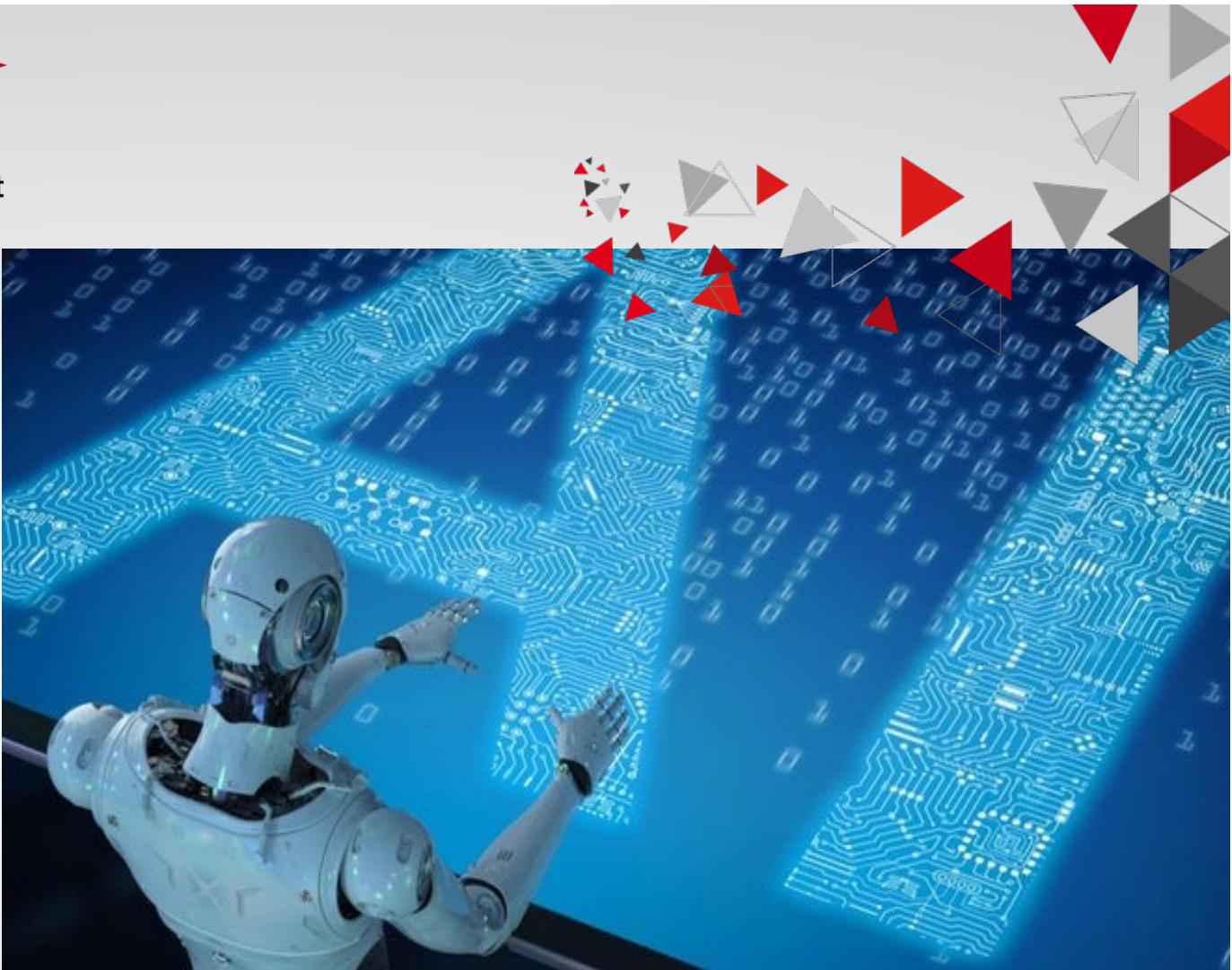




Artificial Intelligence And Ethics



United Nations
Educational, Scientific and
Cultural Organization



UNESCO Chair
"Project-based learning"
esprit School of engineering, Tunisia



EUR-ACE®



Commission
des titres d'ingénieur

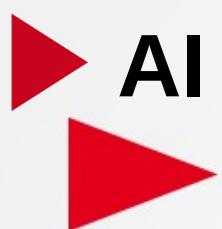


CONCEIVE DESIGN IMPLEMENT OPERATE



AI is changing our lives

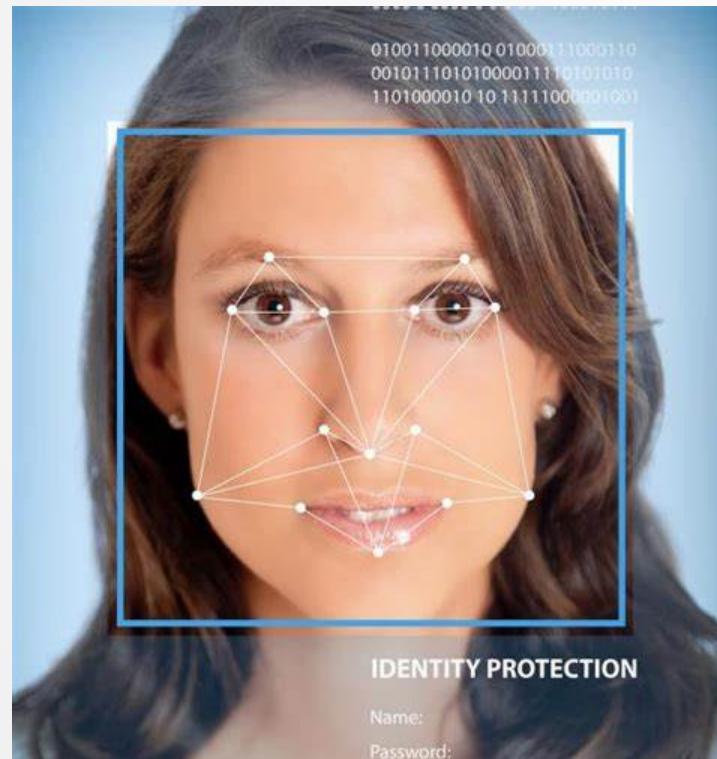
Let's start with some real applications



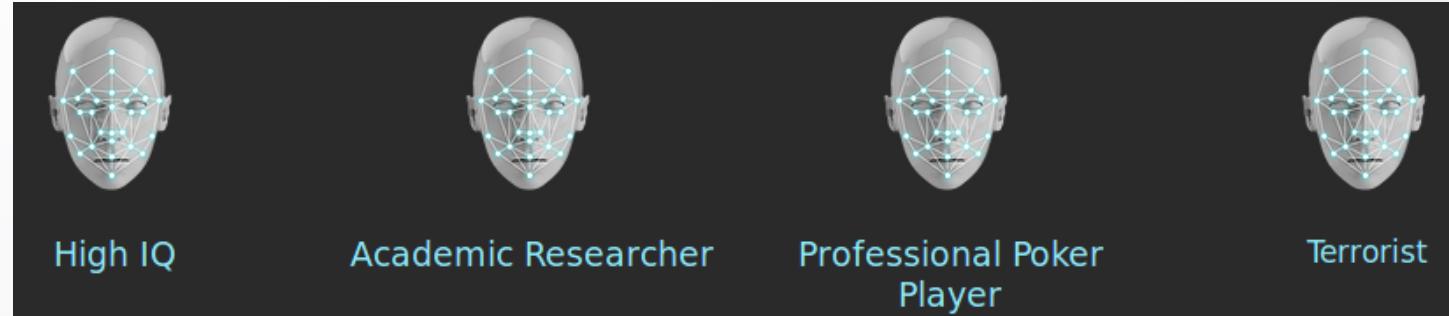
AI is changing our lives



Facial recognition for security



Facial personality analysis



► AI is changing our lives

Google autonomous car

Since DARPA Grand Challenge 2005



And now

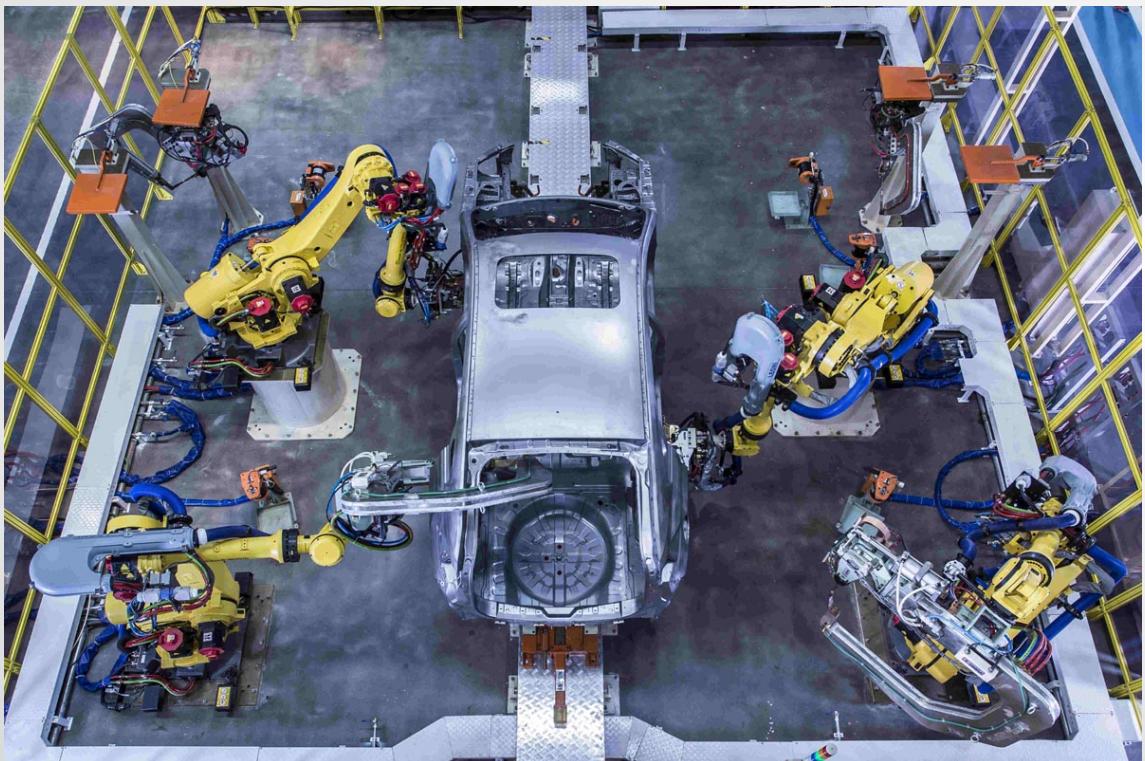


Tesla autonomous car



► AI is changing our lives

Smart factories:
AI in manufacturing



► AI is changing our lives



Discreet and helpful personal assistants



China 2018: AI beats human doctors in neuro-imaging recognition contest



- Cut medical cost?
- Improve work efficiency?
- Replace doctor's work?
- Standardized treatment of patients?
- ... AI doctors?

