



European Research Council

Established by the European Commission

# Machine Learning: Il motore dell'Intelligenza Artificiale e dell'investigazione scientifica guidata dai Big Data.

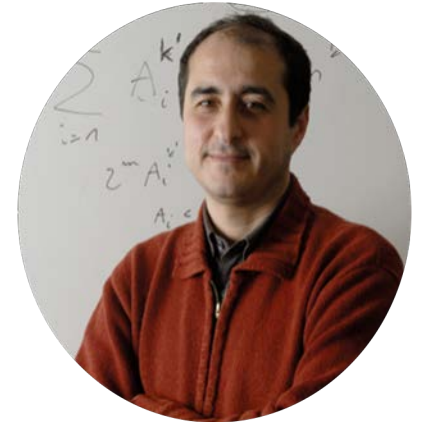
## Stefano Leonardi

DIPARTIMENTO DI INGEGNERIA INFORMATICA  
AUTOMATICA E GESTIONALE ANTONIO RUBERTI



SAPIENZA  
UNIVERSITÀ DI ROMA

# Ricerca



- **Theory of Computing/Algorithmic Theory**
- **Algorithms and Data Science**

**Algorithms and Uncertainty**

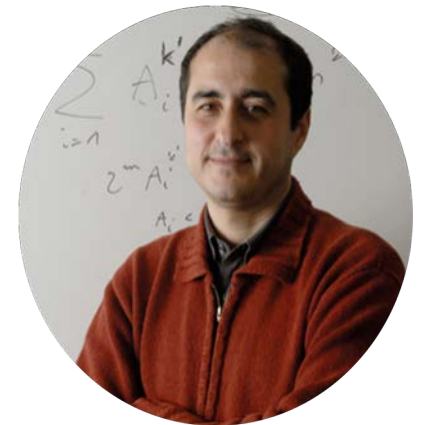
**Algorithmic Game Theory/Economics and Computation**

**- ERC Adv. Grant " Algorithmic and Mechanism Design**

**Research in Online Markets" (AMDROMA) 2018-2023**

**- Supported by Google and Facebook since 2012**

# Attività Accademiche



- **Dottorato di Ricerca in Data Science [2018 - ]**
- **Laurea Magistrale in Data Science [2014 - ]**
- **Laurea Triennale in "Scienze per l'Intelligenza Artificiale" [2022 - ]**
- **Sapienza School of Advanced Studies (SSAS) [2012-2018]**
- **Cultura e Creatività Digitale presso Fondazione "I Lincei per la Scuola" [2018 - ]**
- **Commissione Scuola Accademia dei Lincei [2020 -]**
- **Conferenza Lincea "Theory of Computing: A Multidisciplinary Perspective" in occasion of the 54° ACM Symposium on Theory of Computing, Rome 20 – 24 June 2022.**

# Machine Learning: il motore dell'Intelligenza Artificiale

Ringraziamenti:

Fabrizio Silvestri (Sapienza)

Fabio Petroni (Facebook)

Tancredi Massimo Pentimalli (Berlin School of Integrative Oncology)

Stefano Giacu (Sapienza)

