

PARCOURS IA – 2019/2020

Fondements de l'Intelligence Artificielle et du Machine Learning

SESSION 1 – INTRODUCTION

NICOLAS VAYATIS

école _____
normale _____
supérieure _____
paris—saclay _____

How to start a talk about AI?

Commercial success of AI

Advertising

The screenshot shows the LesEchos.fr website. At the top, there's a navigation bar with 'MENU', 'LesEchos.fr', a search bar, and buttons for 'ABONNEZ-VOUS', 'Newsletters', and 'Mon compte'. Below the navigation bar, there are several promotional banners and news articles. One prominent banner for Honda Civic says 'Découvrez les atouts de la Honda Civic' and 'Le rendez-vous HONDA à retrouver sur lesechos.fr'. Below this, there's a news article titled 'Affaire Khashoggi : Trump menace directement l'Arabie saoudite' with a photo of Donald Trump. To the right of the article, there's a BNP Paribas advertisement with the text 'La Conformité : des métiers incontournables dans le secteur bancaire' and 'imagin8 BRAND STUDIO'. At the bottom, there are two more articles: 'Utilitaires : PSA parie sur l'électrique en centre-ville' and 'Pourquoi le maillot « deux étoiles » reste introuvable'.

Recommender systems

The screenshot shows the Amazon.com homepage. At the top, there's a navigation bar with 'Amazon', 'Rechercher', and 'Panier'. Below the navigation bar, there are several recommendation sections: 'Recommendations for you in CDs & Vinyl', 'Recommended for You in Kindle Books', and 'New for you'. Each section displays a grid of product covers. At the bottom, there's a 'prime' logo and a date '12/03/2018'.



Facebook for the blind

A fact: AI may outperform
Humans

Board games



1996 : Gary Kasparov vs. Deep Blue (IBM)



2016 : Ke Jie vs. AlphaGo (DeepMind)

Supervised learning algorithms outperform human performance in many pattern recognition tasks

- LeCun et al. (1989): Handwritten zip code digit recognition

→ USPS database; about 10,000 digits

→ 10 categories; 7000 training data (16x16 gray level images)

7210414959
0690159784
9665407401
3134727121
1742351244



- Lu and Tang (2015): Face recognition

→ Life Faces in the Wild (LFW) data set

→ 5749 public figures; 13,233 uncontrolled face images

→ Training on 40,000 pairs of images (matched/mismatched)

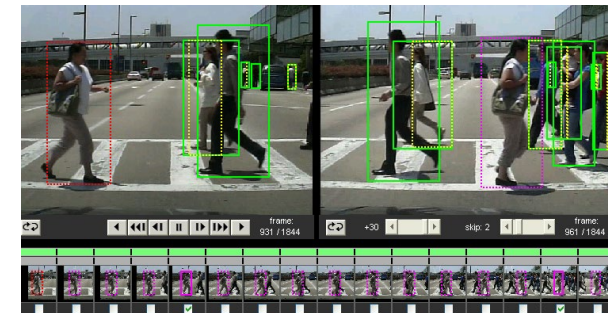
- Zhang et al. (2017): Pedestrian recognition

→ Caltech pedestrian data set

→ 10 hours video at 30Hz; 10^6 frames

→ 10% contain pedestrians; 2300 unique pedestrians

→ Some trouble with partial occlusions...



Questions raised

- What are the drivers of success for AI components and their current limitations?
- Why successful AI applications are typically related to images and text data ?
- Is cross-validated prediction performance the only criterion to adopt AI-driven technologies?
- How does AI rely on heavy (and energy-consuming) computational resources?
- Will there still be a Human in the loop in ten years?

How far can AI go?

(Super)creativity



(Super)creativity



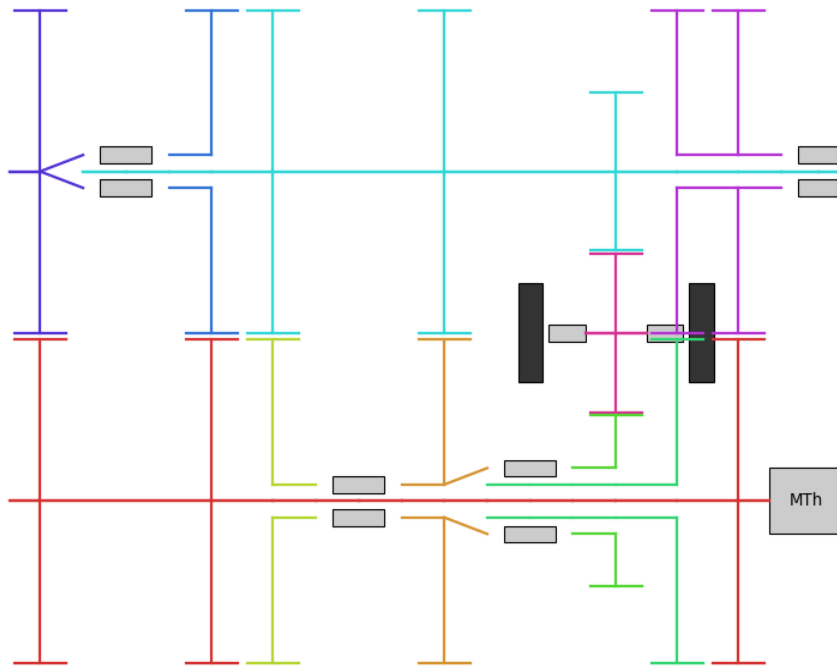
The next Rembrandt



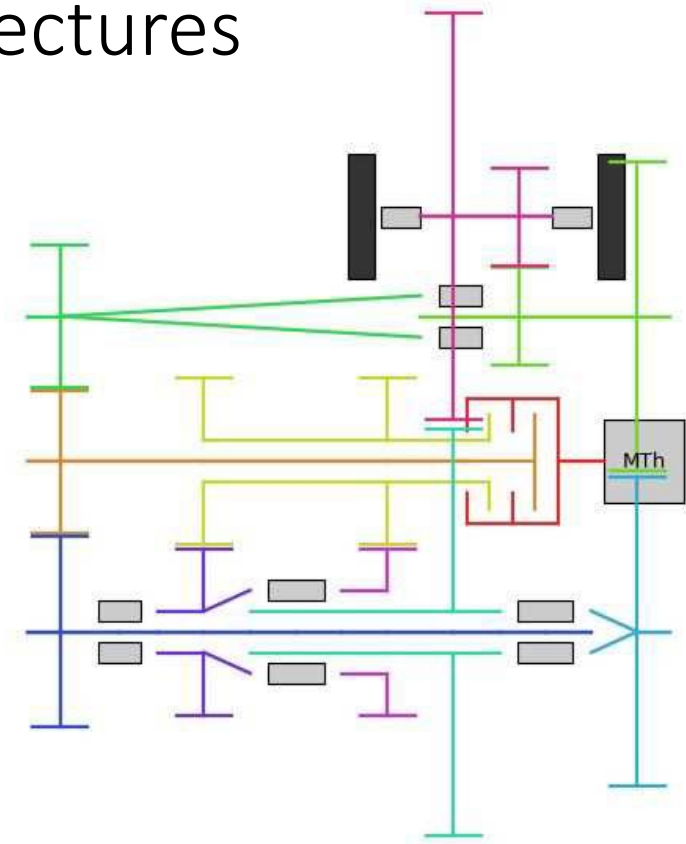
AICAN exhibition, 2018

An example of system design (with HPC)

Exploring the space of gearbox architectures



2D model of a six gearbox ratios
of a manual transmission



2D model of a five gearbox ratios
of a dual-clutch transmission

What does it take to sample and screen the space of gearbox architectures

- 992 architectural schemes scanned
- 1.5×10^9 architectures generated
- 1.5×10^8 architectures tested
- 1,390 viable architectures extracted
- 13,600 CPU-hours on Intel Xeon E5-1620v2
- Further screening based on price and mass constraints
- Expert assessment to evaluate plausibility regarding to volume optimization

2D model of a six gearbox ratios
of a manual transmission

- 142 architectural schemes scanned
- 2.5×10^8 architectures generated
- 2.5×10^7 architectures tested
- 320 viable architectures extracted
- 13,600 CPU-hours on Intel Xeon E5-1620v2
- Further screening based on price and mass constraints
- Expert assessment to evaluate plausibility regarding to volume optimization