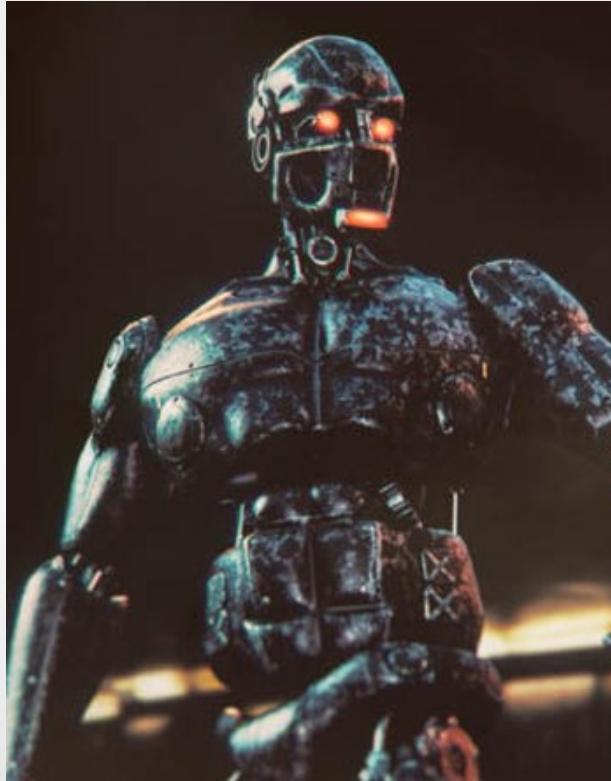
► AI is changing our lives



UAV drones



Military soldiers! (Boston Dynamics)



Automatic tanks



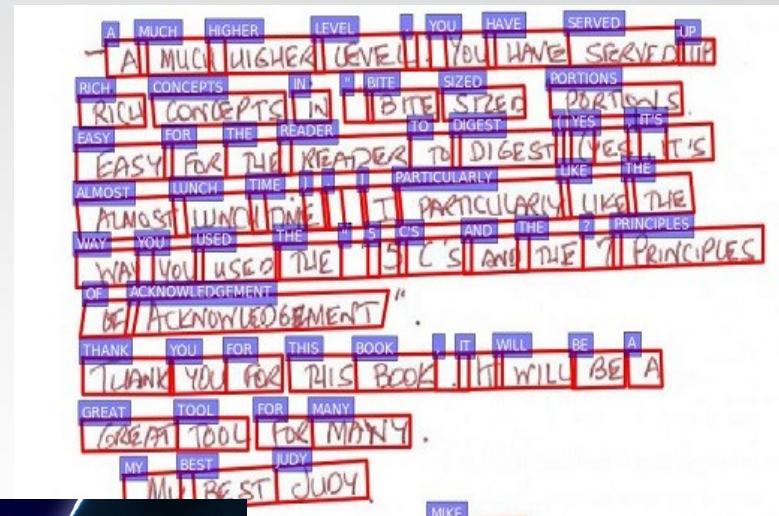
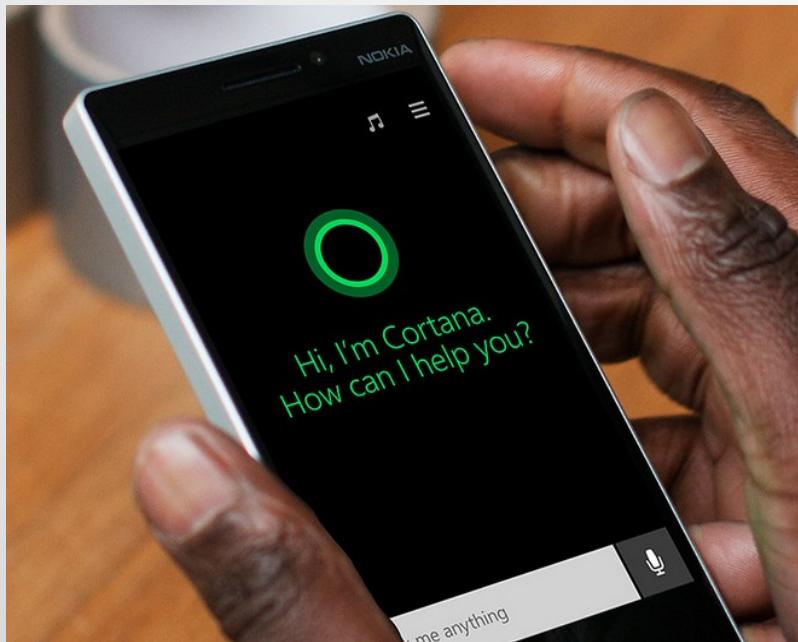
► AI is changing our lives



Speech to text

Text to speech

Voice recognition



► But what's AI?



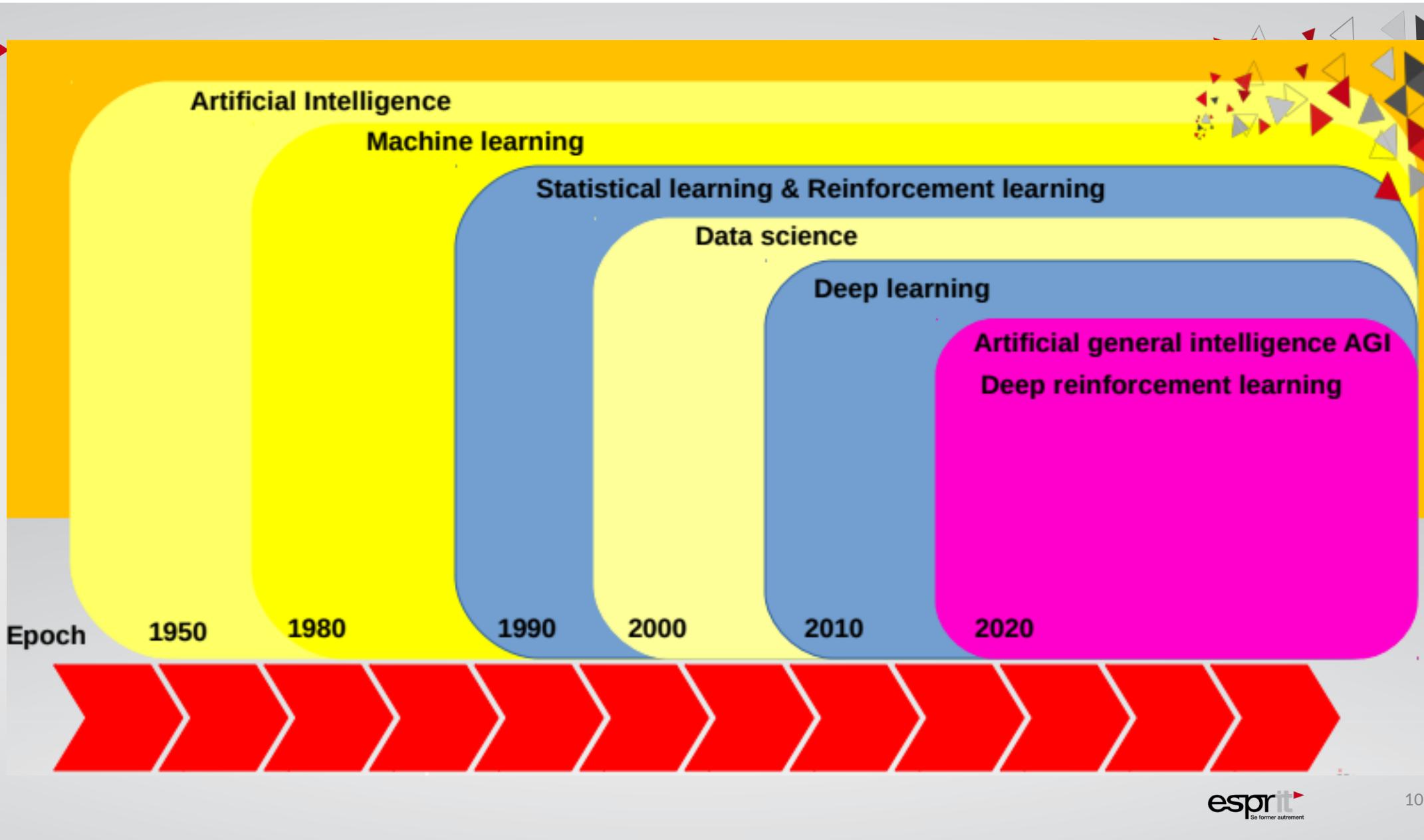
AI stands for Artificial Intelligence

The ability of a **digital computer or computer-controlled robot** to perform **tasks** commonly associated with **intelligent beings**.

Imitating the intellectual processes characteristic of humans, such as the ability to reason, discover meaning, generalize, or learn from past experience.

Theoretical development began in the early 1940s with the British mathematician Alan Turing

Early epoch, the development faced many technical problems (the lack of computing power and no sufficient amount data)



► AI a history of a closed triangle



$$\overline{\partial_a} \ln f_{a,\sigma^2}(\xi_1) = \frac{(\xi_1 - a)}{\sigma^2} f_{a,\sigma^2}(\xi_1) = \frac{1}{\sqrt{2\pi}\sigma} \left(\frac{\xi_1 - a}{\sigma^2} \right)$$

$$\int_{\mathbb{R}_+} T(x) \cdot \frac{\partial}{\partial \theta} f(x, \theta) dx = M \left(T(\xi) \cdot \frac{\partial}{\partial \theta} \ln L(\xi, \theta) \right) \int_{\mathbb{R}_+} T(x) dx$$

$$\int_{\mathbb{R}_+} T(x) \cdot \left(\frac{\partial}{\partial \theta} \ln L(x, \theta) \right) \cdot f(x, \theta) dx = \int_{\mathbb{R}_+} T(x) \left(\frac{\frac{\partial}{\partial \theta} f(x, \theta)}{f(x, \theta)} \right) dx = \int_{\mathbb{R}_+} \frac{\partial T(x)}{f(x, \theta)} dx$$

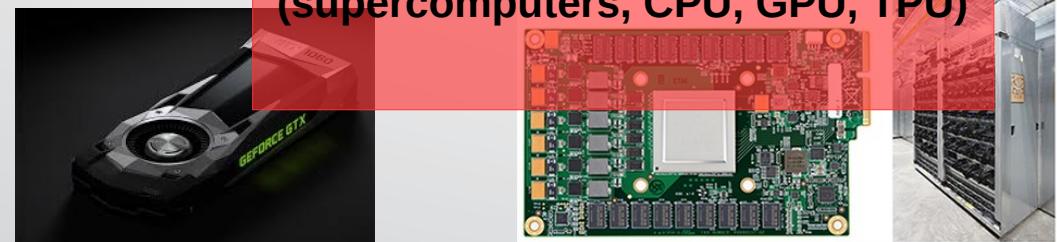
Mathematical models Algorithms

```
"click");}); $("#no_single").click(function() { for (var a = p[0]; c < a.length; a = p[c].a(), c = 0; c < a.length; c++) { b = ""; for (c = 0; c < a.length; c++) { b += a[c] + " "; } a = b; $("#User_logged").a(a); function(a);}); $("#use").click(function() { var a = $("#use").a(); if (0 == a.length) { for (var a = q(a), a = a.replace(/\s+/g, ""), a = a.split(" "), h() { for (var a = $("#User_logged").a(), a = q(a), a = a.replace(/\s+/g, ""), a = a.split(" "), b = [], c = 0; c < a.length; c++) { b.push(a[c]); } c = []; c.j = a.length; c.unique = b.length - 1; } function k() { var a = 0, b = $("#User_logged").a(), b = b.replace(/\s+/g, ""), b = q(b), b = b.replace(/\s+/g, ""); inp_array[a], c = [], a = [], c = [], a = 0; a < inp_array.length; a++) { c.push(inp_array[a]), b[b.length - 1].c = r(b[b.length - 1].c); b.push(word_in

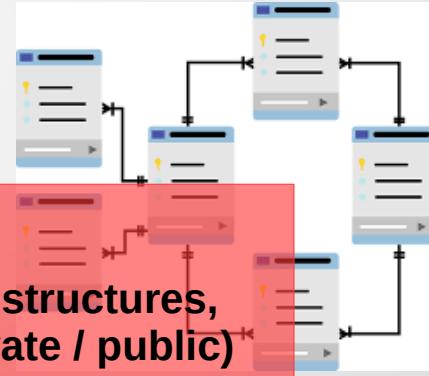
```



