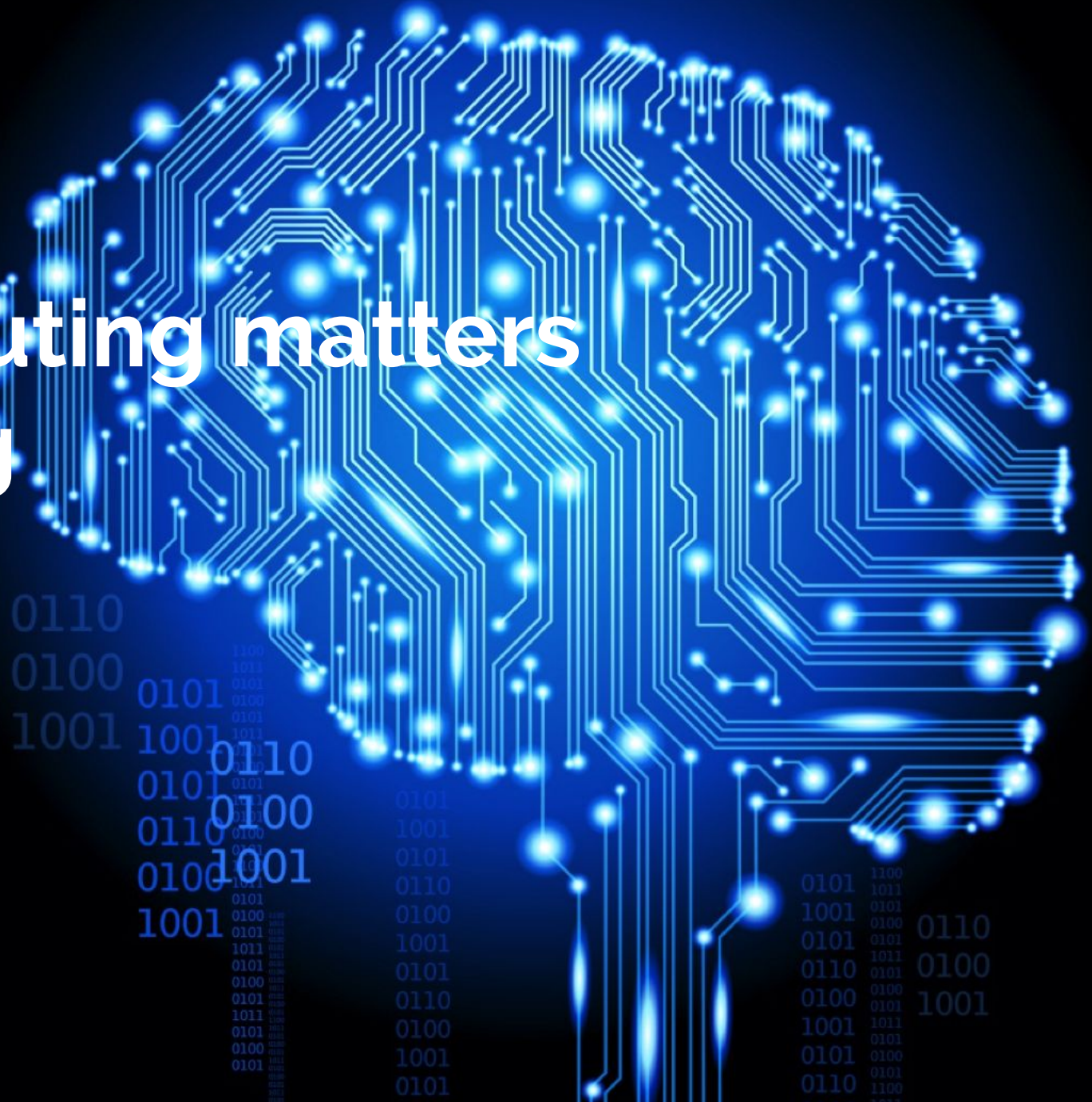# Why Supercomputing matters to Deep Learning