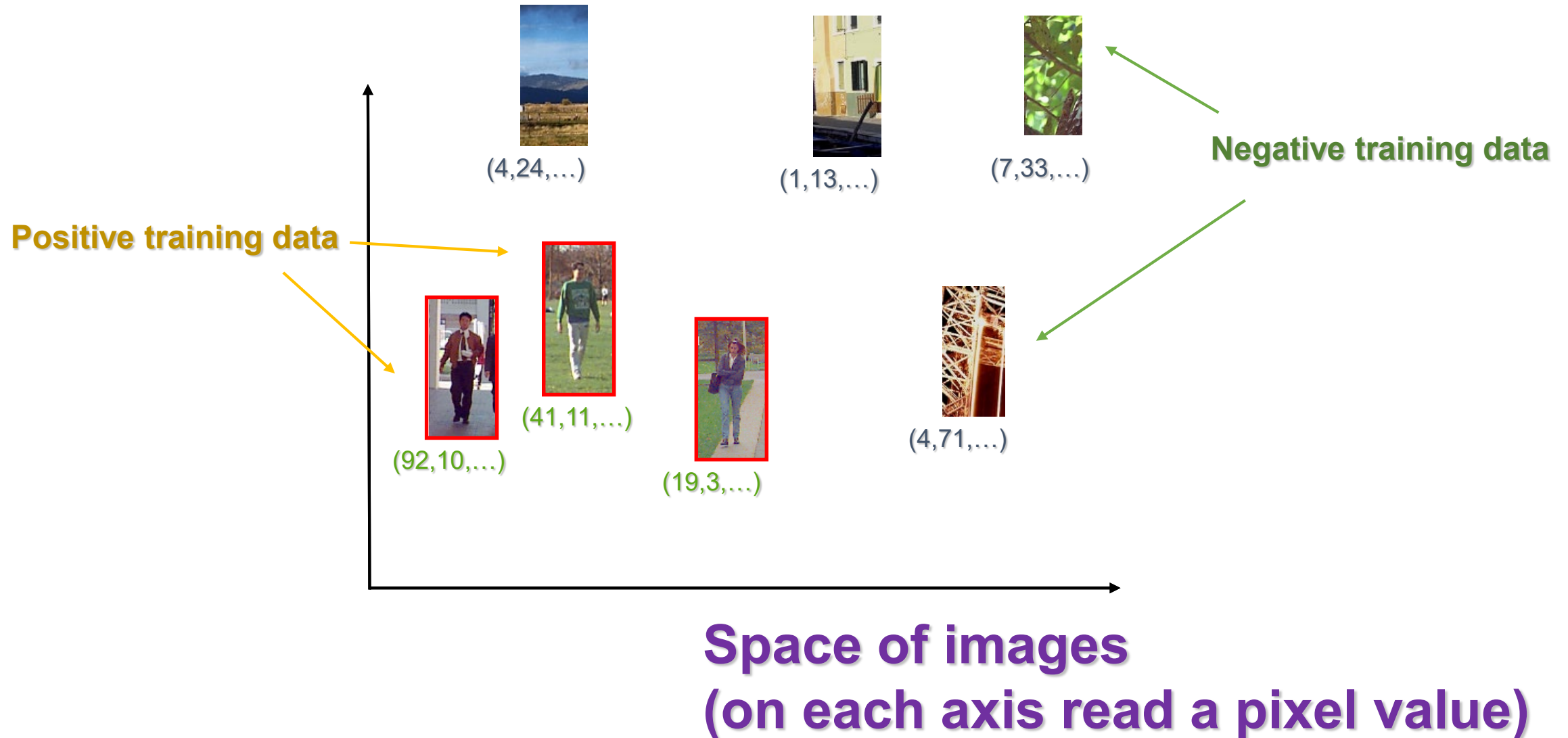
2D model of a five gearbox ratios
of a dual-clutch transmission

What we learned from innovative gearbox design

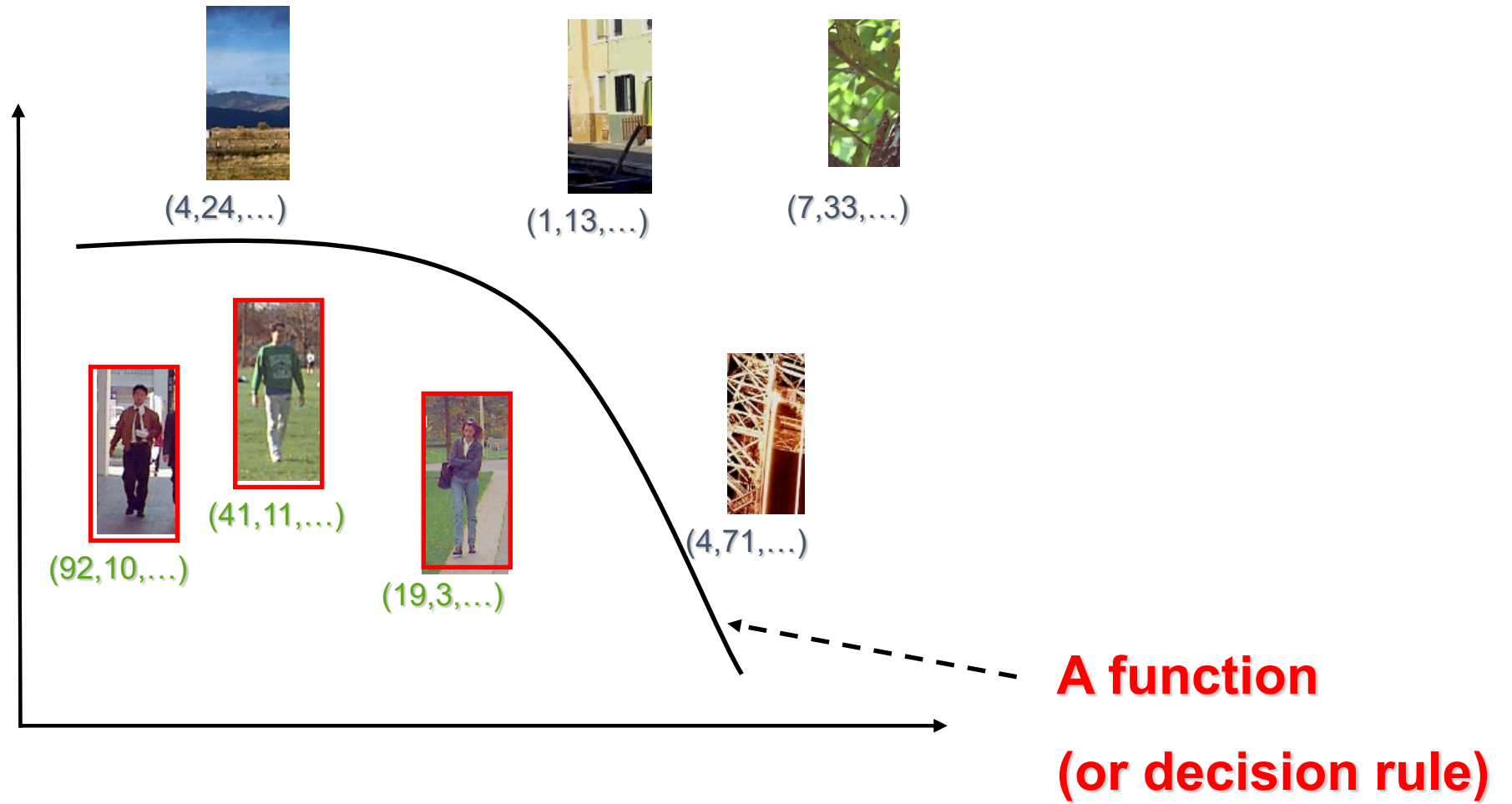
- There might be a DeepBlue for gearbox design
 - Not clear what is the complexity ceiling to extend it to engine design for instance
 - Requires the potential of HPC to sample and screen architectures in order to scale up
- Need to embark field expertise together with modeling ability:
 - Gearbox engineering, mechanical systems, optimization, graph sampling...
- Contribution of machine learning?
 - Not obvious at this stage, but...
 - ... it may help to better select high level design parameters and save brute force exploration time
- How to embrace such a design process disruption?
 - Mindset of the organization
 - Mindset of field experts