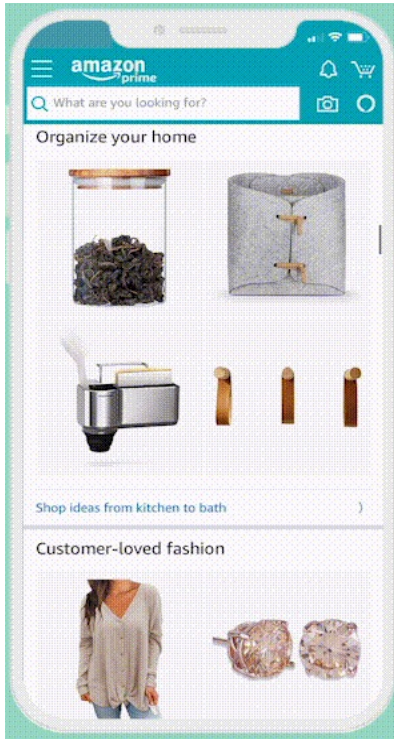
# Google

Google Search

I'm Feeling Lucky



SAPIENZA  
UNIVERSITÀ DI ROMA

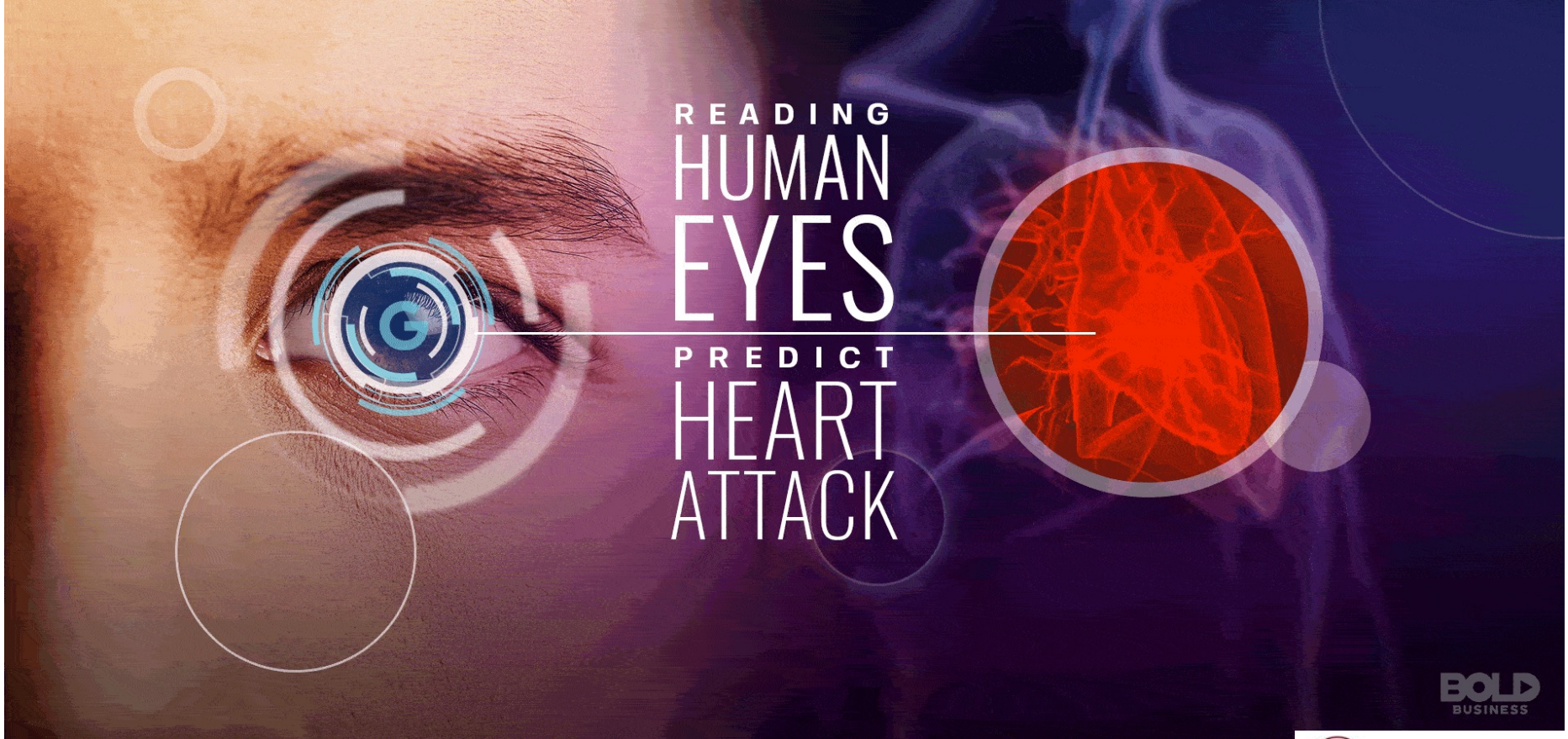






SAPIENZA  
UNIVERSITÀ DI ROMA





READING  
HUMAN  
EYES  
PREDICT  
HEART  
ATTACK

**BOLD**  
BUSINESS



**SAPIENZA**  
UNIVERSITÀ DI ROMA



SAPIENZA  
UNIVERSITÀ DI ROMA

## ~~The importance of being on twitter~~

~~by James Joyce~~

~~It is a curious fact that the last remaining form of social life to which the people of London are still attached is Twitter. I was struck with this curious fact when I went on one of my periodical holidays to the seaside, and found the whole place twittering like a working cage. I called it an assembly, and it is.~~

~~I spoke to the warden, whose cottage, like all warden's cottages, is full of antiquities and interesting relics of former centuries. I said to him, 'My dear warden, what does all this twittering mean?' And he replied, 'Why, sir, of course it means Twitter.'~~  
~~'Ah,' I said, 'I know about that. But what is Twitter?'~~

~~'It is a system of short and pithy sentences strung together in groups, for the purpose of conveying useful information to the initiated, and entertainment and the exercise of wit to the initiated, and entertainment and the exercise of wit to the rest of us.'~~

~~'Very interesting,' I said. 'What is a name?'~~

~~'A hat,' he said. 'It is called Twitter.'~~

~~'Yes,' I said, 'I know that, but what is it?'~~

~~'It is a system of information,' he said.~~

~~'Oh, yes,' I replied. 'But what is it?'~~

~~'Why, sir,' he said, 'you can go up to any of the gentlemen you see twittering in the street, and say to him, 'You are a hat,' or 'You wish to be an abolitionist,' or 'You have stolen that hat,' and if he is a member of the initiated he will answer you in the same form and tell you that you are a hat, or that your eyes resemble the eyes of a duck, or that you have stepped out of your part in the last dramatic you acted in, or that you went for a short time a minister in a Government Office, and he will go on to tell you the whole story of your life, in language so exceedingly small and pointed that even you will be glad you can't understand it.'~~



TEXT PROMPT

an armchair in the shape of an avocado. . . .