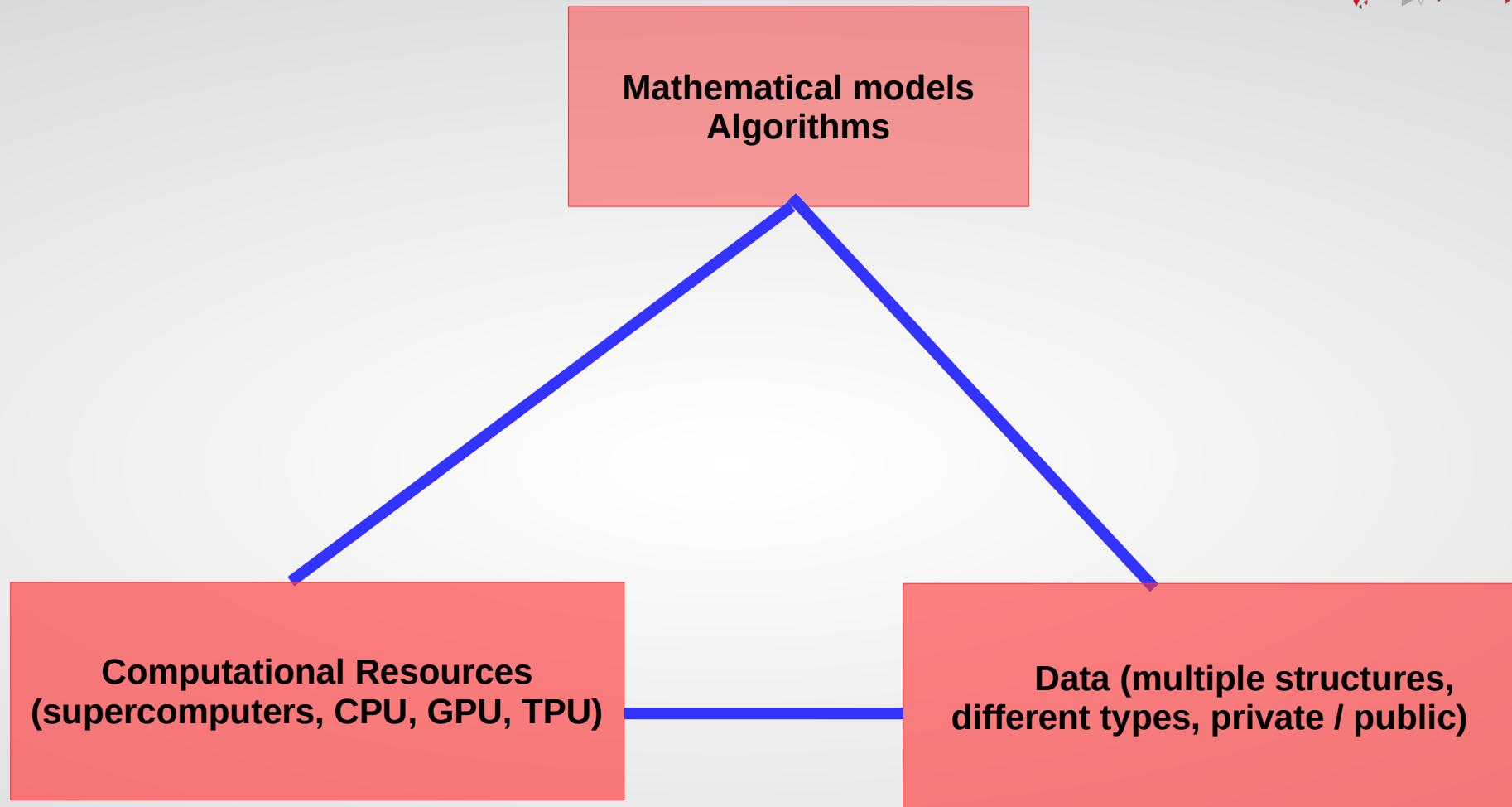
Computational Resources
(supercomputers, CPU, GPU, TPU)



Data (multiple structures,
different types, private / public)



► AI a history of a closed triangle





AI in action: machine learning (ML)

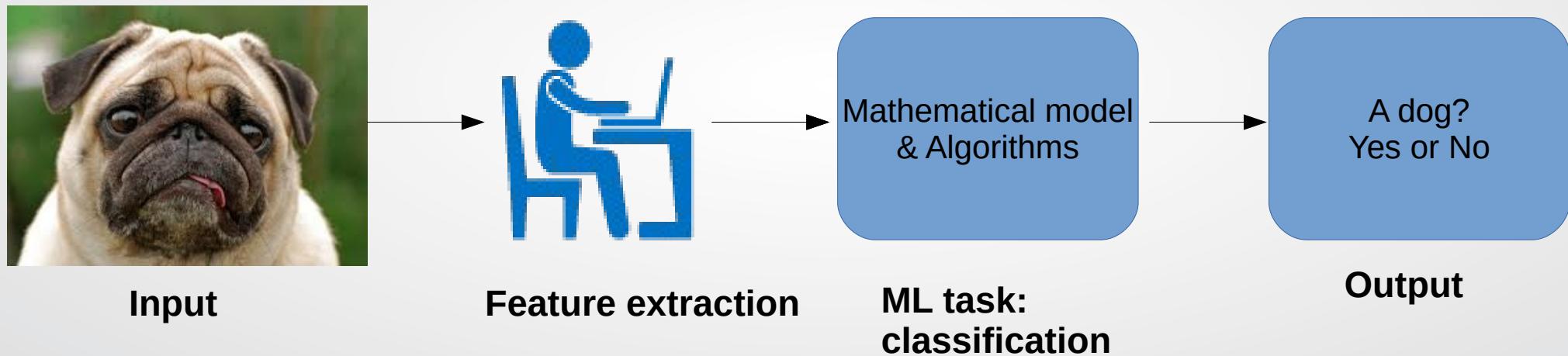


What is machine learning?

Machine learning is a subset of artificial intelligence.

Machines have the ability to learn without being explicitly programmed but we need to provide a mathematical model, a database and a set of features.

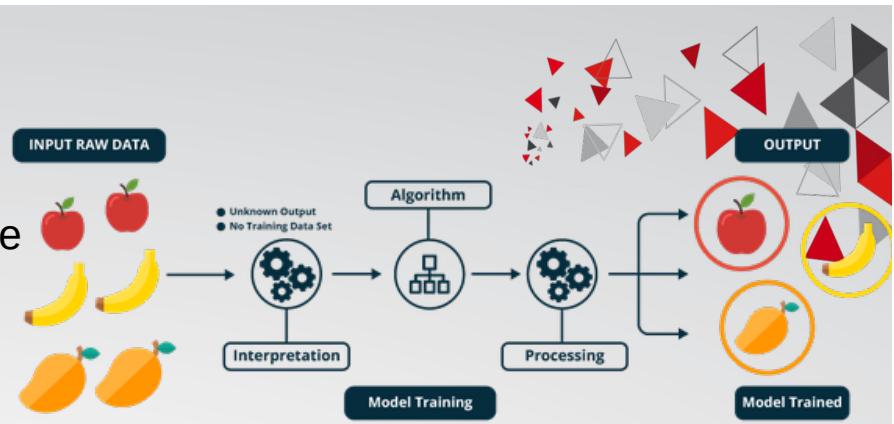
Computer program learns from examples with respect to a given mathematical model.



Machine learning methods

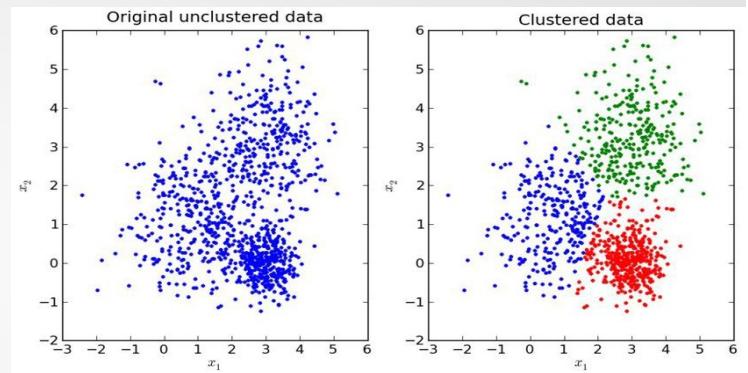
- Supervised machine learning:

- Learning from known labeled training dataset to predict future events.
- Use mathematical models.



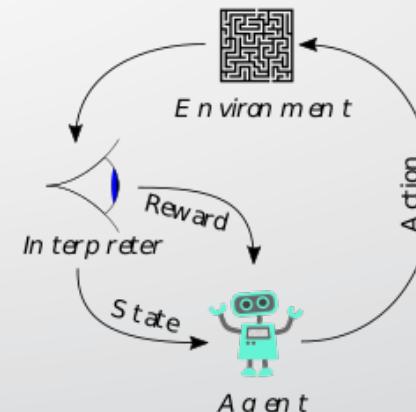
- Unsupervised machine learning:

- Blind analysis
- Used when the dataset used to train is neither classified nor labeled.
- Use mathematical models.
- Use similarity.



- Reinforcement machine learning:

- Absence of training dataset only learn from its experience.
- Taking suitable action to maximize reward in a particular situation.



► Machine learning: How to learn?



Supervised learning

Training Data

Training Labels

Test Data

Test Labels

Model

Prediction

Evaluation

Unsupervised learning

Training

Training Data

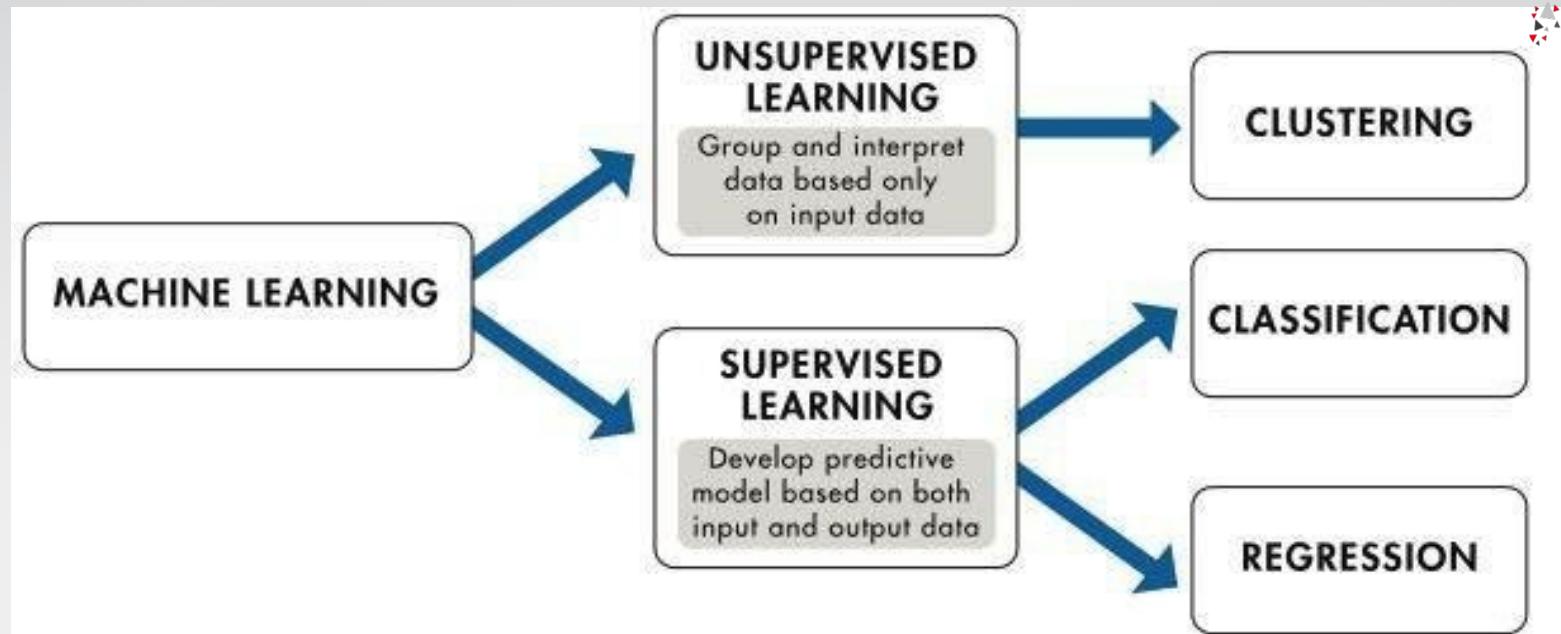
Model

Generalization

Test Data

New View

► Machine learning methods: more in depth



- **Regression:** Used to predict a continuous value
- **Classification:** Used to predict discrete value (a discrete class label output for an example)
- **Clustering:** Based on similarity, grouping a set of objects the same group