**ESADE – MIBA (FALL 2017)**

JORDI **TORRES** | FRANCESC **SASTRE**

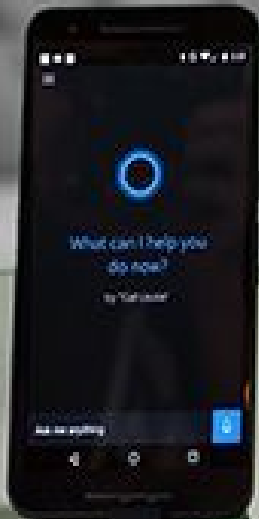Artificial Intelligence is changing our life

JORDI TORRES

Quantum leaps in the quality of a wide range of everyday technologies thanks to Artificial Intelligence

Speech Recognition

We are increasingly interacting with "our" computers by just talking to them

#1. Alexa
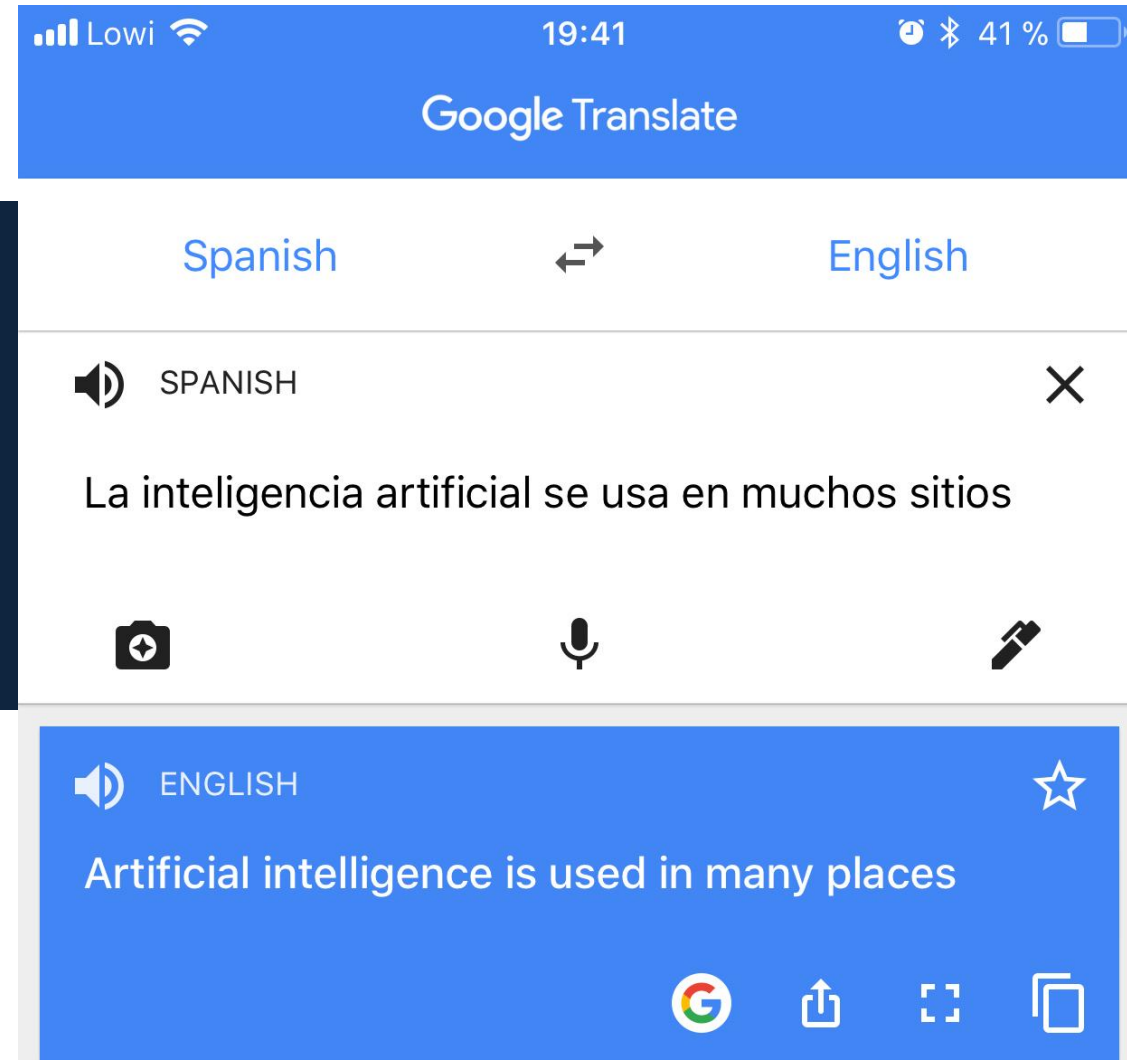(Amazon Echo)