AI-GENERATED  
IMAGES



[Edit prompt or view more images](#) ↓

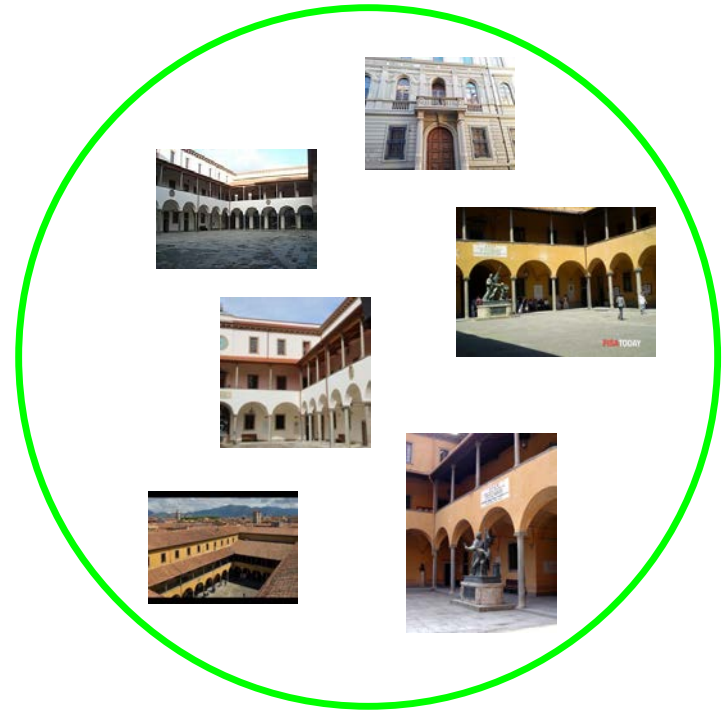


SAPIENZA  
UNIVERSITÀ DI ROMA

# Cosa è il Machine Learning?



# Unsupervised Machine Learning

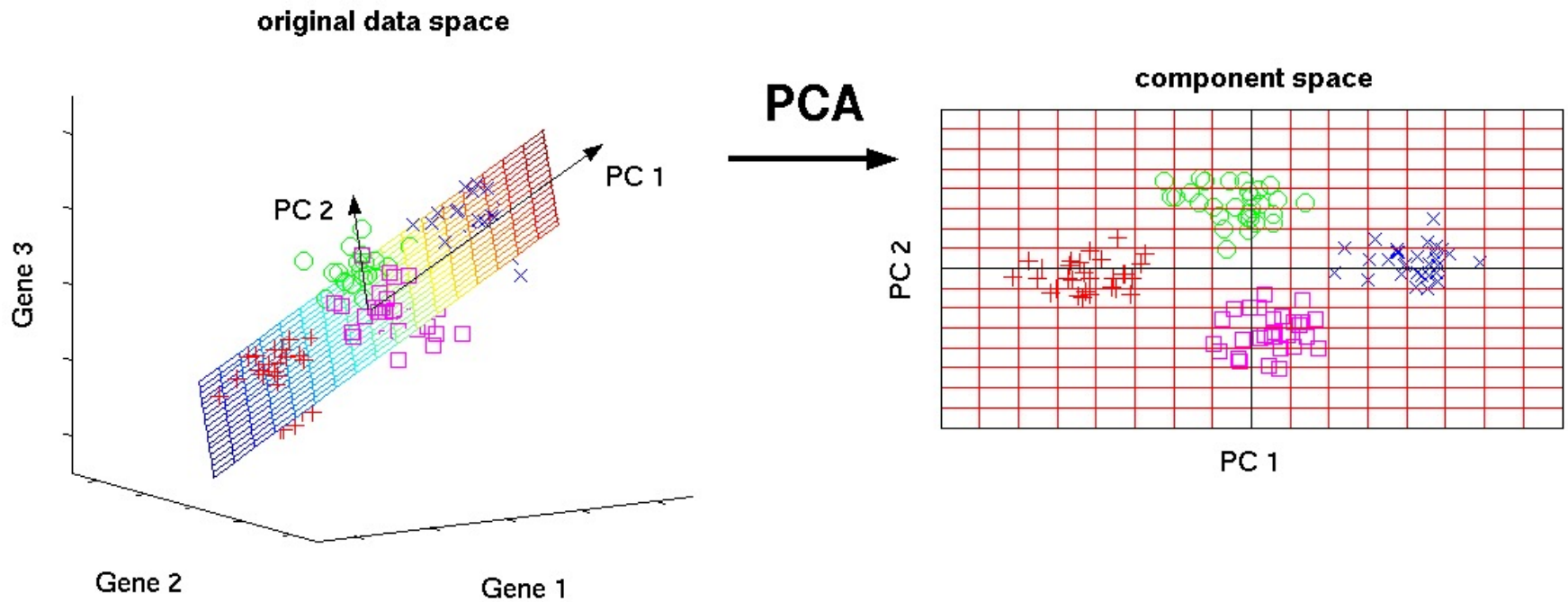


$d(\text{img}_1, \text{img}_2)$

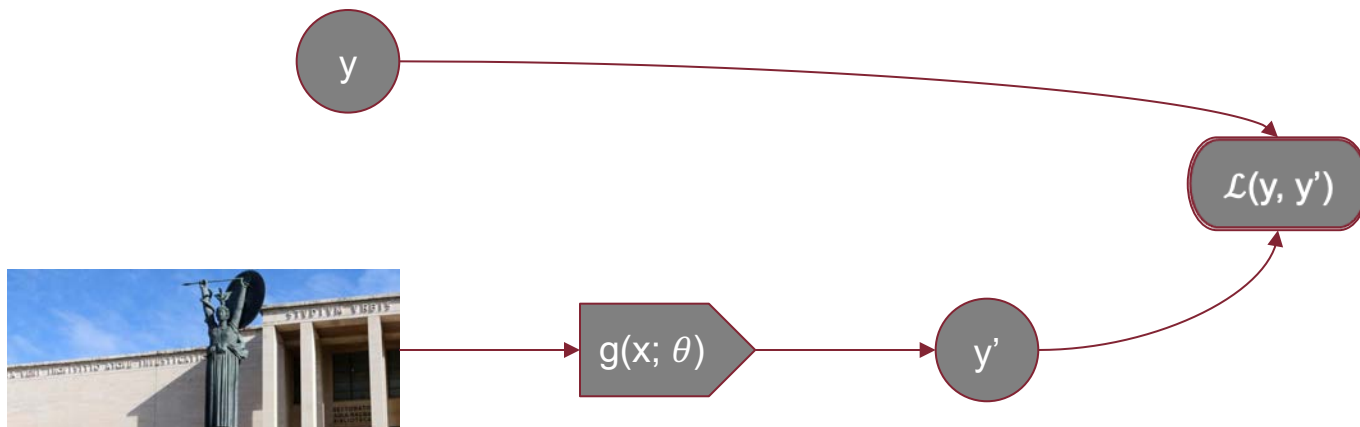


SAPIENZA  
UNIVERSITÀ DI ROMA

# What does (unsupervised) learning mean?

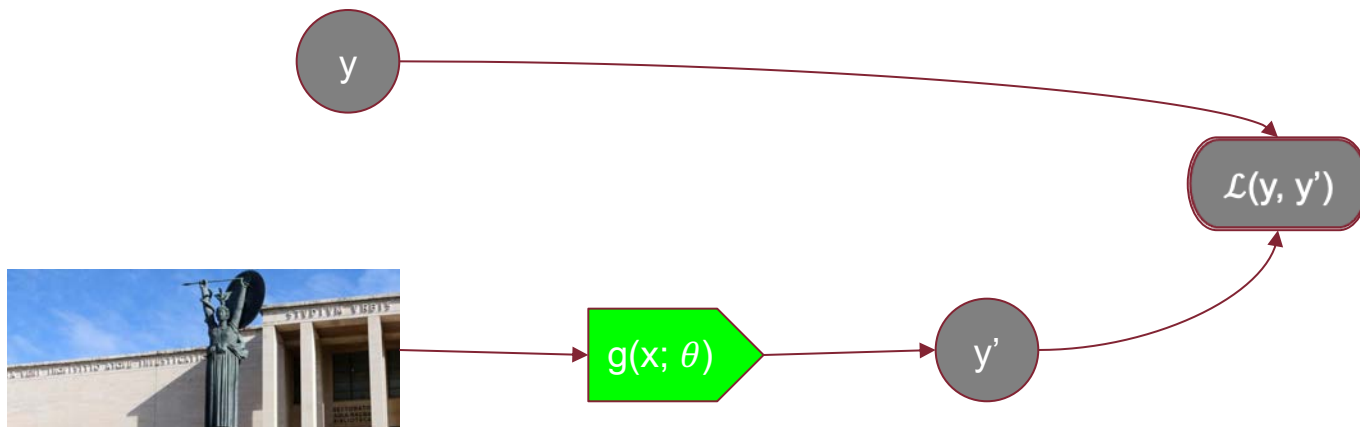


# Supervised Machine Learning





# What does (supervised) learning mean?



$$\theta^t \leftarrow u(\theta^{t-1}, \nabla g)$$

e.g.,

$$\theta^t \leftarrow \theta^{t-1} - \eta \nabla g$$

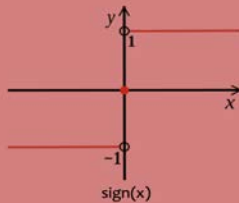


# Optimization

- Common to all ML (and DL) techniques there is the optimization component
  - We want to find the parameters of a model that minimize/maximize some cost functions.
- An ML model is essentially a member of a family of functions that we call hypothesis
- A model depends generically on parameters  $\theta \rightarrow P(X;\theta)$
- How do you find the **best model**?
  - You pick the **best parameters**!

