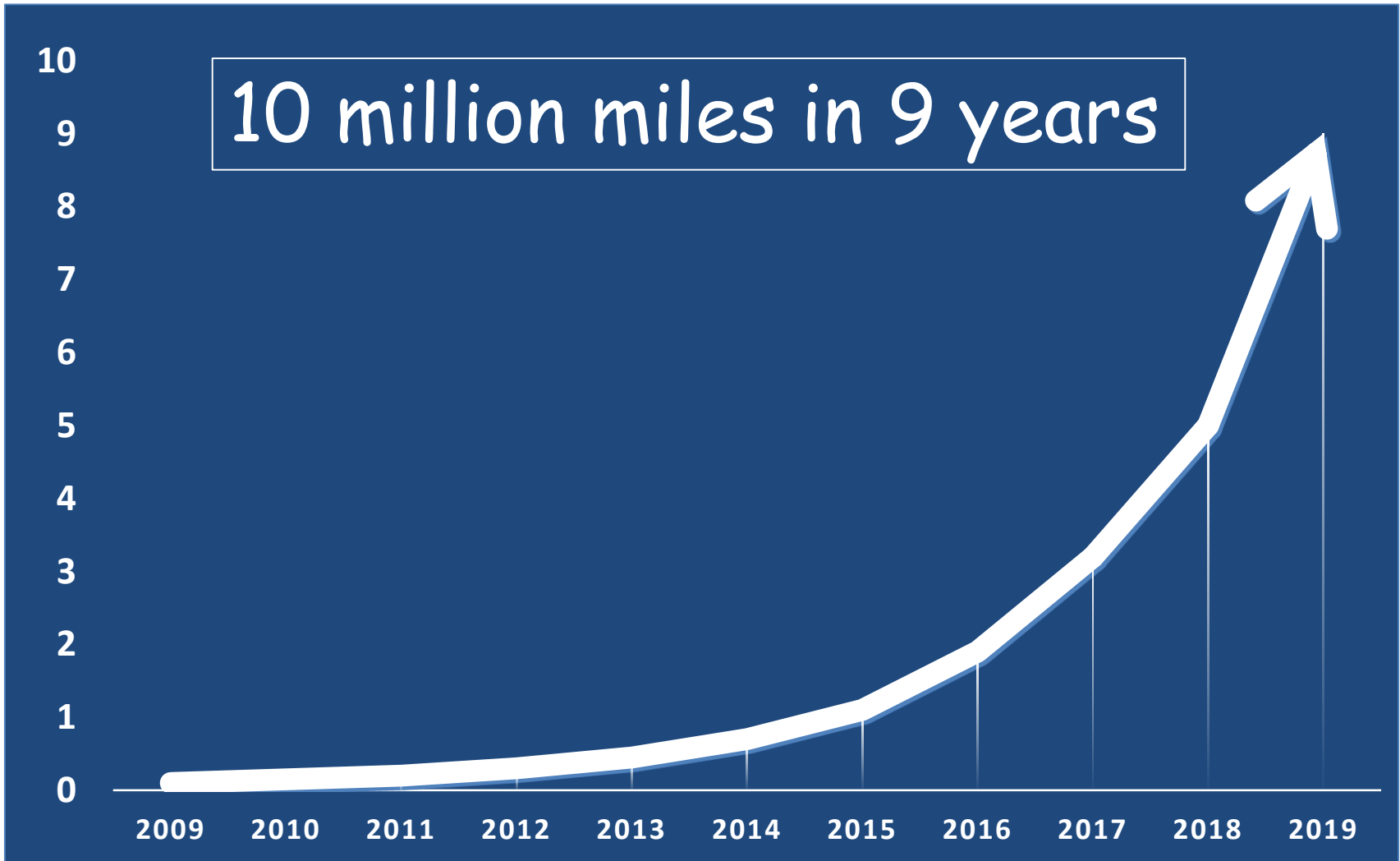
'Neuron' Coverage



Full 'neuron' coverage shows that every neuron is 'activated' (value above zero) at least once (but - basic coverage - typically finds no adversarial examples)