Demystification of AI/machine learning

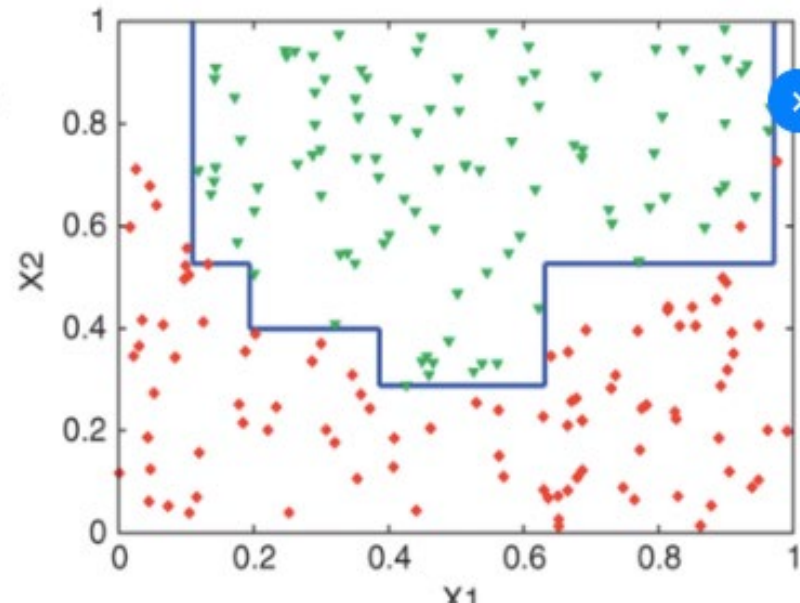
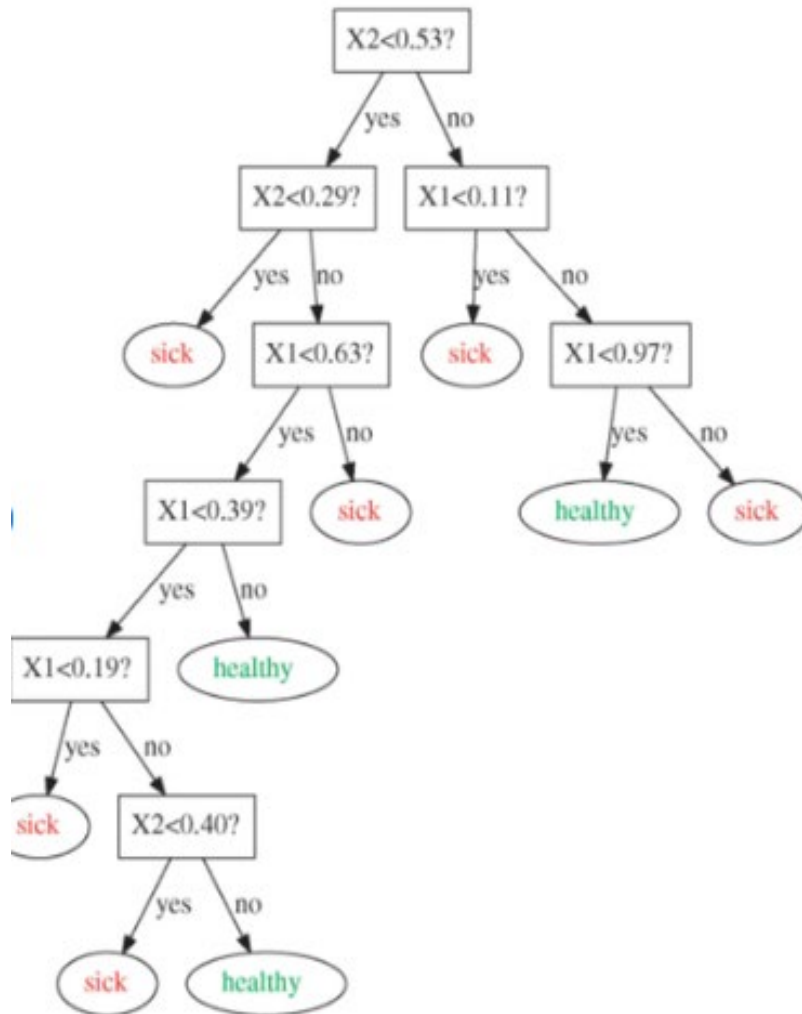
Supervised training data for pedestrian recognition



How does the machine “represent what it learned”?

