



Introduction to Artificial Intelligence for CyberSecurity Applications

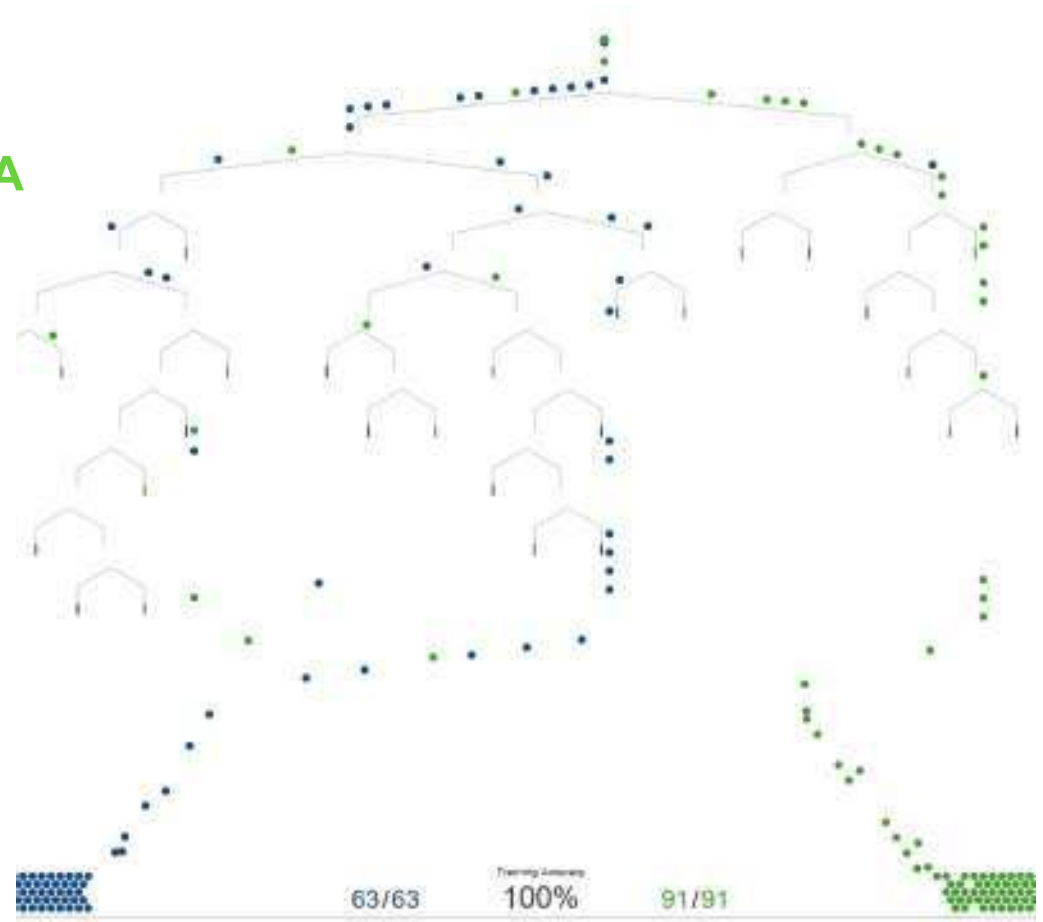
Rob Collins
Director – Sales Engineering, APAC

