

# Data Science & ML Course

## Lesson #1 - Outline & Directions

Ivanovitch Silva  
September, 2018





# Introduction

---





Ivanovitch Silva ([ivan@imd.ufrn.br](mailto:ivan@imd.ufrn.br))



# Roadmap history of research activities



# Deployments of research products



# Data Science & Innovation Lab

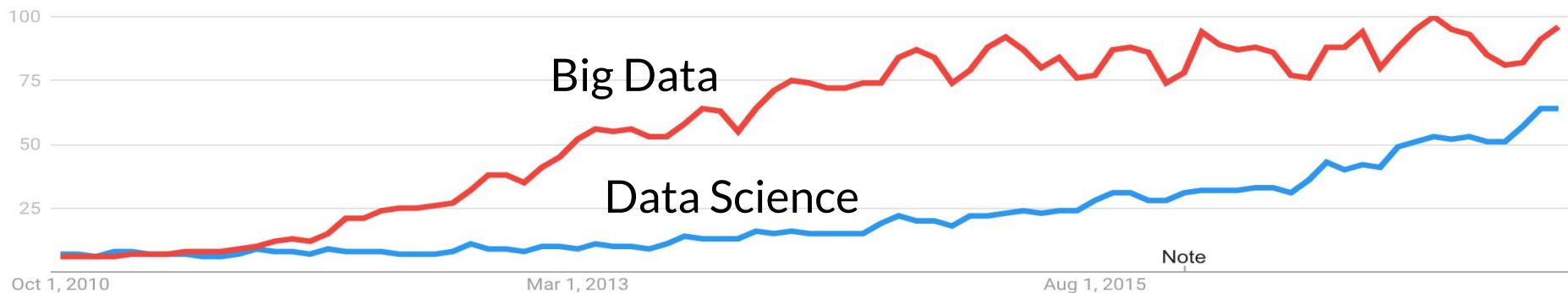


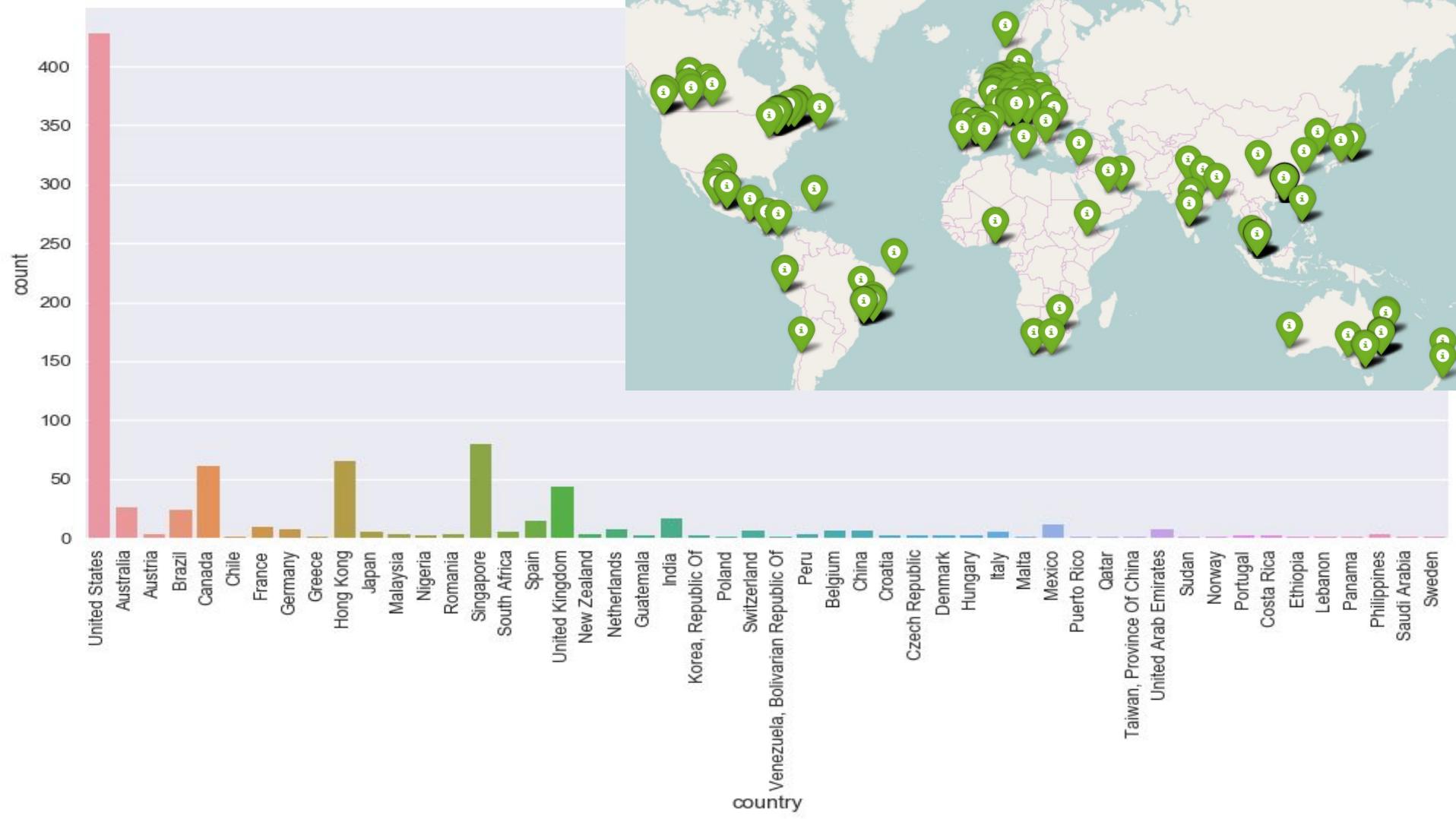
# Big Data and Social Analytics certificate course

2017 DATES TO BE CONFIRMED

[DOWNLOAD COURSE PROSPECTUS](#)

*Discover a new way to think about big data analysis when you explore the theory behind "social analytics", and practically apply that knowledge as you learn pioneering data analytics techniques from the creators of those very tools and methods.*







## 2016/2017 - Specialization course in Big Data

### Undergraduate

- 2017.1 - IMD0105 Introduction to Data Science
- 2017.2 - IMD0252 Learning Analytics
- 2017.2 - DCA0046 Data Science
- 2018.1 - DCA0132 Data Engineering
- 2018.2 - IMD0905 Data Science
- 2018.2 - DCA0131 Data Science

### Graduate

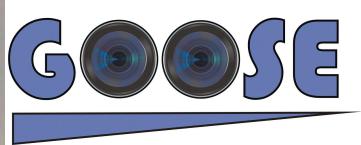
- 2017.2 - EEC2006 Data Science Foundations
- 2017.2 - ITE0021 Learning Analytics
- 2018.2 - EEC1509 Machine Learning

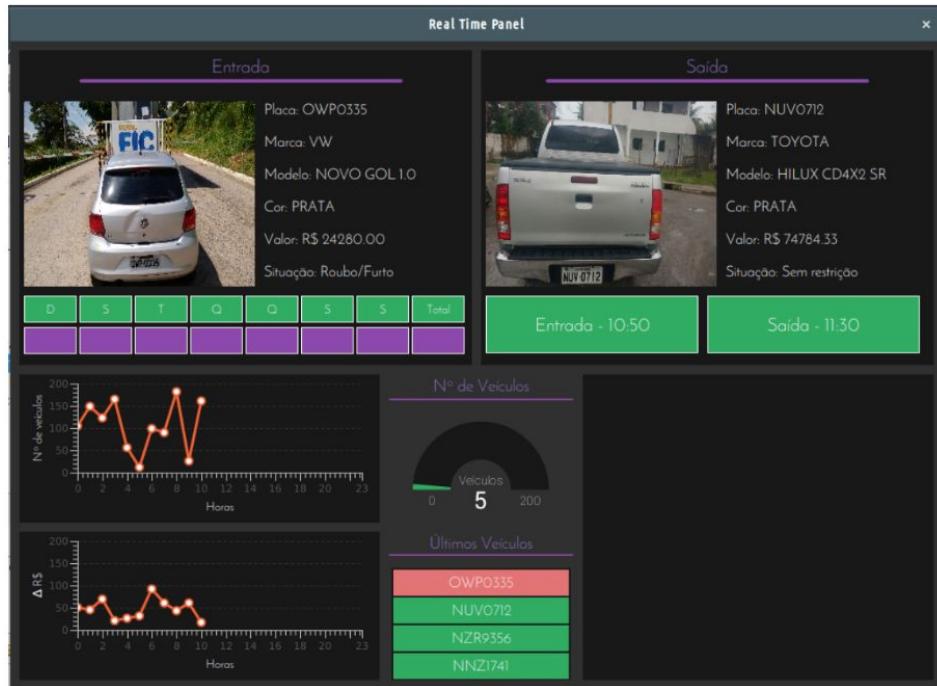
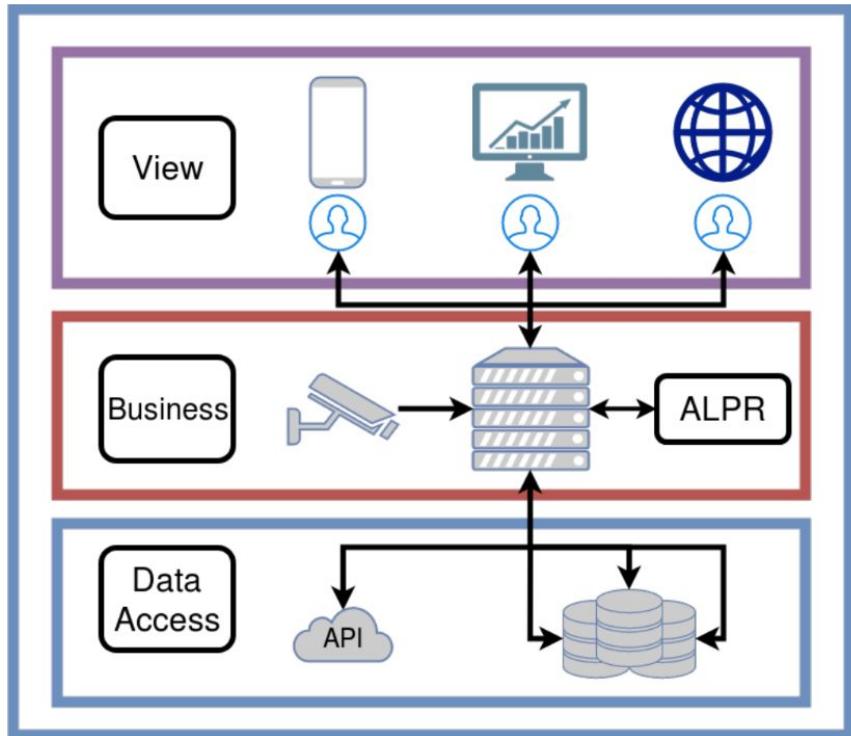