#2. Cortana
(Windows 10 Phone)

#3. Siri
(iPhone)

#4. Google Now
(Android)

JORDI TORRES

*Google Translate* now renders spoken sentences in one language into spoken sentences in another, for **32 pairs** of languages and offers text translation for **100+ languages**.

**Natural Language Processing**

*Google Translate* now renders spoken sentences in one language into spoken sentences in another, for **32 pairs** of languages and offers text translation for **100+ languages**.

**Natural Language Processing**

Now our computers can recognize images and generate descriptions for photos in seconds.

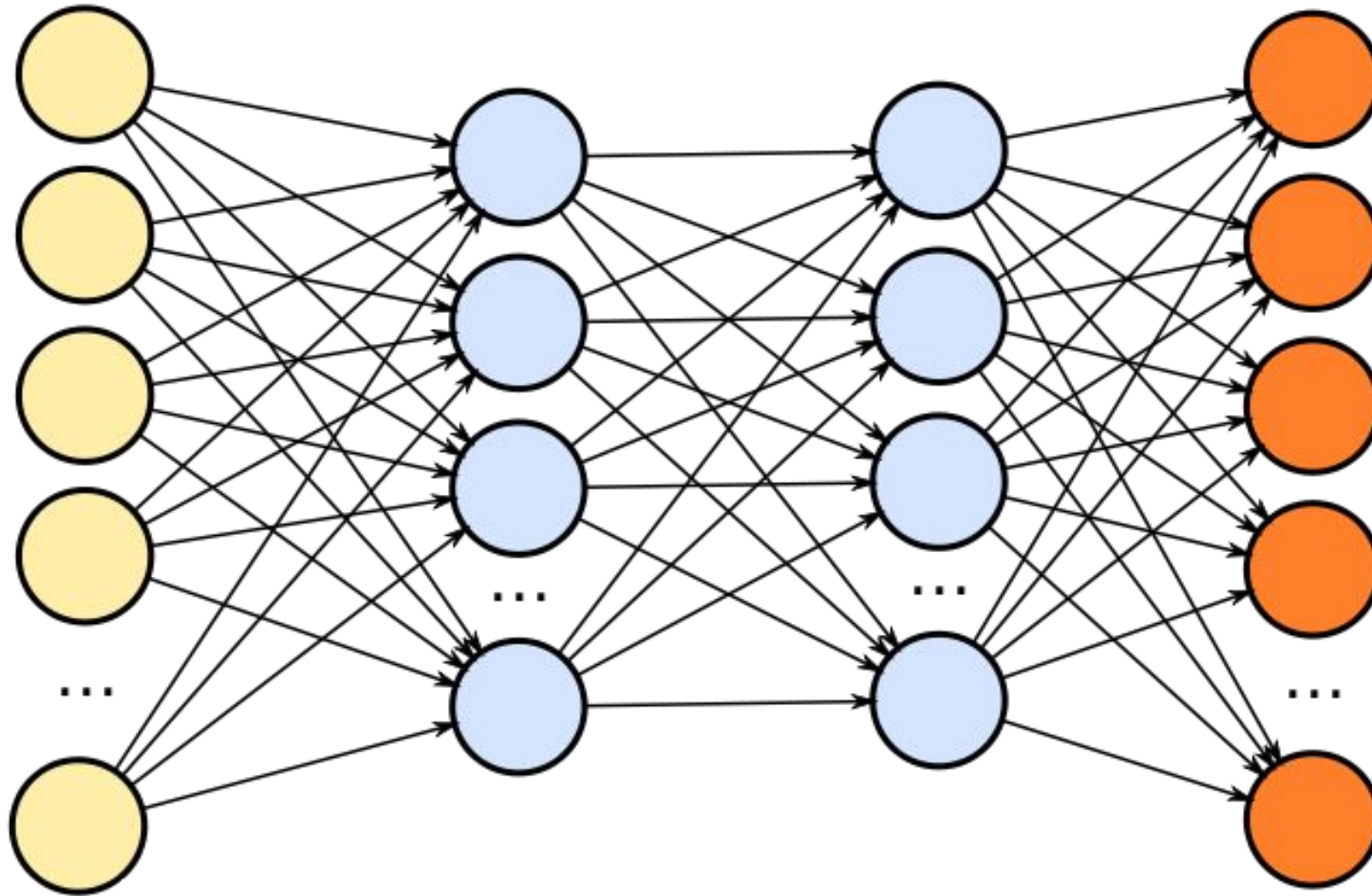All these three areas are crucial to unleashing improvements in **robotics**, drones, self-driving cars, etc.

