



**Barcelona  
Supercomputing  
Center**  
Centro Nacional de Supercomputación



# Big data e inteligencia artificial y HPC en salud y biomedicina

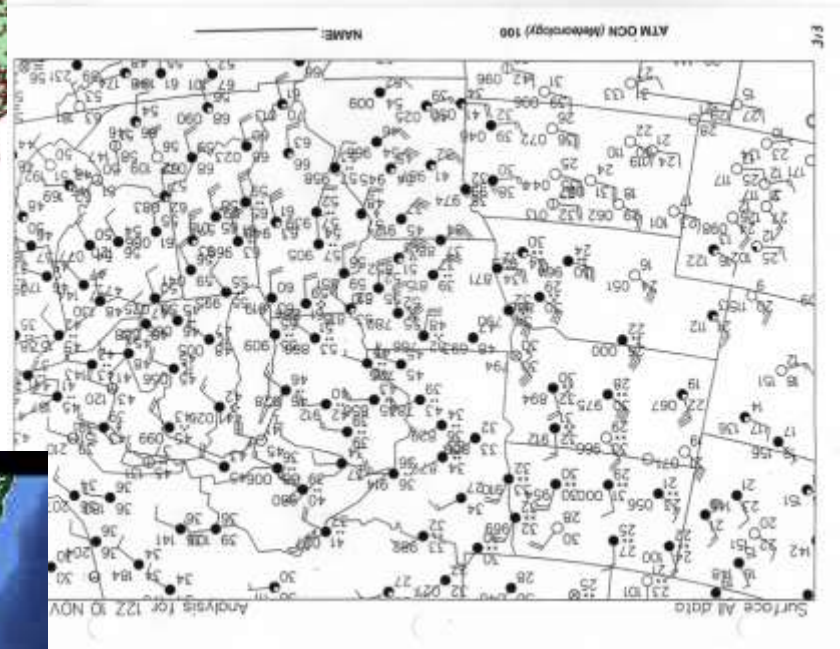
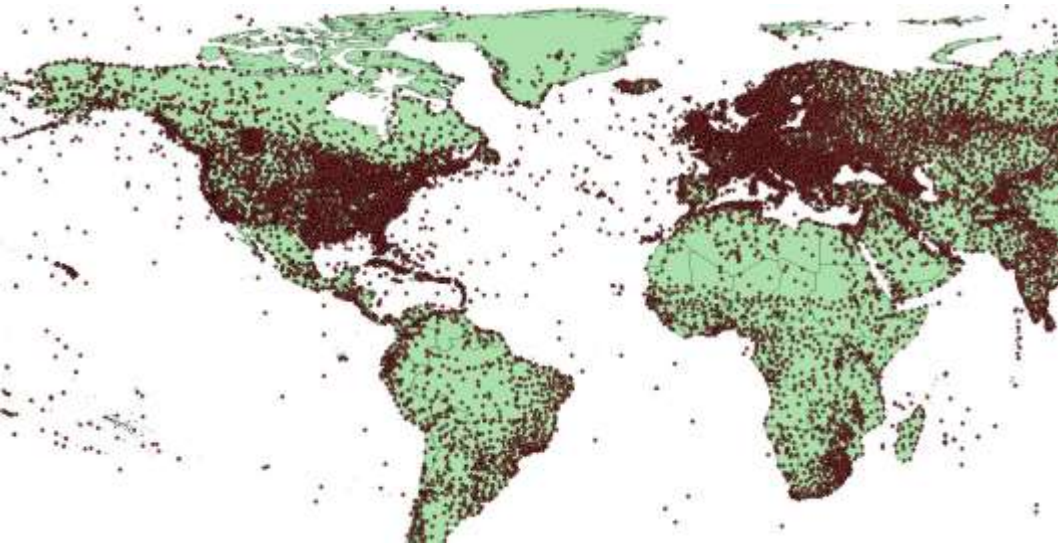
**Alfonso Valencia. Ph.D.**  
***ICREA Professor***  
***Director Life Sciences Dept. BSC***  
***Director Spanish Bioinformatics Institute***  
***INB-ISCIII ELIXIR-ES***