DataCamp

# ALPR Mobile









Multispectral Satellite  
Landsat 5

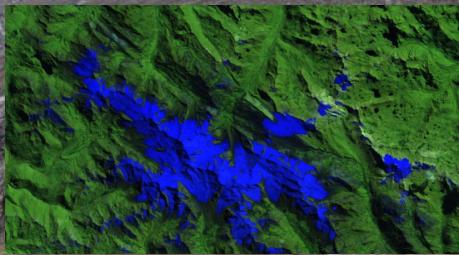
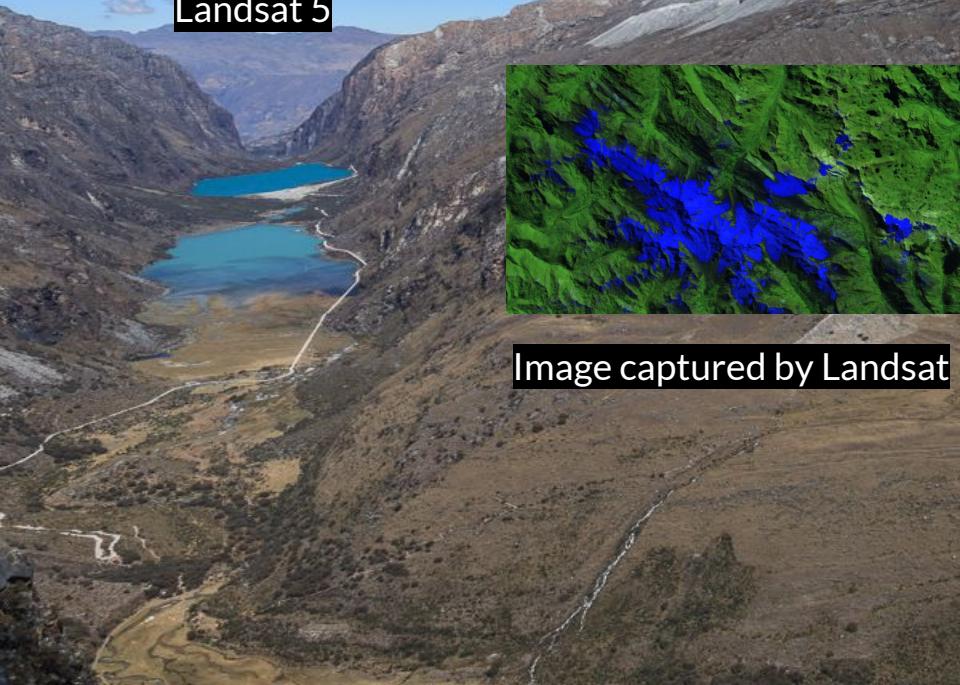
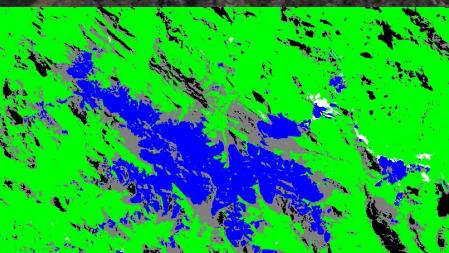
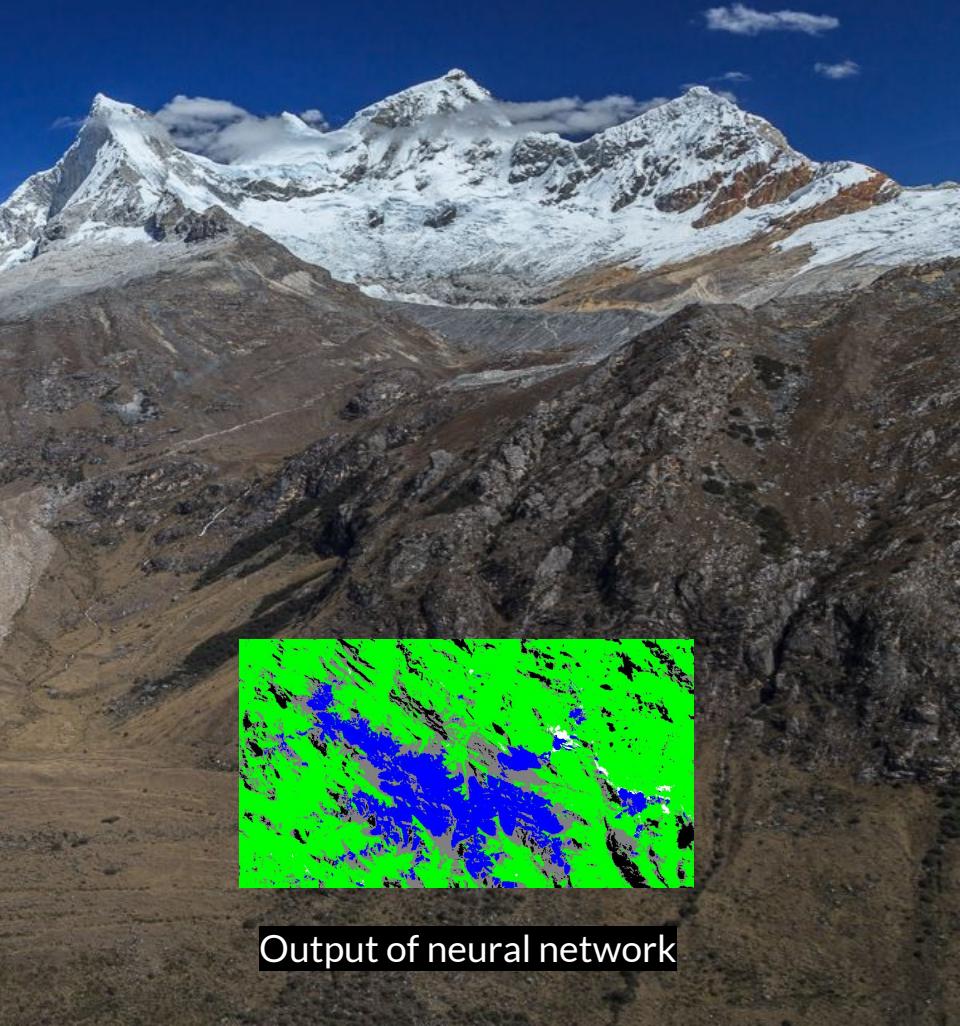


Image captured by Landsat



Output of neural network

Protocolo de Autorização  
06/02/2015 10:58:18

Consulta pelo número da licença  
1515 0004 8000 no site da Cetesb

Cetesb - Agência Estadual do Meio Ambiente

Consulta



Protocolo de Autorização  
06/02/2015 10:58:18

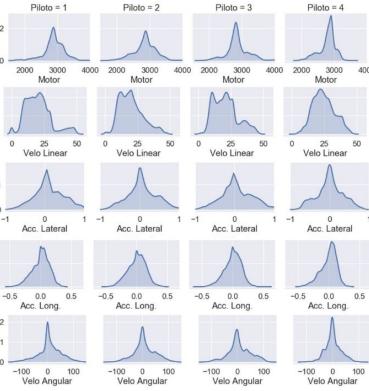
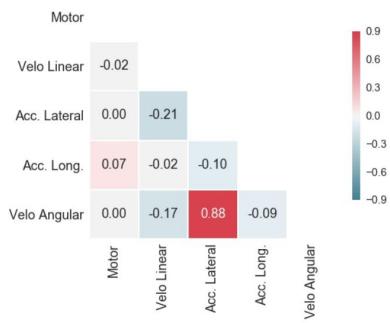
OBD-II