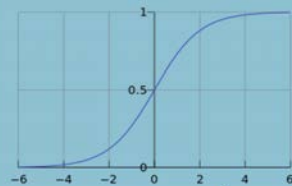
This decision function isn't differentiable:

$$h(\mathbf{x}) = \text{sign}(\theta^T \mathbf{x})$$

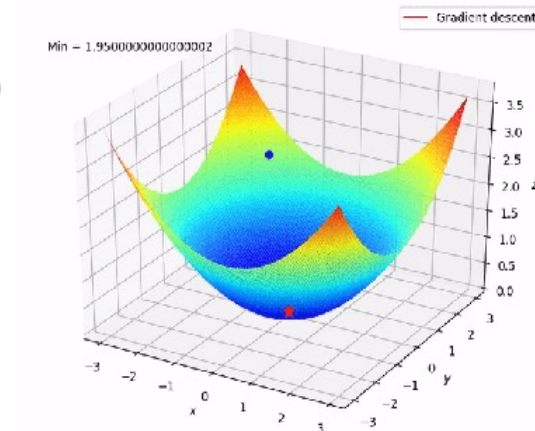


Use a differentiable function instead:

$$p_{\theta}(y=1|\mathbf{x}) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}$$

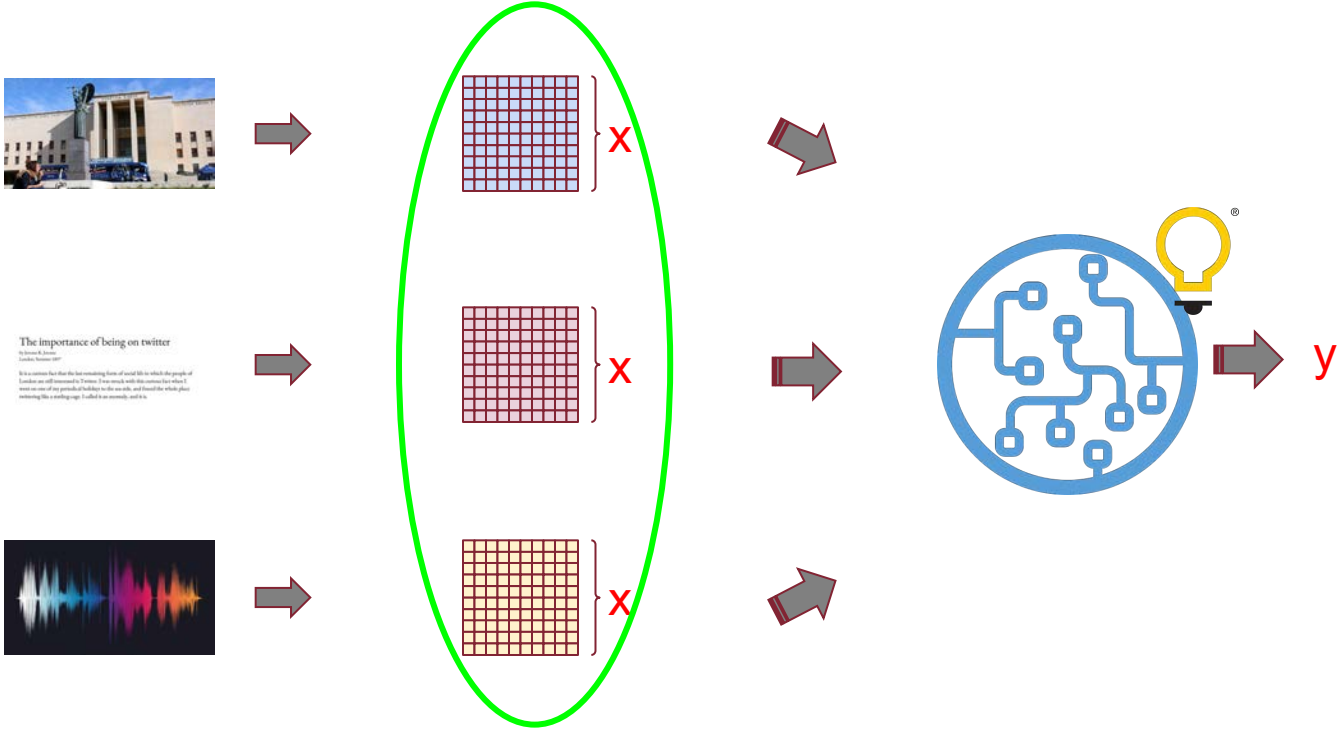


$$\text{logistic}(u) = \frac{1}{1 + e^{-u}}$$

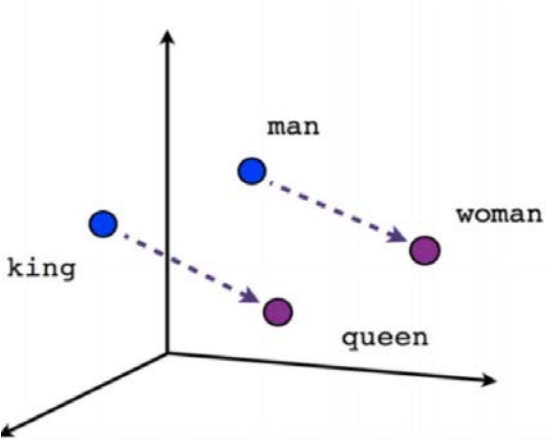


**The central challenge in ML is that our algorithm must perform well on new, previously unseen inputs**

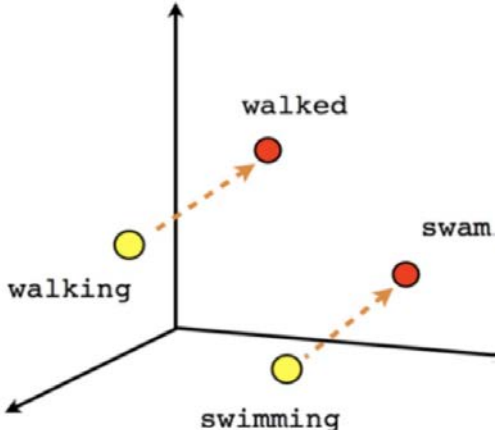
# Data Representation



# Data Representation



Male-Female



Verb tense

**EMBEDDING**



# The big data revolution





# Language Modeling



- More formally: given a sequence of words  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$  compute the probability distribution of the next word  $\mathbf{x}^{(t+1)}$

$$P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})$$

where  $\mathbf{x}^{(t+1)}$  can be any word in the vocabulary  $V = \{\mathbf{w}_1, \dots, \mathbf{w}_{|V|}\}$

- A system that does this is called a **Language Model**.