***Plan TL InfoDay***

21 October 19

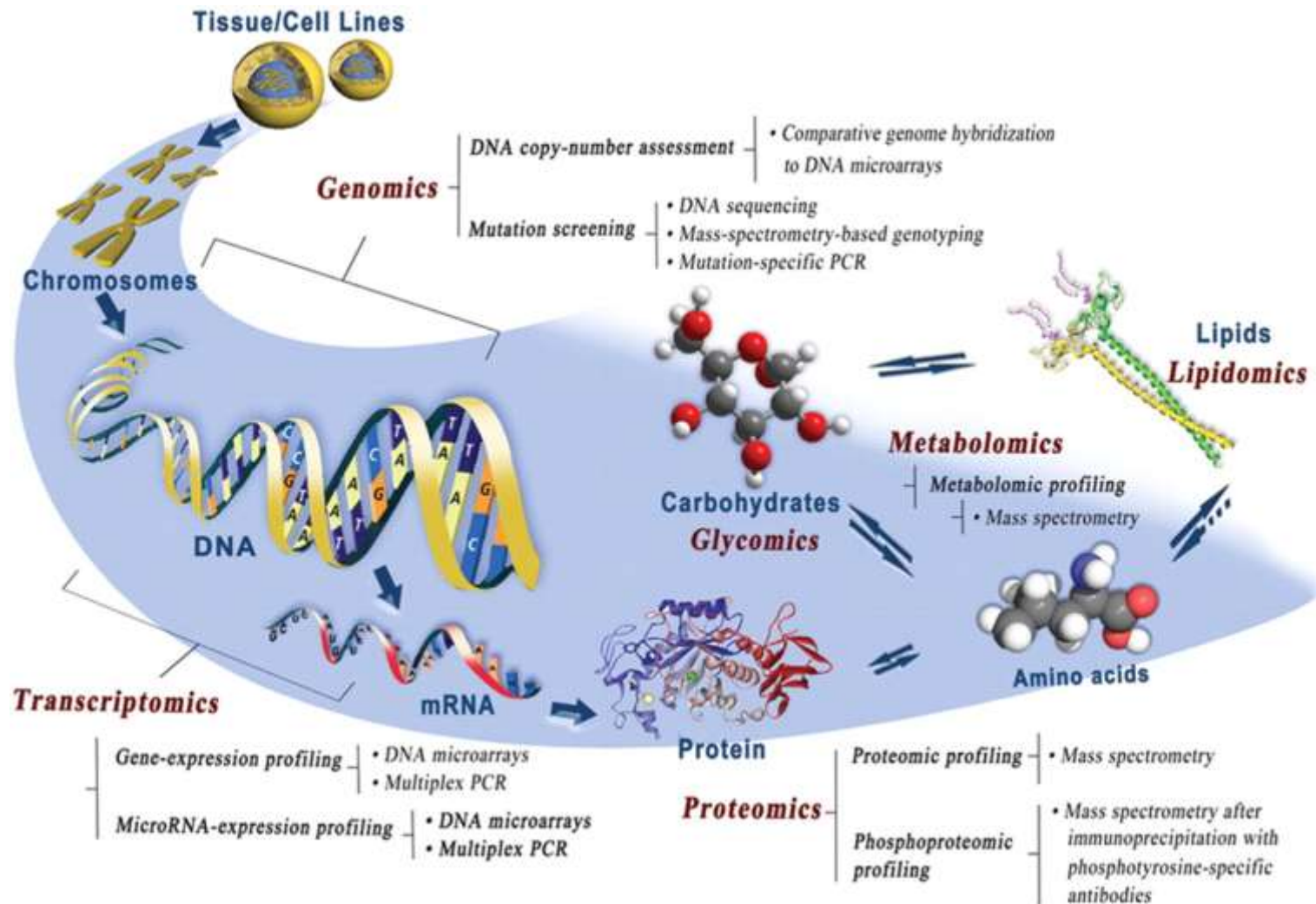
- **DATOS**

# Datos sobre el tiempo

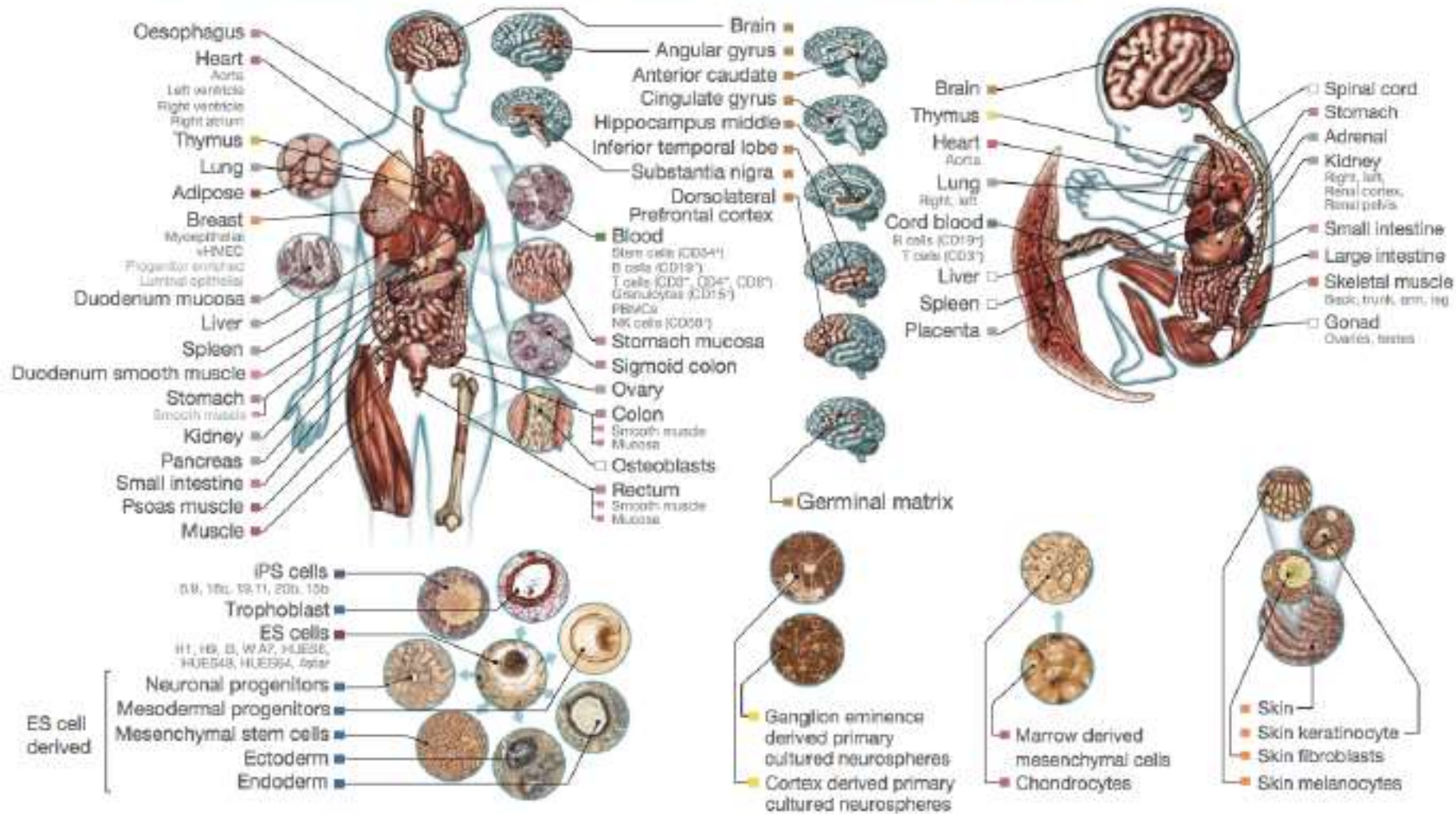




# Datos sobre Biología: Genómica, Transcriptómica, Epigenómica, Proteómica, Metabolómica, Lipidómica, ...



# Tissues and cell types profiled in the Roadmap Epigenomics



**Figure 1 | Tissues and cell types profiled in the Roadmap Epigenomics Consortium.** Primary tissues and cell types representative of all major lineages

immune lineages, ES cells and iPS cells, and differentiated lineages derived from ES cells. Box colours match groups shown in Fig. 2b. Epigenome identifiers:



# We will have 40 ZettaBytes of Healthdata by 2020

## What are we really saying here?

### Volume SCALE OF DATA

40 ZETTAYTES (40,000,000,000,000,000,000 BYTES) of data will be created by 2020, an increase of 3,000 times from 2005.

It's estimated that 2.5 QUINTILLION BYTES (2,500,000,000,000,000,000 BYTES) of data are created each day.

Most companies in the U.S. have at least 100 TERABYTES (100,000,000,000,000 BYTES) of data stored.

6 BILLION PEOPLE have cell phones.

WORLD POPULATION 7 BILLION.

## The FOUR V's of Big Data

From health information, downloads to web, mining and social networks, data is becoming more and more varied to make the technology and services that the world uses on a daily basis work better. It's big data, and it's not the massive amounts of data to do!

As a result in the world, 100% data scientists speak big data into four dimensions: Volume, Variety, Velocity and Veracity.

Expanding on the industry and expanding the data professionals' information from multiple personal and professional sources such as Facebook, social media, wireless content, sensors and mobile devices. Companies that leverage data to make their products and services to better meet customer needs, optimize operations and relationships, and find new sources of revenue.

By 2015, 4.4 MILLION IT JOBS will be created globally to support big data, with 1.9 million in the United States.

### Variety DIFFERENT FORMS OF DATA

As of 2013, the global size of data in healthcare was estimated to be 150 EXABYTES (150,000,000,000,000,000 BYTES).

30 BILLION PIECES OF CONTENT are shared on Facebook every month.

4 BILLION+ HOURS OF VIDEO are watched on YouTube each month.

400 MILLION TWEETS are sent per day by about 200 million monthly active users.

By 2014, it's anticipated there will be 420 MILLION WEARABLE, WIRELESS HEALTH MONITORS.

### Velocity ANALYSIS OF STREAMING DATA

The New York Stock Exchange captures 1 TB OF TRADE INFORMATION during each trading session.

Modern cars have close to 100 SENSORS that monitor items such as fuel level and tire pressure.

By 2016, it is projected there will be 18.9 BILLION NETWORK CONNECTIONS - almost 2.5 connections per person on earth.

### Veracity UNCERTAINTY OF DATA

1 IN 2 BUSINESS LEADERS don't trust the information they use to make decisions.

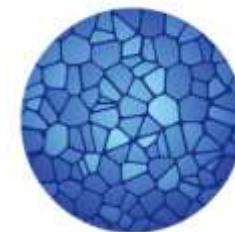
27% OF RESPONDENTS in one survey were unsure of how much of their data was inaccurate.

Four data quality costs the US economy around \$3.1 TRILLION A YEAR.

Source: McKinsey Global Institute, "Smarter, Faster, Stronger: Emerging IT and Analytics Trends, 2013-2015", 2013

Diez mil millones de cell phones





Single-cell multi-omics



Machine learning



Personalized organoid disease  
models

About 500 scientists, clinicians, patients, and stakeholders from 20 European countries and beyond have gathered in Berlin at the LifeTime Opening Conference. The pan-European initiative aims to revolutionize healthcare. It applies breakthrough technologies to the progression of human diseases and intends to find and implement new methods for personalized prevention, early diagnosis and treatment.