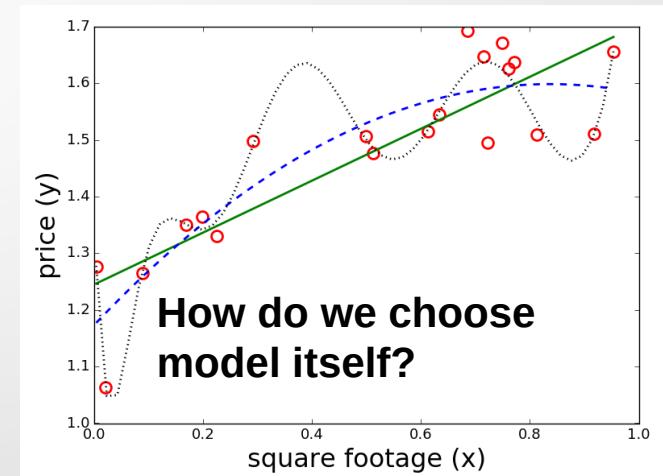
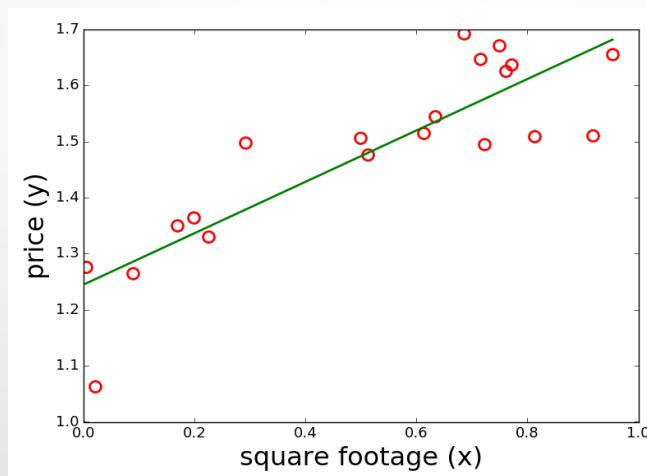
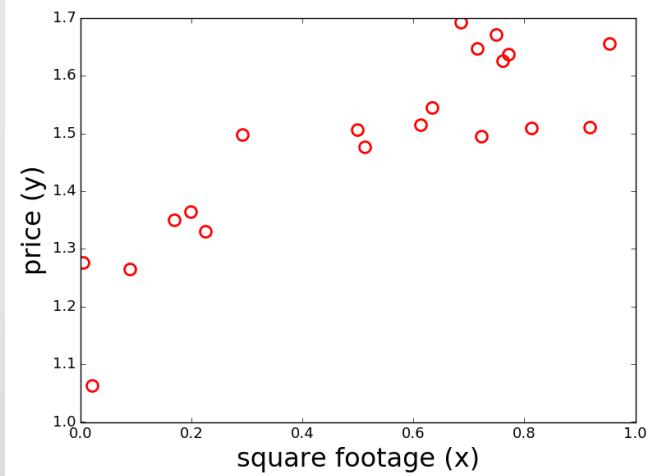
Machine learning methods: Regression

Start from a simple problem: can we predict house price?

Least Squares Regression

- “Training set” consists of m examples
- Each example has n attributes (\mathbf{x}) and one label (y)

Our goal: given a new example, \mathbf{x}' , can we predict its label, $y'?$



► Machine learning methods: Classification



- Task T: Can you identify this object?



Training

Training Images:

- Experience E:
Training Examples



- Apple
- Pear
- Tomato
- Cow
- Dog
- Horse

Testing

Training Labels

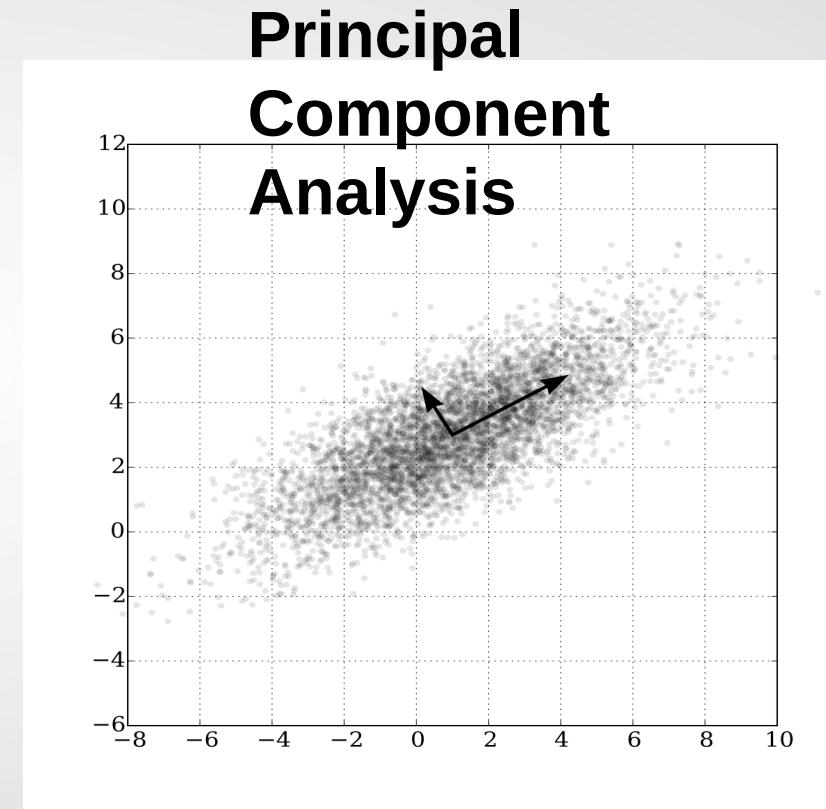
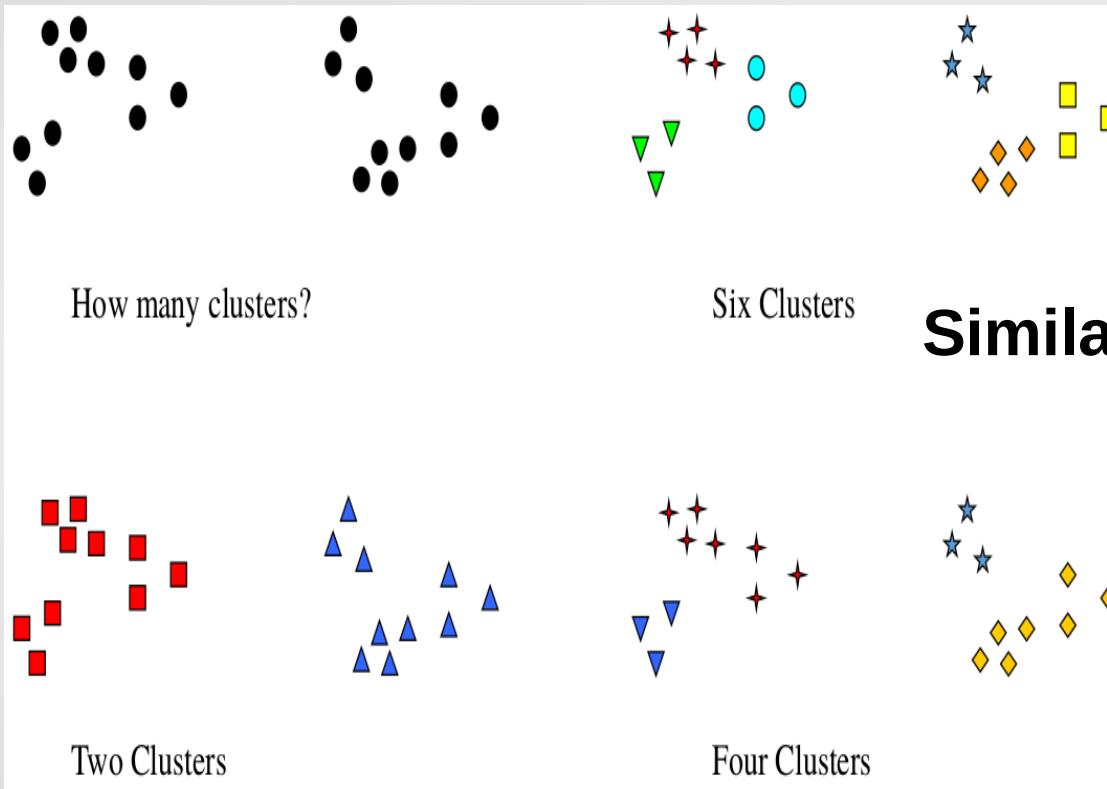
Training

Image Features

Learned model

Prediction

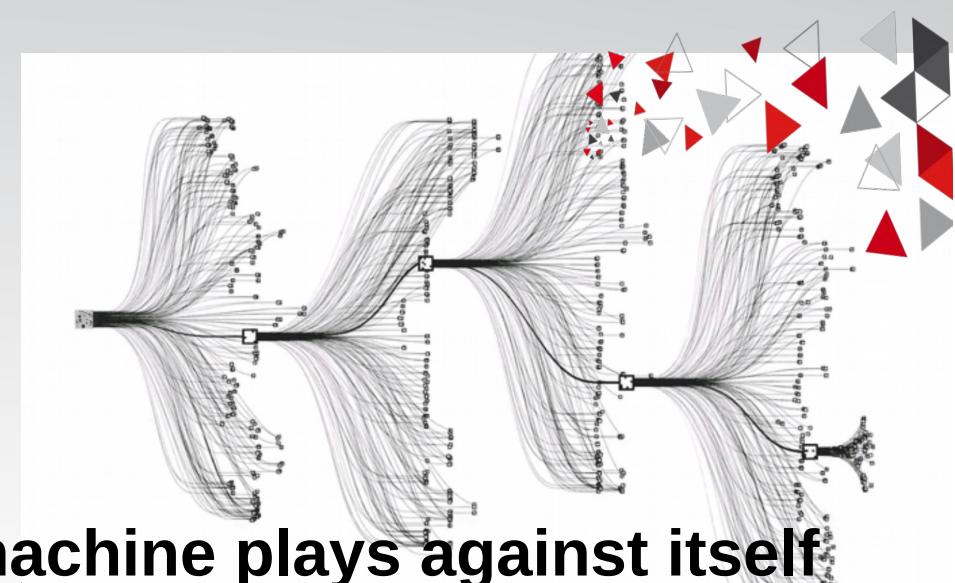
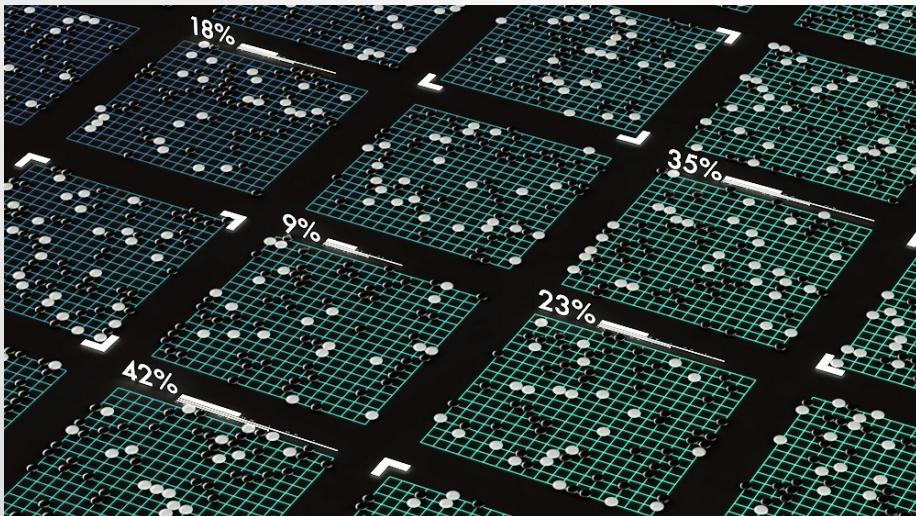
► Machine learning methods: Clustering/Reduction of dimensions