***The Necessity of
Virtual
Test Environments***

Waymo On-Road Test Miles (millions)

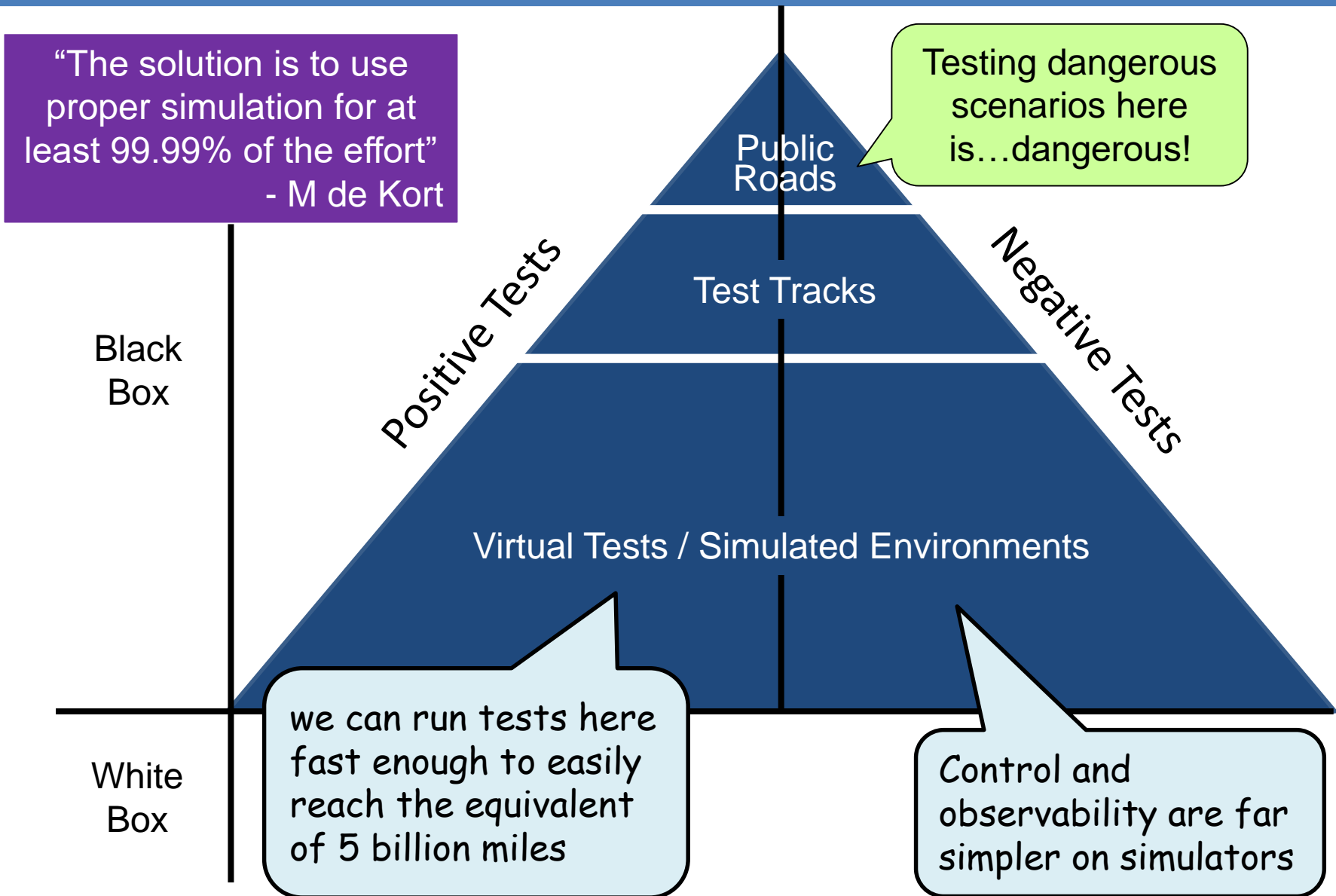


20% Better (than human drivers)