## Unravelling tumour heterogeneity by single-cell profiling of circulating tumour cells

Laura Keller & Klaus Pantel

Nature Reviews Cancer (2019) | Download Citation

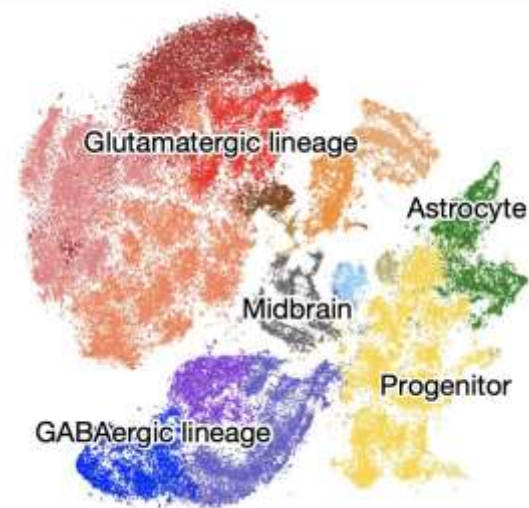
## ARTICLE

NATURE | [www.nature.com/nature](http://www.nature.com/nature)

<https://doi.org/10.1038/s41586-019-1434-6>

# Resolving medulloblastoma cellular architecture by single-cell genomics

Received: 13 September 2018; Accepted: 21 June 2019;  
Published online: 24 July 2019



# ELIXIR:

European infrastructure for biological information



## Spanish National Bioinformatics Institute (INB) Spanish Node of ELIXIR (ELIXIR-ES)



**TransBioNet**  
network of Bioinformaticians in  
Research Institutes of Spanish  
Hospitals  
(INB hosted)





- **DATOS**
- **INTELIGENCIA ARTIFICIAL**

# High-performance medicine: the convergence of human and artificial intelligence

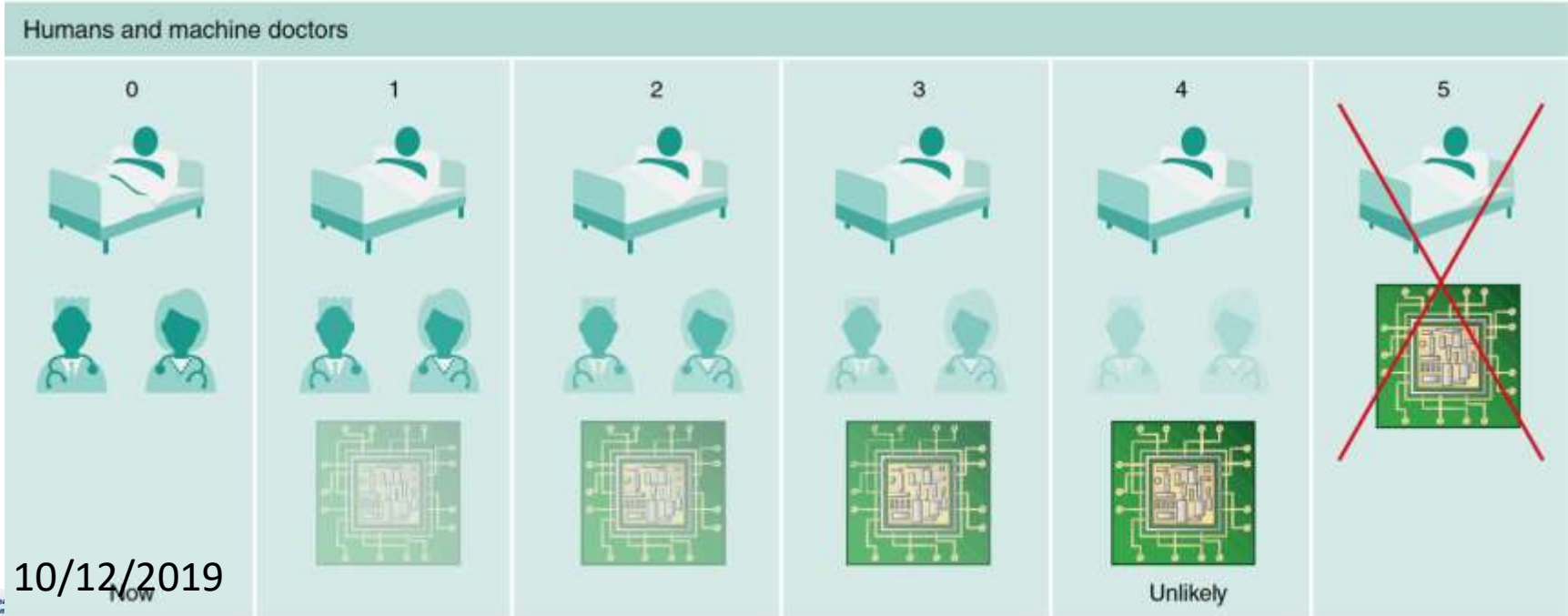


Eric J. Topol

NATURE MEDICINE | VOL 25 | JANUARY 2019 | 44-56 |

Human driver monitors environment		
0 No automation	1 Driver assistance	2 Partial automation
The absence of any assistive features such as adaptive cruise control.	Systems that help drivers maintain speed or stay in lane but leave the driver in control.	The combination of automatic speed and steering control—for example, cruise control and lane keeping.

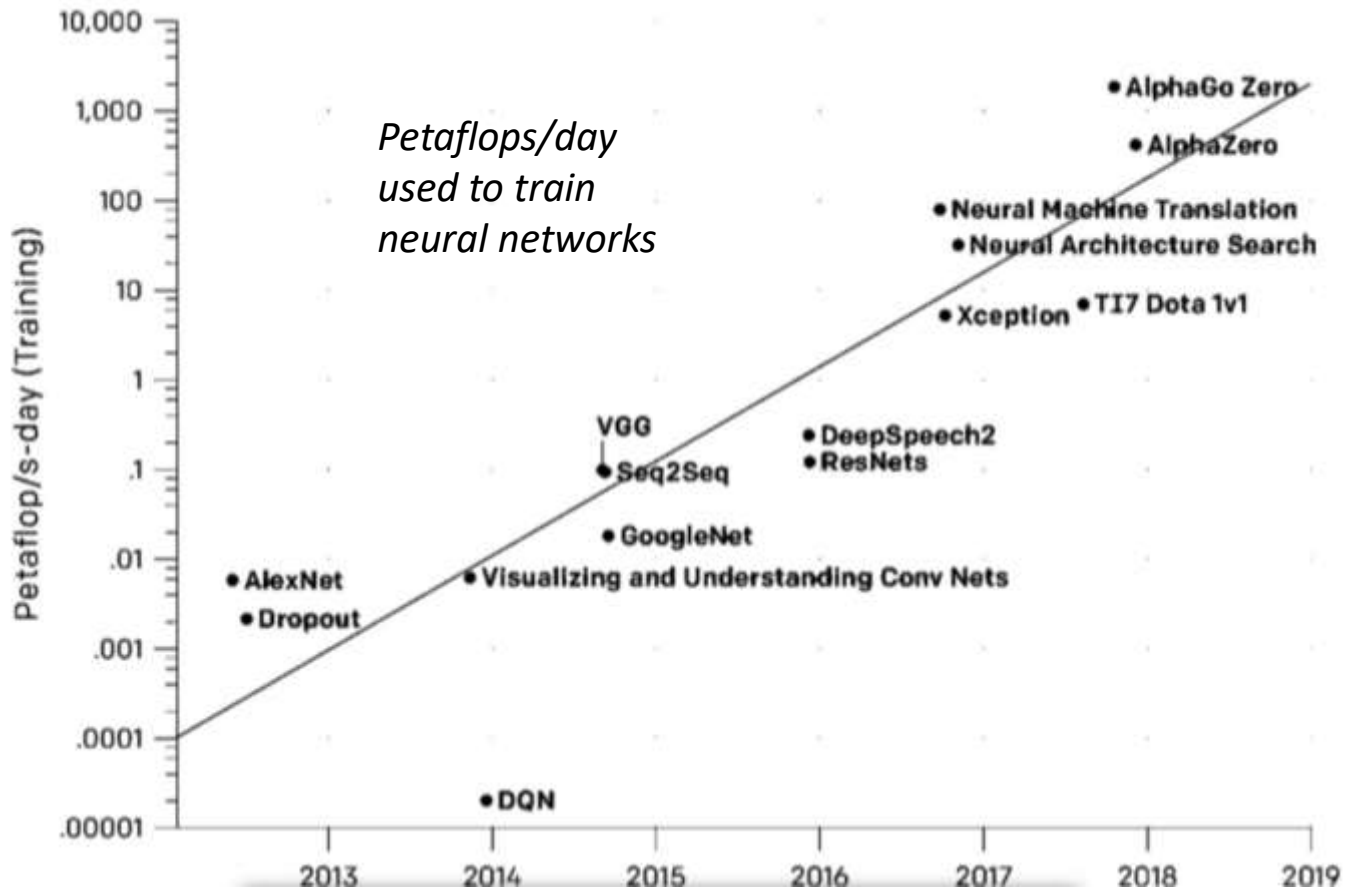
System monitors environment		
3 Conditional automation	4 High automation	5 Full automation
Automated systems that drive and monitor the environment but rely on a human driver for backup.	Automated systems that do everything—no human backup required—but only in limited circumstances.	The true electronic chauffeur: retains full vehicle control, needs no human backup, and drives in all conditions.



10/12/2019

# Advances in AI and HPC go hand by hand

Since GPUs were first used in AI (2012), **computing power** available to generate AI models has increased exponentially – and improvements in computing power has been key for **AI progress**.