GLOSSARY

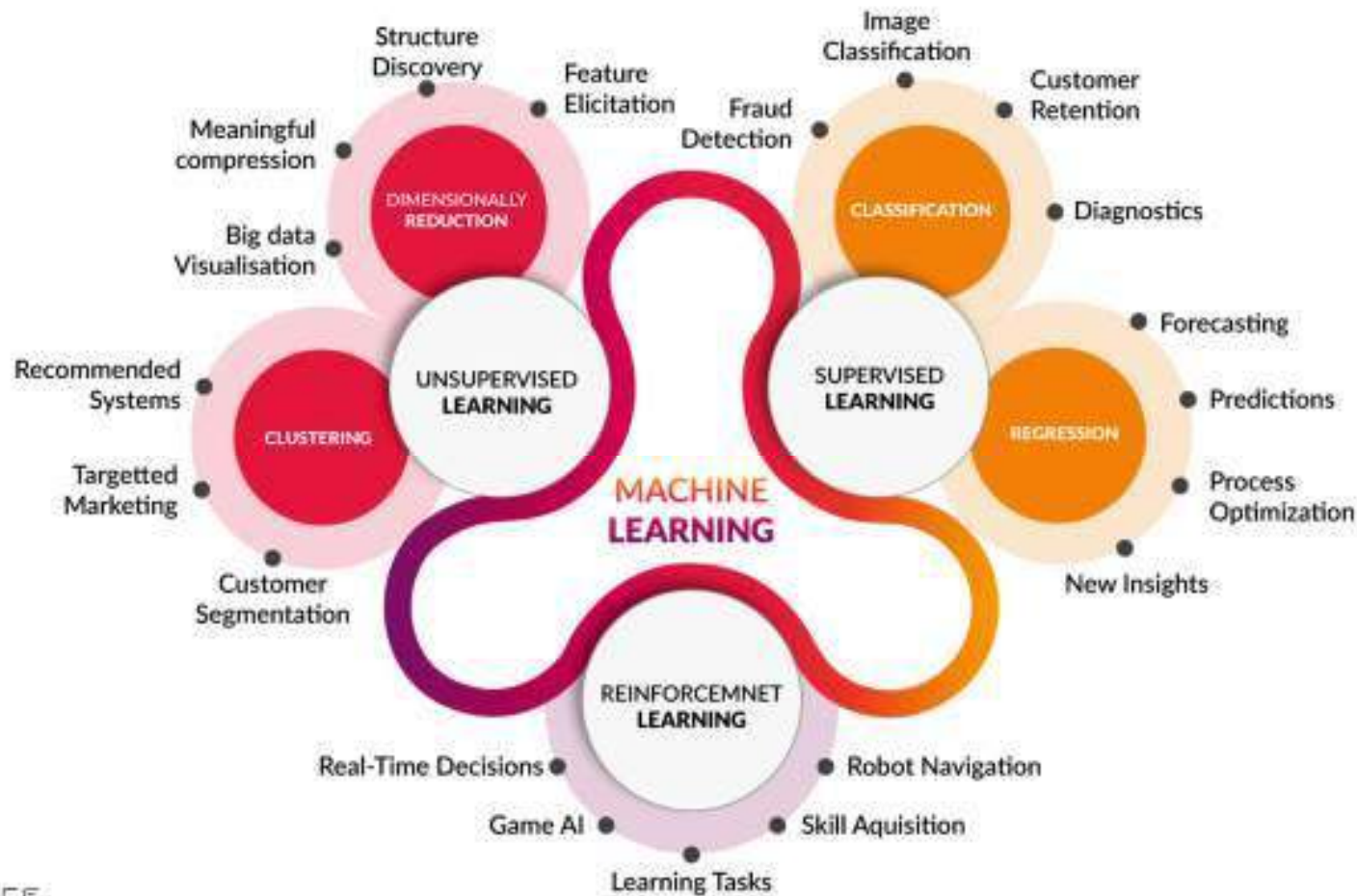
- k-means – simple clustering algorithm
- DBSCAN – more advanced clustering algorithm
- NB – Naïve Bayes classifier model
- GMM – Gaussian Mixture Model clustering algorithm
- LSTM – Long Short-Term Memory Neural Network algorithm
- CNN – Convolutional Neural Network
- RNN – Recurrent Neural Network
- LR – Logistic Regression classifier
- DT – Decision Tree

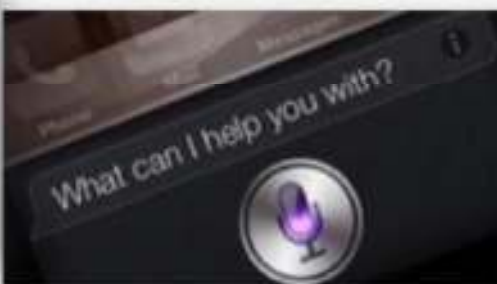
**MACHINE LEARNING IS A
FIELD OF STUDY THAT
GIVES COMPUTERS THE
ABILITY TO LEARN
WITHOUT EXPLICITLY
BEING PROGRAMMED**

- Arthur Samuel, 1959



MACHINE LEARNING WILL BE EVERYWHERE





MACHINE LEARNING WILL BE EVERYWHERE

The trend to incorporate ML capabilities into new and existing security products will continue apace. According to an April 2016 Gartner report:

- By 2018, 25% of security products used for detection will have some form of machine learning built into them.
- By 2018, prescriptive analytics will be deployed in at least 10% of UEBA products to automate response to incidents, up from zero today.

Gartner Core Security, *The Fast-Evolving State of Security Analytics*, April, 2016, Report ID: G00298030 accessed at <https://hs.coresecurity.com/gartnerreprint-2017>