JORDI TORRES

Source: http://edition.cnn.com/2013/05/16/tech/innovation/robot-bartender-mit-google-makr-shakr/

All these three areas are crucial to unleashing improvements in robotics, **drones**, self-driving cars, etc.

JORDI **TORRES**

Source: http://axisphilly.org/article/military-drones-philadelphia-base-control/

All these three areas are crucial to unleashing improvements in robotics, drones, **self-driving cars**, etc.

JORDI TORRES

AI is at the heart of today's technological innovation.

image source: http://acuvate.com/blog/22-experts-predict-artificial-intelligence-impact-enterprise-workplace

# Many of these breakthroughs have been made possible by a family of AI known as Neural Networks

Neural networks, also known as a **Deep Learning,** enables a computer to learn from observational data

Although the greatest impacts of deep learning may be obtained when it is integrated into the whole toolbox of other AI techniques

Universitat Politècnica de Barcelona

FACULTAT D'INFORMÀTICA
GUIA DOCENT

Curs 1982/83

John McCarthy coined the term Artificial Intelligence in the 1950s

In 1958 Frank Rosenblatt built a prototype neural net, which he called the Perceptron

JORDI TORRES

One of the key drivers:
The data deluge

# One of the key drivers: The data deluge

Thanks to the advent of Big Data AI models can be "trained" by exposing them to large data sets **that were previously unavailable**.

# Training DL neural nets has an insatiable demand for Computing



16X Model

**2012 AlexNet**
8 layers
1.4 GFLOP
~16% Error

**2015 ResNet**
152 layers
22.6 GFLOP
~3.5% error

source: cs231n.stanford.edu/slides/2017/cs231n_2017_lecture15.pdf

Thanks to advances in Computer Architecture, nowadays we can solve problems that would have been intractable some years ago.

# 1982

## FACOM  230 – Fujitsu
Instructions per second:     few Mips * (M = 1.000.000)
Processors : 1

JORDI TORRES

# 2012

## MARENOSTRUM III - IBM
Instructions per second: 1.000.000.000 MFlops
Processors :   6046 (48448 cores)

JORDI TORRES

# 2012

**MARENOSTRUM III - IBM**
Instructions per second: 1.000.000.000 MFlops
Processors : 6046 (48448 cores)

**only 1.000.000.000 times faster**

JORDI TORRES

# CPU improvements!

Until then, the increase in computational power every decade of "my" computer, was mainly thanks to CPU

# CPU improvements!

Until then, the increase in computational power every decade of "my" computer, was mainly thanks to CPU

Since then, the increase in computational power for Deep Learning has not only been from CPU improvements . . .

but also from the realization that GPUs (NVIDIA) were 20 to 50 times more efficient than traditional CPUs.

# Deep Learning requires computer architecture advancements

Fast tightly coupled network interfaces



Dense computer hardware

AI specific processors



Optimized libraries and kernels

**COMPUTING POWER**
is the real enabler!

What if I do not have this hardware?

Now we are entering into an era of computation democratization for companies !

And what is "my/your" computer like now?

And what is "my/your" computer like now?



Source: http://www.google.com/about/datacenters/gallery/images

And what is "my/your" computer like now?