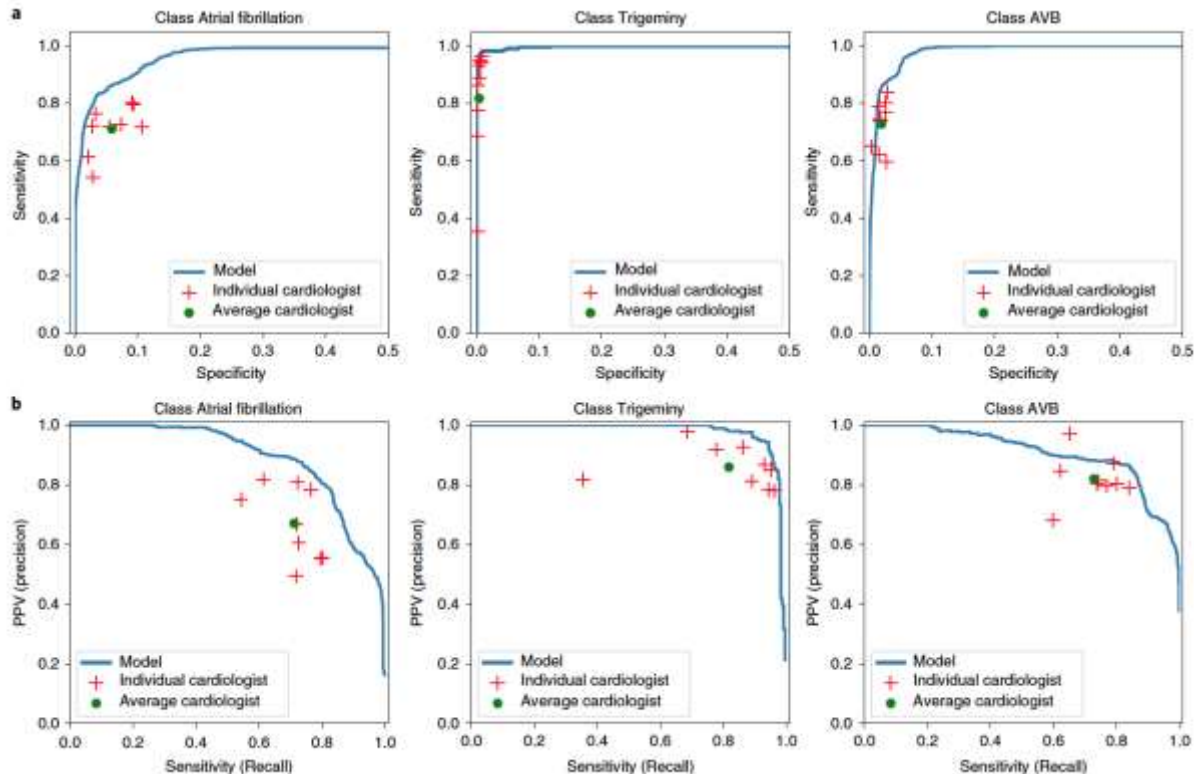


# Cardiologist-level arrhythmia detection and classification in ambulatory electrocardiograms using a deep neural network



Awni Y. Hannun<sup>1,6\*</sup>, Pranav Rajpurkar<sup>1,6</sup>, Masoumeh Haghpanahi<sup>2,6</sup>, Geoffrey H. Tison<sup>3,6</sup>,  
Codie Bourm<sup>2</sup>, Mintu P. Turakhia<sup>4,5</sup> and Andrew Y. Ng<sup>1</sup>

NATURE MEDICINE | VOL 25 | JANUARY 2019 | 65–69 |



10/12/2019

# Dermatologist-level classification of skin cancer with deep neural networks

Andre Esteve, Brett Kuprel, Roberto A. Novoa, Justin Ko, Susan M. Swetter, Helen M. Blau & Sebastian Thrun

*Nature* volume 542, pages 115–118 (02 February 2017) | Download Citation

A Corrigendum to this article was published on 28 June 2017

## AbstractAbstract

Skin cancer, the most common human malignancy<sup>1,2,3</sup>, is primarily diagnosed visually, beginning with an initial clinical screening and followed potentially by dermoscopic analysis, a biopsy and histopathological examination. Automated classification of skin lesions using images is a challenging task owing to the fine-grained variability in the appearance of skin lesions. Deep convolutional neural networks (CNNs)<sup>4,5</sup> show potential for general and highly variable tasks across many fine-grained object categories<sup>6,7,8,9,10,11</sup>. Here we demonstrate classification of skin lesions using a single CNN, trained end-to-end from images directly, using only pixels and disease labels as inputs. We train a CNN using a dataset of 129,450 clinical images—two orders of magnitude larger than previous datasets<sup>12</sup>—consisting of 2,032 different diseases. We test its performance against 21 board-certified dermatologists on biopsy-proven clinical images with two critical binary classification use cases: keratinocyte carcinomas versus benign seborrheic keratoses; and malignant melanomas versus benign nevi. The first case represents the identification of the most common cancers, the second represents the identification of the deadliest skin cancer. The CNN achieves performance on par with all tested experts across both tasks, demonstrating an

## Access options

```
<html><head></head><body></body></html>
```

THE LANCET  
Oncology

Log in 🔍 ☰

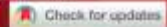
ARTICLES | ONLINE FIRST

### Comparison of the accuracy of human readers versus machine-learning algorithms for pigmented skin lesion classification: an open, web-based, international, diagnostic study

Philipp Tschandl, PhD + Noel Codella, PhD + Bengü Nisa Akay, MD + Prof Giuseppe Argenziano, PhD +

Ralph P Braun, MD + Prof Horacio Cabo, MD + et al. et al. Show all authors + Show all authors

Published: June 11, 2019 + DOI: [https://doi.org/10.1016/S1470-2045\(19\)30333-X](https://doi.org/10.1016/S1470-2045(19)30333-X) +



PlumX Metrics

### Summary

#### Background

Whether machine-learning algorithms can diagnose all pigmented skin lesions as accurately as human experts is unclear. The aim of this study was to compare the diagnostic accuracy of state-of-the-art machine-learning algorithms with human readers for all clinically relevant types of benign and malignant pigmented skin lesions.

#### Methods

For this open, web-based, international, diagnostic study, human readers were asked to diagnose dermoscopic images selected randomly in 30-image batches from a test set of 1511 images. The diagnoses from human readers were compared with those of 139 algorithms created by 77 machine-learning labs, who participated in the International Skin Imaging Collaboration 2018 challenge and received a training set of 10 015 images in advance. The ground truth of each lesion fell into one of

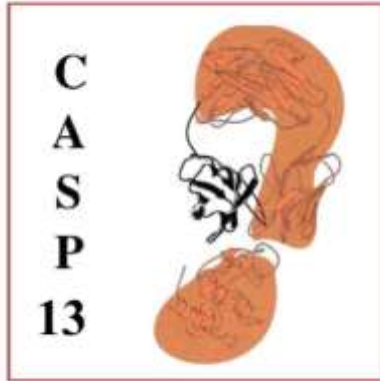
10/12/2019

# BSC works on Medical Imaging

- Detecting retina pathologies
  - Trained models competitive with ophthalmologists
  - With Lenovo & Hospital Vall Hebron
- Learning from liver conditions
  - Learning about rare diseases
  - With Hospital Clinic
- Predicting and guiding in-vitro success
  - Finding the best embryo ASAP
  - With Hospital Clinic
- Supporting medical doctors on Rx review
  - Aid for Dr. in rural areas
  - With Asepeyo and ICS







Thirteenth meeting  
Riviera Maya, Mexico  
DECEMBER 1-4, 2018

Predicting protein structure from the sequence is one of the fundamental problems in molecular biology.