TWO THINGS CAME TOGETHER TO ENABLE AI

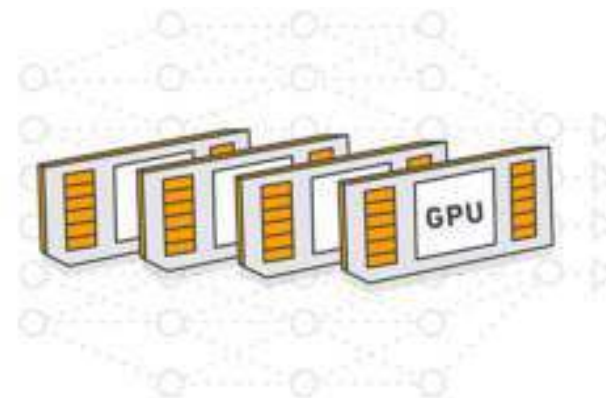
- **Big Data**

- Large collections of Spam, malware, exploits, network traffic, user behaviors

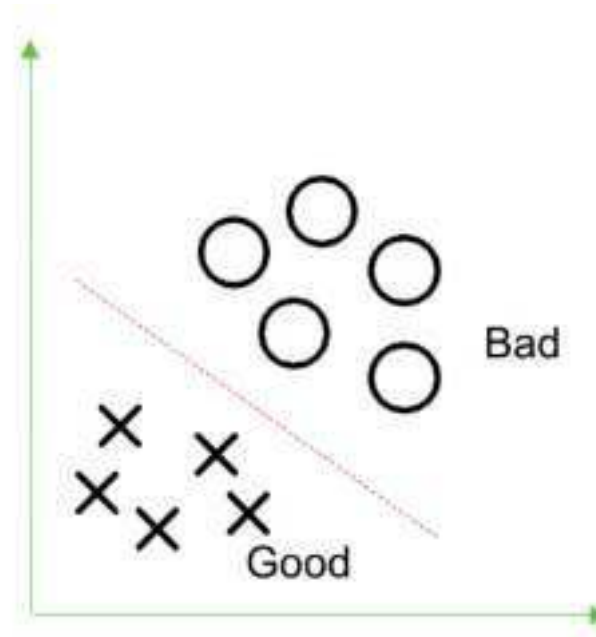


- **Cloud Computing Power**

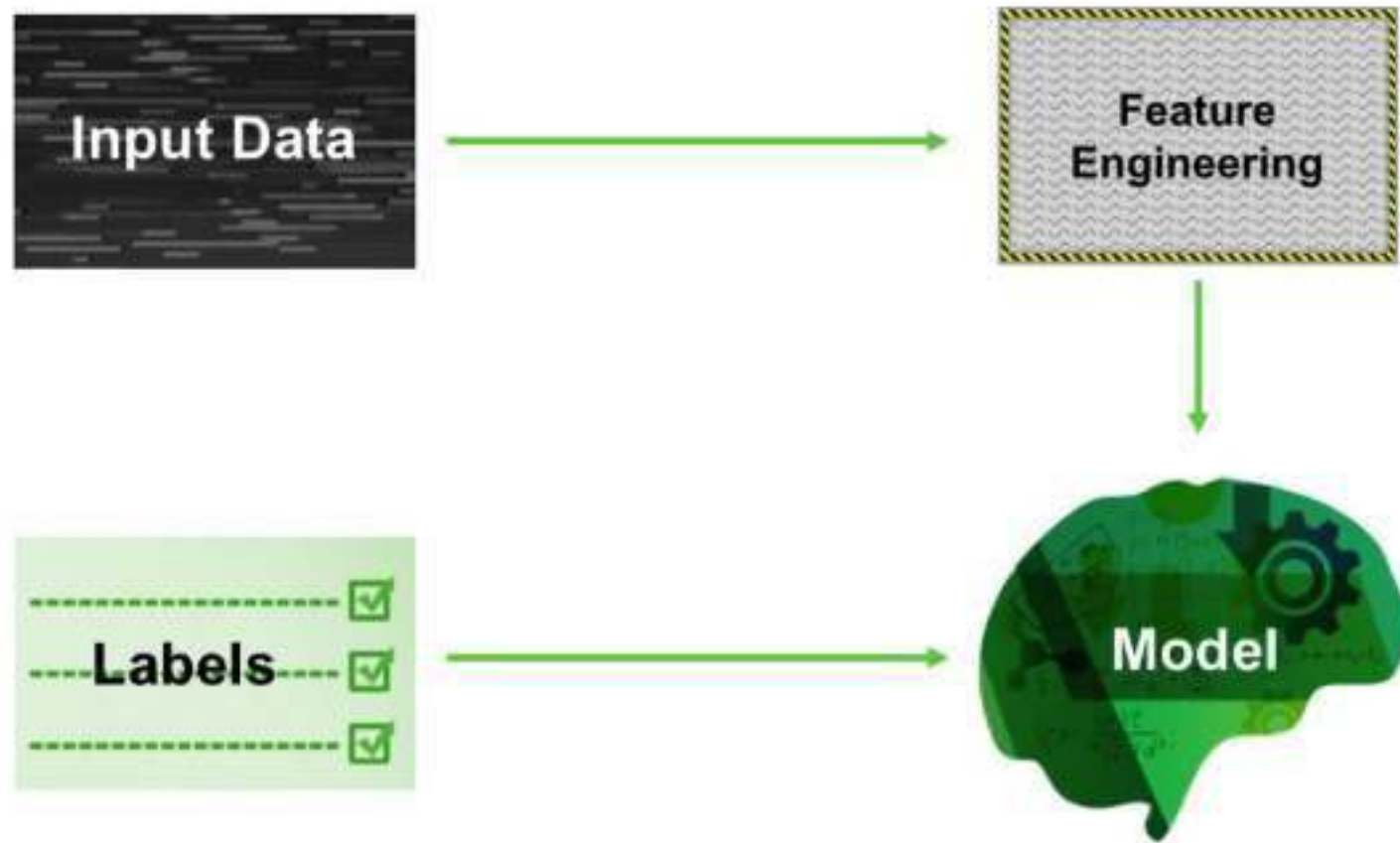
- Possible to consume over 100,000 CPU/GPU cores