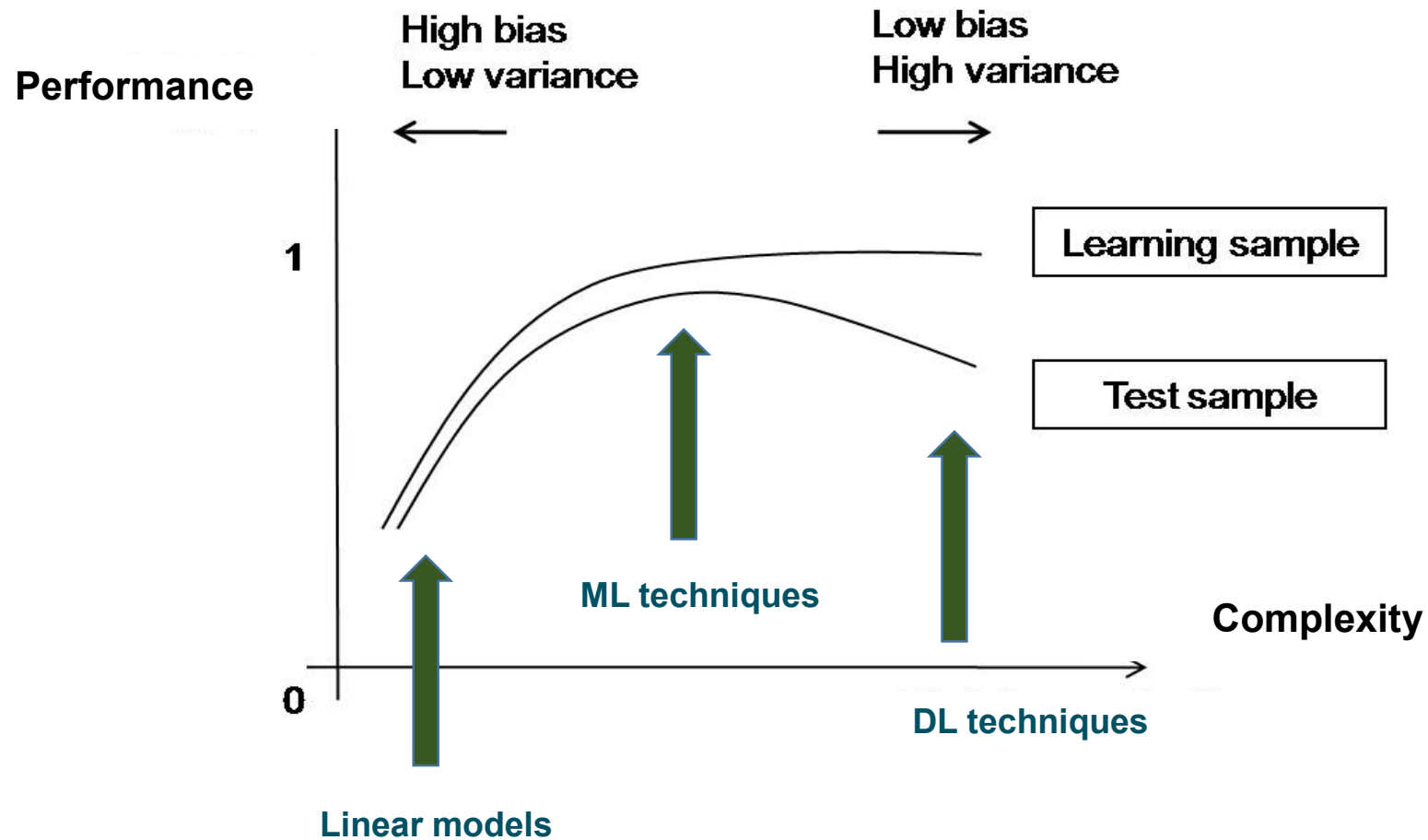


Machine learning is mainly about
function estimation /approximation



ALL OF MACHINE LEARNING IN ONE FIGURE

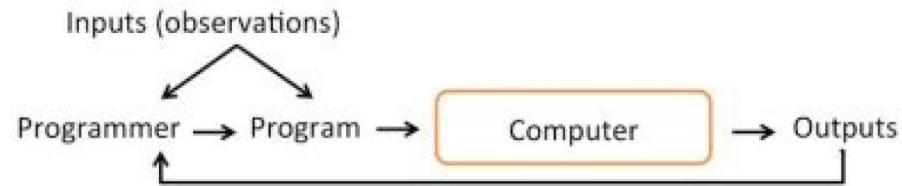
UNDERFITTING AND OVERFITTING, IT'S ALL ABOUT TRADE-OFF



Historical perspective on AI

From Symbolic AI to Machine Learning

The Traditional Programming Paradigm

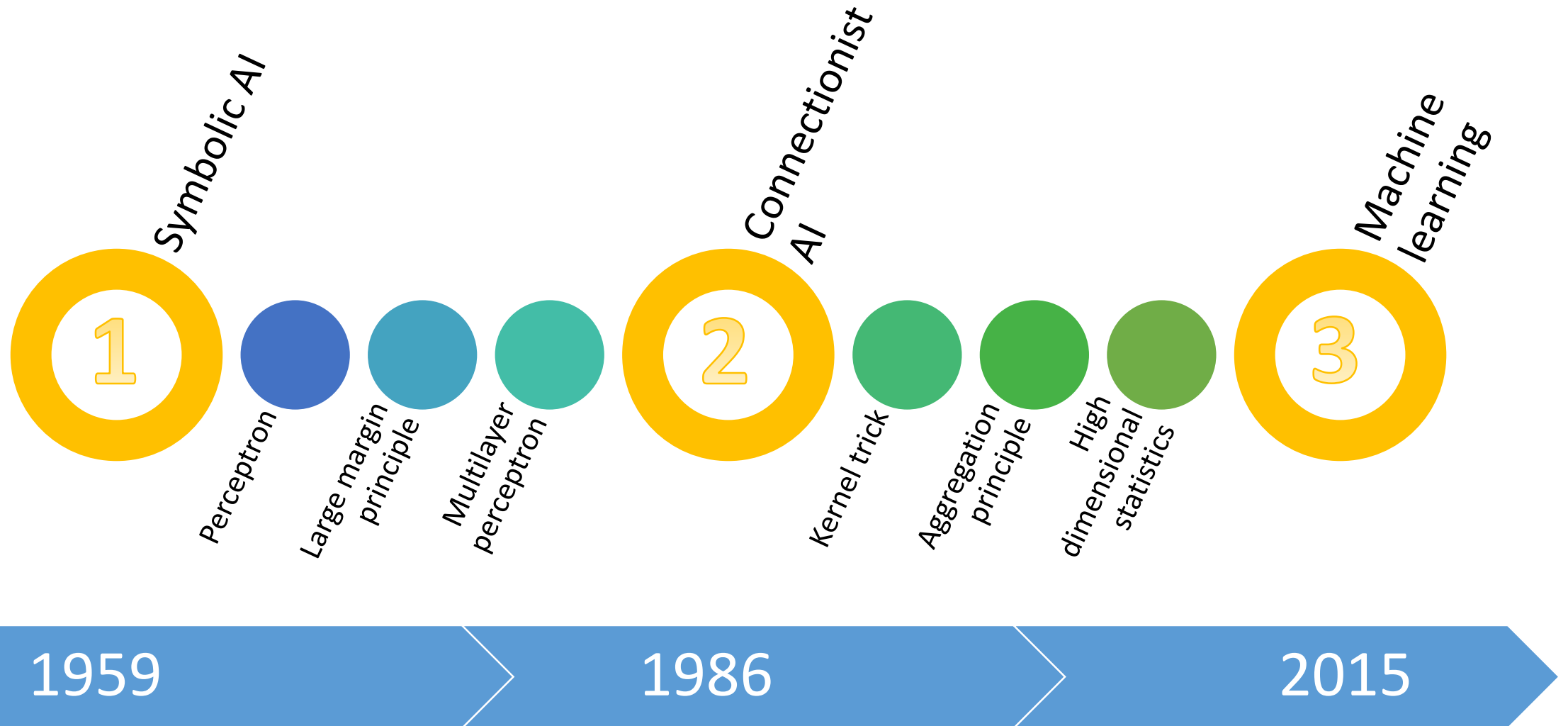


Machine Learning is the field of study that gives computers the ability to learn without being explicitly programmed
– Arthur Samuel (1959)

Machine Learning



Three AI waves... and two AI winters



Winters explained

1

Symbolic AI

Only able to capture explicit knowledge

Lack of robustness

Computational burden

1959

2

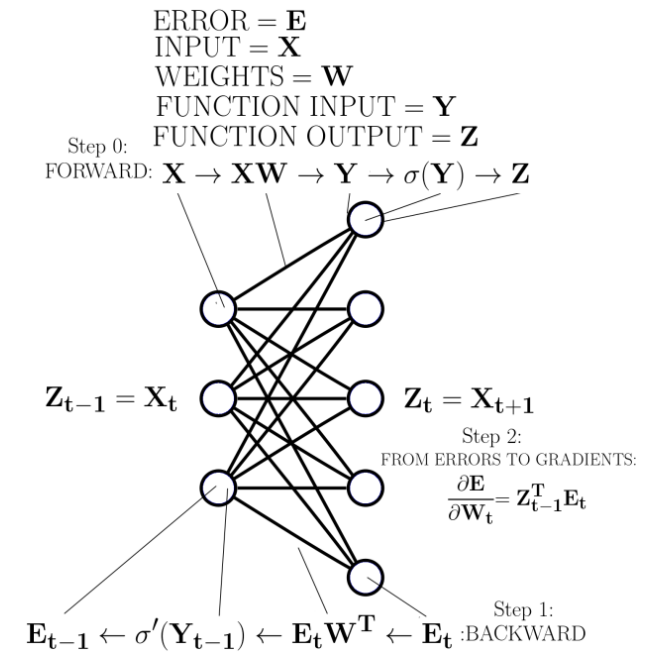
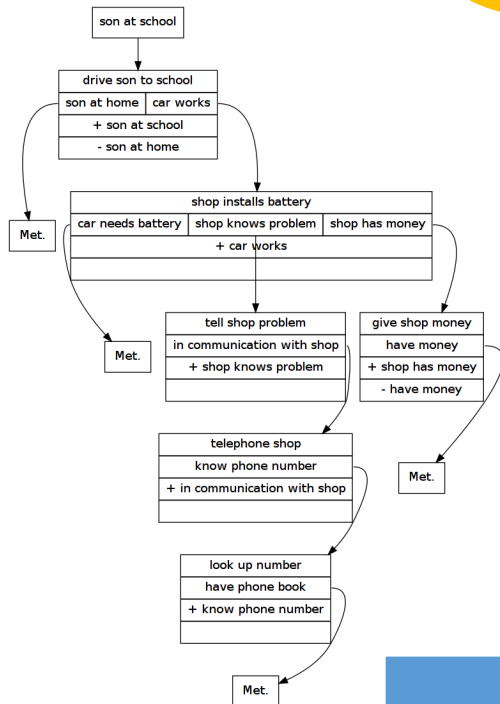
Connectionist AI