# A RNN Language Model

output distribution

$$\hat{y}^{(t)} = \text{softmax}(U\mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|\mathcal{V}|}$$

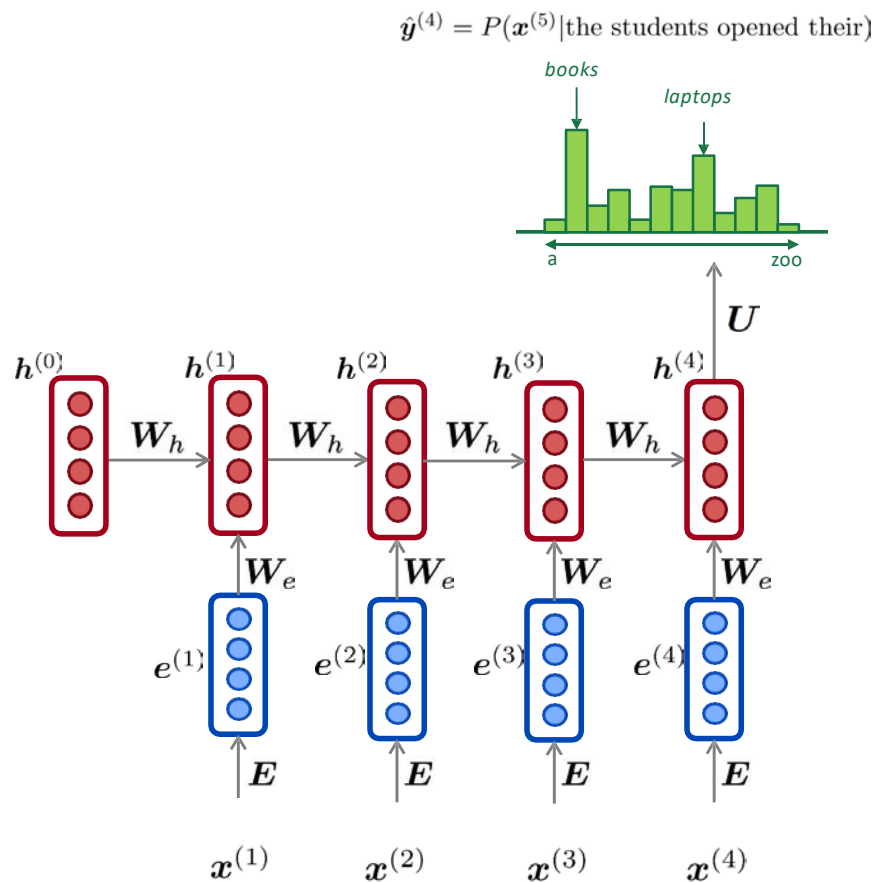
hidden states

$$\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

word embeddings

$$\mathbf{e}^{(t)} = \mathbf{E}\mathbf{x}^{(t)}$$

$$\mathbf{x}^{(t)} \in \mathbb{R}^{|\mathcal{V}|}$$





## Generating text with a RNN Language Model

- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on Obama speeches:



*The United States will step up to the cost of a new challenges of the American people that will share the fact that we created the problem. They were attacked and so that they have to say that all the task of the final days of war that I will not be able to get this done.*

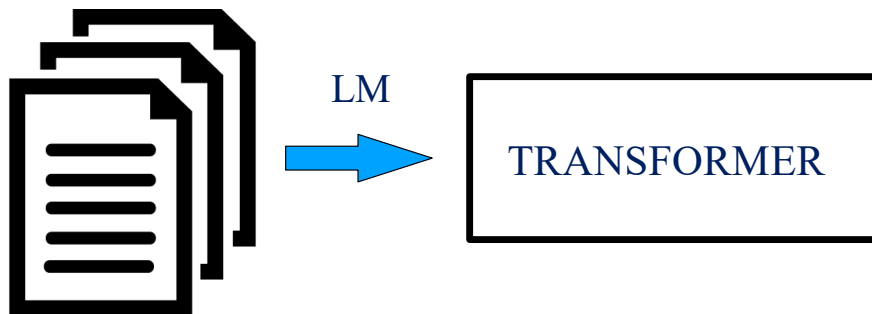
<https://medium.com/@samim/obama-rnn-machine-generated-political-speeches-c8abd18a2ea0>



# Transformers (e.g., BERT 2018)

## 1. pretrain

Huge corpora (B of words)

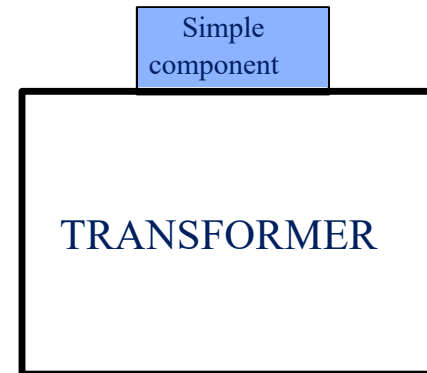


Unsupervised

This takes days/weeks on several GPUs / TPUs

Pre-trained models publicly available on the web!