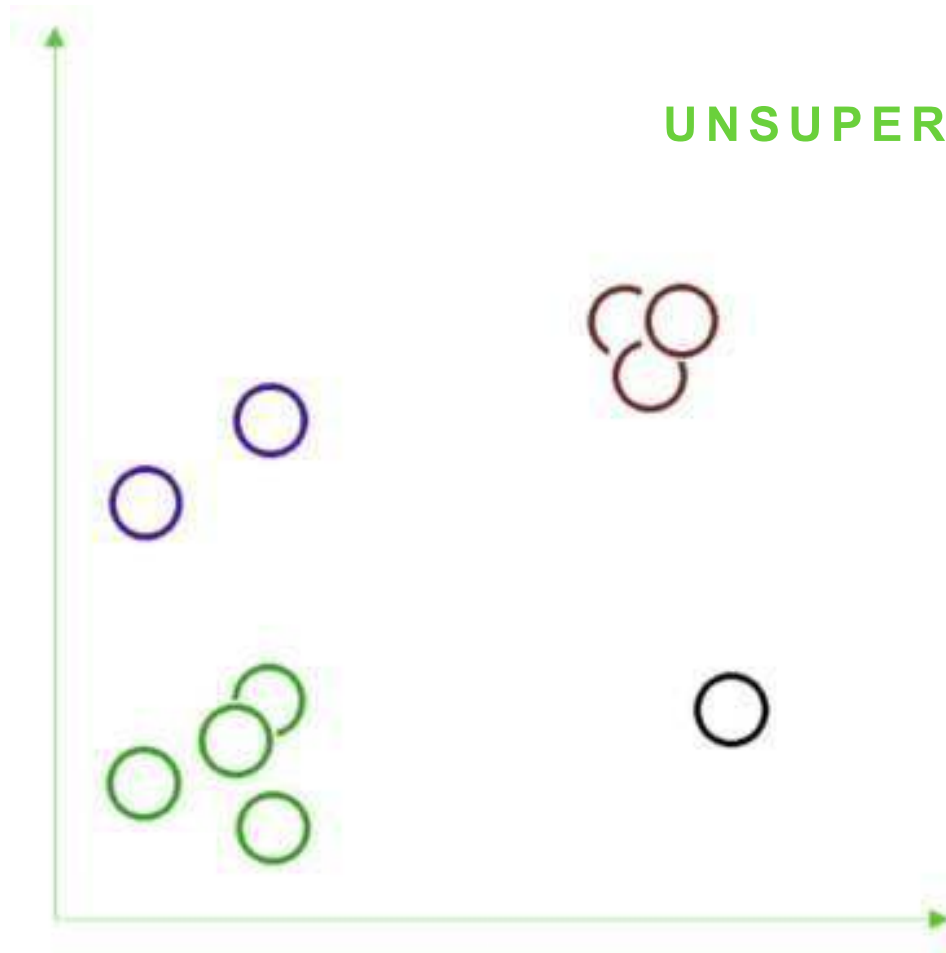
SUPERVISED



SUPERVISED PROCESS



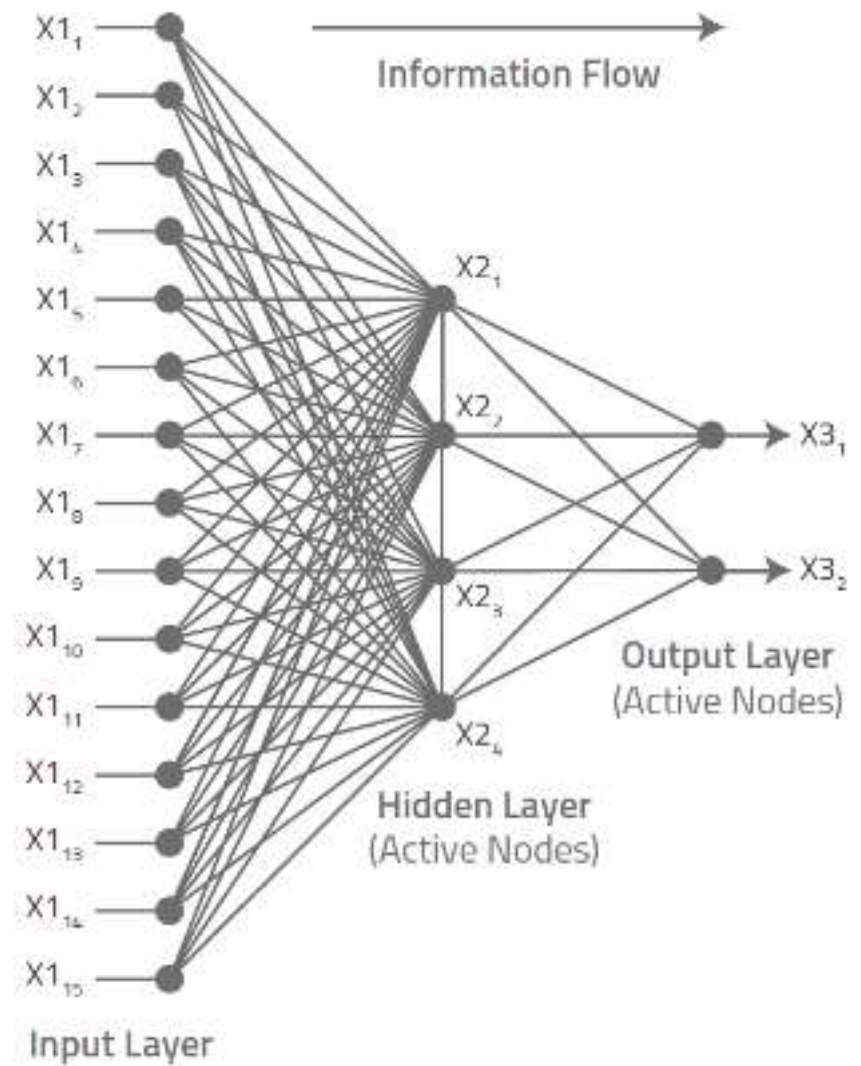
UNSUPERVISED

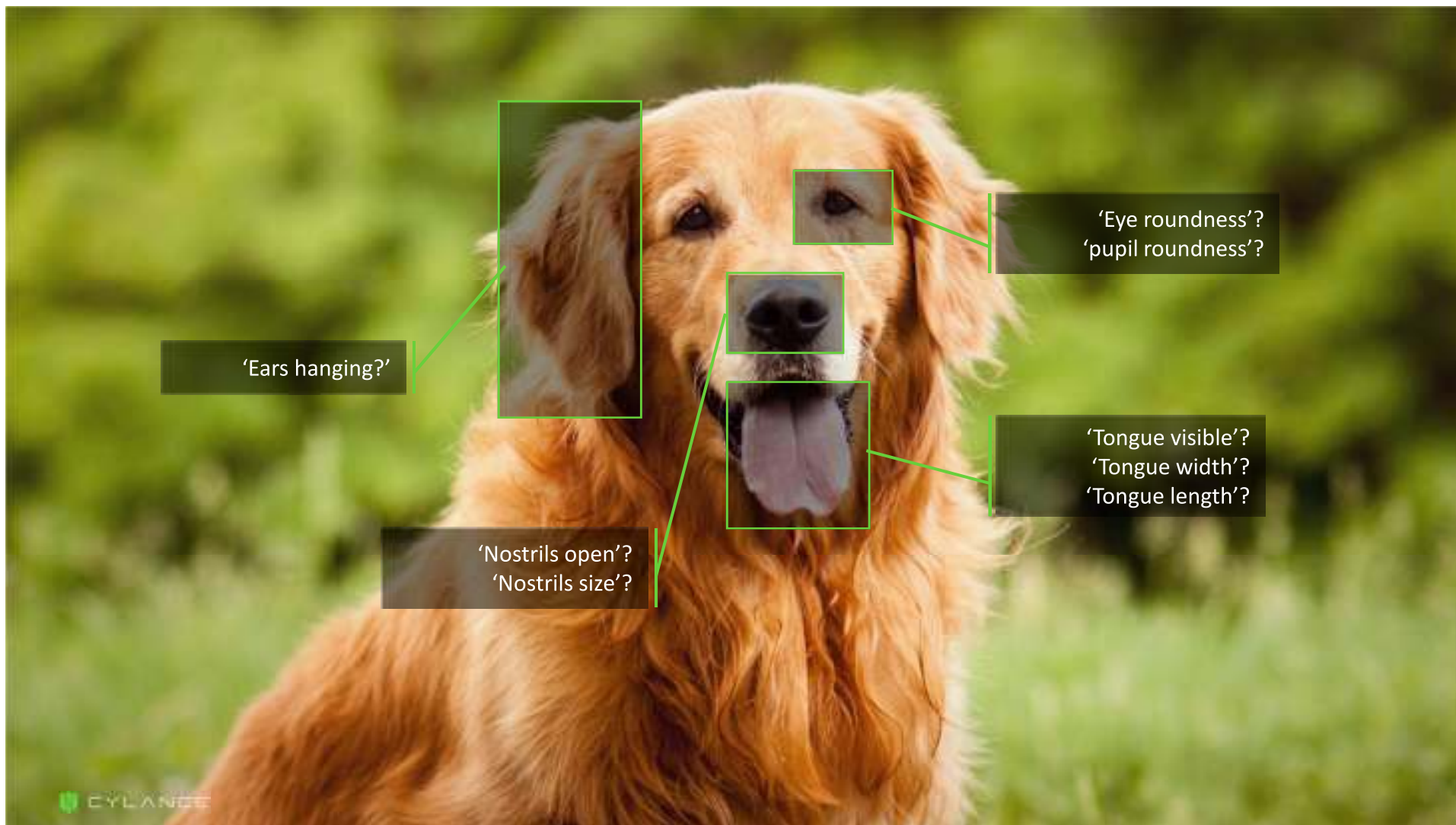