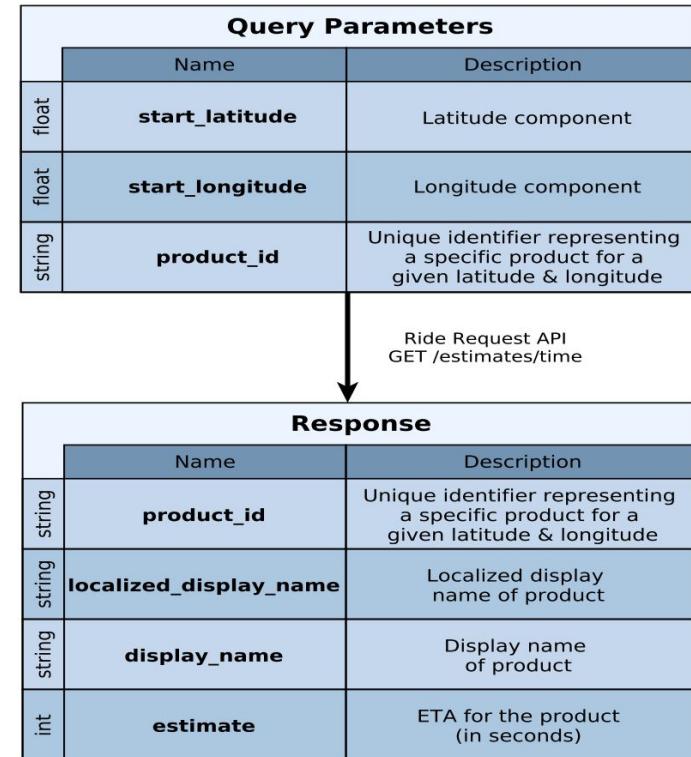
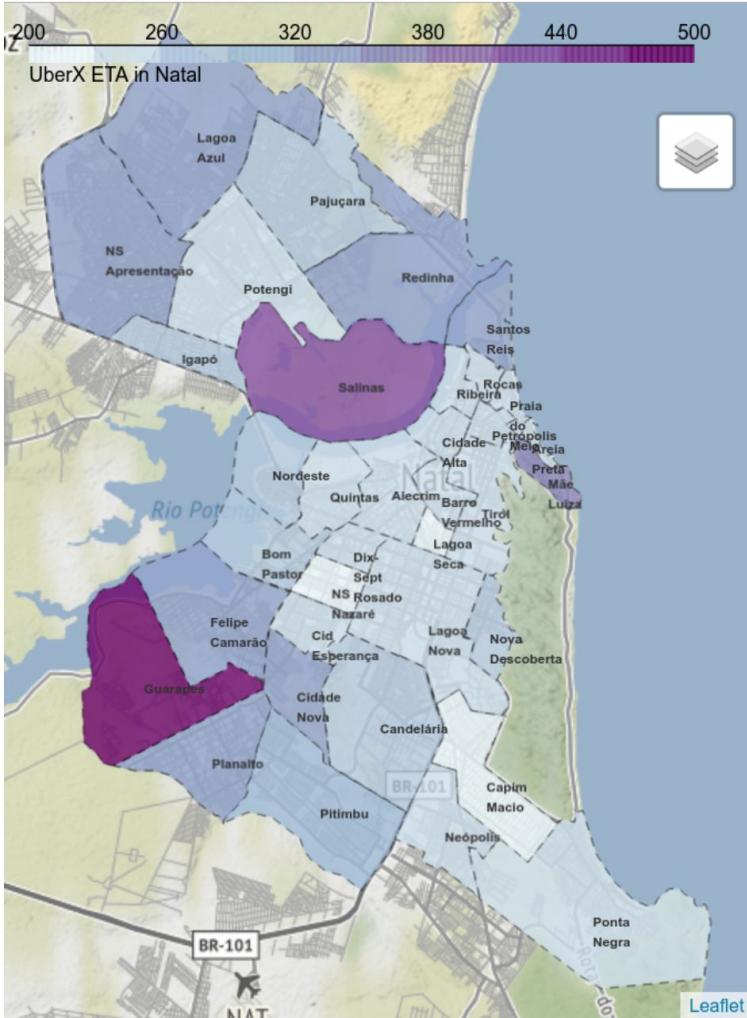




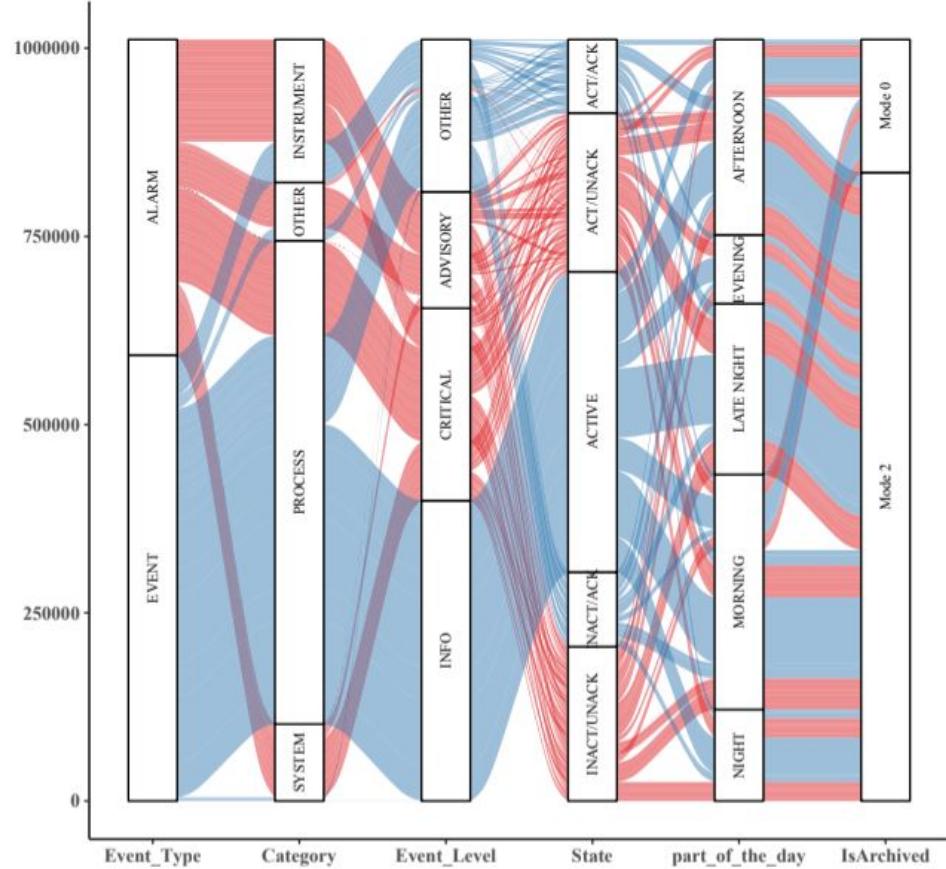
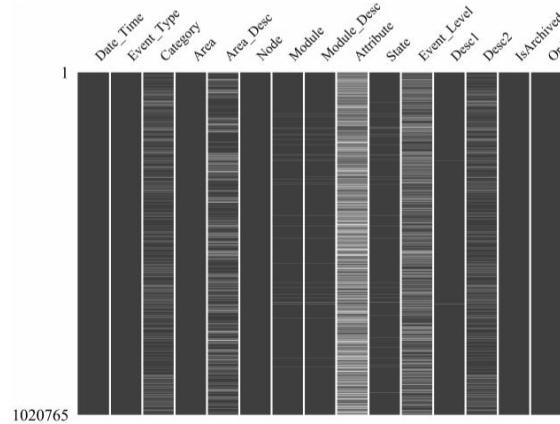
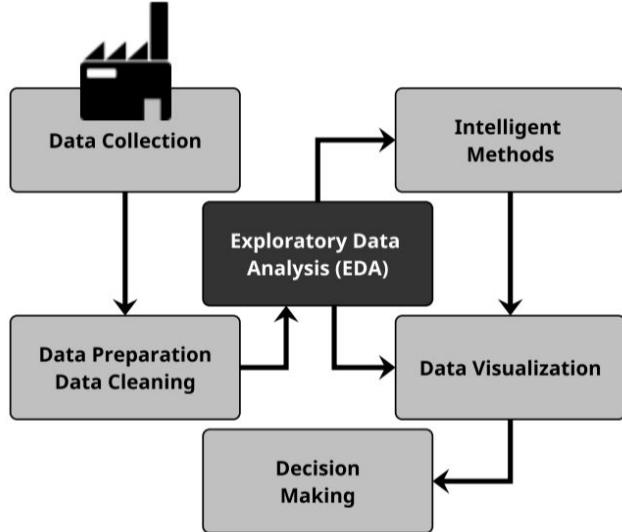
UNIVERSITY  
OF BRESCIA

**UFRN**  
UNIVERSIDADE FEDERAL DO RIO GRANDE DO NORTE





# UBER



# UFRN Analytics

---

Este repositório é dedicado à análise de dados do [Portal de Dados Abertos da UFRN](#), que reune dados de diversas áreas da Universidade com o intuito de prover transparência às ações da instituição e incentivar o desenvolvimento de aplicações que contribuam com a gestão universitária.

[Servidores](#)

[Ensino](#)

[Pesquisa](#)

[Extensão](#)