It is the key to the prediction of the consequences of mutations in human diseases and to drug design

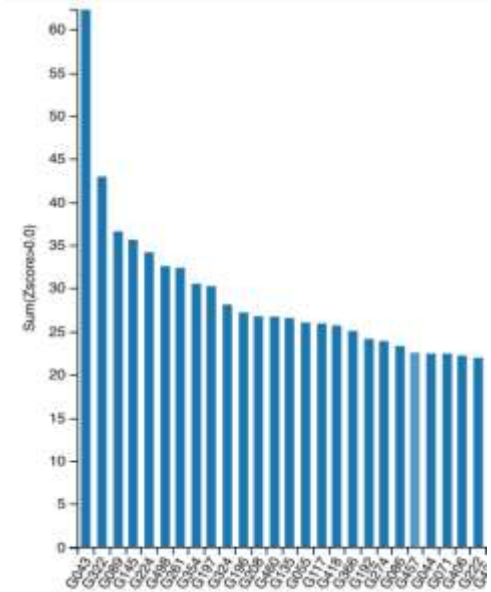


World UK **Science** Cities Global development Football Tech Business Environment More

Science

# Google's DeepMind predicts 3D shapes of proteins

AI program's understanding of proteins could usher in new era of medical progress



On its first foray into the competition, AlphaFold topped a table of 98 entrants, predicting the most accurate structure for 25 out of 43 proteins, compared with three out of 43 for the second placed team in the same category.



Nico Callewaert @NicoCallewaert · 11h

Probably my nomination for basic molecular science advance of 2018, need to see a bit more methods details but results in blinded CASP13 test clearly impressive. [deepmind.com/blog/alphafold/](https://deepmind.com/blog/alphafold/)

- **DATOS**
- **INTELIGENCIA ARTIFICIAL**
- **HPC**

# The Evolution of the Research Paradigm

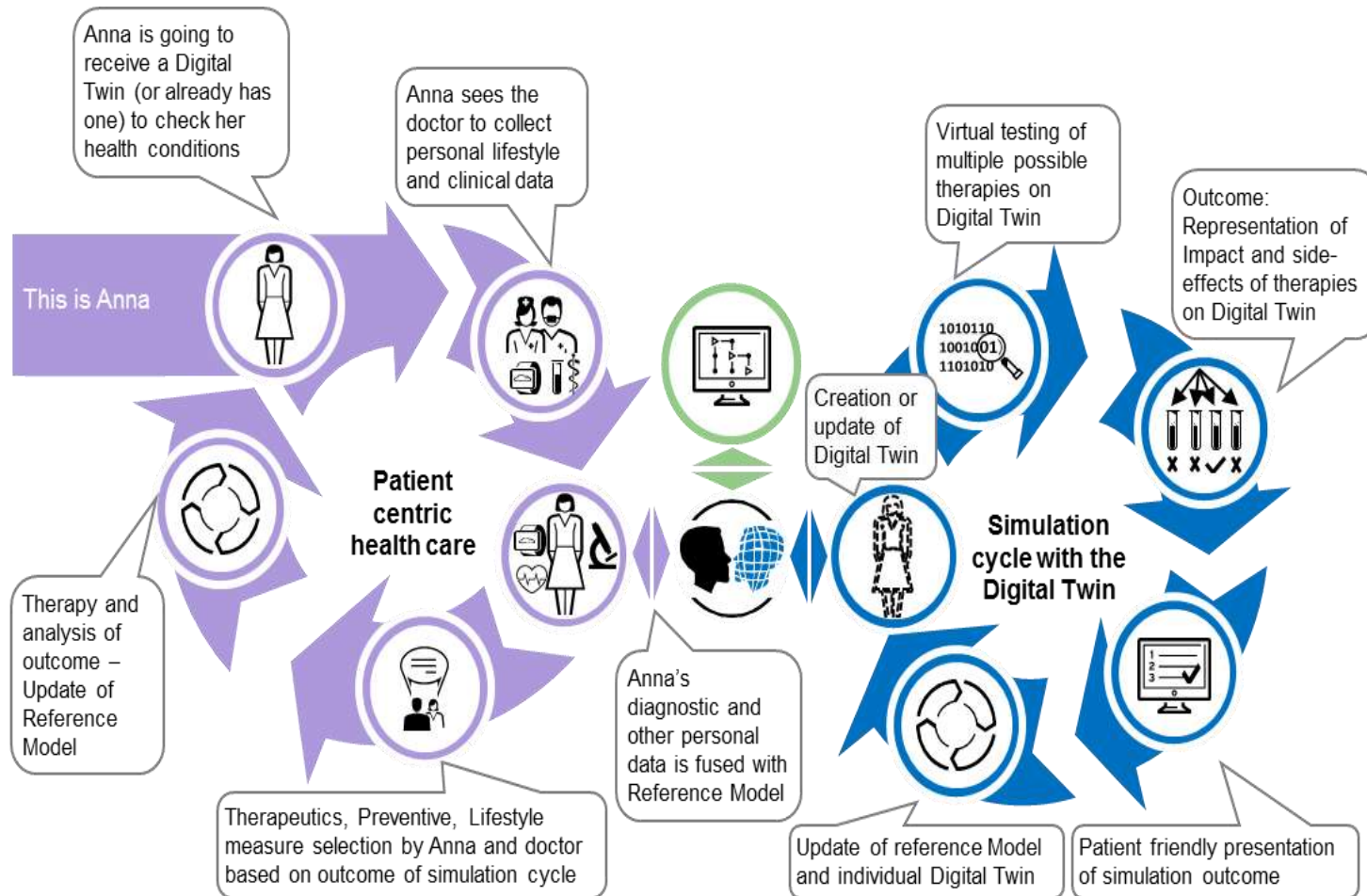


## Numerical Simulation and Big Data Analysis

- Reduce expense
- Avoid suffering
- Help to build knowledge where experiments are impossible or not affordable



# Digital Twin for Future Medicine



**Is this scenario possible?  
When?**

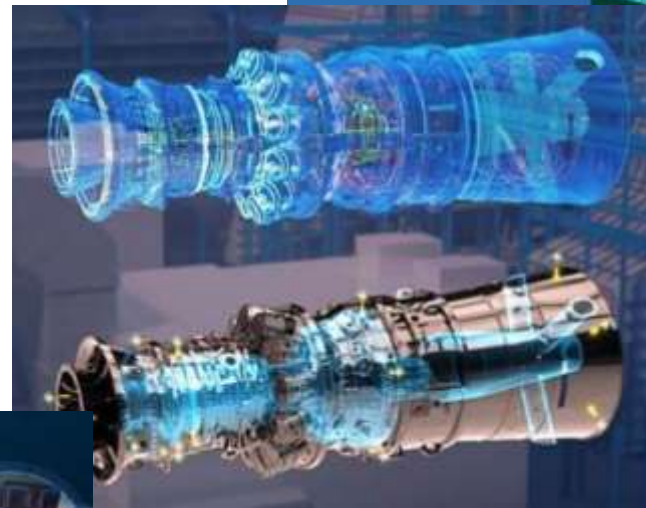
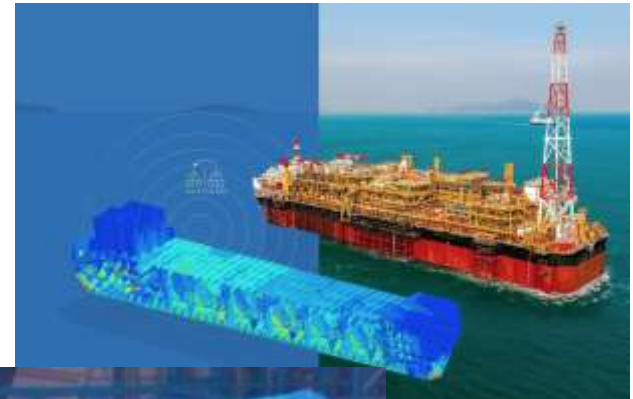
# HORIZON