No labels, asking for clusters

NEURAL NETWORK

- Works like a human brain – useful connections remain, others dropped



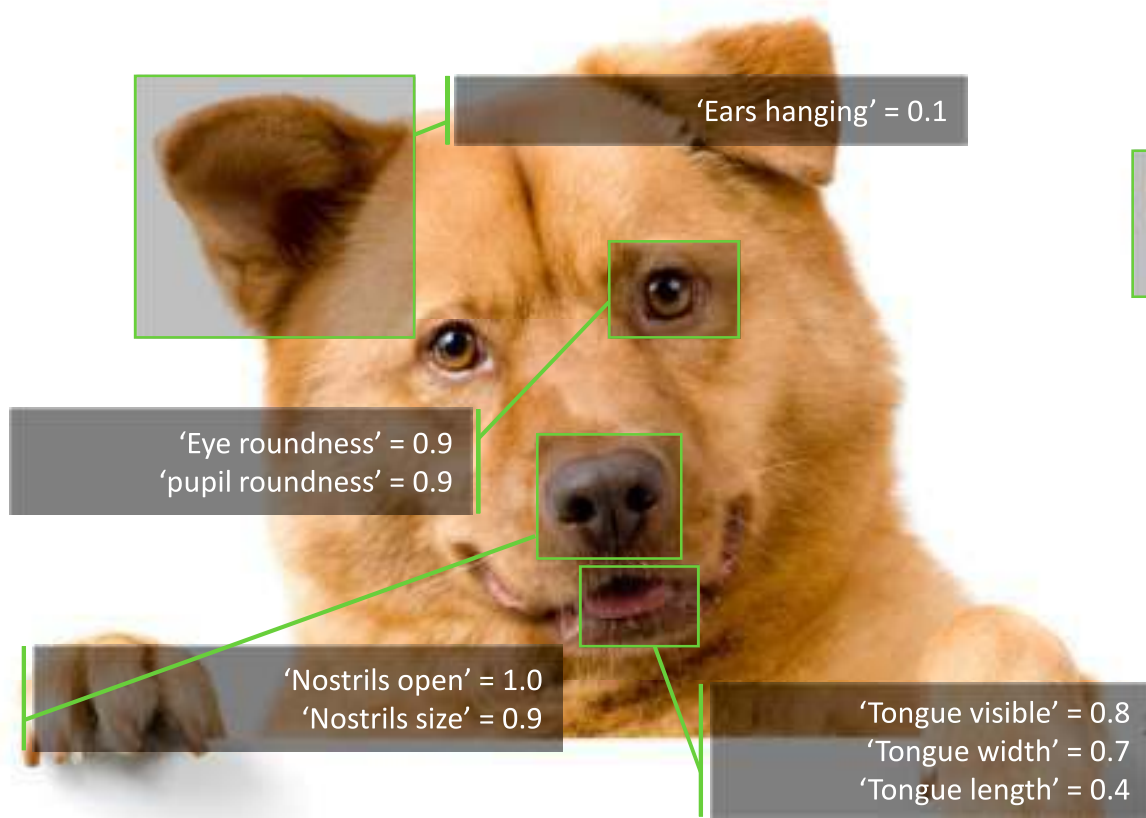