<https://github.com/ycaroravel/UFRN-Analytics>





## Group Survey

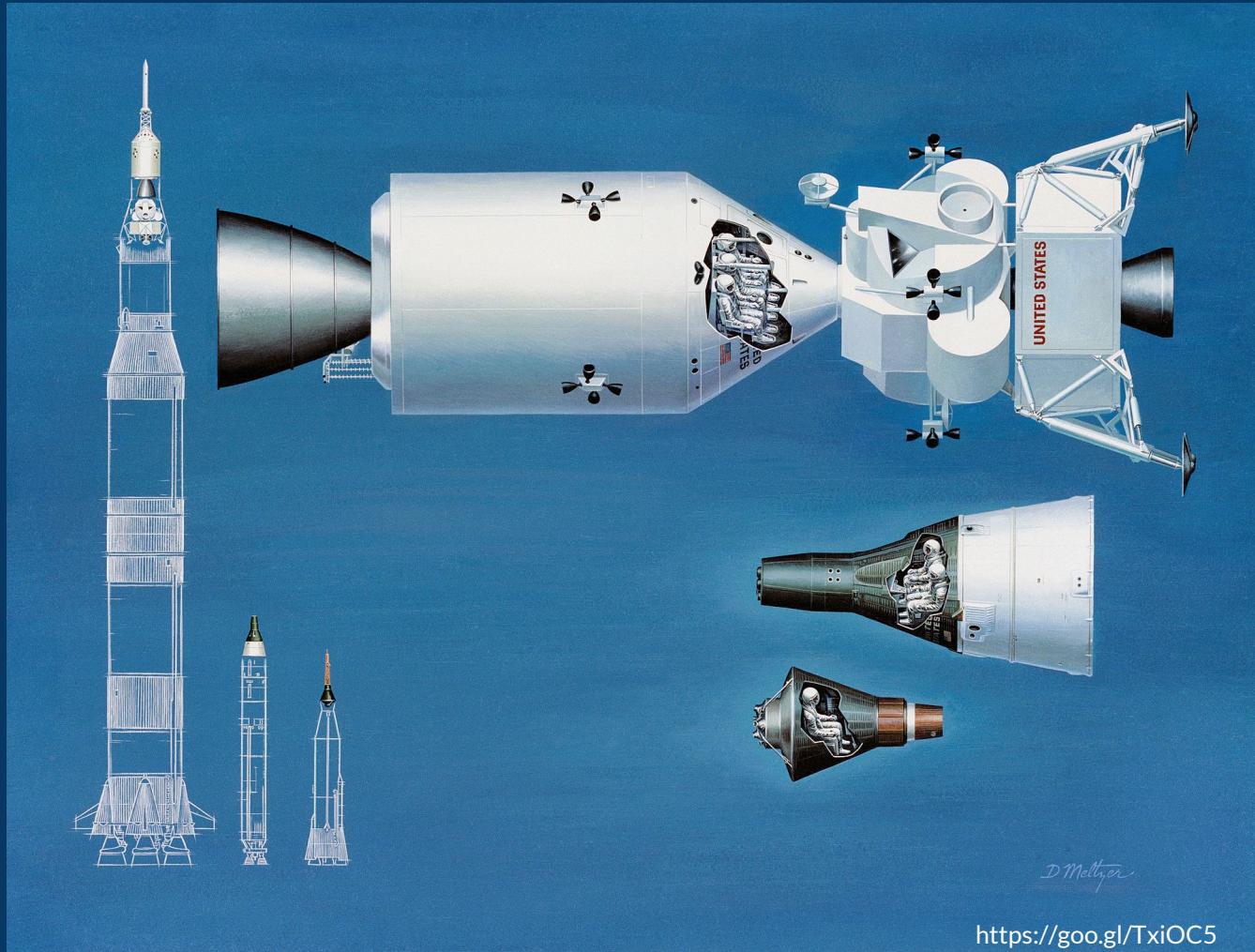


# Provocation #1

hardware

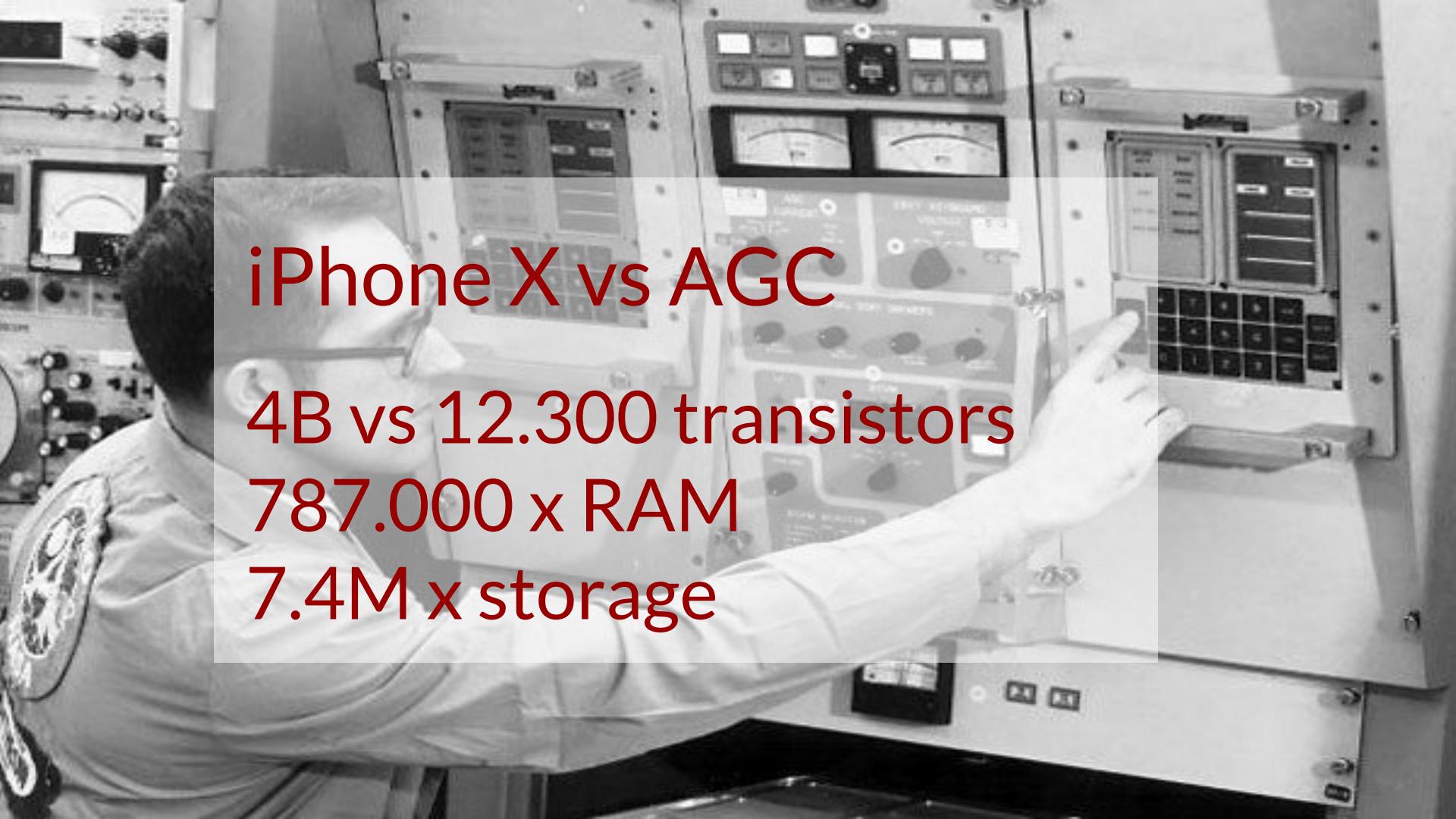


1969



D. Meltzer

<https://goo.gl/TxiOC5>



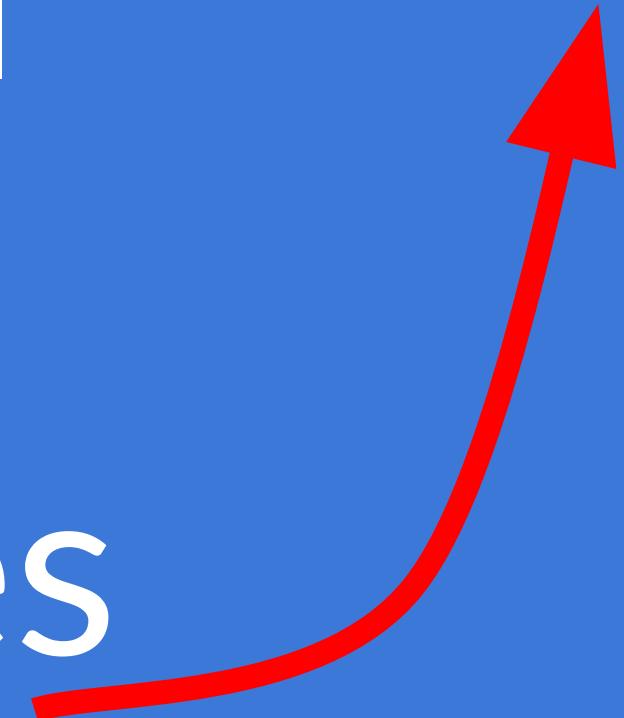
# iPhone X vs AGC

4B vs 12.300 transistors

787.000 x RAM

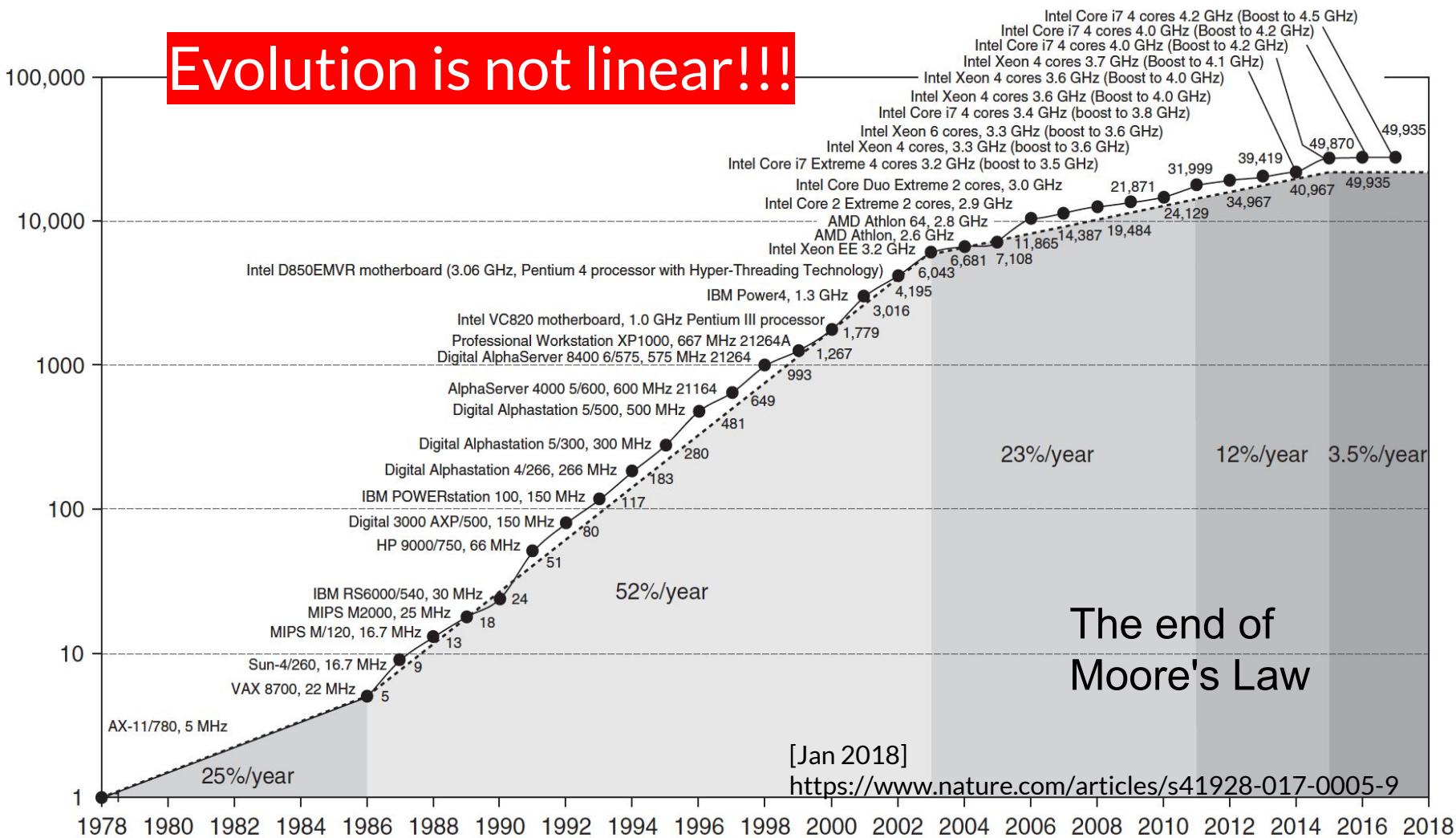
7.4M x storage

exponential  
growth of  
technologies



# Evolution is not linear!!!

Performance (vs. VAX-11/780)