Lack of data

Lack of interpretability

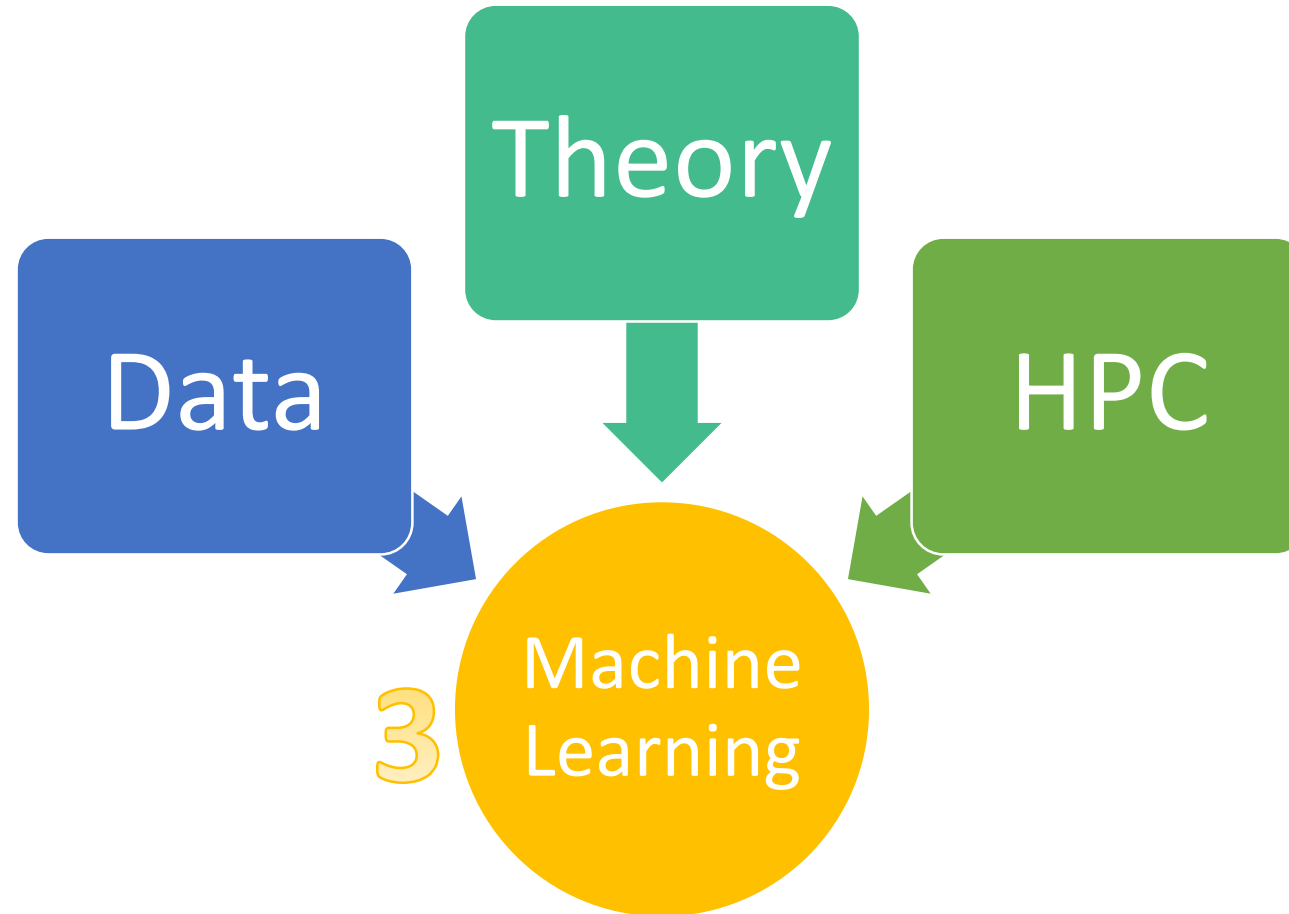
Hard to maintain

1986

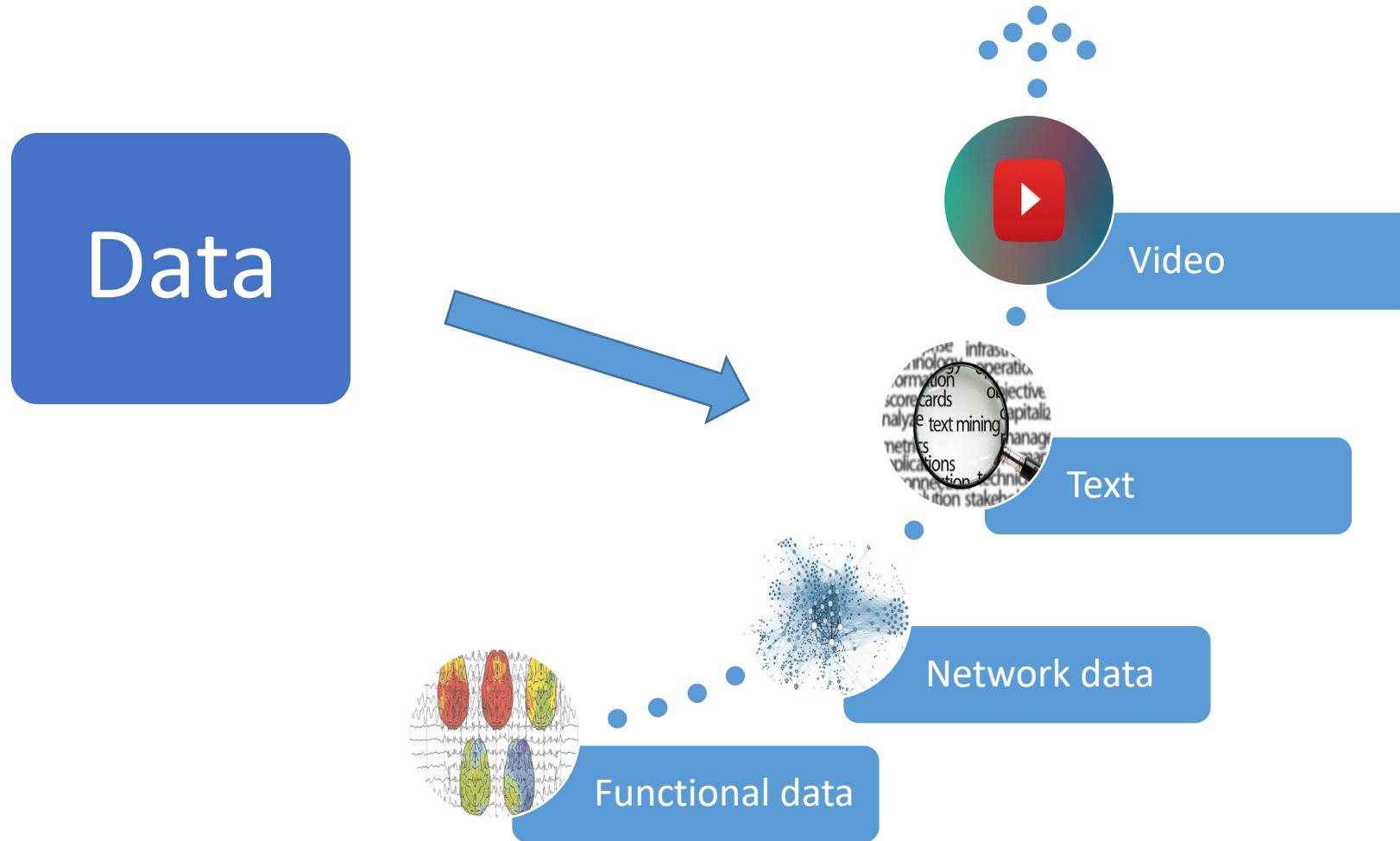


The three drivers
of the third AI wave

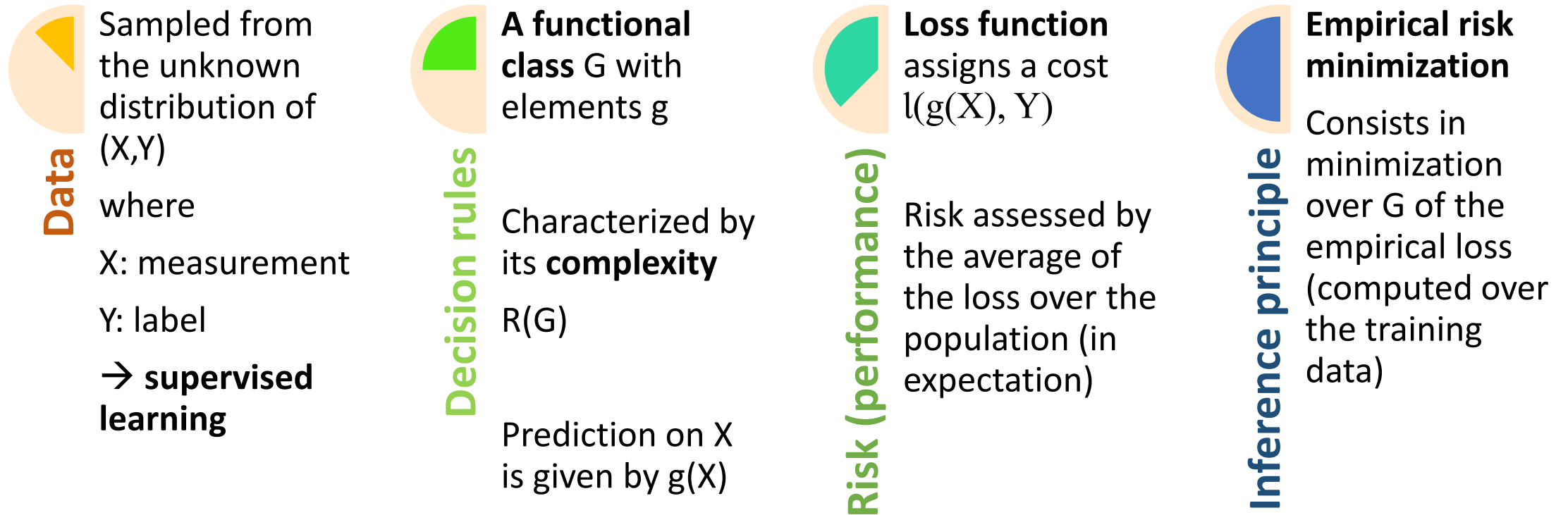
What is different now



All is data



Statistical Learning Theory: setup



(Machine) learning amounts to functional optimization

Statistical Learning Theory: main ingredients

$$\hat{R}_n(\mathcal{F}) = \mathbb{E} \left(\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(Z_i) \middle| D_n \right)$$

Assume h is a function with bounded differences and denote by $c_1, \dots, c_n > 0$ the upper bounds on its componentwise variations

We have, for any $t > 0$

$$\mathbb{P}(h(Z_1, \dots, Z_n) - \mathbb{E}(h(Z_1, \dots, Z_n)) > t) \leq \exp \left(-\frac{2t^2}{\sum_{i=1}^n c_i^2} \right)$$

Under (A1-A2), we have, for some $s > 1$, that any real-valued measurable f satisfies :

$$L(g_f) - L^* \leq 2c(A(f) - A^*)^{1/s}$$

Rademacher
complexity

Concentration
inequality

Risk
communication

Statistical Learning Theory: typical guarantee

$$L(\hat{g}_n) \leq \inf_{g \in \mathcal{G}} L(g) + \hat{R}_n(\mathcal{G}) + 3\sqrt{\frac{\log(2/\delta)}{2n}}$$

with probability at least $1 - \delta$

True loss of the ERM

\leq Minimal loss in the class

+ Complexity

+ Precision

Is HPC a necessary tool for Machine Learning?

- Big data + Deep Learning + Real-time training → definitely needs HPC
- What if:
 - Real-time decisions but not necessarily real-time training
 - Not so big data
 - Satisfied with other (shallow) Machine Learning algorithms
(e.g. Random Forests, Boosting, SVM)

?