'Ears hanging?'

'Eye roundness'?
'pupil roundness'?

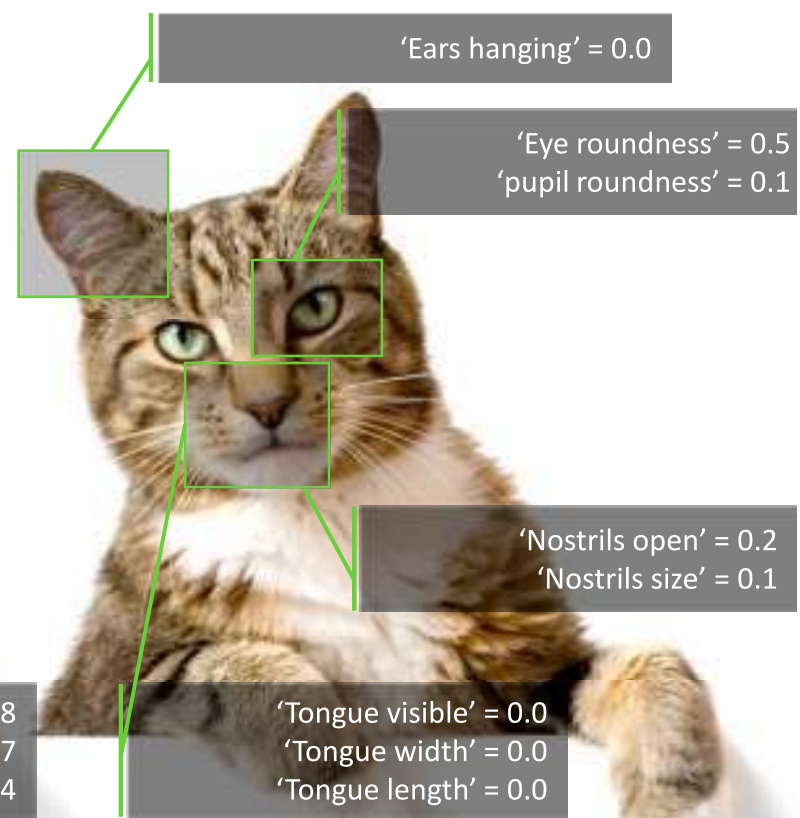
'Tongue visible'?
'Tongue width'?
'Tongue length'?