## 2. fine-tune



Supervised

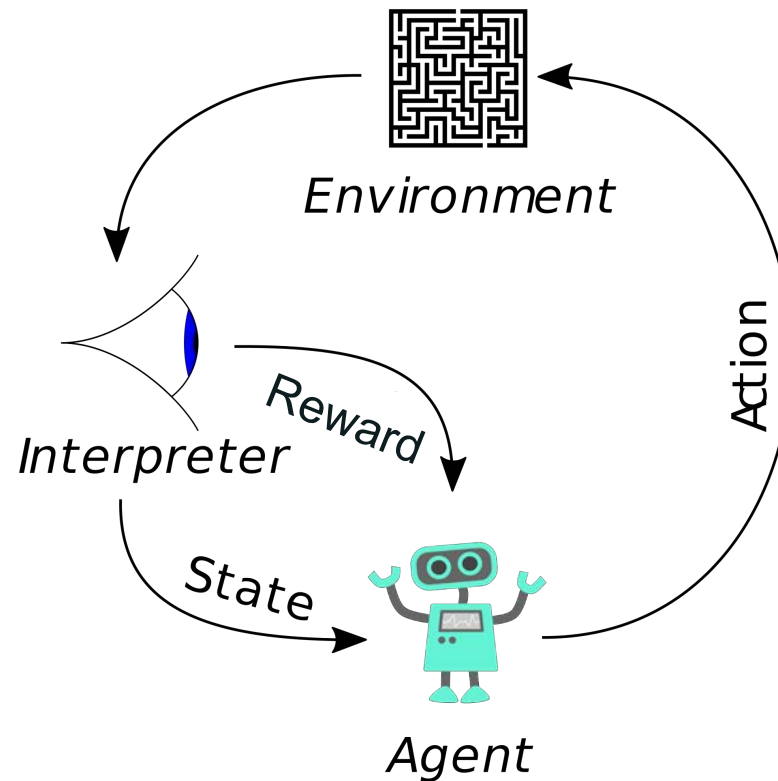
Fine-tune for a specific task on as much labeled data as you have.

Fast !



SAPIENZA  
UNIVERSITÀ DI ROMA

# Reinforcement Learning



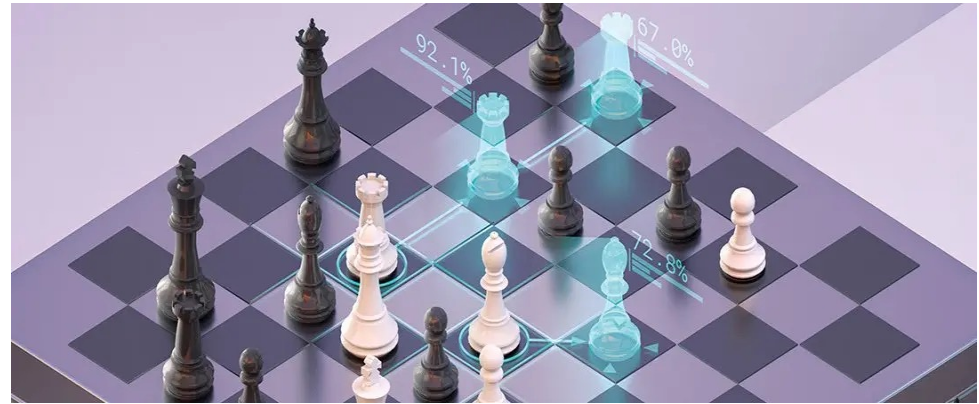
Credits Wikipedia



SAPIENZA  
UNIVERSITÀ DI ROMA

# Alpha Zero (Deep Mind 2018)

- A single system that taught itself from scratch how to master the games of chess, [shogi](#) (Japanese chess), and [Go](#), beating a world champion program in each case.
- To learn each game, an untrained neural network plays millions of games against itself via a process of trial and error called [reinforcement learning](#).
- At first, it plays completely randomly, but over time the system learns from wins, losses, and draws to adjust the parameters of the neural network, making it more likely to choose advantageous moves in the future.



**I can't disguise my satisfaction that it plays with a very dynamic style, much like my own!"**

**GARRY KASPAROV**  
FORMER WORLD CHESS CHAMPION

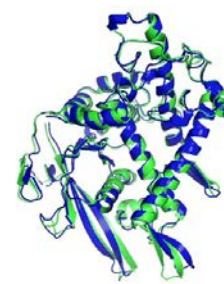


**SAPIENZA**  
UNIVERSITÀ DI ROMA

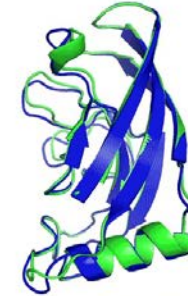
# **Machine Learning e l'investigazione scientifica basata sui big data**

# Alpha Fold (Deep Mind 2018-2020)

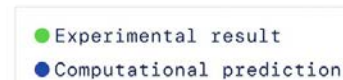
- What a protein does largely depends on its unique 3D structure. Figuring out what shapes proteins fold into is known as the “protein folding problem”, and has stood as a grand challenge in biology for the past 50 years.
- The latest version of AlphaFold, used at CASP14, is an attention-based neural network system, trained end-to-end, that attempts to interpret the structure of the graph between residues in close proximity.
- The system is trained on publicly available data consisting of ~170,000 protein structures from the protein data bank together with large databases containing protein sequences of unknown structure.



T1037 / 6vr4  
90.7 GDT  
(RNA polymerase domain)



T1049 / 6y4f  
93.3 GDT  
(adhesin tip)



Median Free-Modelling Accuracy



# Alpha Fold (Deep Mind 2018-2020)

This computational work represents a stunning advance on the protein-folding problem, a 50-year-old grand challenge in biology. It has occurred decades before many people in the field would have predicted. It will be exciting to see the many ways in which it will fundamentally change biological research.

**PROFESSOR VENKI RAMAKRISHNAN**  
NOBEL LAUREATE AND PRESIDENT OF THE  
ROYAL SOCIETY

We have been stuck on this one problem – how do proteins fold up – for nearly 50 years. To see DeepMind produce a solution for this, having worked personally on this problem for so long and after so many stops and starts, wondering if we'd ever get there, is a very special moment.

**PROFESSOR JOHN MOULT**  
CO-FOUNDER AND CHAIR OF CASP, UNIVERSITY OF  
MARYLAND

In partnership with EMBL-EBI, it has been launched the AlphaFold Protein Structure Database.



**SAPIENZA**  
UNIVERSITÀ DI ROMA

# ML and disease classification: brain tumors

## Problem:

- How many subtypes of brain tumors exist?
- How to identify (i.e. classify and Dx) them?
- Tumors are rare and Dx is not standardized

## Relevance:

- Completely changes prognosis and treatment
- A process combining an advanced imaging technology and artificial intelligence (AI) can accurately diagnose brain tumors in fewer than 3 minutes during surgery.

## Available data:

- Tumor methylation profiles  
(Epigenetic marks of active and inactive DNA)  
(What are the tumor cells doing? Which genes are active?)



<https://www.nytimes.com/2020/01/06/health/artificial-intelligence-brain-cancer.html>



# CNNs and digital pathology: AI-assisted Dx

## Problem:

- How to reliably identify skin cancer?
- Clinicians have different levels of experience

## Relevance:

- Skin lesions extremely frequent (most common cancer!)
- Early detection dramatically changes treatment (i.e simple surgical excision)
- Low-cost, widespread access (e.g. smartphones!)