Demystification of big data

First issue of big data: Sampling bias

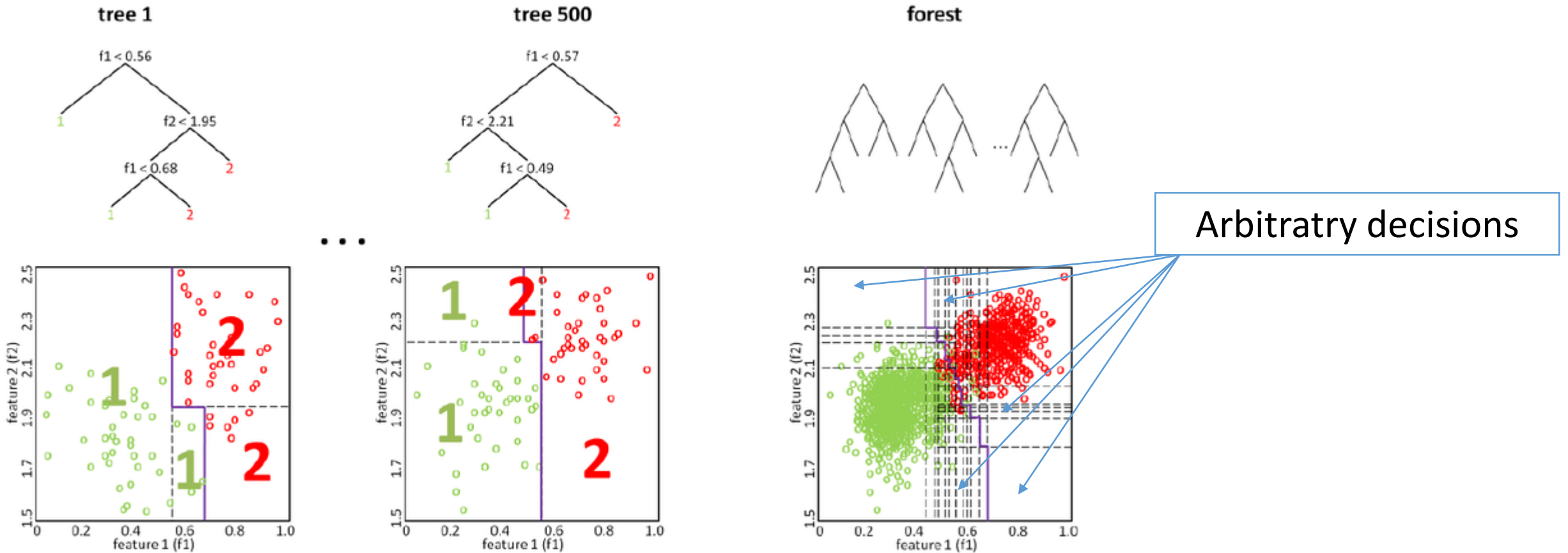


TRAIN



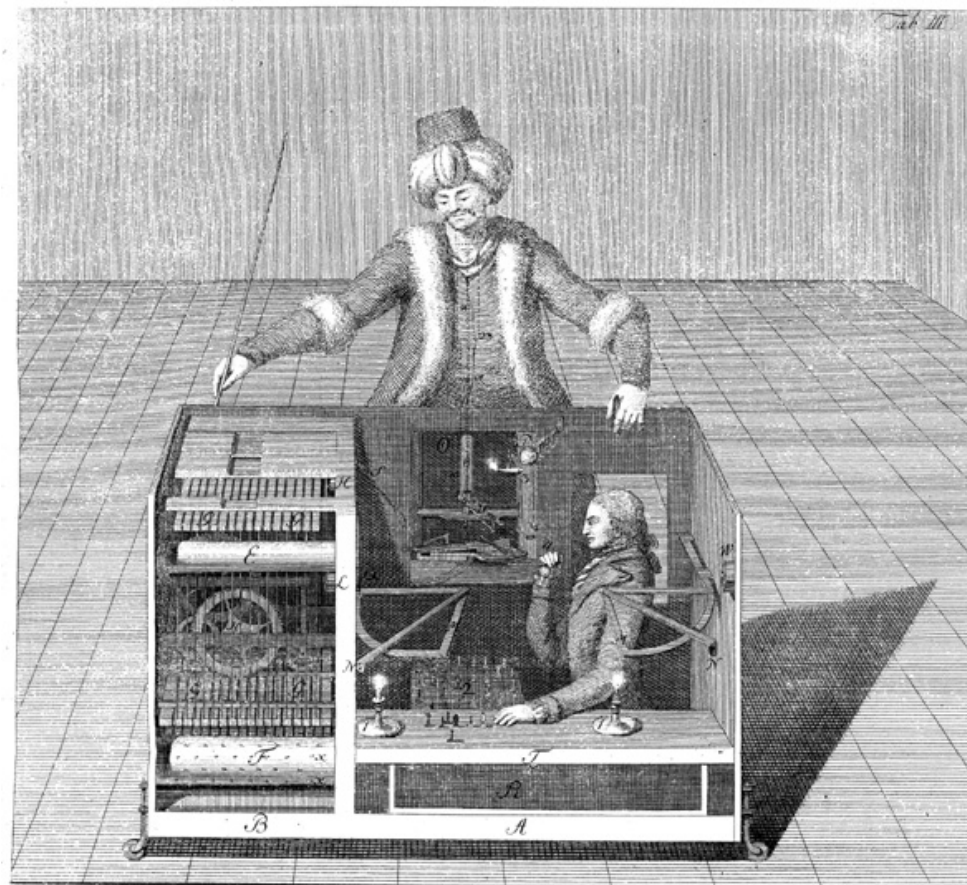
TEST

Biases and artefacts in data-driven partitions



→ Need to learn with a reject option (see work by Marten Wegkamp, 2005-...)

Second issue with big data: Need a Turk!



The cost of data labeling

Discover, preview and complete HITs on the new Worker website. Try it out Today!

amazonmechanical turk
Artificial Intelligence

[Your Account](#) [HITs](#) [Qualifications](#)

[Introduction](#) | [Dashboard](#) | [Status](#) | [Account Settings](#)

Mechanical Turk is a marketplace for work.
We give businesses and developers access to an on-demand, scalable workforce.
Workers select from thousands of tasks and work whenever it's convenient.

311,197 HITs available. [View them now.](#)

Make Money by working on HITs

HITs - *Human Intelligence Tasks* - are individual tasks that you work on. [Find HITs now.](#)

As a Mechanical Turk Worker you:

- Can work from home
- Choose your own work hours
- Get paid for doing good work

Find an interesting task → **Work** → **Earn money**

Find HITs Now

Get Results from Mechanical Turk Workers

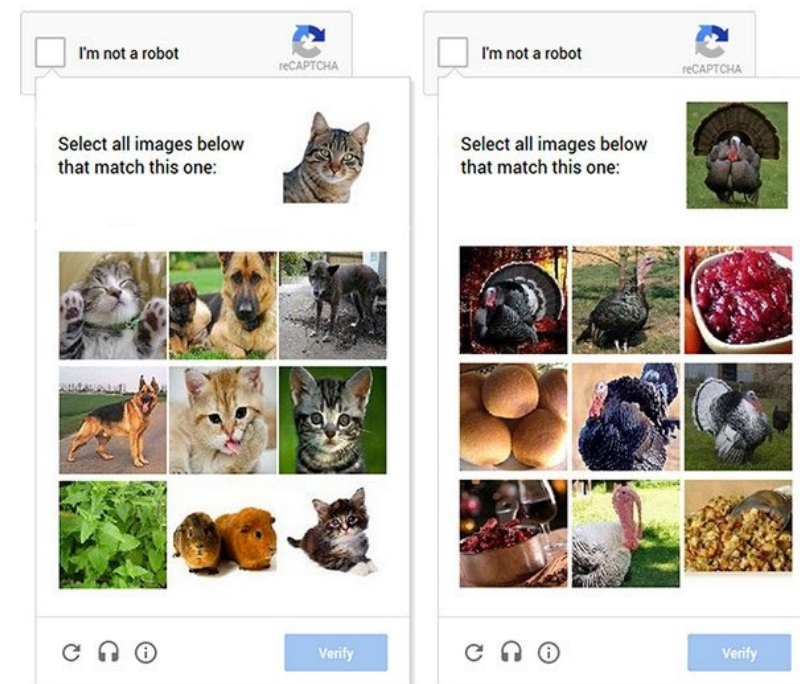
Ask workers to complete HITs - *Human Intelligence Tasks* - and get results using Mechanical Turk. [Get Started.](#)

As a Mechanical Turk Requester you:

- Have access to a global, on-demand, 24 x 7 workforce
- Get thousands of HITs completed in minutes
- Pay only when you're satisfied with the results

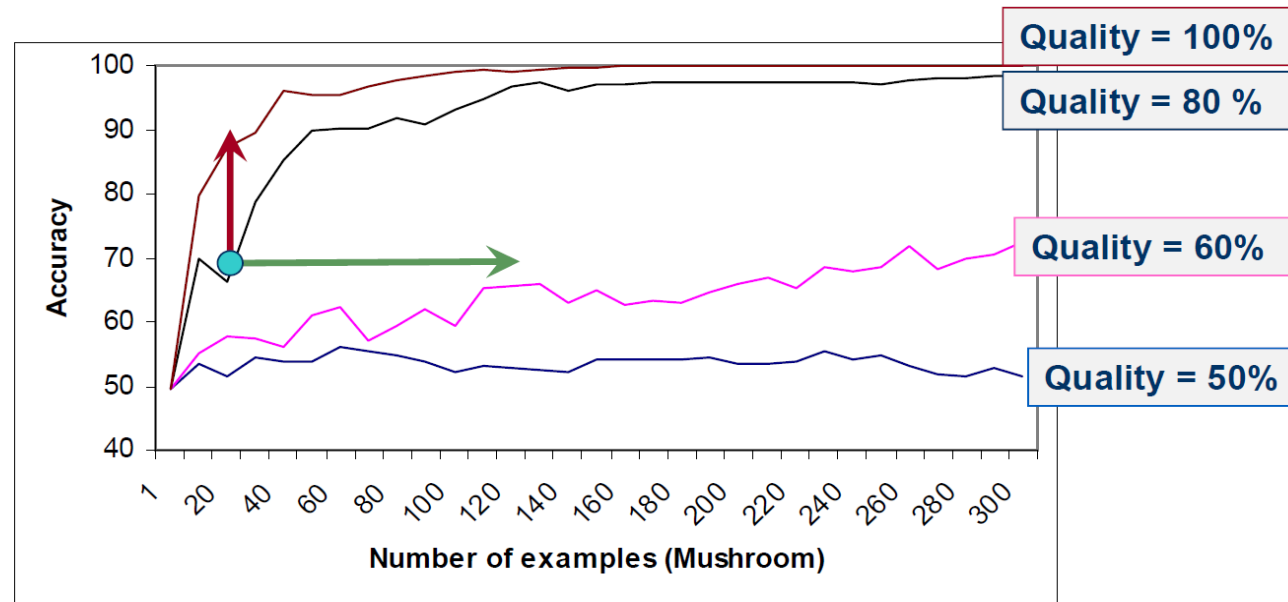
Fund your account → **Load your tasks** → **Get results**