The Cloud

# Huge data centers!



28.000 m2

28.000 m2

Foto: Google

28.000 m2

Foto: Google

28.000 m2

For those (experts) who want to develop their own software, cloud services like Amazon Web Services provide GPU-driven deep-learning computation services

New P2 Instance Type for Ama ✕

🔒 https://aws.amazon.com/blogs/aws/new-p2-instance-type-for-amazon-ec2-up-to-16-gpus/

☰ Menu    **amazon** web services    Live AWS re:Invent    Products ▾    Solutions    Pricing    Software    More ▾

AWS Blog

# New P2 Instance Type for Amazon EC2 – Up to 16 GPUs

by Jeff Barr | on 29 SEP 2016 | in Amazon EC2, Launch | Permalink | 💬 Comments

| Instance Name | GPU Count | vCPU Count | Memory | Parallel Processing Cores | GPU Memory | Network Performance |
|---|---|---|---|---|---|---|
| p2.xlarge | 1 | 4 | 61 GiB | 2,496 | 12 GiB | High |
| p2.8xlarge | 8 | 32 | 488 GiB | 19,968 | 96 GiB | 10 Gigabit |
| p2.16xlarge | 16 | 64 | 732 GiB | 39,936 | 192 GiB | 20 Gigabit |

# And Google ...

And Google ...

And all major cloud platforms...

Microsoft Azure
IBM Cloud
Aliyun
Cirrascale
NIMBIX
Outscale

. . .

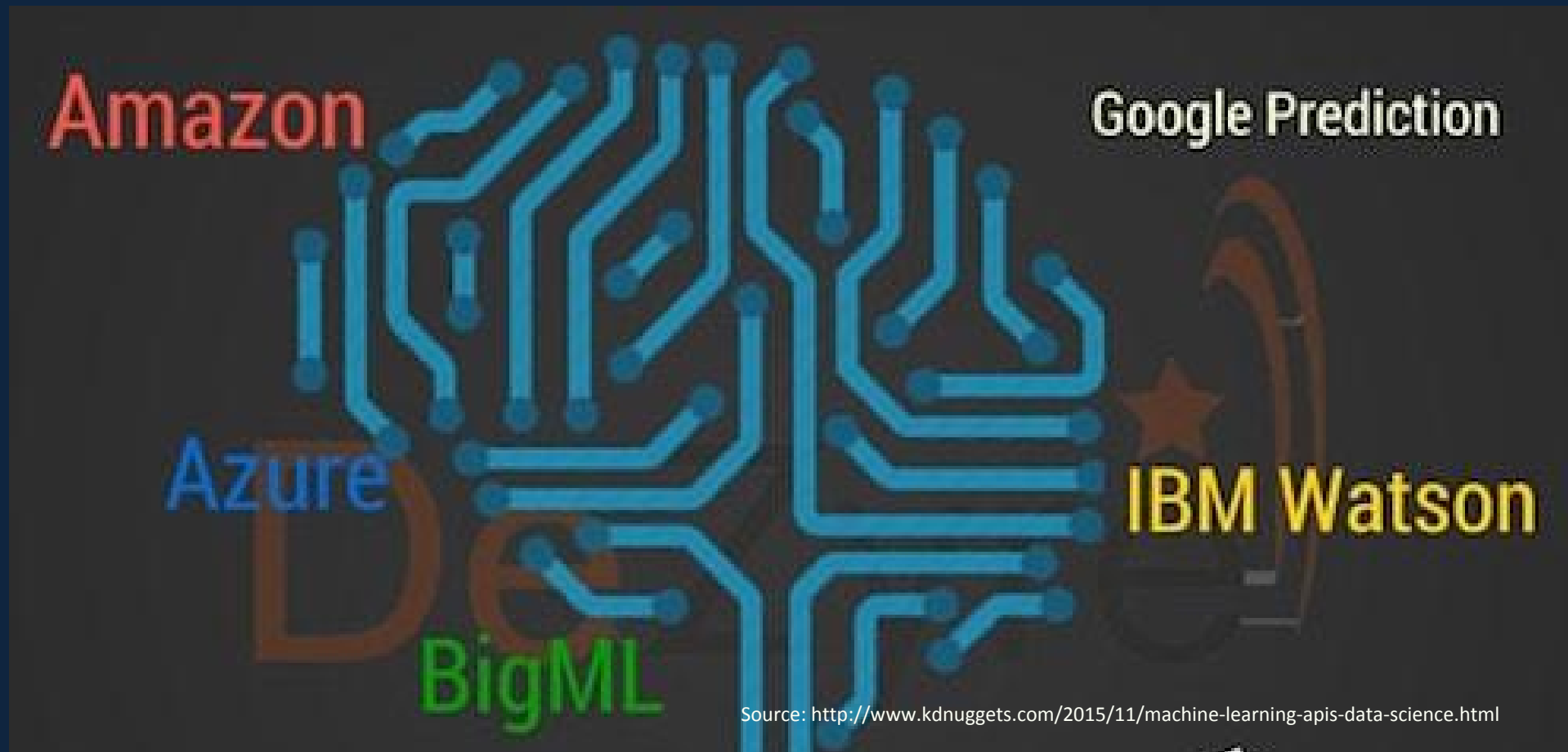Cogeco Peer 1
Penguin Computing
RapidSwitch
Rescale
SkyScale
SoftLayer