\$ 2,999.<sup>00</sup>

Architecture

Frame Buffer

Boost Clock

Tensor Cores

CUDA Cores

## NVIDIA TITAN V

The Most Powerful PC GPU Ever Created

NVIDIA TITAN V is the most powerful graphics card ever created for the PC, driven by the world's most advanced architecture—NVIDIA Volta. NVIDIA's supercomputing GPU architecture is now here for your PC, and fueling breakthroughs in every industry.

[LEARN MORE](#)

**NVIDIA TITAN V**

**NVIDIA Volta**

**12 GB HBM2**

**1455 MHz**

**640**

**5120**

**1. COTAÇÕES**

	USD - Dólar:	3.58
	EUR - Euro:	4.25
	GBP - Libra Esterlina:	4.84

**2. ALÍQUOTA DE ICMS**

Tipo de envio:	Courier
Unidade federativa:	RN
Valor da alíquota:	17%

**3. TAXAS**

Conversão monetária:	3.58
Imposto de importação (%):	60
ICMS (%):	0
IOF (%):	6.38

**4. IMPOSTO DE IMPORTAÇÃO****VALORES**

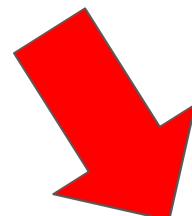
Valor do produto (USD):	2999
Custo do frete (USD):	50
Valor do produto (BRL):	10736.42
Custo do frete (BRL):	179.00
<b>TOTAL DA COMPRA (BRL):</b>	<b>10915.42</b>

**TRIBUTOS**

<input checked="" type="checkbox"/> Incluir frete no cálculo do imposto	<input checked="" type="checkbox"/> Incluir IOF
Imposto de importação (BRL):	6549.25
ICMS (BRL):	0.00
IOF (BRL):	696.40
<b>TOTAL DE TRIBUTOS (BRL):</b>	<b>7245.66</b>

**TOTAIS**

**Valor final com impostos (BRL):** 18161.08



# Lara Croft has changed over 21 years

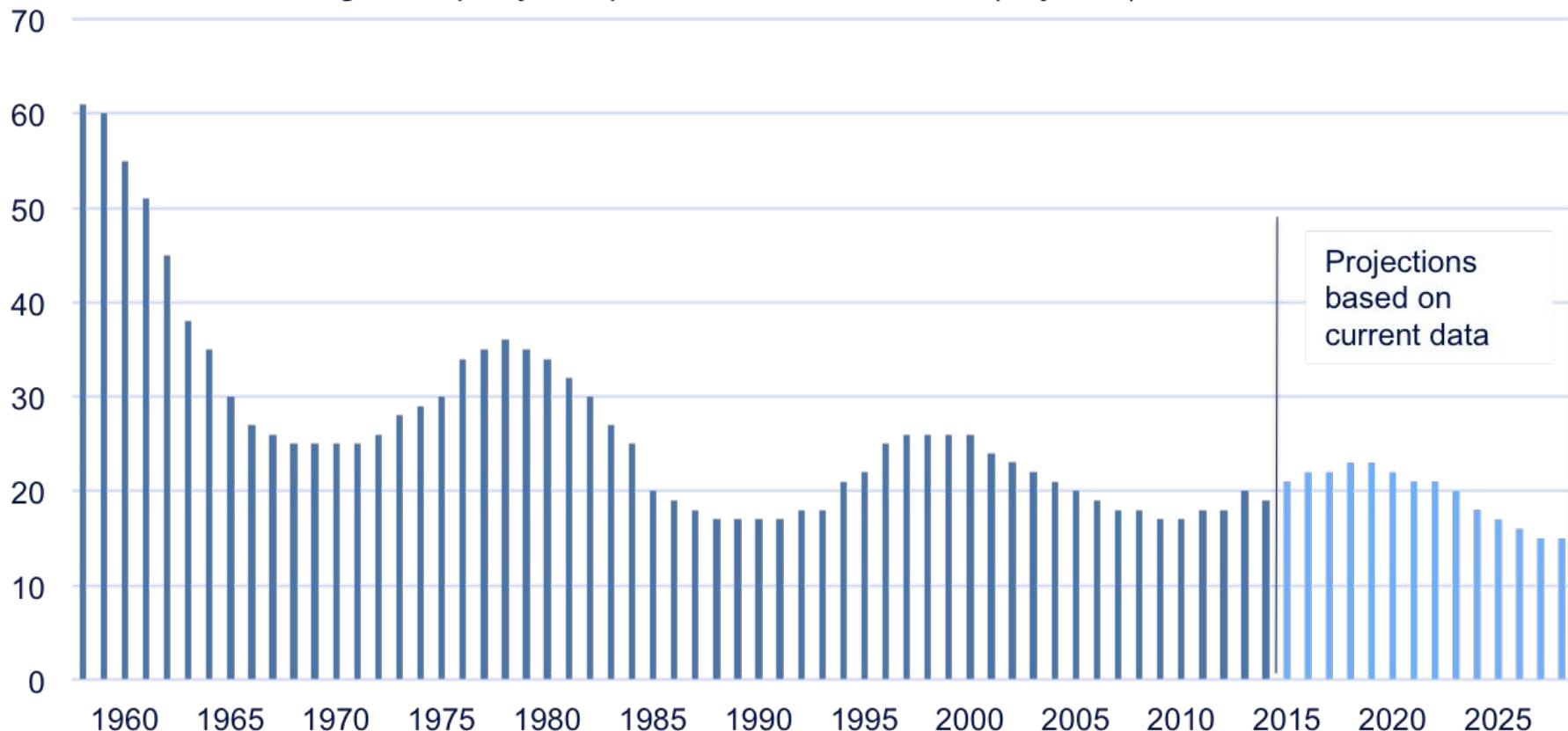


1996 - 2014

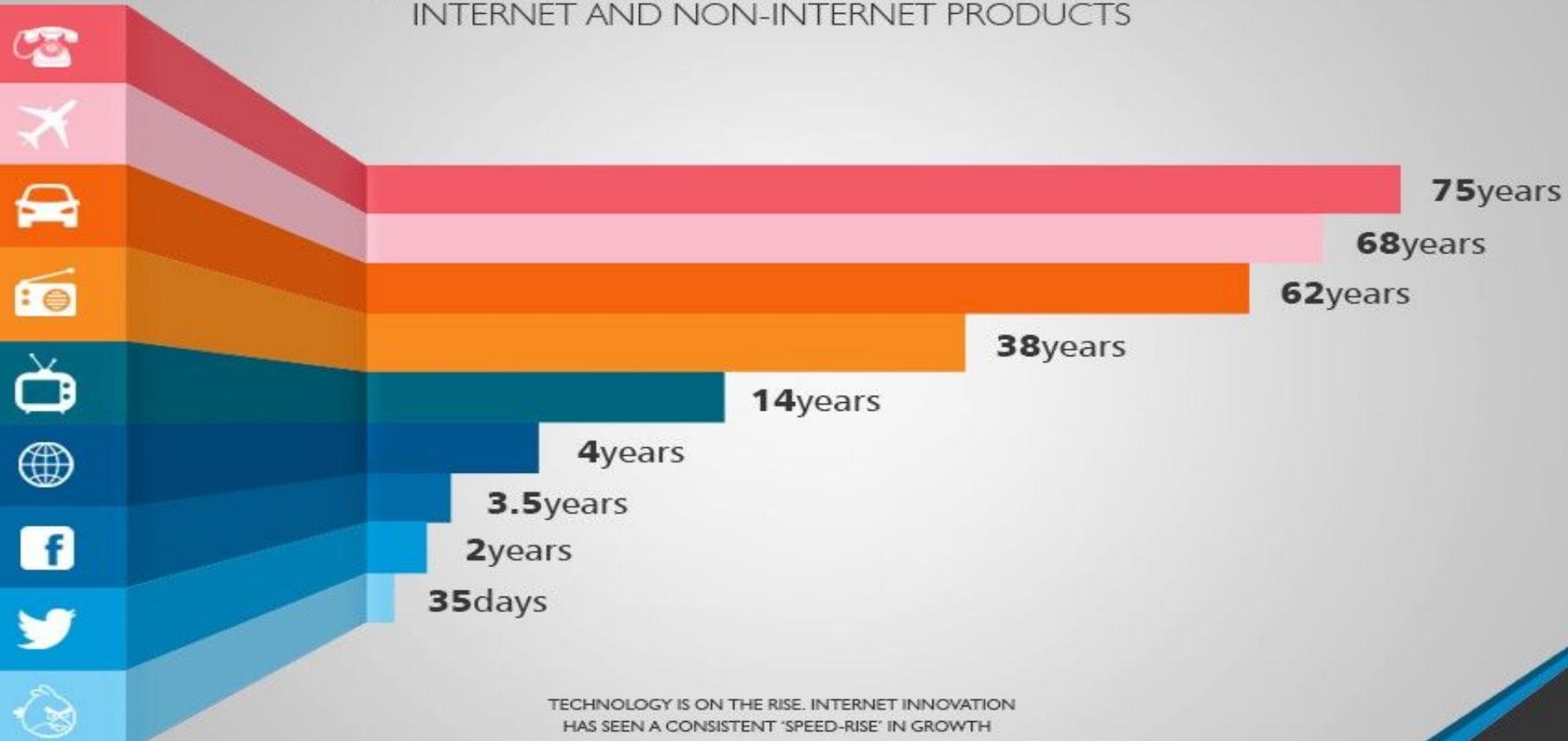
1/6 cost | 1/20 power | 4 hacks in a box



## Average company lifespan on S&P 500 Index (in years)



## REACHING 50 MILLION USERS: THE JOURNEY OF INTERNET AND NON-INTERNET PRODUCTS



# Provocation #2

data & internet & AI

# 90's

"Read"  
Search Engine  
Google

# Update

# 00's

"Write"  
Social Networks  
Facebook

## Participate

# 10's

"Act"  
App  
Uber

Act

# 20's

"Change"  
AI  
?