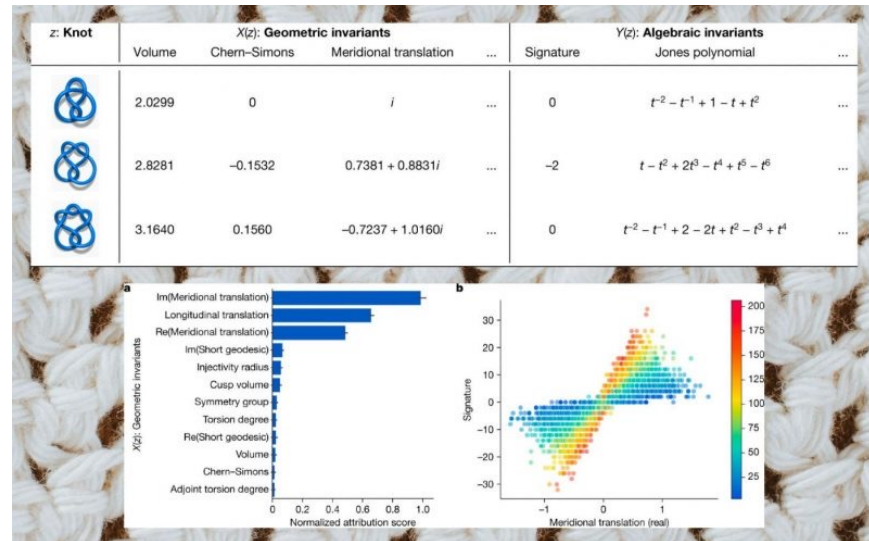
- CNN end-to-end training  
Google's Inception CNN  
130'000 images for training  
In: pixels and labels (biopsy results, 2'000 diseases)  
Out: Benign vs Malignant
- Compared to 21 dermatologists, CNN outperformed humans!

<https://mashable.com/article/ibm-research-melanoma-testing-with-phone-camera>

# Math Conjectures (Deep Mind 2021)

- Use machine learning to make significant new discoveries in pure mathematics
- Conjecture a new connection between the algebraic and geometric structure of knots, and a candidate algorithm predicted by the combinatorial invariance conjecture for symmetric groups
- A supervised learning model was able to detect the existence of a pattern between a large set of geometric invariants and the signature  $\sigma(K)$  of a knot that was not previously known to be related to the hyperbolic geometry.

[Nature](#) volume 600, pages70–74 (2021)



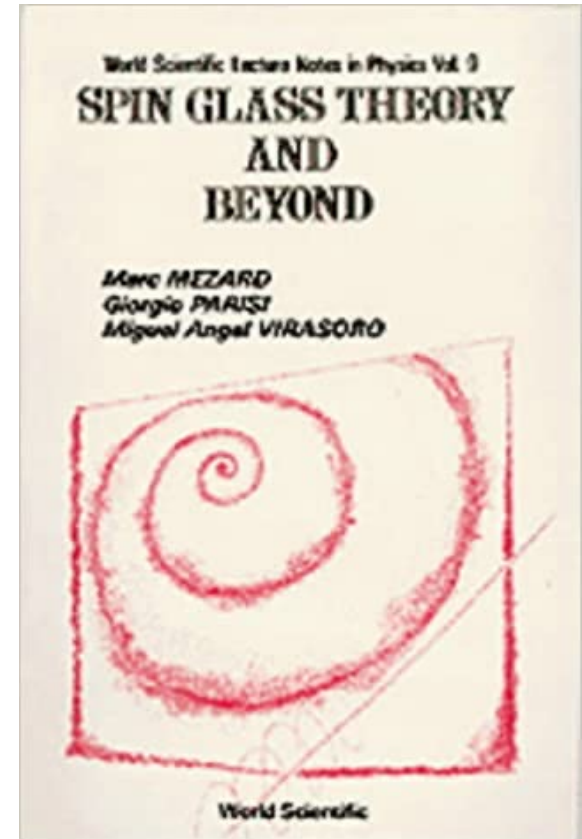
- Neural networks pre-trained on text and fine tuned on code perfectly solves university-level problems.

Drori et al (PNAS 2022)



# Spin Glass Theory

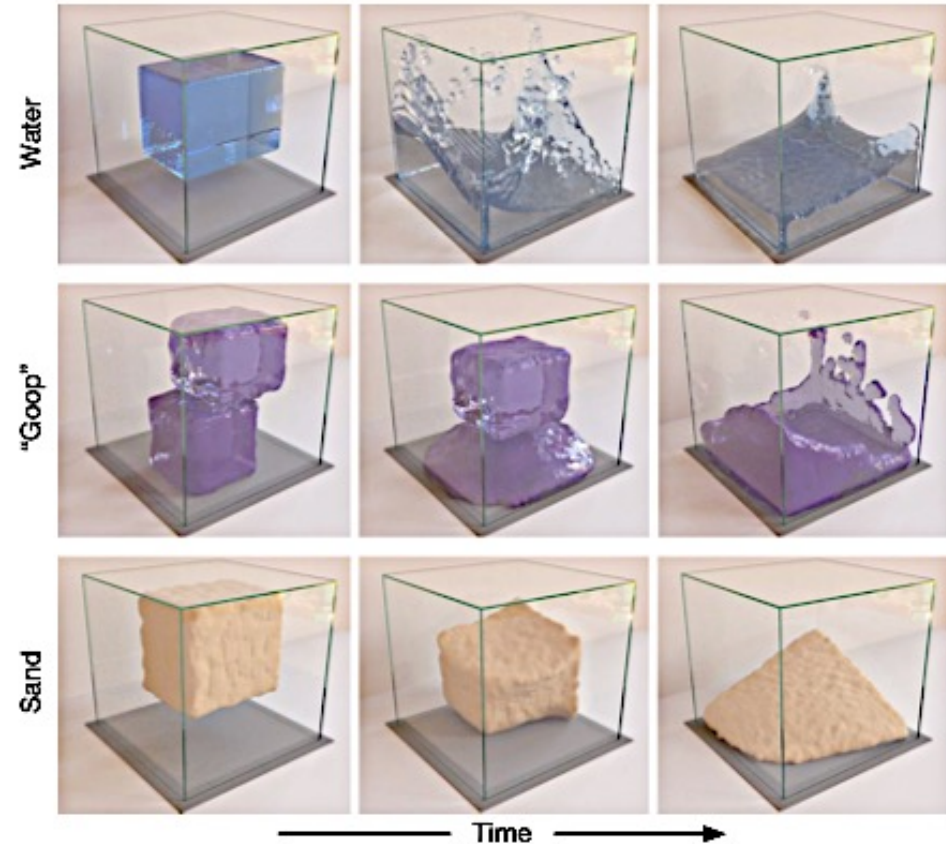
- Finding the ground states of spin glasses is an essential computationally hard problem for the understanding of the nature of disordered magnetic and other physical systems. DIRAC Fan et al [2021] uses Deep Neural Networks trained purely on small-scale spin glass instances and then applied to arbitrarily large ones.
- Use methods developed in statistical physics of glassy systems to analyze numerically the training dynamics of deep neural networks. Baity-Jesi et al. [ICML 2018]