. . .

And for "less expert" people, various companies are providing a working scalable implementation of ML/AI algorithms as a Service (**AI-as-a-Service**)
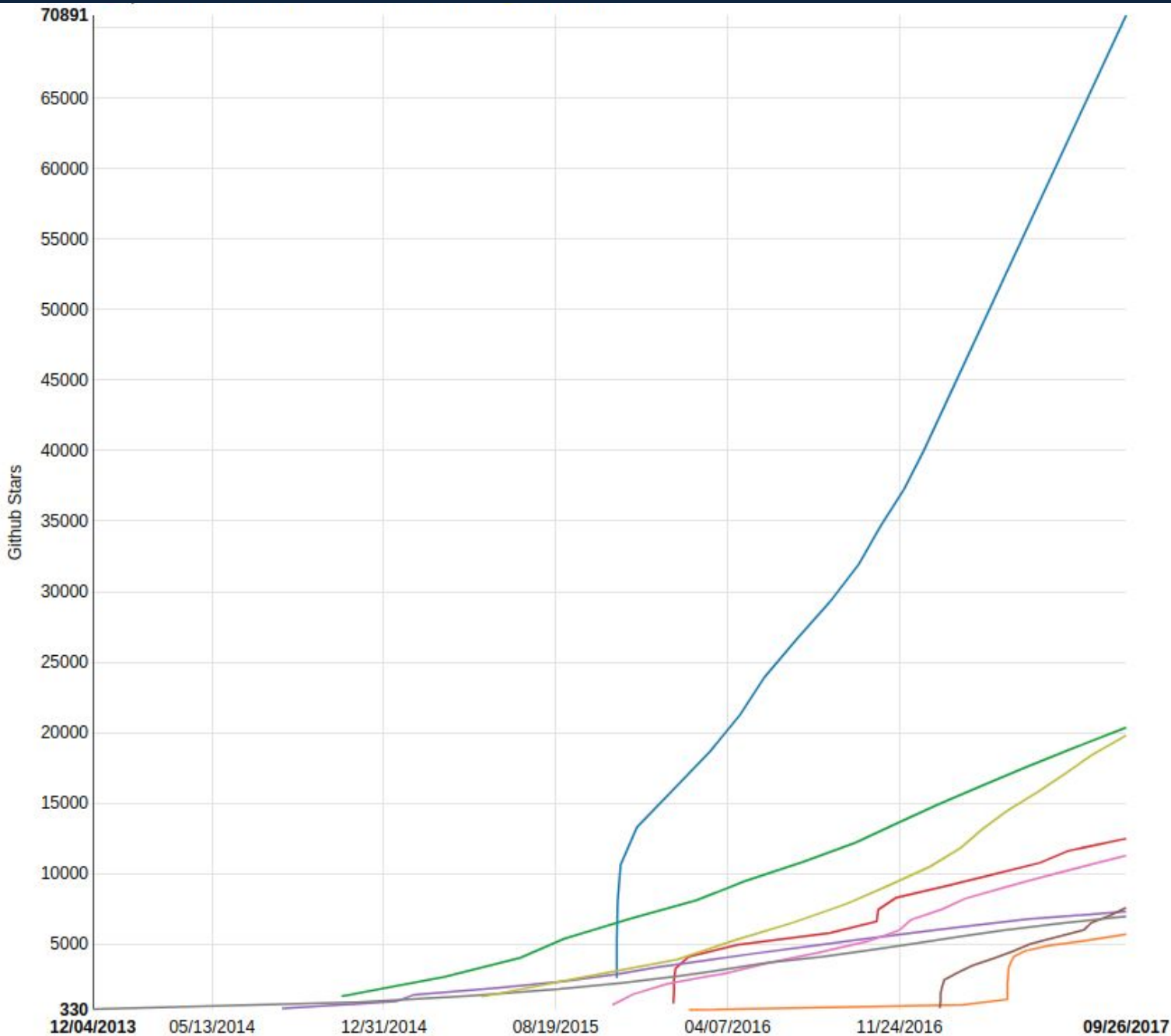
# An open-source world for the Deep Learning community
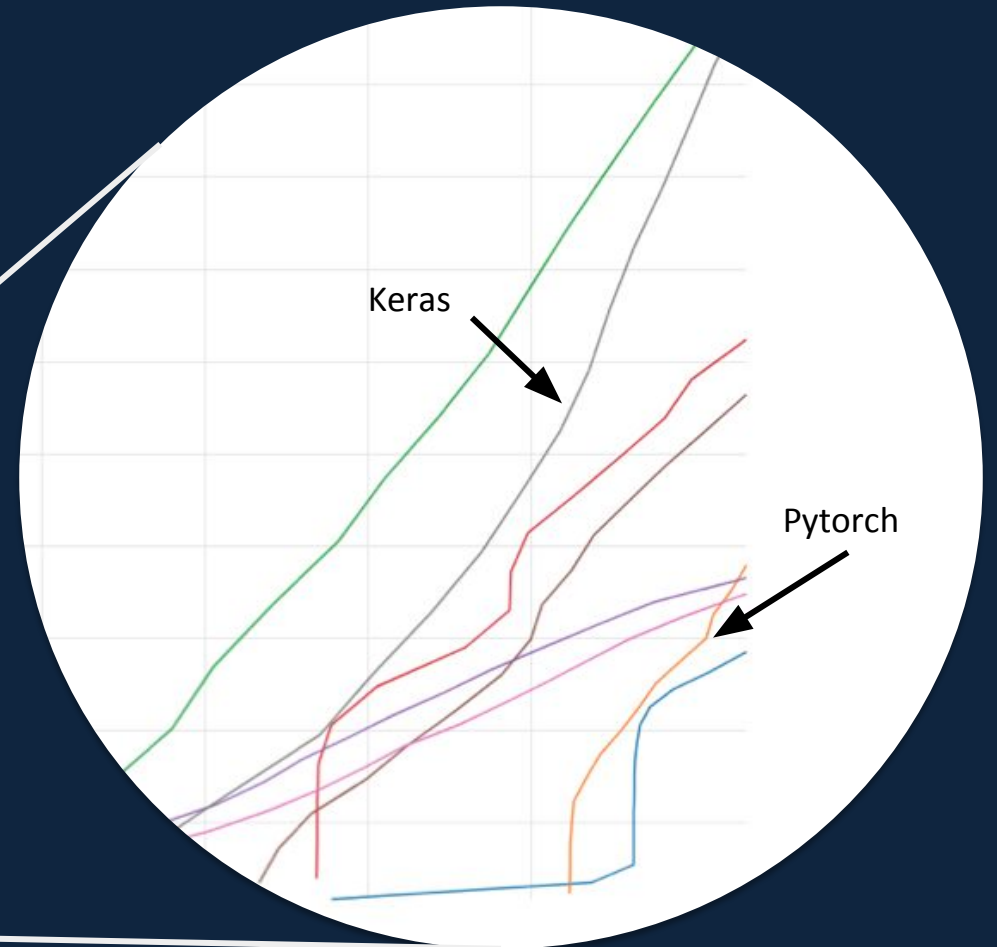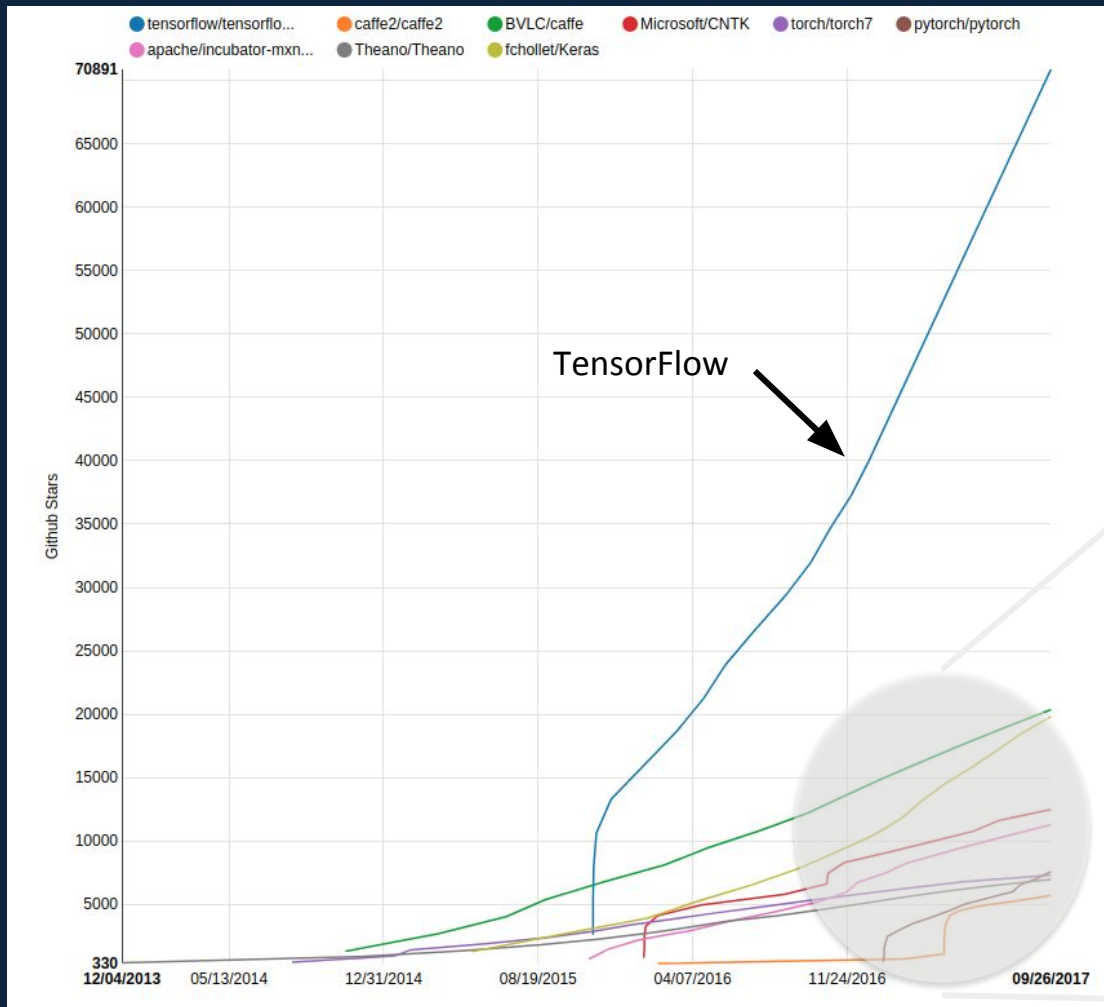
# Many **open-source DL software** have greased the innovation process

# Github Stars



| | |
|---|---|
| ● | tensorflow/tensorflo... |
| ● | apache/incubator-mxn... |
| ● | caffe2/caffe2 |
| ● | Theano/Theano |
| ● | Microsoft/CNTK |
| ● | pytorch/pytorch |
| ● | torch/torch7 |

source: Francesc Sastre

# In this course: **we will consider Keras**



frameworks with more slope

and no less important, **an open-publication ethic**, whereby many researchers publish their results immediately on a database without awaiting peer-review approval.

JORDI **TORRES** | FRANCESC **SASTRE**