The EU Research & Innovation Magazine

INDUSTRY SCIENCE IN SOCIETY ICT

## How digital 'twins' are guiding the future of maintenance and manufacturing

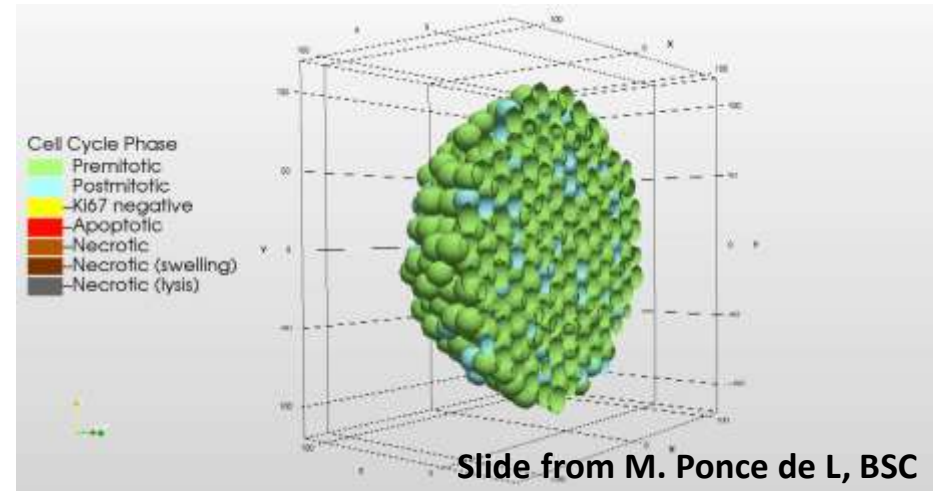
15 November 2019



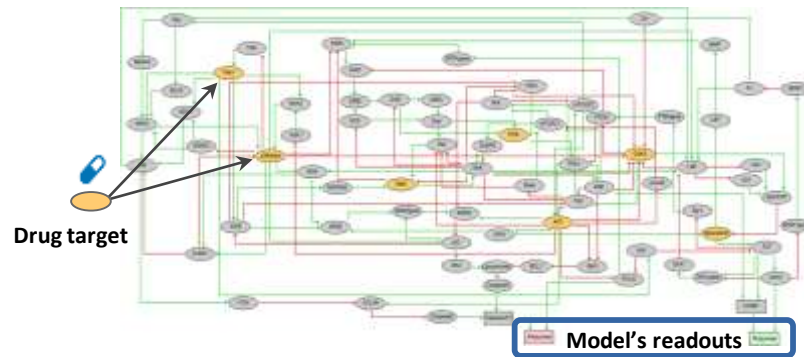
# Simulations of biological systems at different levels



By *Mariano Vazquez, CASE - BSC*



~48 h simulation time, 30 min wall time  
~2500 cells



By *Victor Guallar ICREA & BSC*

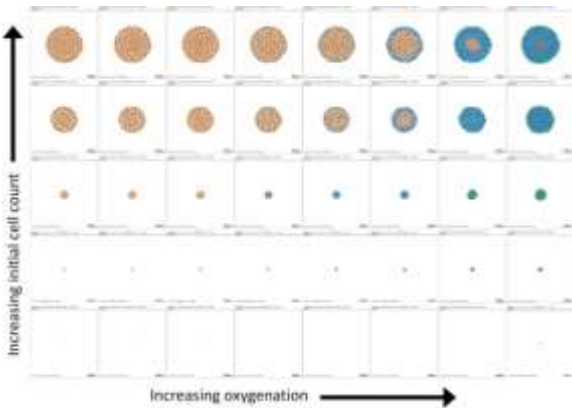
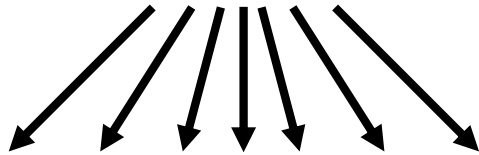


# Large Scale Simulations



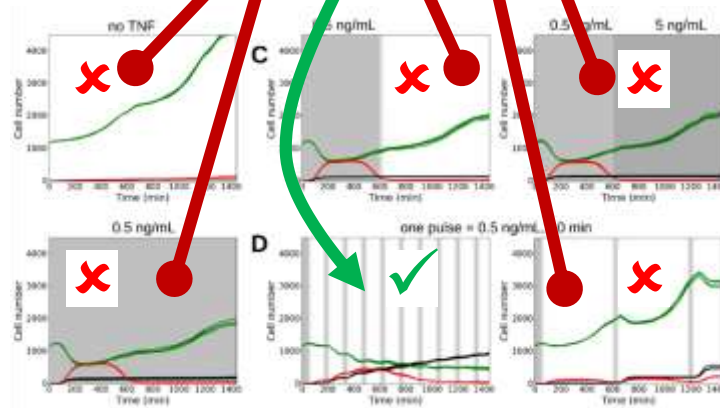
## AI / ML systems

Exploration of the parameter space

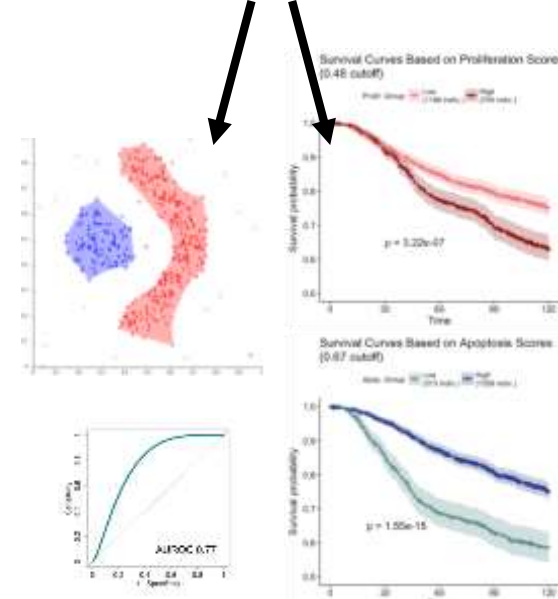


Monitoring and tailoring simulations during execution time

Event Recognition System

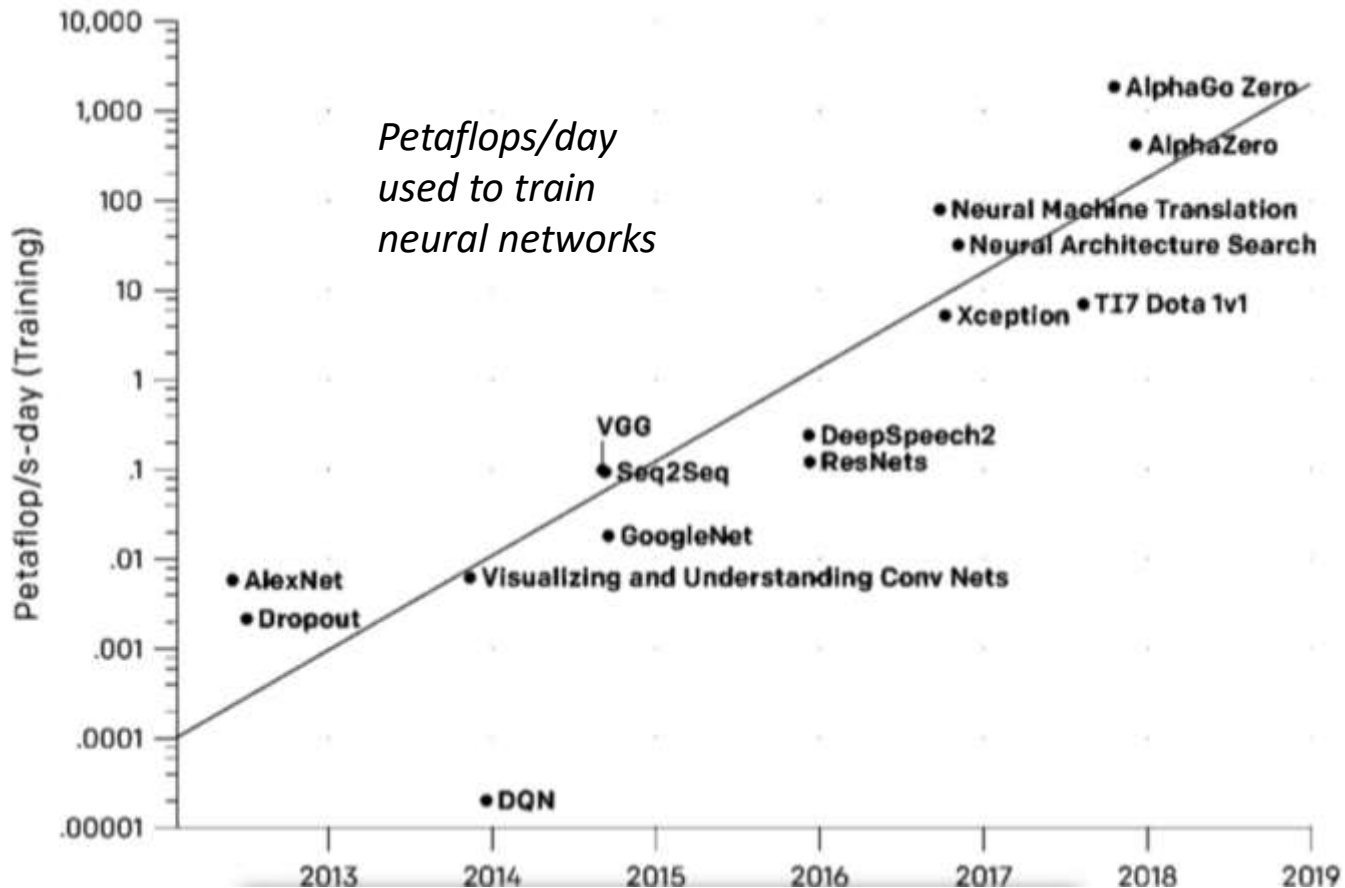


Analysing the results of the simulations



# Advances in AI and HPC go hand by hand

Since GPUs were first used in AI (2012), **computing power** available to generate AI models has increased exponentially – and improvements in computing power has been key for **AI progress**.