Get Started



Quality labels more powerful than big data

- Get more examples → Improve classification
- Get more labels → Improve label quality → Improve classification



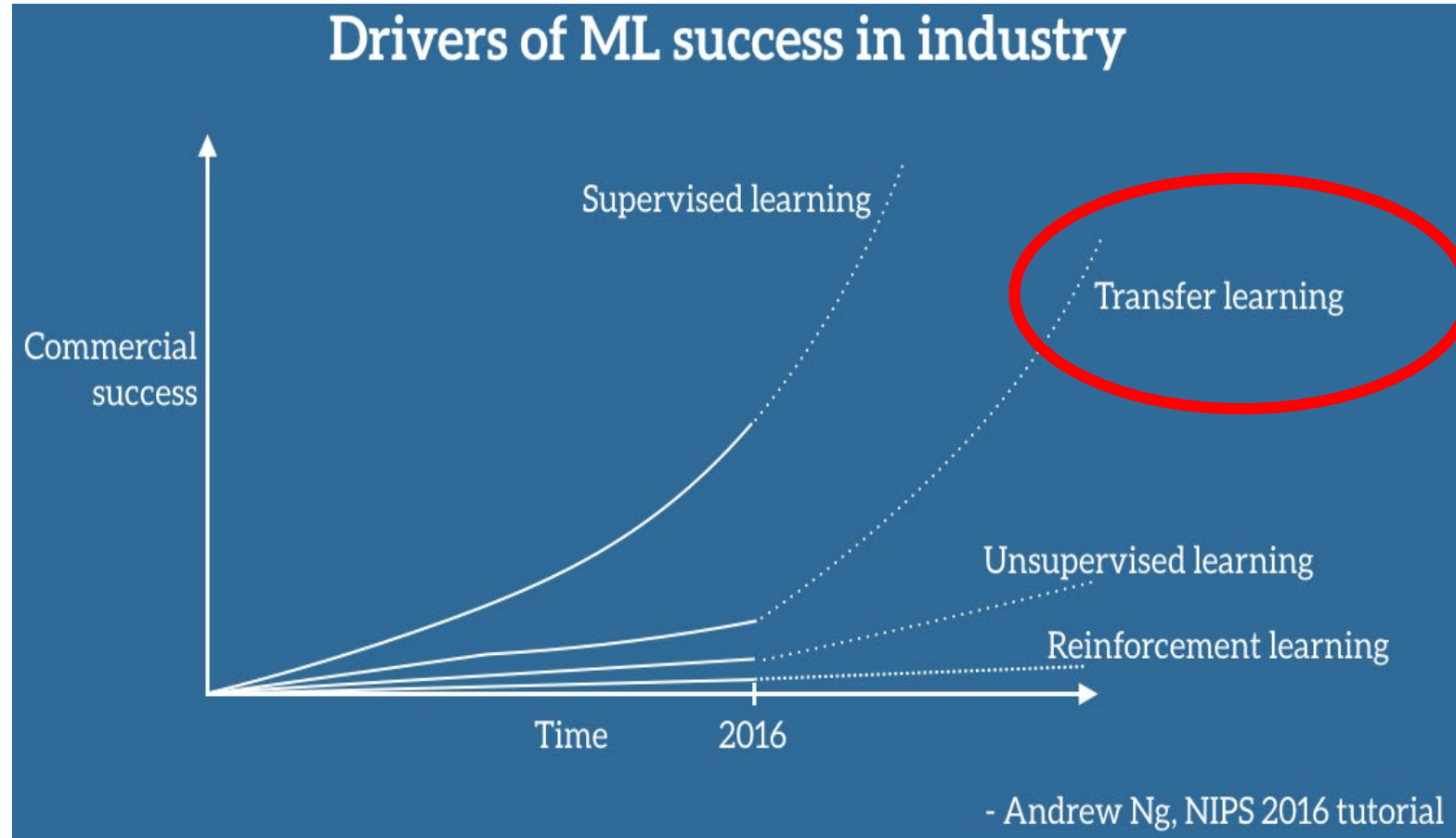
Source: Get Another Label? Improving Data Quality and Data Mining Using Multiple, Noisy Labelers. Proceedings of KDD-2008 by V. S. Sheng, F. Provost, P. G. Ipeirotis.

Machine Learning for Industry or Science:
a different story...

Industrial requirements heavier than for marketing

- Industrial processes go under continuous improvement → **Sampling bias is the rule!**
- Labeling training data relies on **field expertise** → Turks are expensive and unwilling!
- Expectations for **performance** are at a different scale when comparing decisions for critical systems or clinical applications to advertising or book recommendation

Expected impact of machine learning in the industry

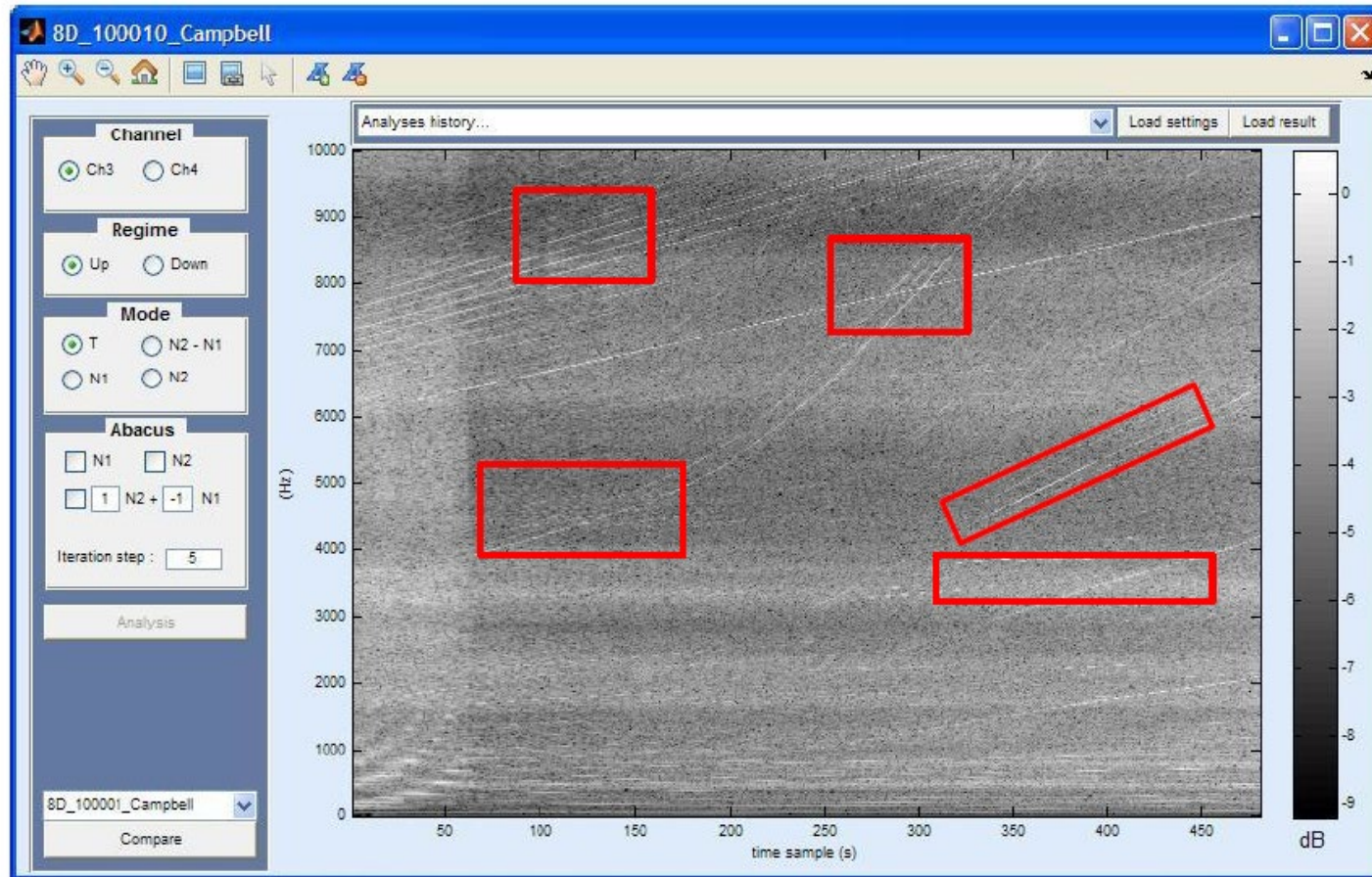


NB: Andrew Ng is VP & Chief Scientist of Baidu, Co-Chairman and Co-Founder of Coursera, and an Adjunct Professor at Stanford University.