## Transform

Change



Act

UBER

Write

facebook

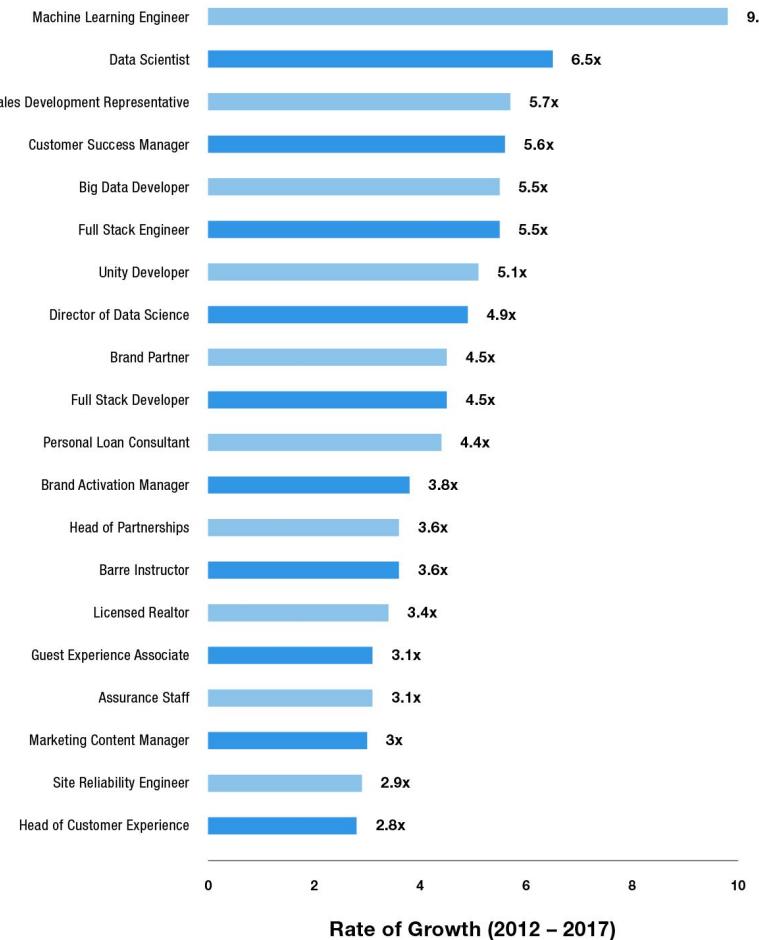
Read

Google



# THE YEAR OF INTELLIGENCE

## Top 20 Emerging Jobs



There are **9.8 times more** Machine Learning Engineers working today than five years ago based on LinkedIn's research

[Dec. 2017]  
<http://bit.do/forbesjobs>



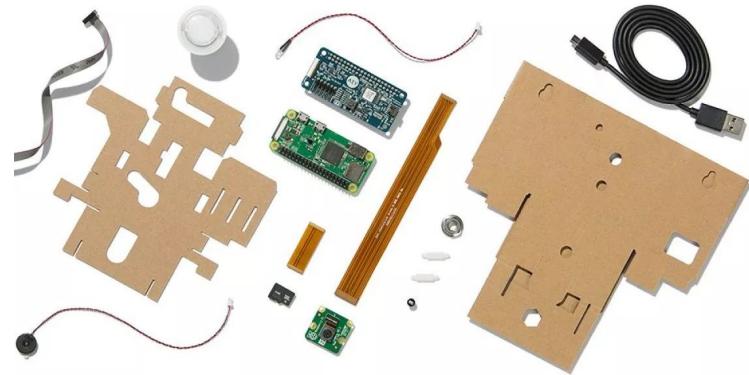
Do-it-yourself artificial intelligence

We want to put AI into the maker toolkit, to help you solve real problems that matter to you and your communities. These kits will get you started by adding natural human interaction to your maker projects.



79.20%  
Fuji  
Category: Edible Fruit

20.80%  
Orange  
Category: Edible Fruit

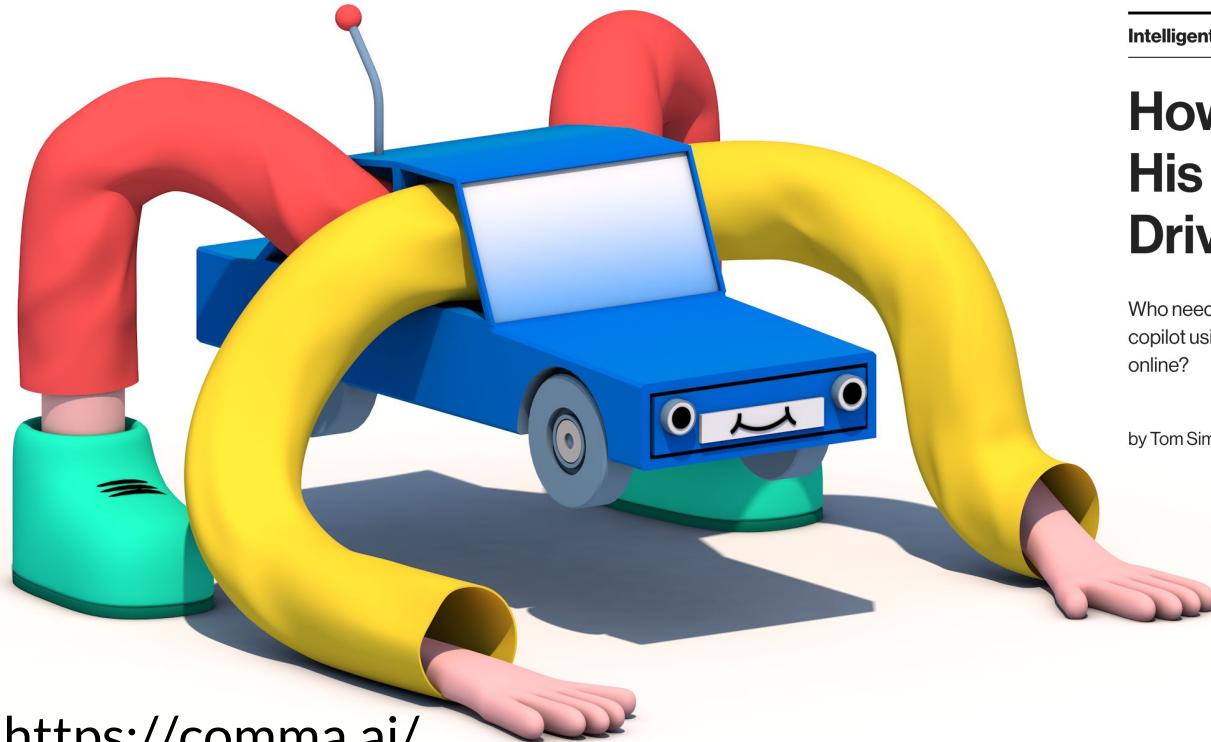


What is the speed of light



The speed of light is  
299,792,458 meters  
per second.