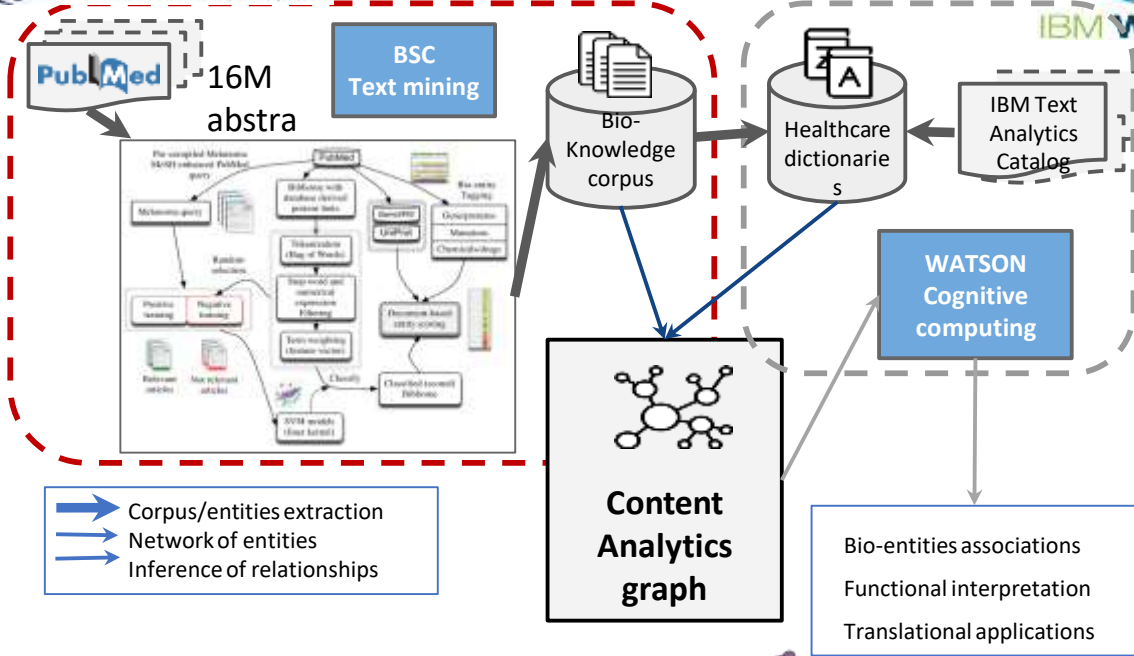


# Text mining & Cognitive computing for melanoma research

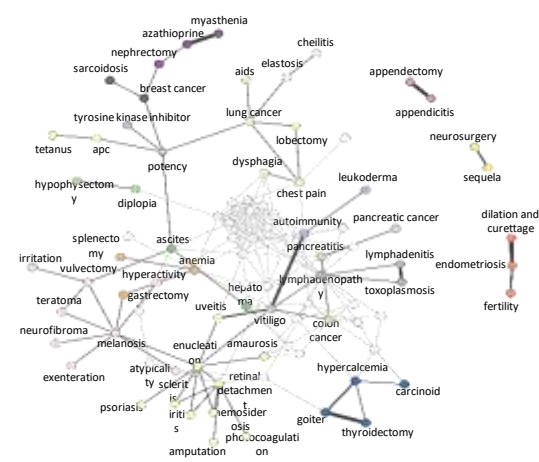
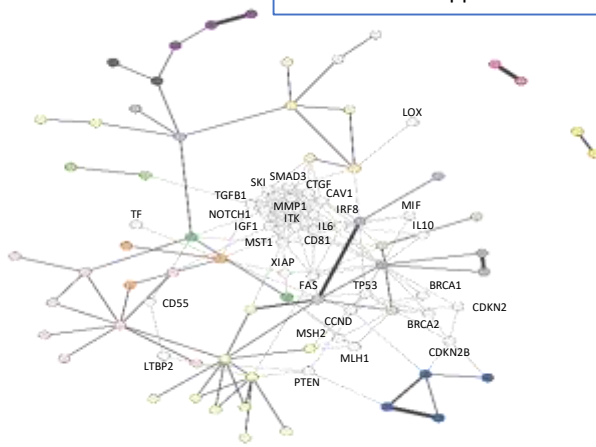


Specialized info

Watson General info



- ➔ Corpus/entities extraction
- ➔ Network of entities
- ➔ Inference of relationships



In collaboration with IBM Spain



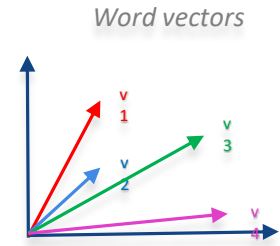
# Word embeddings



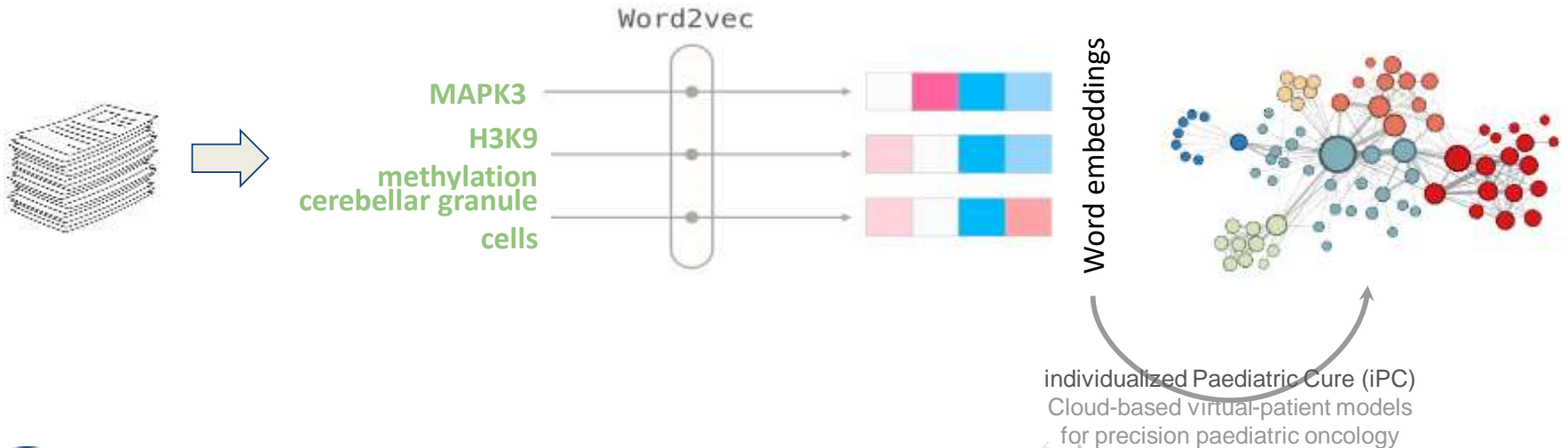
Word2Vec  
fastText

→ |V|

v1	0,78	0,65	0,98
v2	0,23	0,12	0,32
v3	0,90	0,32	0,56
v4	0,08	0,43	0,65
v5	0,77	0,88	0,77



Evaluación intrínseca: cálculo similitud entre términos (sinónimos en SNOMED)  
Evaluación extrínseca: comprobar su utilidad en otras tareas PLN (neuroNER)