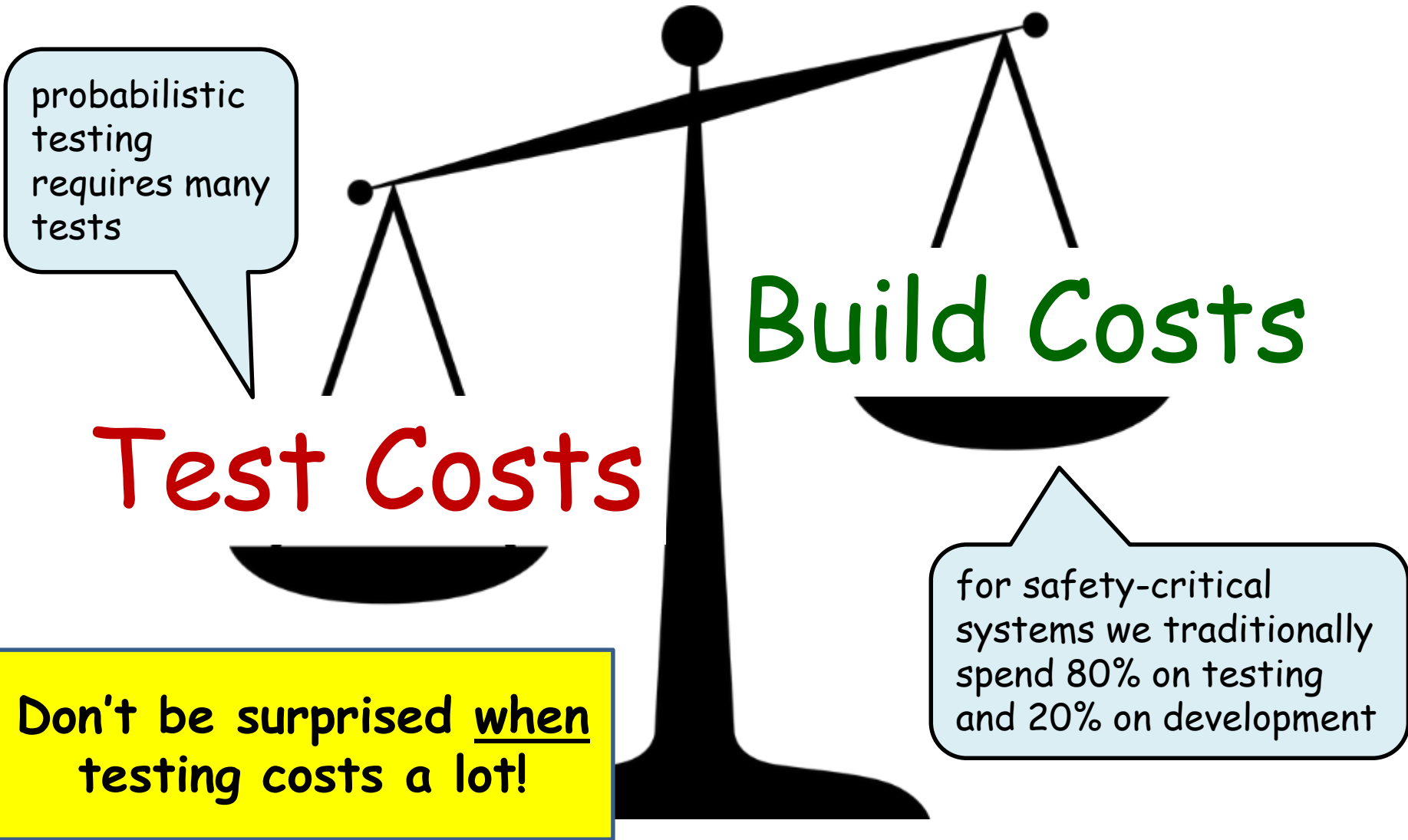
500x what has
gone before =
5 Billion Miles