'Nostrils open'?
'Nostrils size'?



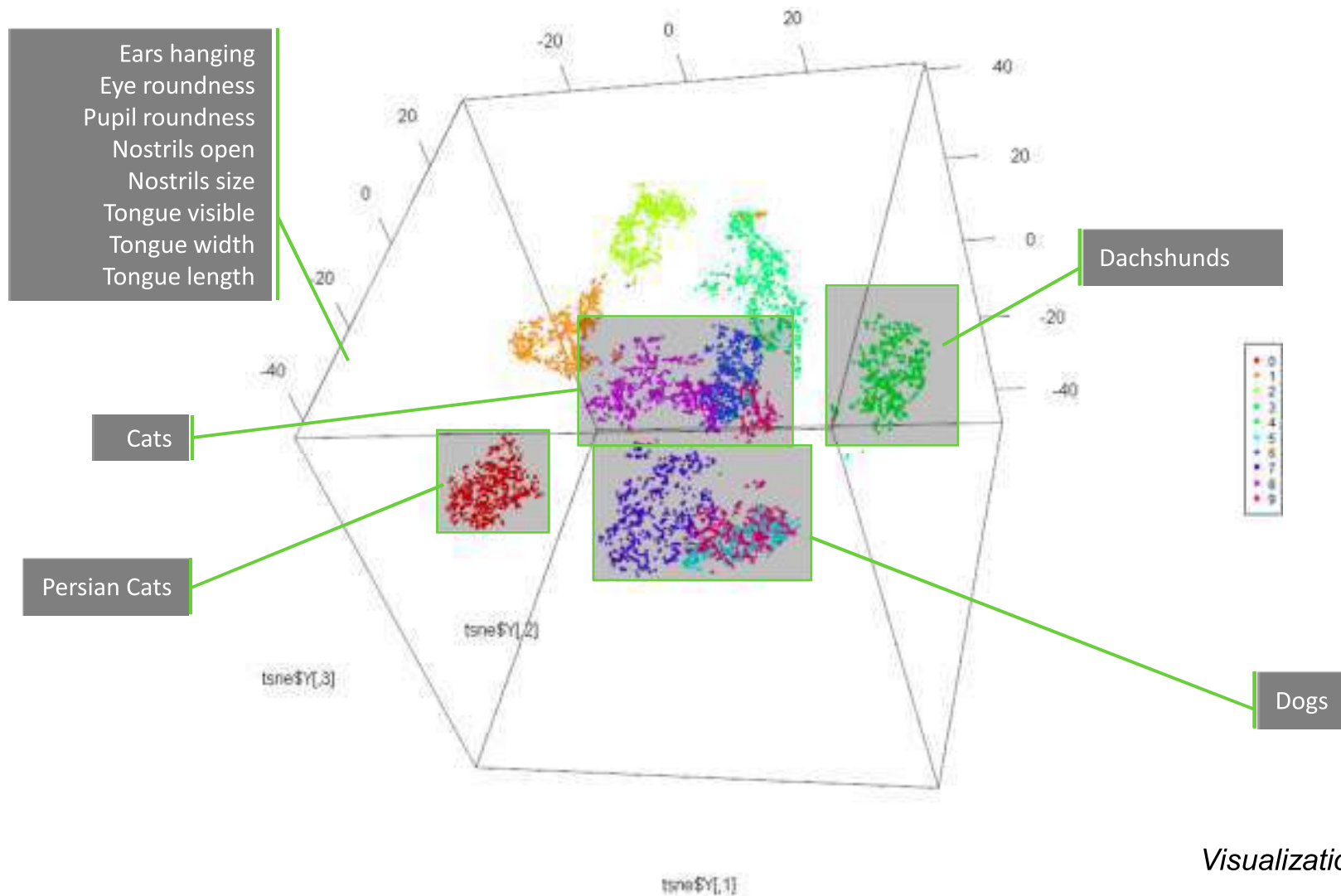
(0.1, 0.9, 0.9, 1.0, 0.9, 0.8, 0.7, 0.4)