# LEARNING TO SIMULATE COMPLEX PHYSICS WITH GRAPH NETWORKS

[arXiv:2002.09405](https://arxiv.org/abs/2002.09405)  
[\[cs.LG\]](#)

- **Graph Network-based Simulators** that can learn to simulate challenging physical domains: fluids, rigid solids, deformable materials interacting with one another.
- State of a physical system represented with particles, expressed as nodes in a graph.
- Computes dynamics via learned message-passing.
- Shown that it can generalize from single-timestep predictions with thousands of particles during training, to different initial conditions, thousands of timesteps, and at least an order of magnitude more particles at test time.



# Forecast Earthquakes on Tectonic Faults

**Machine Learning Predicts Laboratory Earthquakes** GRL 2017  
Bertrand Rouet-Leduc<sup>1,2</sup>, Claudia Hulbert<sup>1</sup>, Nicholas Lubbers<sup>1,3</sup>, Kipton Barros<sup>1</sup>, Colin J. Humphreys<sup>2</sup>, and Paul A. Johnson<sup>4</sup>

<sup>1</sup>Theoretical Division and CNLS, Los Alamos National Laboratory, Los Alamos, NM, USA, <sup>2</sup>Department of Materials Science and Metallurgy, University of Cambridge, Cambridge, UK, <sup>3</sup>Department of Physics, Boston University, Boston, MA, USA, <sup>4</sup>Geophysics Group, Los Alamos National Laboratory, Los Alamos, NM, USA

**Estimating Fault Friction From Seismic Signals in the Laboratory** GRL 2018  
Bertrand Rouet-Leduc<sup>1</sup>, Claudia Hulbert<sup>1</sup>, David C. Bolton<sup>2</sup>, Christopher X. Ren<sup>3</sup>, Jacques Riviere<sup>2,4</sup>, Chris Marone<sup>2</sup>, Robert A. Guyer<sup>1</sup>, and Paul A. Johnson<sup>1</sup>

**Earthquake Catalog-Based Machine Learning Identification of Laboratory Fault States and the Effects of Magnitude of Completeness** GRL 2018  
Nicholas Lubbers<sup>1</sup>, David C. Bolton<sup>2</sup>, Jamaludin Mohd-Yusof<sup>3</sup>, Chris Marone<sup>2</sup>, Kipton Barros<sup>1</sup>, and Paul A. Johnson<sup>1</sup>

**Lab Earthquake Prediction**

**Geophysical Research Letters** GRL 2019  
**RESEARCH LETTER** 10.1029/2019GL081251  
**Machine Learning Can Predict the Timing and Size of Analog Earthquakes**  
F. Corbi<sup>1</sup>, L. Sandri<sup>2</sup>, J. Bedford<sup>3</sup>, F. Funicello<sup>4</sup>, S. Brizzi<sup>1,4</sup>, M. Rosenau<sup>1</sup>, and S. Lallemand<sup>5</sup>

**Geophysical Research Letters** GRL 2020  
**RESEARCH LETTER** 10.1029/2020GL090255  
**Deformation Precursors to Catastrophic Failure in Rocks**  
J. A. McKeck<sup>1</sup>, J. M. Aiken<sup>2</sup>, J. Mathiesen<sup>3</sup>, Y. Ben-Zion<sup>4</sup>, and F. Renard<sup>1,4</sup>

**Precursors:**  
Consistent changes prior to failure  
(elastic wave speed, amplitude, fault  
dilation, fluid pressure/chemistry, etc.)

<https://www.lanl.gov/discover/news-release-archive/2018/December/1217-machine-learning.php>



SAPIENZA  
UNIVERSITÀ DI ROMA



# In Codice Ratio – ML based transcription of the 85 km shelves of the Vatican Secret Archive

**Automatically transcribe** the contents of the manuscripts. Follow a novel approach, based on character segmentation.

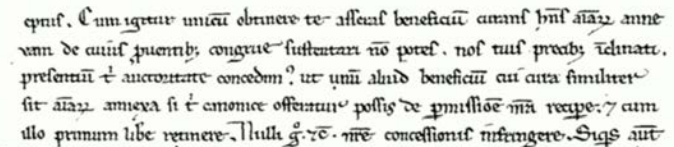
Govern imprecise character segmentation by considering that correct segments are those that give rise to a sequence of characters that more likely compose a Latin word.

Designed a principled solution that relies on **convolutional neural networks** and **statistical language models**.

High School students will be involved in collecting and labelling training data through the Program of digital of Fondazione "I Lincei per La Scuola" e Ministero dell'Istruzione

## Transcription

Input:  
image



-eptus. Cum igitur unicu(m) obtinere te asseras beneficiu(m) curans h(abe)ns a(n)i(m)ar(um) annexam de cuius p(ro)ventib(us) congrue sustentari no(n) potes nos tuis precib(us) i(n)clinati presentiu(m) t(ibi)i auctoritate concedim(us) ut unu(m) aliud beneficiu(m) cui cura similiter sit a(n)i(m)ar(um) annexa si t(ibi)i canonice offeratur possis de p(er)missio(n)e n(ost)ra recip(er)e et cum illo primum lib(er)e retinere. Nulli(er)go et c(etera) n(ost)re concessionis infringere. Siq(ui)s aut(em)

Output:  
text



<http://www.inf.uniroma3.it/db/icr/>



SAPIENZA  
UNIVERSITÀ DI ROMA

**“I don’t know how long it’s going to be before  
the name of our field is changed from  
computer science to machine learning...”**

**Donald Knuth on Machine Learning and the  
Meaning of Life  
(2021)**