<https://comma.ai/>

<https://www.youtube.com/watch?v=3jstaBeXgAs>

---

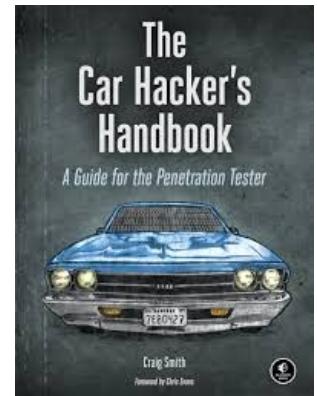
Intelligent Machines

---

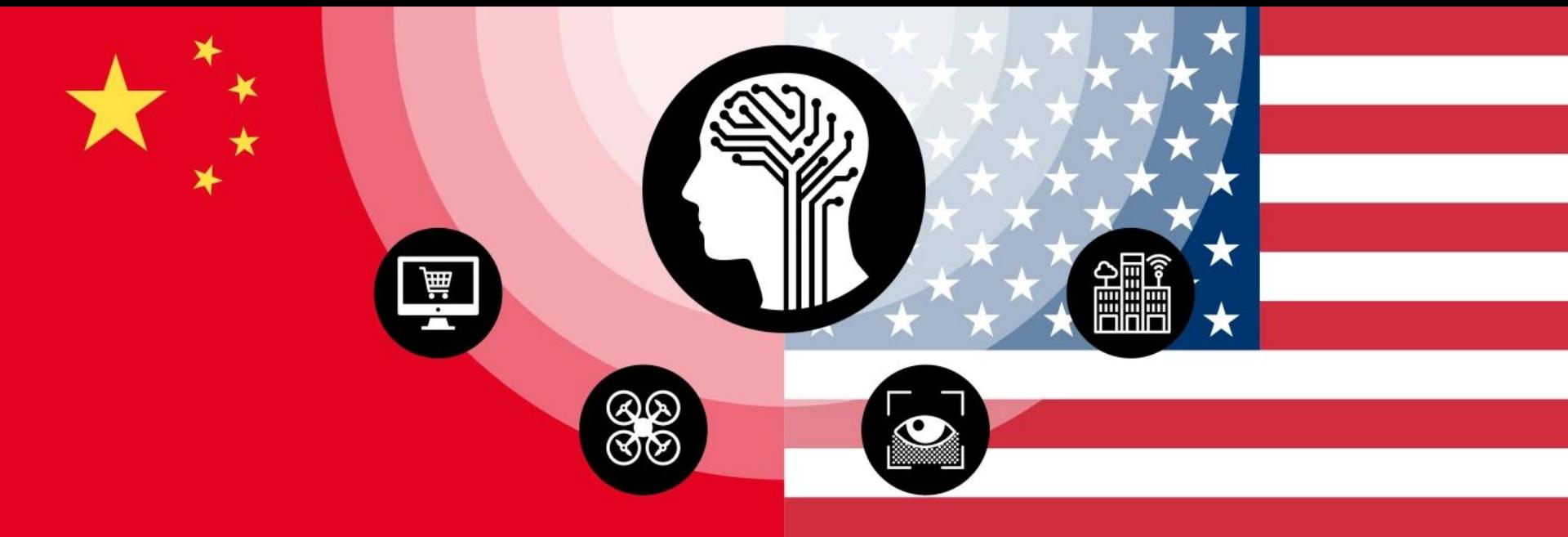
## How a College Kid Made His Honda Civic Self-Driving for \$700

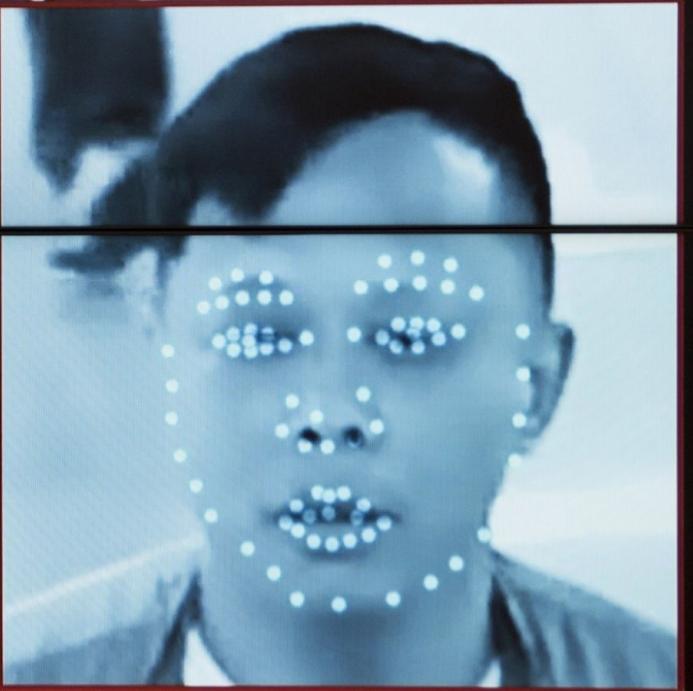
Who needs a Tesla when you can build your own automated copilot using free hardware designs and software available online?

by Tom Simonite February 21, 2017



# AI arms race





相似度  
SIMILARITY

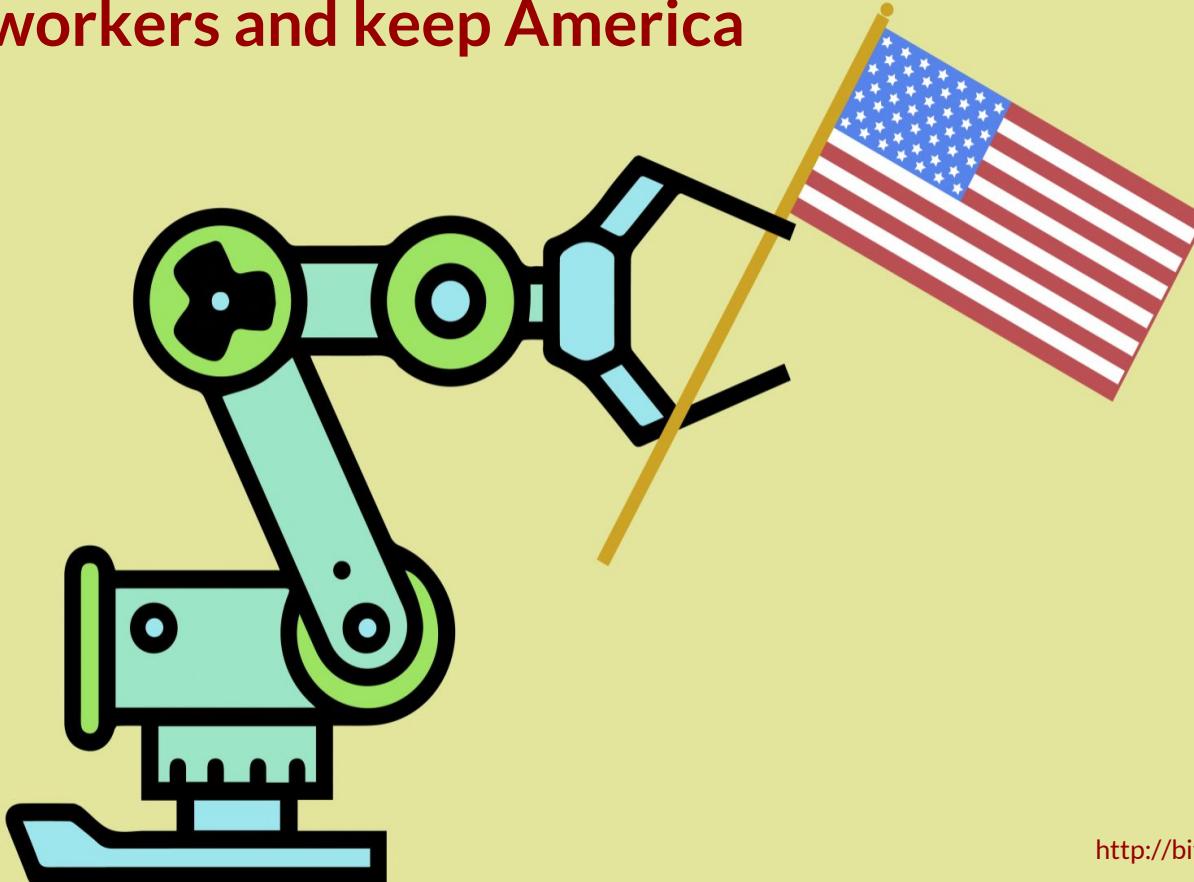
67.2%

<https://www.ft.com/content/e33a6994-447e-11e8-93cf-67ac3a6482fd>

China's  
watchful eye

The White House says a new AI task force  
will protect workers and keep America  
first.

May 10, 2018



[http://bit.do/trump\\_aiprogram](http://bit.do/trump_aiprogram)

# HUMANITY

FRENCH STRATEGY  
FOR ARTIFICIAL  
INTELLIGENCE

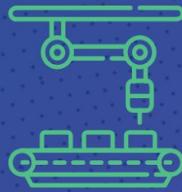
The President of the French Republic presented **his vision and strategy** to make France a **leader in artificial intelligence (AI)** at the Collège de France on **29 March 2018**.



Download the  
Villani Report

<https://www.aiforhumanity.fr/en/>





# *Agenda brasileira para a Indústria 4.0*

O Brasil preparado para os desafios do futuro

[CONHEÇA A AGENDA](#)

## Estudo “Internet das Coisas: um plano de ação para o Brasil”



# Internet das coisas: um plano de ação para o Brasil



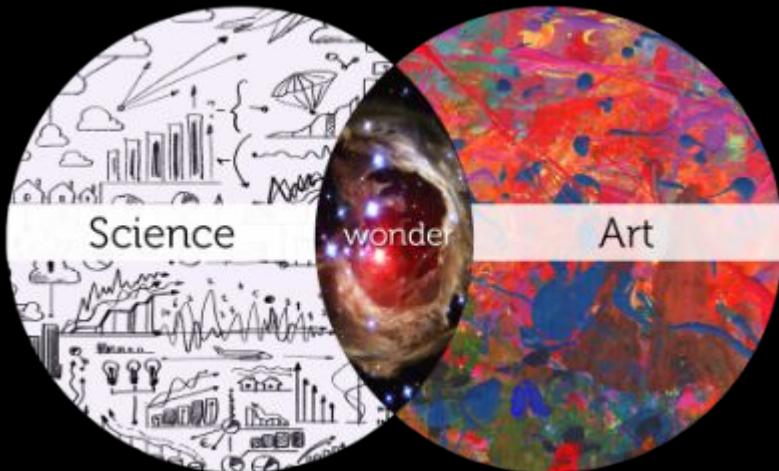
Baixe o Relatório de Plano de Ação (PDF - 3,3 MB) para o desenvolvimento de IoT no Brasil.