this dog's feature vector

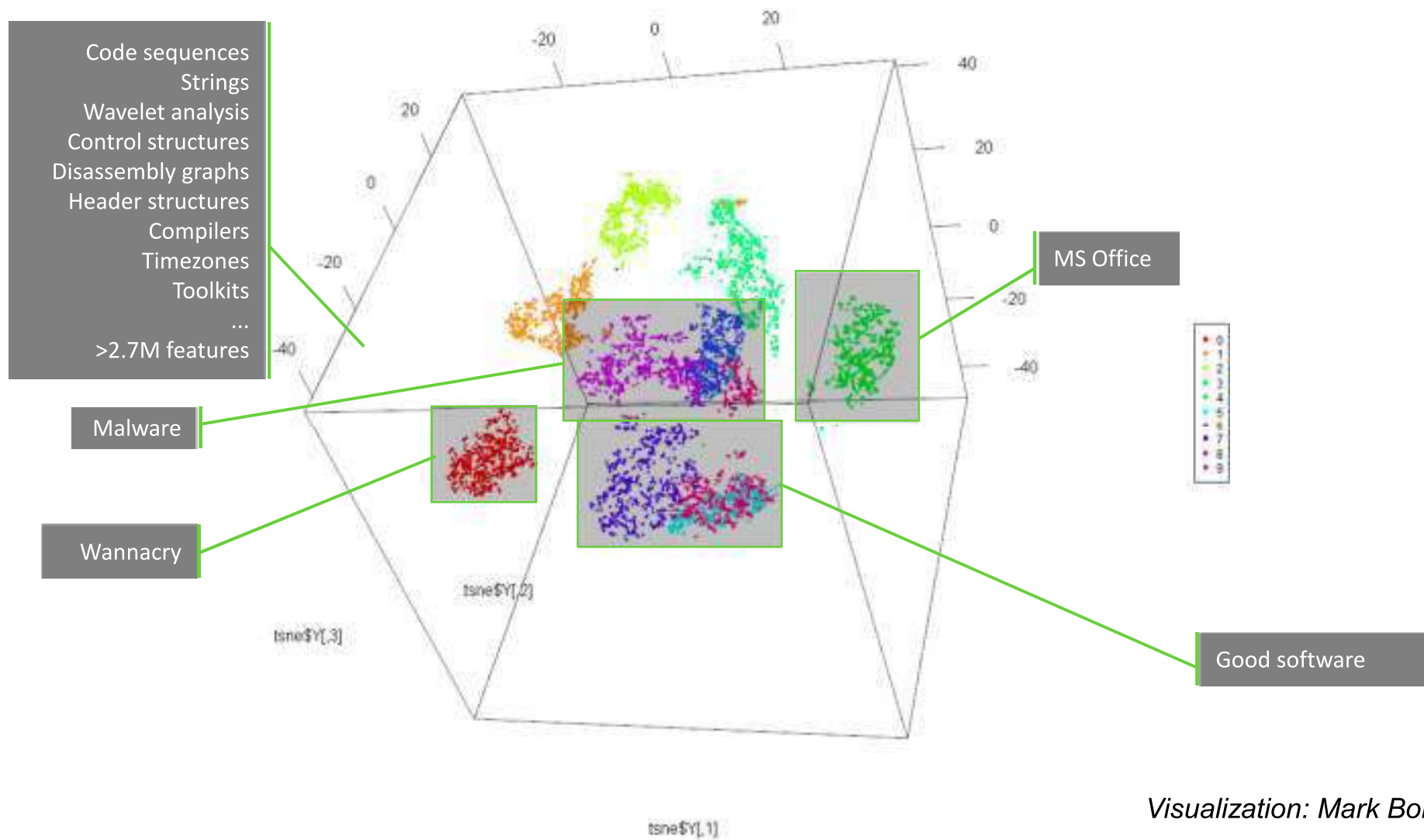


(0.0, 0.5, 0.1, 0.2, 0.1, 0.0, 0.0, 0.0)

this cat's feature vector



Visualization: Mark Borg



Visualization: Mark Borg

5 GENERATIONS OF ML FOR CYBERSECURITY



Generational Factors

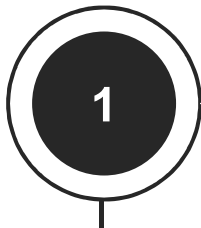
- Runtime
- Features
- Datasets
- Human Interaction
- Goodness of Fit

DARPA's Three AI Waves:

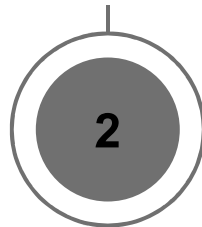
DESCRIBE

CATEGORIZE

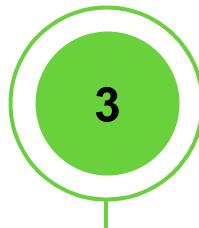
EXPLAIN



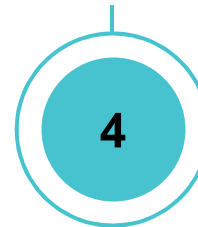
R: Cloud only training / prediction
F: Small features (~1,000)
D: Small samples (~1M)
D: Hand picked and human labeled
H: Easily interpretable
G: High FPs / Underfit / Easy to bypass



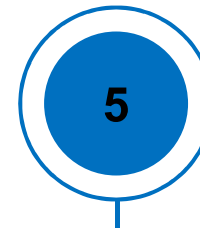
R: Cloud training / local prediction
F: Medium features (~100,000)
D: Medium samples (~100M)
D: Mostly human labeled / some heuristic
H: Largely uninterpretable
G: Misleading FP rate / Overfit



R: Cloud enhanced models
F: Large features (~3M)
D: Large samples (~1B)
D: Largely heuristic labeled
H: Some interpretability with visualization
G: Fit appropriately / accuracy metrics generalize



R: Models learn from local training
F: Large features (>3M)
D: Online learning
H: Model explains strategy & gets feedback
G: Model fits current and future inputs



R: Unsupervised local training
F: Unlimited with semi-supervised discovery and data collection
D: Active learning
H: Human input optional
G: Model identifies and adapts to concept drift



QUESTIONS — AND — ANSWERS



THANK YOU