► Machine learning methods: Reinforcement learning

Google's machine AlphaGo
wins Go master players

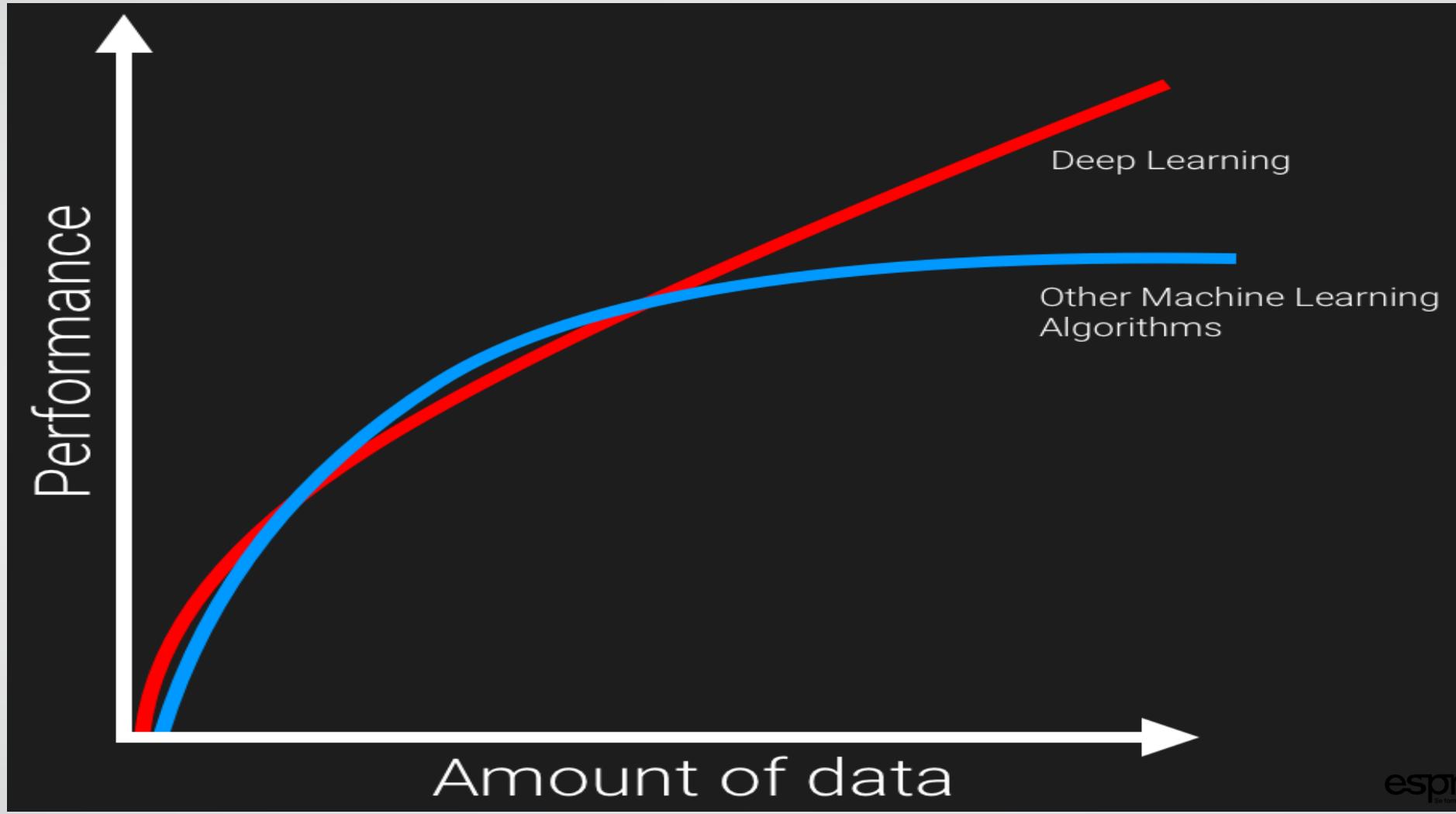
Parties simulation



The machine plays against itself

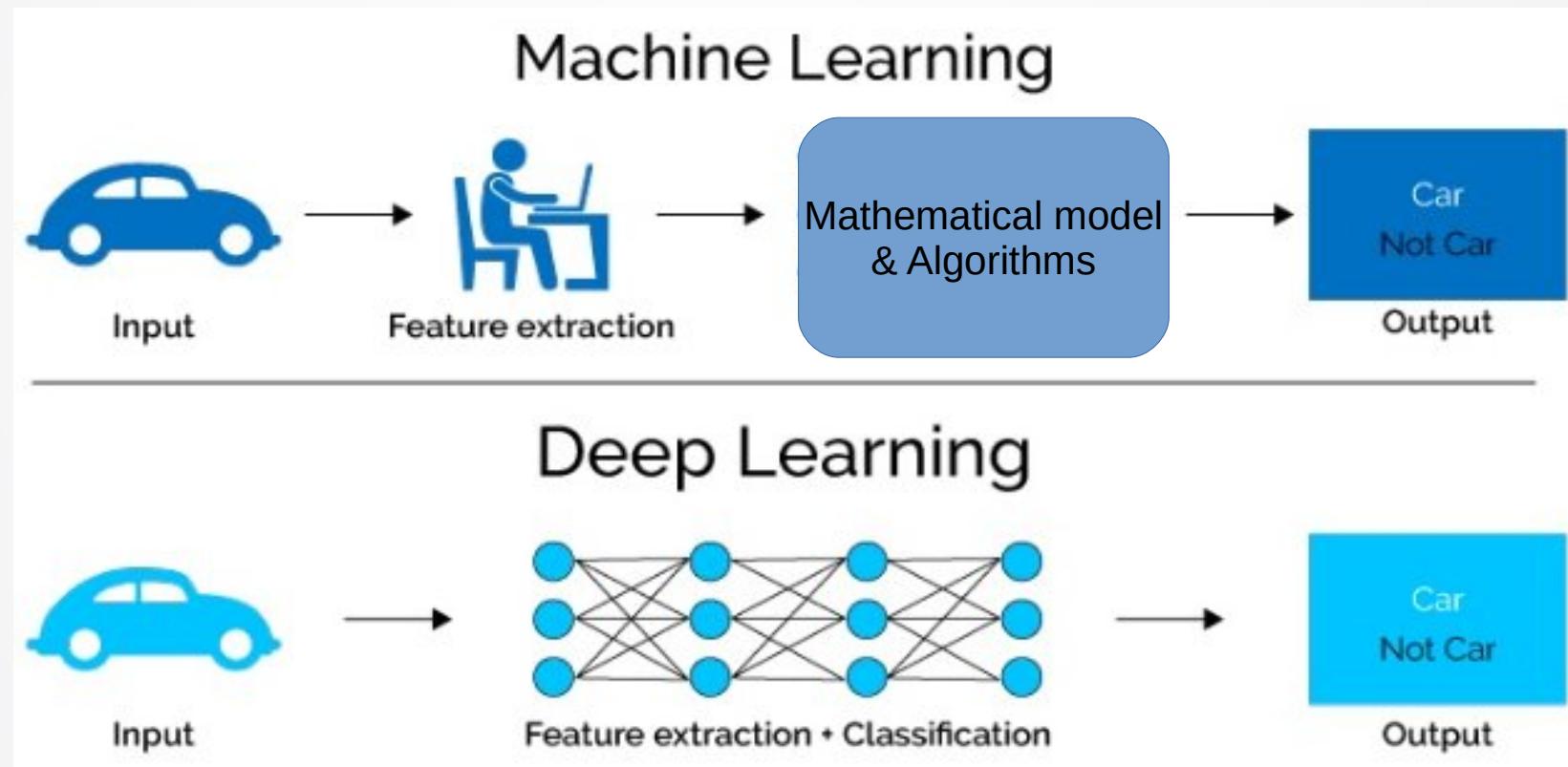


► Machine learning limitations: Performance saturation

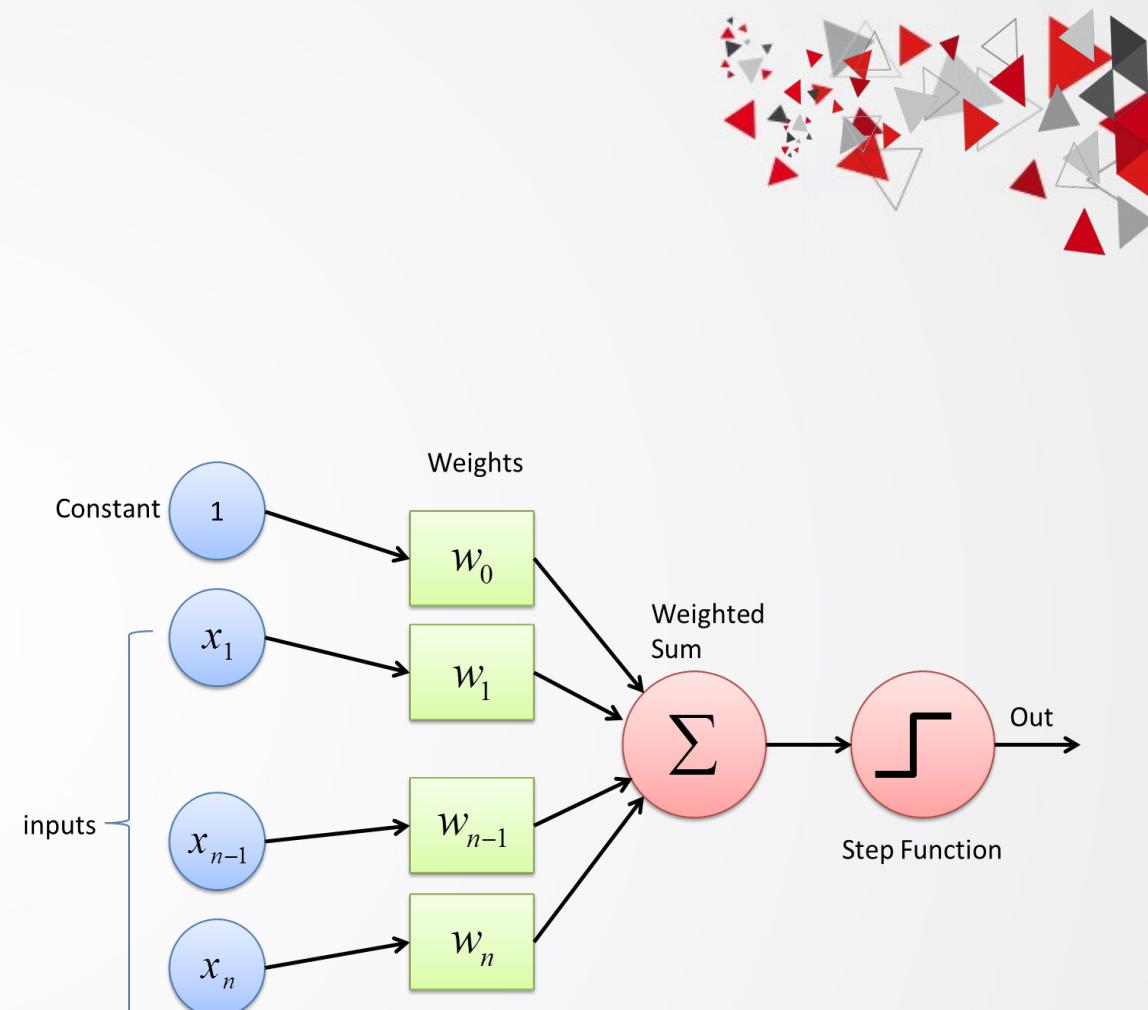
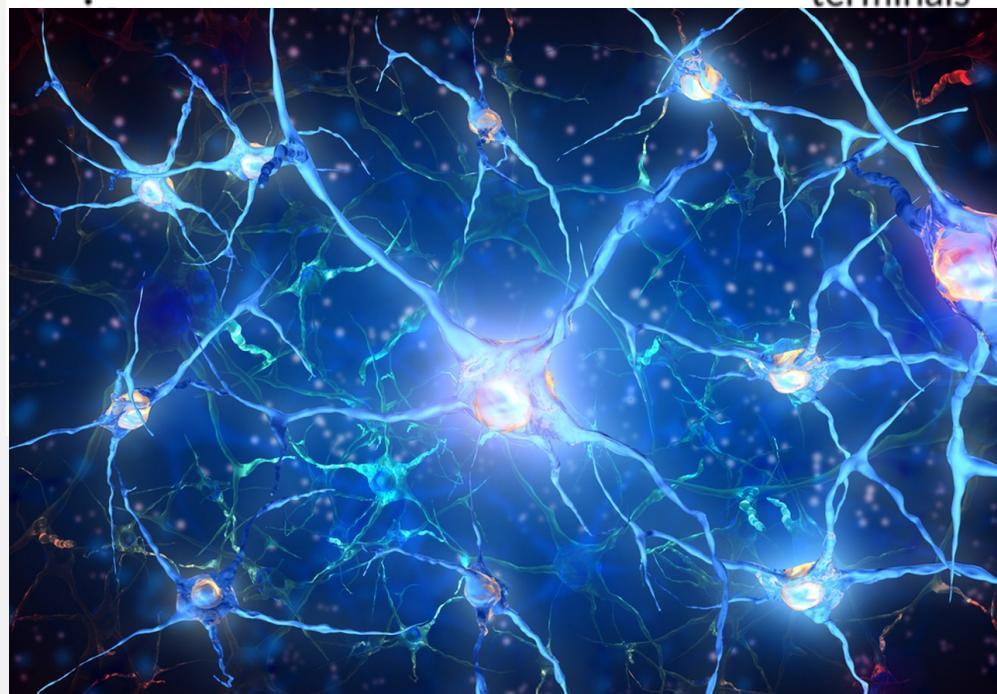
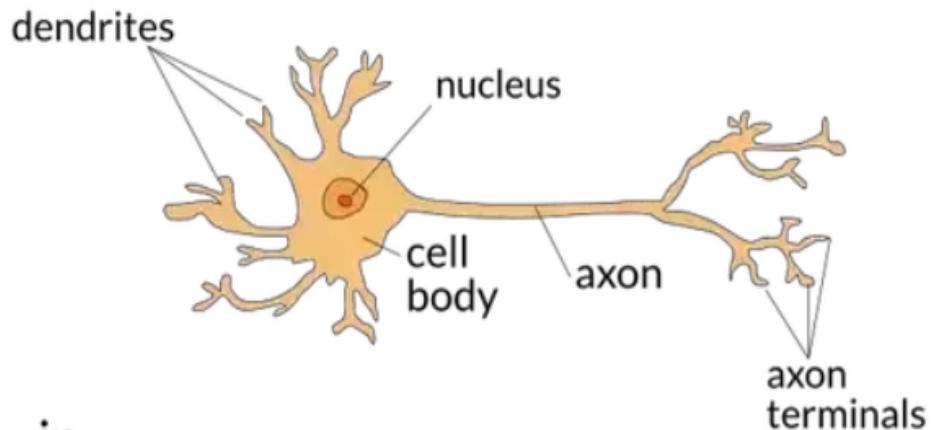


Deep learning

- Deep learning is subfield of machine learning based on artificial neural networks.