# Gender and other biases ...

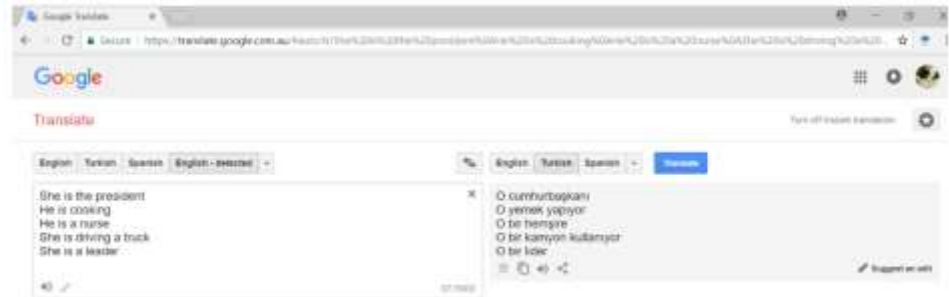
## How AI systems amplify bias



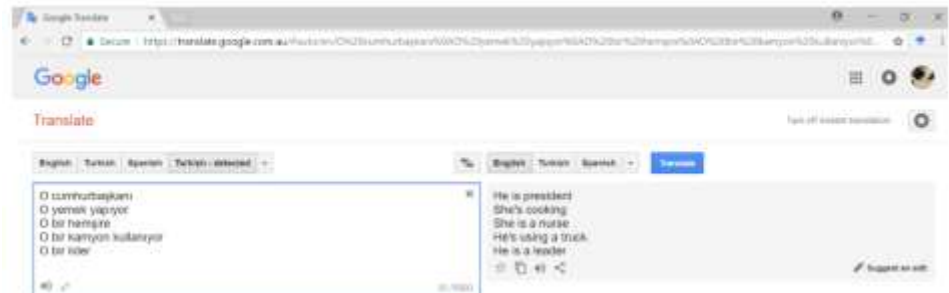
Image recognition systems that use biased machine learning data sets will inadvertently magnify that bias. Researchers are examining ways to reduce the effects.



In this example of gender bias, adapted from a report published by researchers from the University of Virginia and the University of Washington, a visual semantic role labeling system has learned to identify a person cooking as female, even when the image is male.



Google image search



- Research activities in five key interconnected AI scientific areas (Explainable AI, Physical AI, Verifiable AI, Collaborative AI, Integrative AI), which arise from the application of AI in real-world scenarios;
- The creation of a European Ethical Observatory to ensure that European AI projects adhere to high ethical, legal, and socio-economical standards;

## AI4EU ETHICAL OBSERVATORY



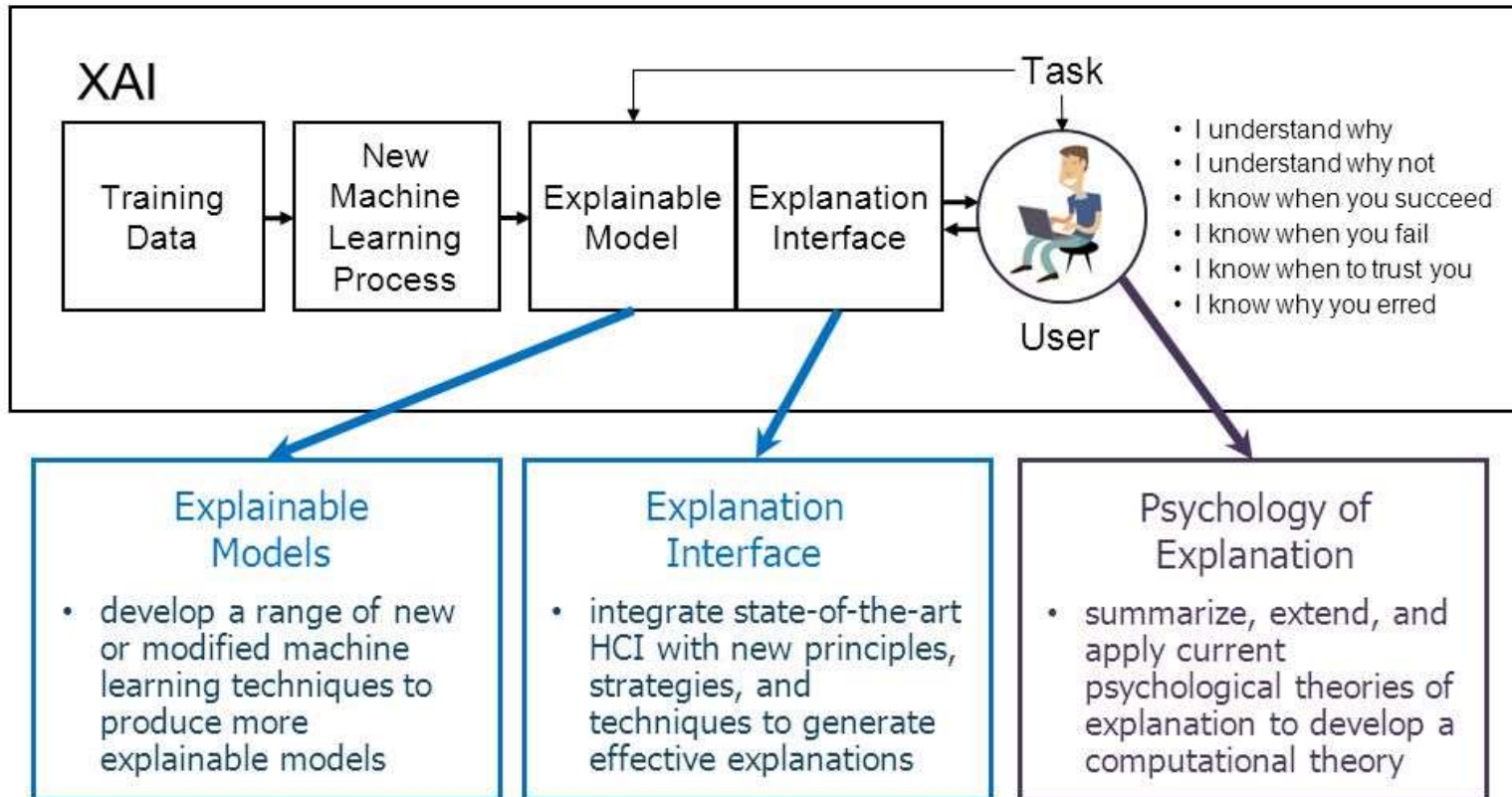
The main goal of the [AI4EU Observatory on Society and Artificial Intelligence \(OSAI\)](#) is to support the distribution and the discussion of knowledge about the Ethical, Legal, Socio-Economic and Cultural Issues of AI (ELSEC-AI) within Europe.



# Explainable Artificial Intelligence



## B. Program Scope – XAI Development Challenges



## El nuevo miedo de Europa se llama China

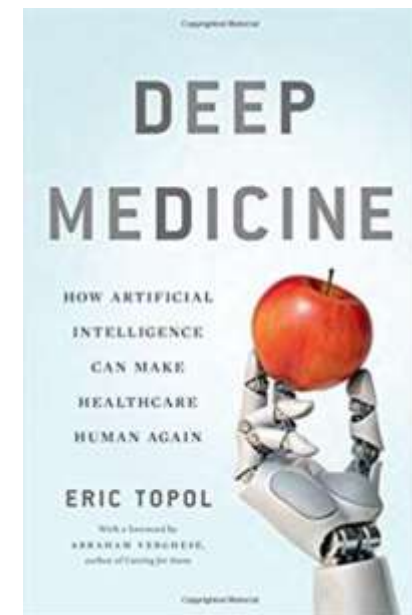
Mientras el Gobierno de Merkel ha decidido invertir 3.000 millones de euros en inteligencia artificial hasta 2025, China invertirá 130.000 millones hasta 2030", señala Müller-Markus. La reducción de

The United Kingdom is also betting big on AI's future and emphasizing healthcare. When the UK government issued four Grand Challenges, one centered on medicine, Theresa May declared, "The development of smart technologies to analyse great quantities of data quickly and with a higher degree of accuracy than is possible by human beings, opens up a whole new field of medical research and gives us a new weapon in our armory in the fight against disease."<sup>66</sup> In 2018, I was commissioned by the UK to work with the National Health Service on planning the future of its healthcare, particularly leading a review on the impact of AI and other medical technologies on its workforce over the next two decades.<sup>67</sup> The opportunity to

*AI and Health Systems*

257

work with leaders of AI, digital medicine, genomics, and robotics, along with ethicists, economists, and educators was an extraordinary experience in the context of a single-payer healthcare system with a palpable will to change and adapt. The full report was issued in 2019, where we project major impacts at every level—the patient, the clinicians, and the health systems throughout the country.



10/12/2019