Autonomous Cars – Test Environments



Conclusions

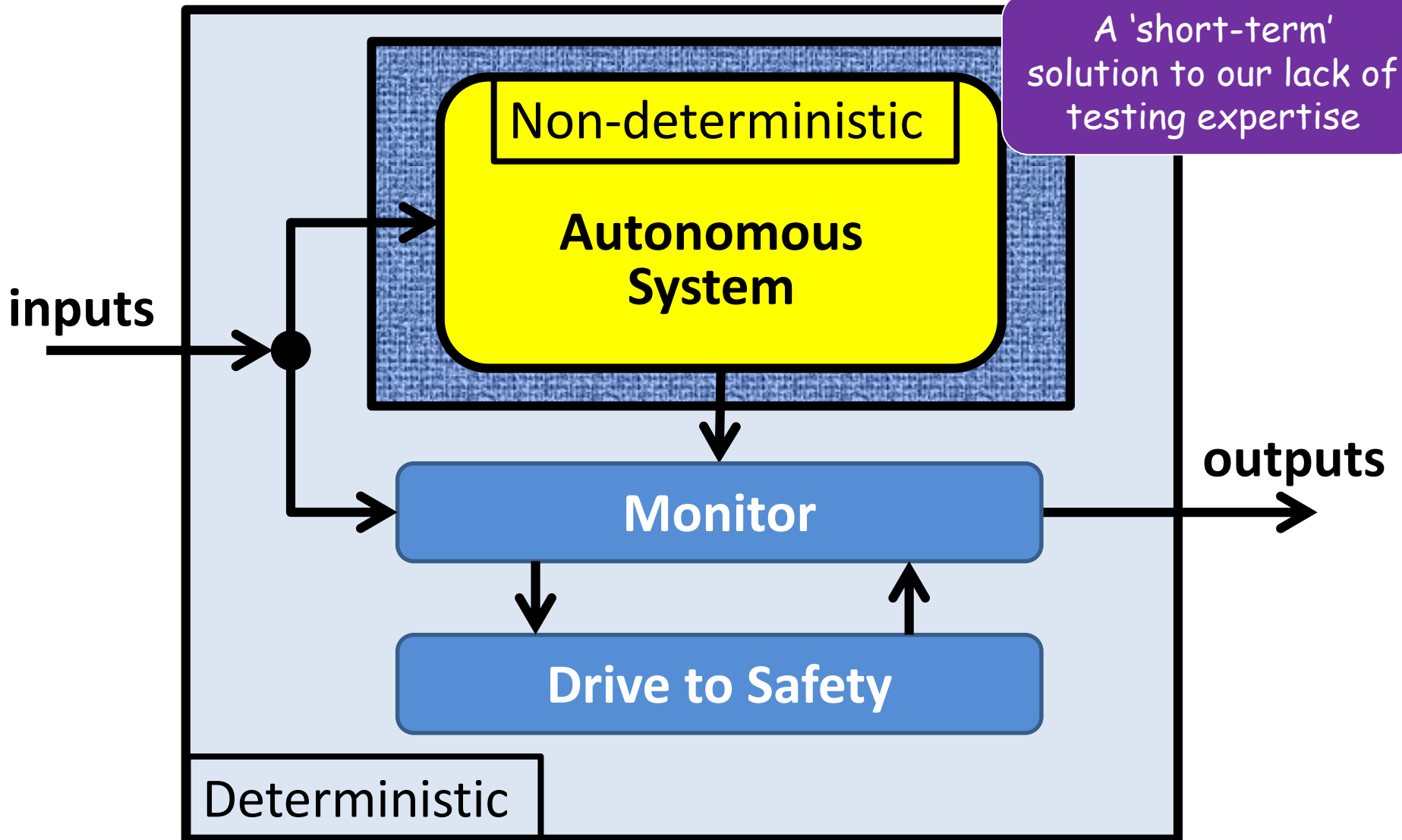
Autonomous System Costs



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?