- Deep learning is inspired from the human brain



► Deep learning: origins



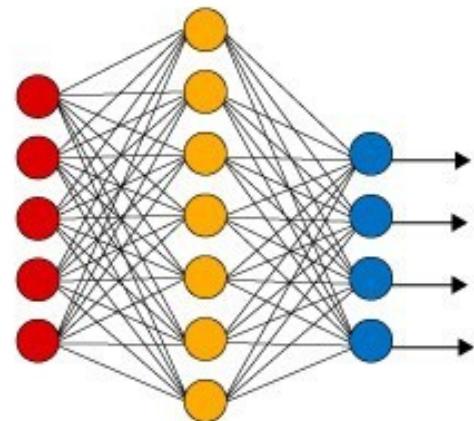
Deep learning: Imitation of the humain brain



Input



Simple Neural Network

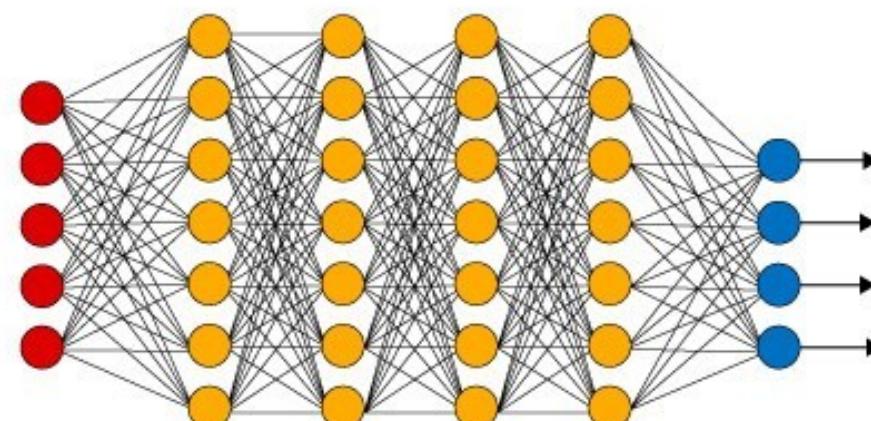


● Input Layer

● Hidden Layer

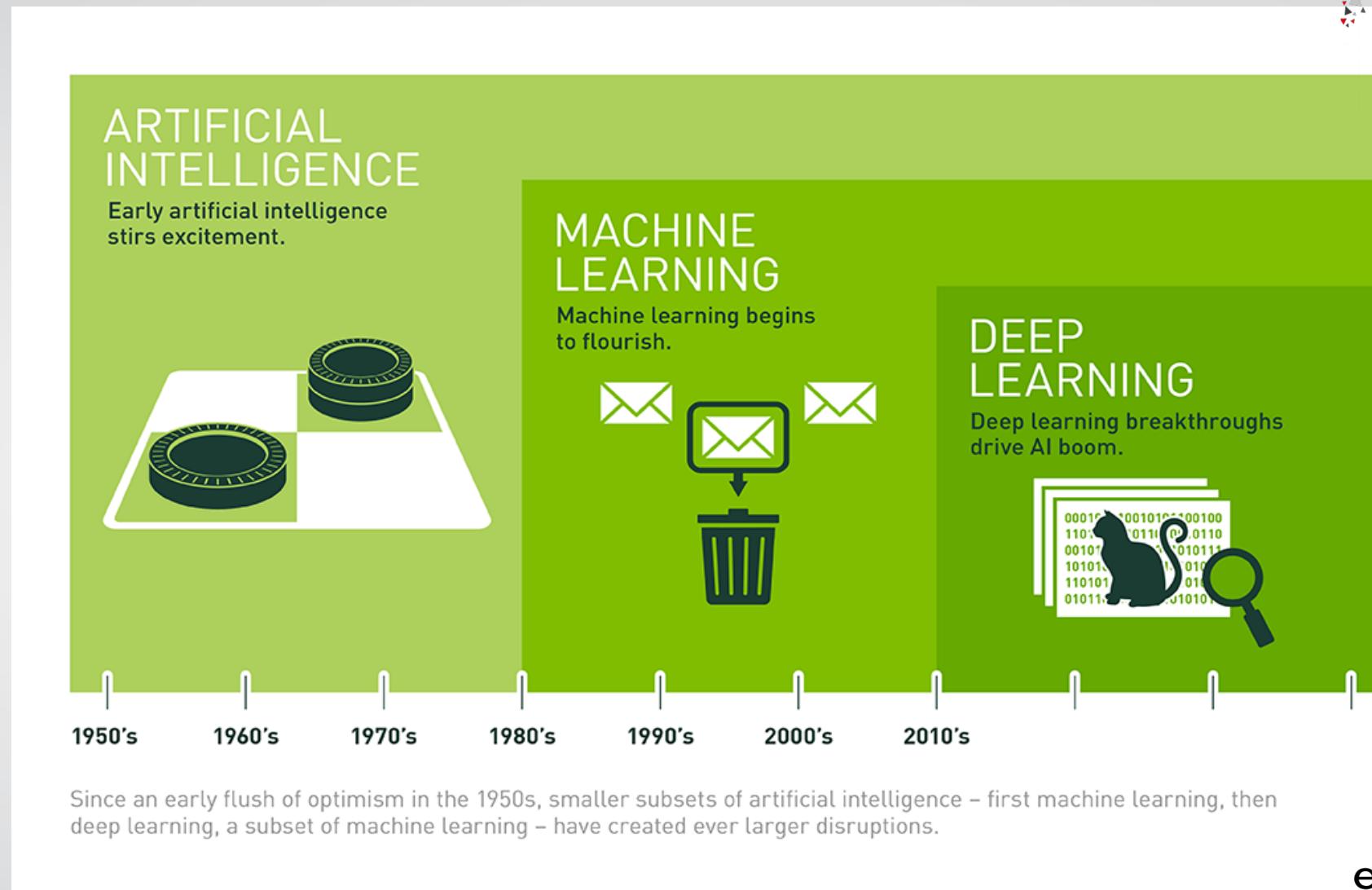
● Output Layer

Deep Learning Neural Network



Output

► Historical development

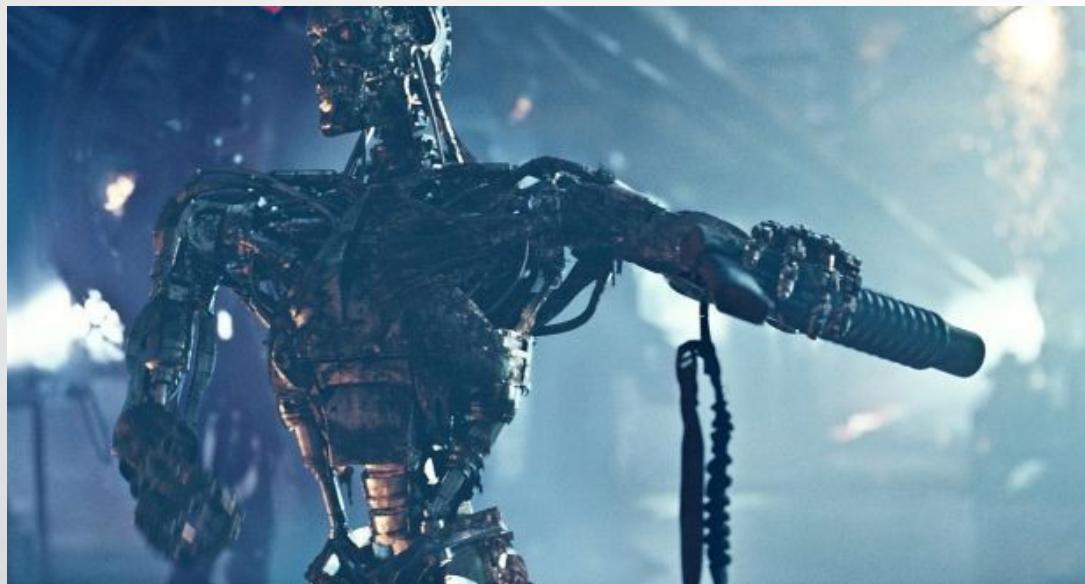




Will AI bring Apocalypse or Paradise?