Machine learning to support
scientific computing and
simulation projects

An example of anomaly detection objective

Benchmark assessment of aircraft engine

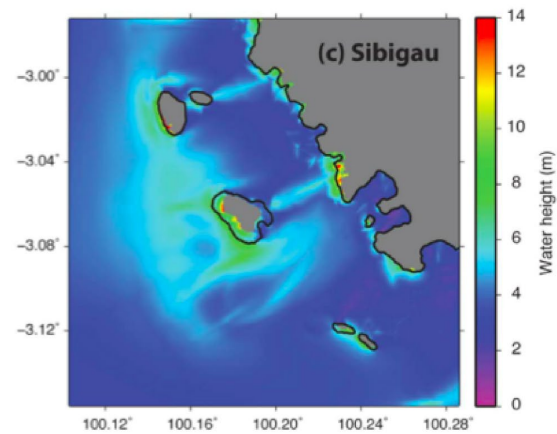


Source: Confidential report (2012) - Mathilde Mougeot, NV

- What we see?
Time-frequency representation of vibration signals (Campbell diagram) wrt to speed during acceleration and deceleration regimes.
- Nature of anomalies
Tiny details in those images. Require a lot of expertise to tag.
- Databases are small
Only a few hundreds engines have been recorded with a very limited number of anomalies reported.
- But image structure helps!
Anomaly detectors can be built using adapted representations of such signals and basic nearest neighbors in feature space.

Example of project (1/2)

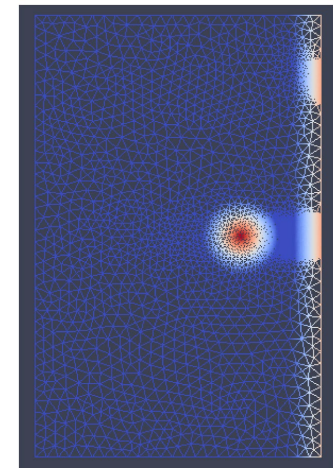
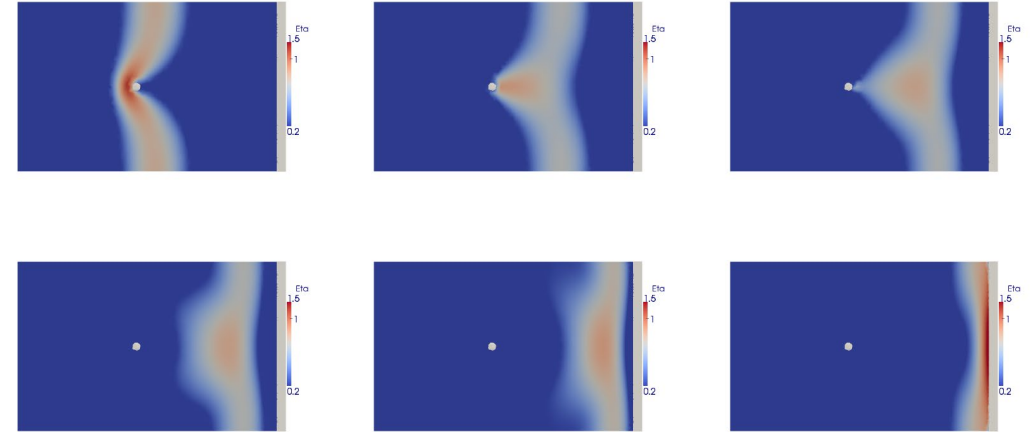
Tsunami run-up amplification



2010 Sumatra tsunami and the Mentawai Islands (Hill et al., 2012)

From:

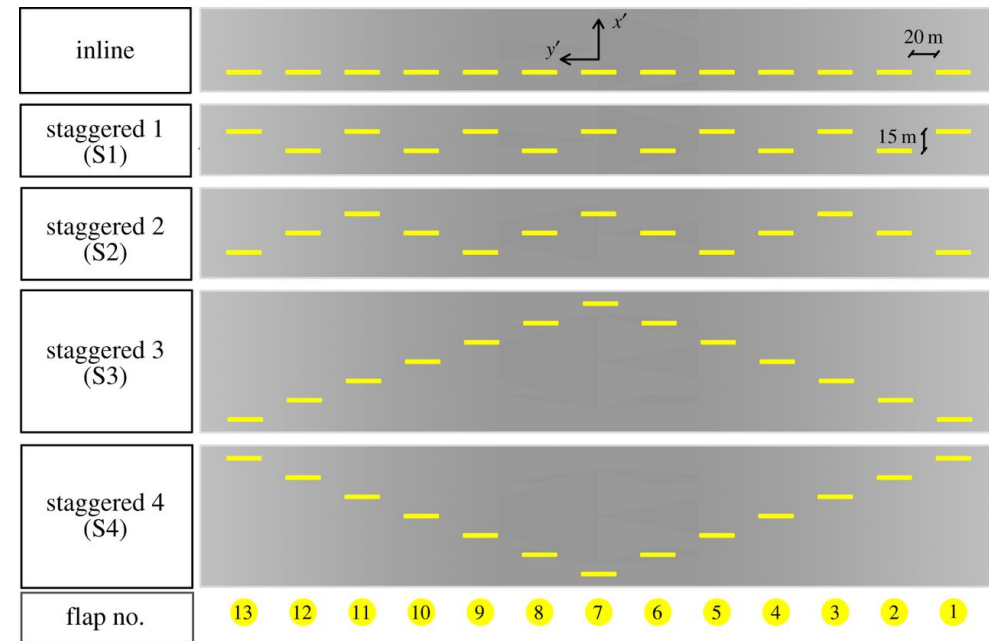
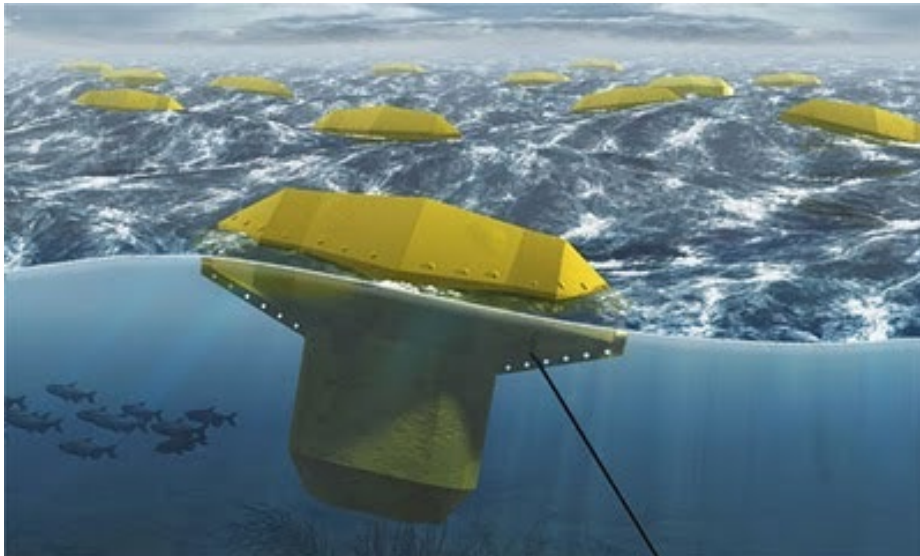
Themistoklis S. Stefanakis, Emile Contal, Nicolas Vayatis, Frédéric Dias, and Costas E. Synolakis (2014).
Can Small Islands Protect Nearby Coasts From Tsunamis? An Active Experimental Design Approach .
Proceedings of the Royal Society-A, 470: 20140575.



Adaptive mesh grid of the VOLNA solver

Example of project (2/2)

System design for WEC farms



From:

Dripta Sarkar, Emile Contal, Nicolas Vayatis, Frederic Dias (2015).

A Machine Learning Approach to the Analysis of Wave Energy Converters.

Proceedings of OMAE 2015.

Take-home messages

Machine learning achieves some kind of regression in high dimensional (or structured) spaces

- Heavily relies on mathematics to model complex data and formulate the task-related optimization problem
- AI-based technologies may outperform humans in certain *well-defined* prediction tasks: detection, recognition, planning, etc.
- Missing piece: few studies on control actions (after prediction)
- Strong AI not for tomorrow... still need to define the search space and the objective...

Scaling up and industrialization of AI modules in science and industry raises scientific challenges

- Sciences (life sciences, engineering sciences, social sciences, physics...) and Industry (energy, healthcare, banking, defense...) will not benefit of supervised learning 'as-is'
- Naive implementation of AI has/may/will lead to industrial disasters
- The main risk: to be driven by a method and not by the problem to be solved in *its* context
- Secondary risk: believe too much in training data and Proof-of-Concept
- Eventually: for the previous reasons, the risk to be out of the game, and there will be a game!