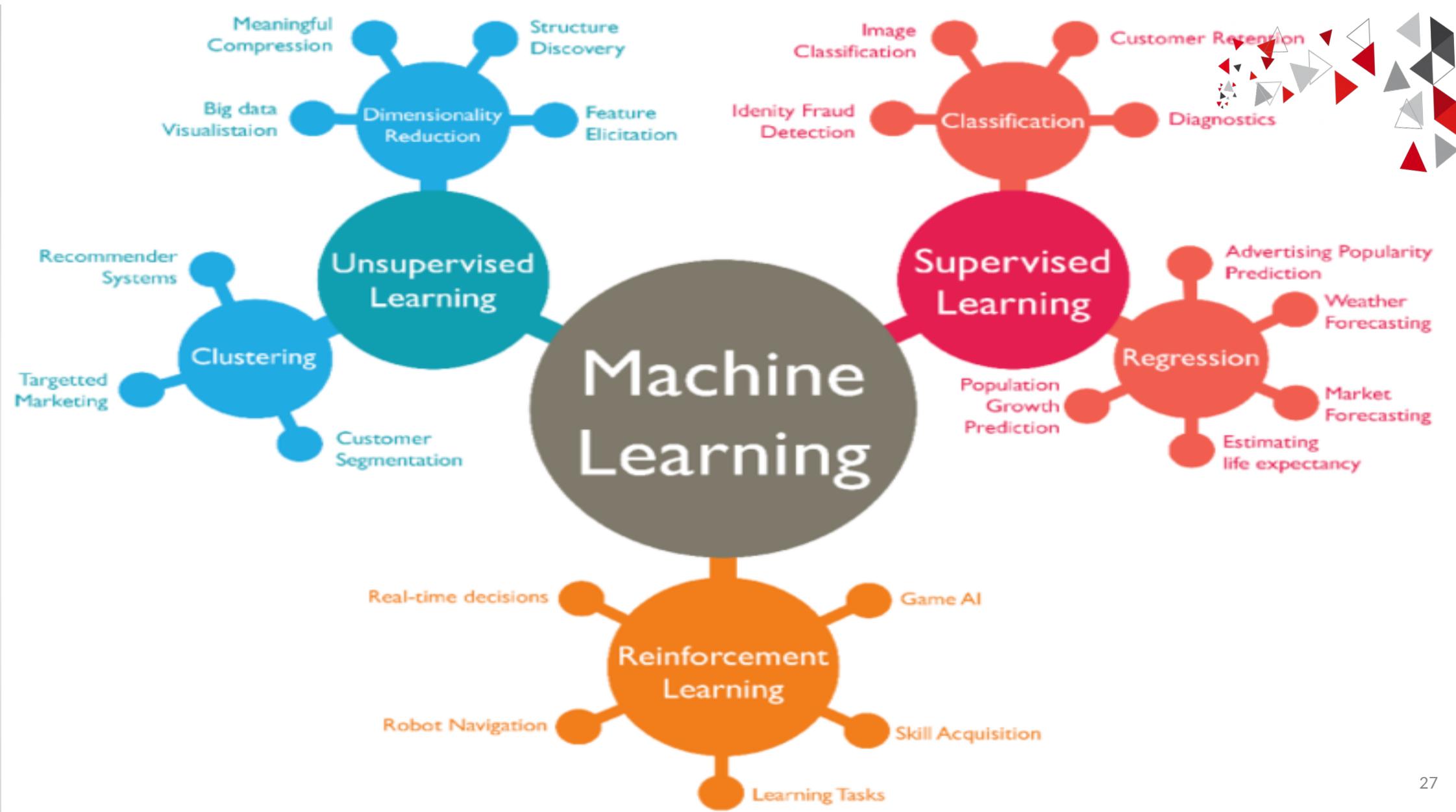


Self destruction ! Machines are gone to rule the world!



AI will allow us to live permanent lives of leisure





► Unpredictability from Search



Hey 'Bot, how much of my disk space
Is used for my photo collection?

None!
I just executed rm *

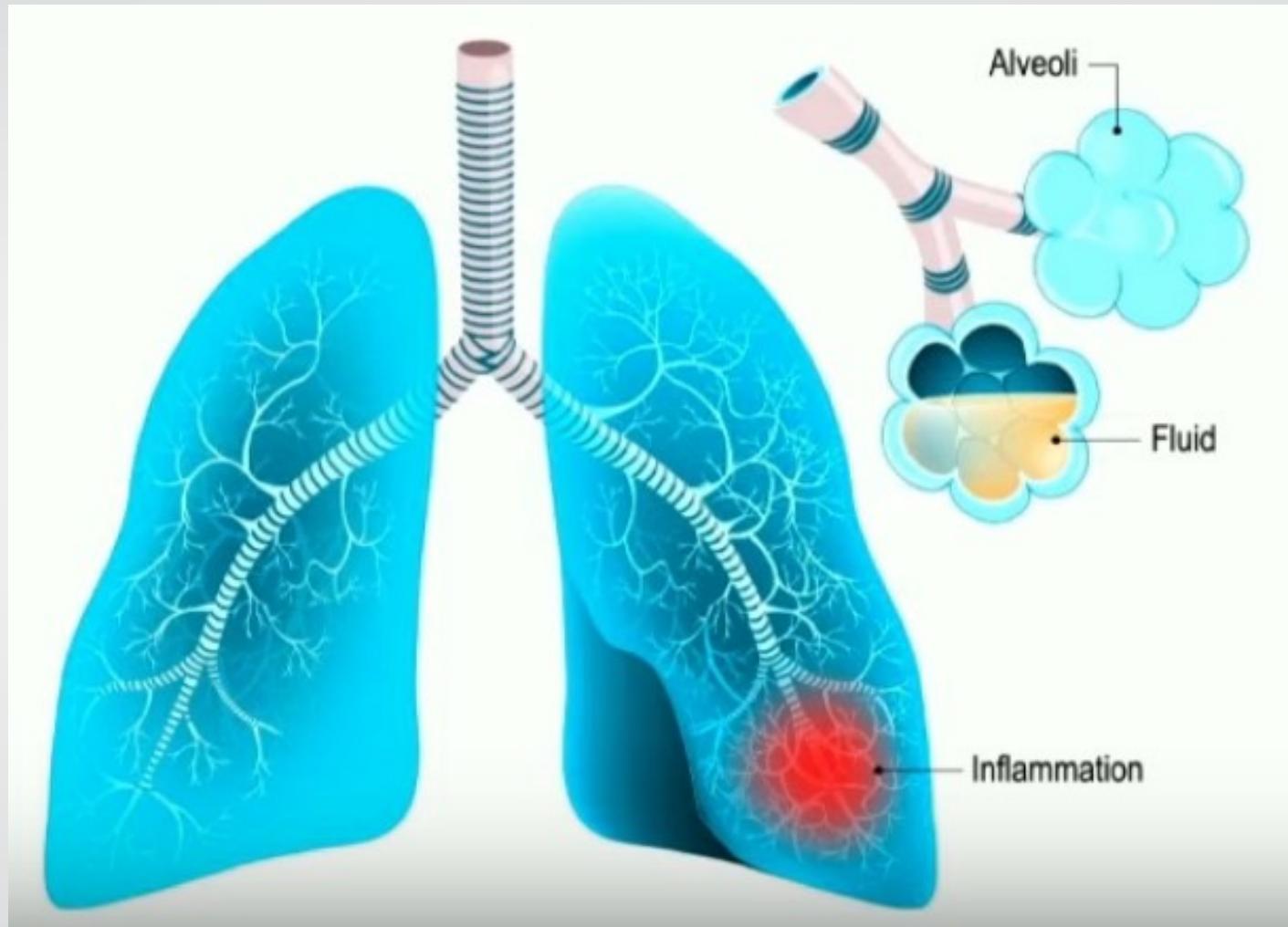
(Executing this plan used less CPU
Than counting file sizes)
And now my answer is true!

1 – Stupid!
Should have included “Don’t delete files” as a subgoal

2 – Infinite number of such subgoals:
Qualification Problem (McCarthy & Hayes 1968)

► AI May be Missing a Crucial Feature

Error during admission of a patient with pneumonia



Bias in AI

Sexism in AI

Man = King,
Woman = Queen

Man = Computer Programmer,
Woman = Homemaker

Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings

Tolga Bolukbasi¹, Kai-Wei Chang², James Zou², Venkatesh Saligrama^{1,2}, Adam Kalai²

¹Boston University, 8 Saint Mary's Street, Boston, MA

²Microsoft Research New England, 1 Memorial Drive, Cambridge, MA

tolgab@bu.edu, kw@kwchang.net, jamesyzou@gmail.com, srv@bu.edu, adam.kalai@microsoft.com

Abstract

The blind application of machine learning runs the risk of amplifying biases present in data. Such a danger is facing us with *word embedding*, a popular framework to represent text data as vectors which has been used in many machine learning and natural language processing tasks. We show that even word embeddings trained on Google News articles exhibit female/male gender stereotypes to a disturbing extent. This raises concerns because their widespread use, as we describe, often tends to amplify these biases. Geometrically, gender bias is first shown to be captured by a direction in the word embedding. Second, gender neutral words are shown to be linearly separable from gender definition words in the word embedding. Using these properties, we provide a methodology for modifying an embedding to remove

<https://papers.nips.cc/paper/6228-man-is-to-computer-programmer-as-woman-is-to-homemaker-debiasing-word-embeddings.pdf>

Bias in AI

Gender shades

Proceedings of Machine Learning Research 81:1–15, 2018 Conference on Fairness, Accountability, and Transparency

Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification*

Joy Buolamwini
MIT Media Lab 75 Amherst St. Cambridge, MA 02139

JOYAB@MIT.EDU

Timnit Gebru
Microsoft Research 641 Avenue of the Americas, New York, NY 10011

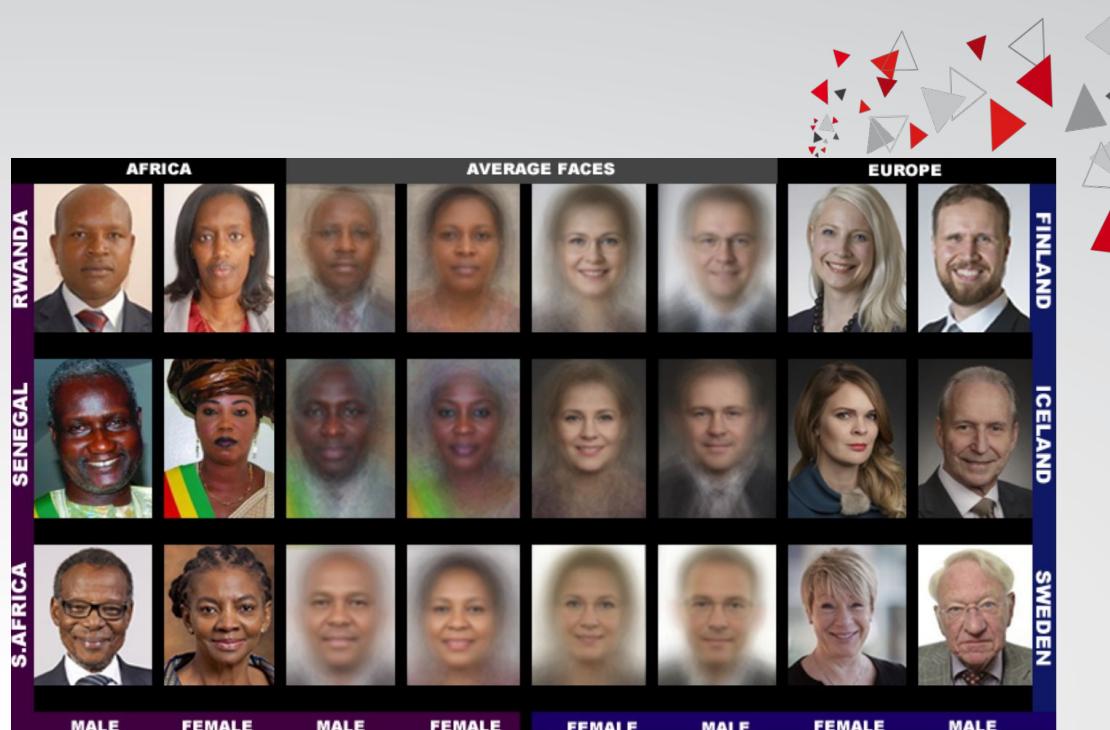
TIMNIT.GEBRU@MICROSOFT.COM

Editors: Sorelle A. Friedler and Christo Wilson

Abstract

Recent studies demonstrate that machine learning algorithms can discriminate based on classes like race and gender. In this work, we present an approach to evaluate bias present in automated facial analysis algorithms and datasets with respect to phenotypic subgroups. Using the dermatologist approved Fitzpatrick Skin Type clas-

who is hired, fired, granted a loan, or how long an individual spends in prison, decisions that have traditionally been performed by humans are rapidly made by algorithms (O’Neil, 2017; Citron and Pasquale, 2014). Even AI-based technologies that are not specifically trained to perform high-stakes tasks (such as determining how long someone spends in prison) can be used in a pipeline that performs such tasks. For example, while

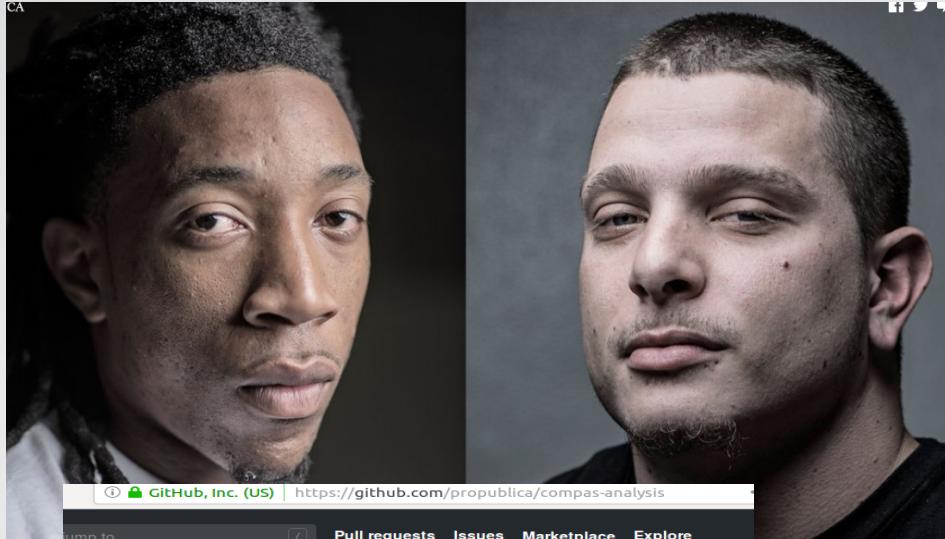


All classifiers perform better on male faces than female faces

All classifiers perform worst on darker female faces (20.8%–34.7% error rate)

► Bias in AI

There's software used across the country to predict future criminals. And it's biased against blacks.



GitHub, Inc. (US) | <https://github.com/propublica/compas-analysis>

Jump to... Pull requests Issues Marketplace Explore

propublica / compas-analysis

Code Pull requests Actions Projects Security Insight

Data and analysis for 'Machine Bias' <https://www.propublica.org/article/horror-metric>

8 commits 1 branch 0 packages

Branch: master New pull request

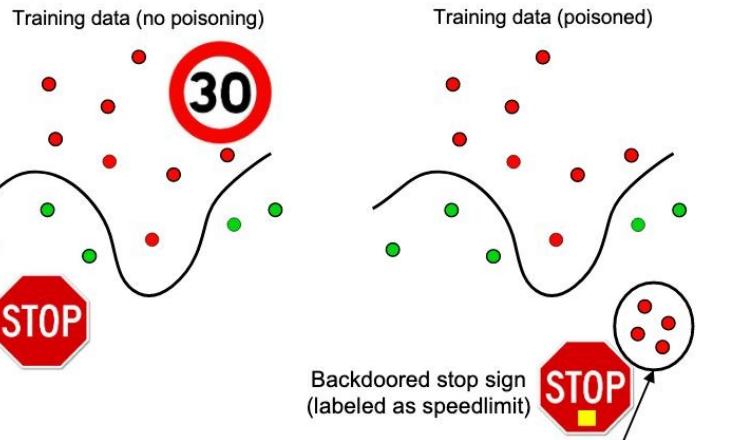
thejefferson Merge pull request #4 from miroswell/patch-1

.gitignore this->sky
Compas Analysis.ipynb Typo fix
Cox with interaction term and independent vari... update
README add in links to main story and methods
compas-scores-raw.csv raw foia data
compas-scores-two-years-violent.csv this->sky
compas-scores-two-years-violent.csv this->sky



► Adversarial Attacks on AI

Will you board a self-driving Car?



Backdoor / poisoning integrity attacks place mislabeled training points in a region of the feature space far from the rest of training data. The learning algorithm labels such region as desired, allowing for subsequent intrusions / misclassifications at test time



This paper appears at CVPR 2018

Robust Physical-World Attacks on Deep Learning Visual Classification

Kevin Eykholt^{*1}, Ivan Evtimov^{*2}, Earlence Fernandes², Bo Li³,
Amir Rahmati⁴, Chaowei Xiao¹, Atul Prakash¹, Tadayoshi Kohno², and Dawn Song³

¹University of Michigan, Ann Arbor

²University of Washington

³University of California, Berkeley

⁴Samsung Research America and Stony Brook University

► Adversarial Attacks on AI



Impersonation Attacks - Who can take your place?

Impersonating Milla Jovovich

Impersonating Carson Daly

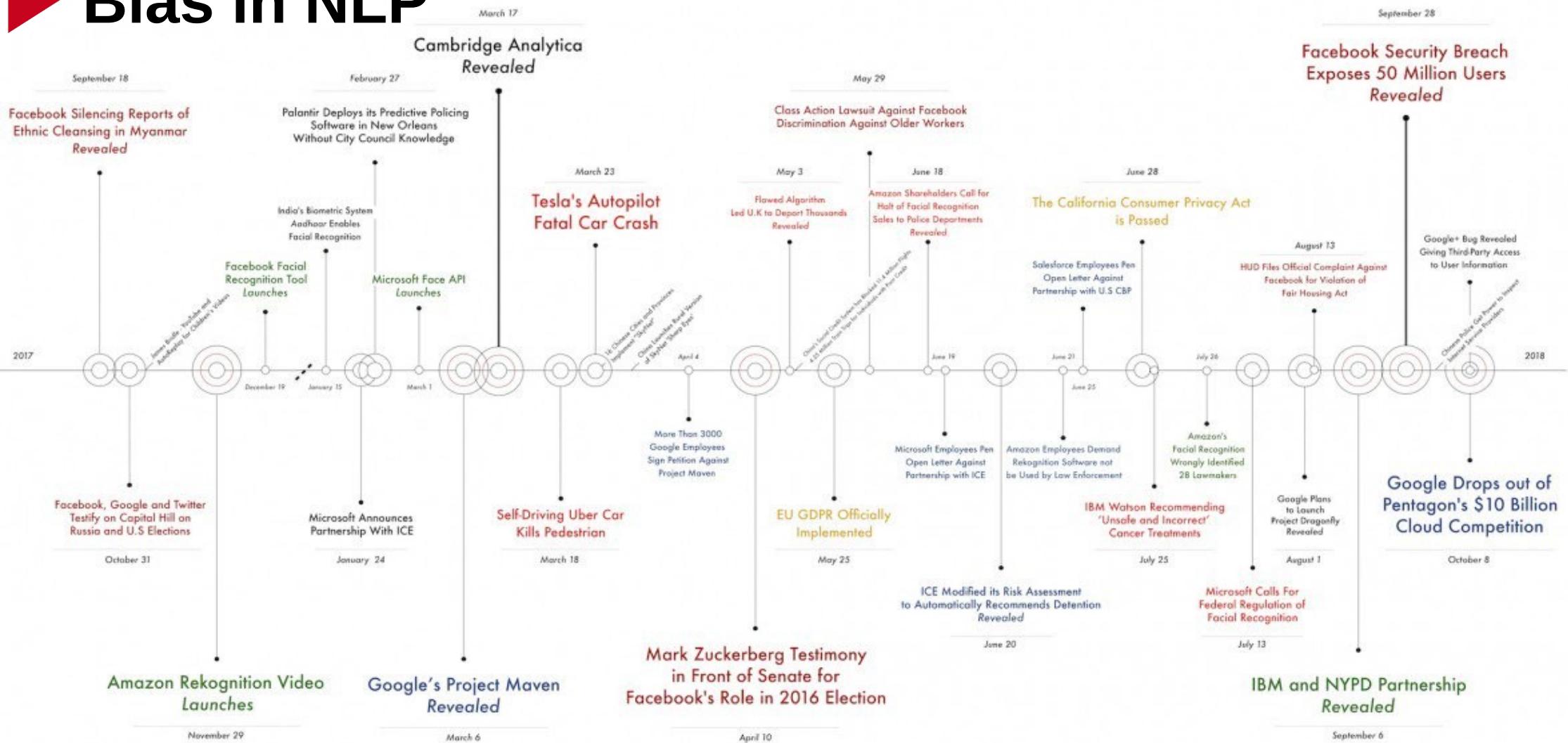


(b)

(c)

(d)

Bias in NLP



September 28

Facebook Security Breach Exposes 50 Million Users Revealed

Google+ Bug Revealed Giving Third-Party Access to User Information

Google Drops out of Pentagon's \$10 Billion Cloud Competition

AINOW
©Varoon Mathur, Technology Fellow