**NEURAL**



**SABOTAGE**



<https://www.youtube.com/watch?v=SOtm7vylwxc>





# Spotify MACHINE LEARNING DAY

MONDAY, JULY THE 9<sup>TH</sup>  
STOCKHOLM, SWEDEN



# HOSPITAL SÍRIO-LIBANÊS

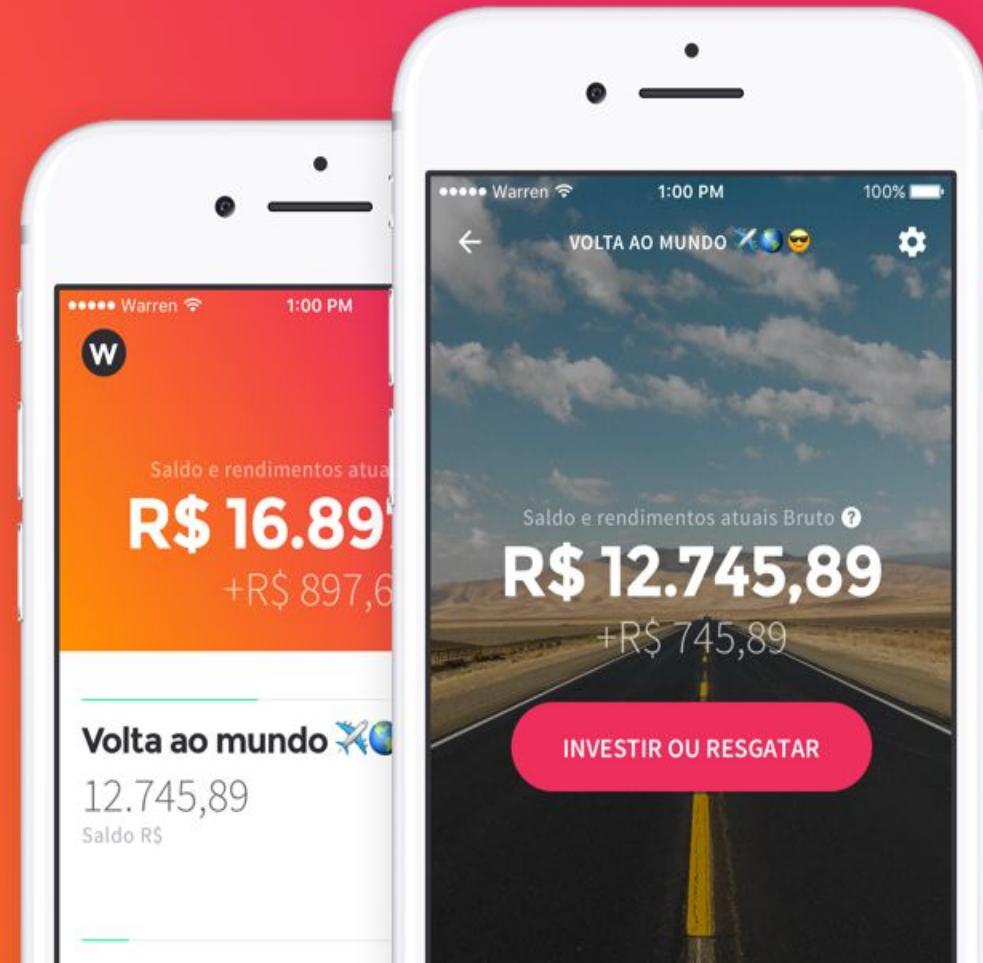
KUNUMI

<http://bhetc.org.br/empresas-do-bhetc/kunumi/>



<https://www.ufmg.br/online/radio/arquivos/046358.shtml>

# Uma nova forma de investir.



# Niantic is going to crowdsource AR maps



**90% of the data in the world today  
has been created in the last two years alone**

# This volume is continuing to grow

Data volume will grow

**800%**

over the next five years

gartner

By 2022

**93%**

of the data will be **free-form**

IDC

# data comes from everywhere



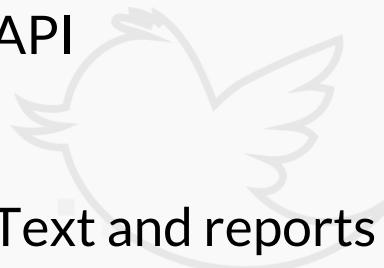
Data Files  
(XML, CSV, Excel, JSON, ...)



Database  
(MySQL, Oracle, ...)



API



Text and reports



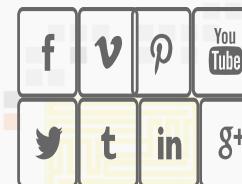
Image and videos



Sites



Maps



Social Media

BIG

CHALLENGES

# THE COMING FLOOD OF DATA IN AUTONOMOUS VEHICLES

RADAR

~10-100 KB  
PER SECOND

SONAR

~10-100 KB  
PER SECOND

GPS

~50KB  
PER SECOND

CAMERAS

~20-40 MB  
PER SECOND

LIDAR

~10-70 MB  
PER SECOND

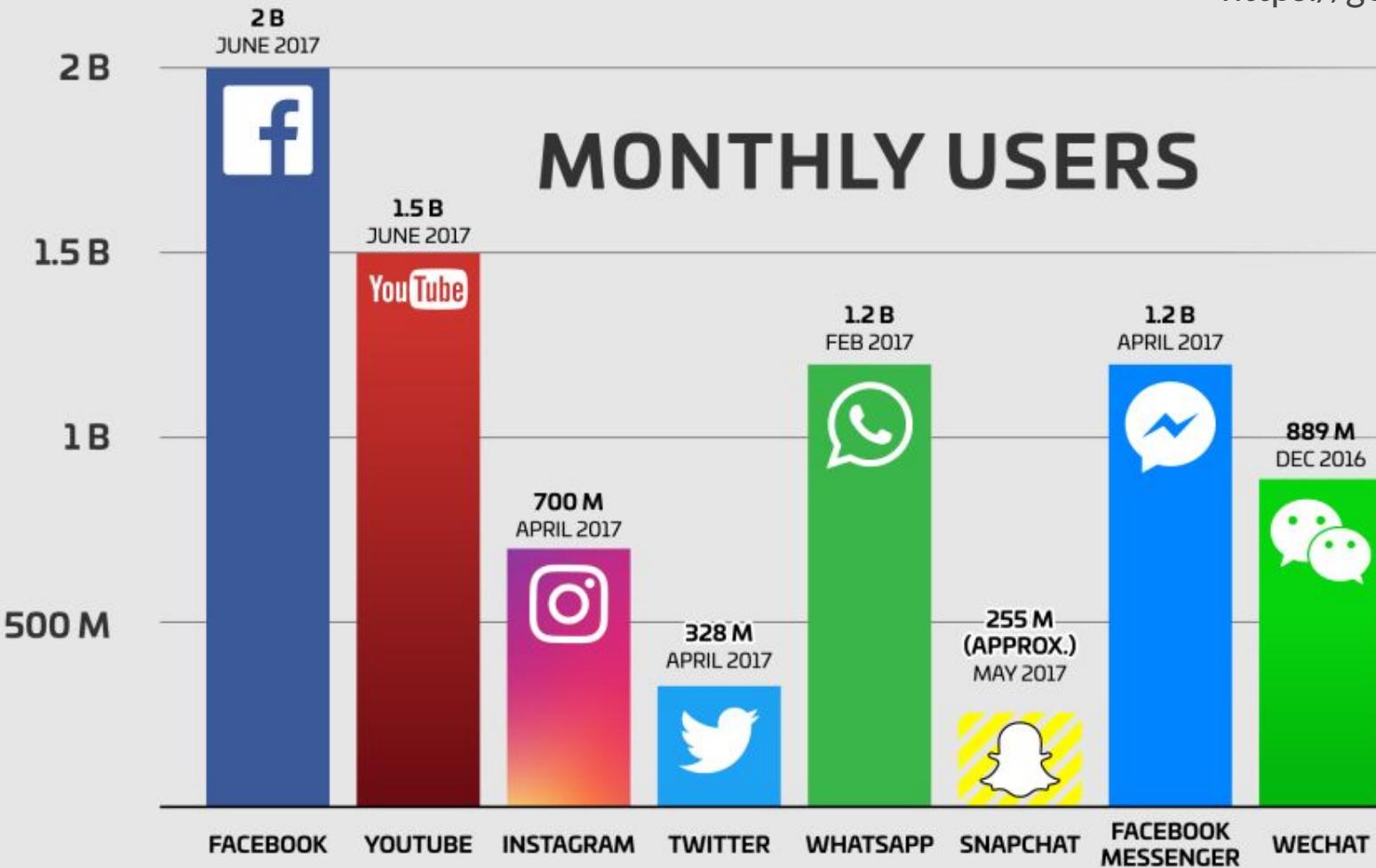
AUTONOMOUS VEHICLES

4,000 GB  
PER DAY... EACH DAY



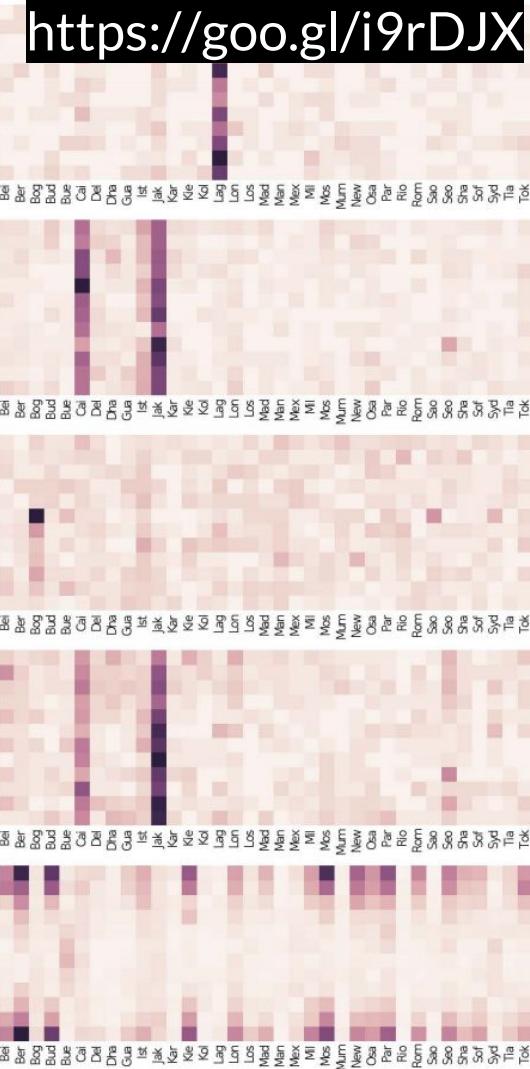
BIG

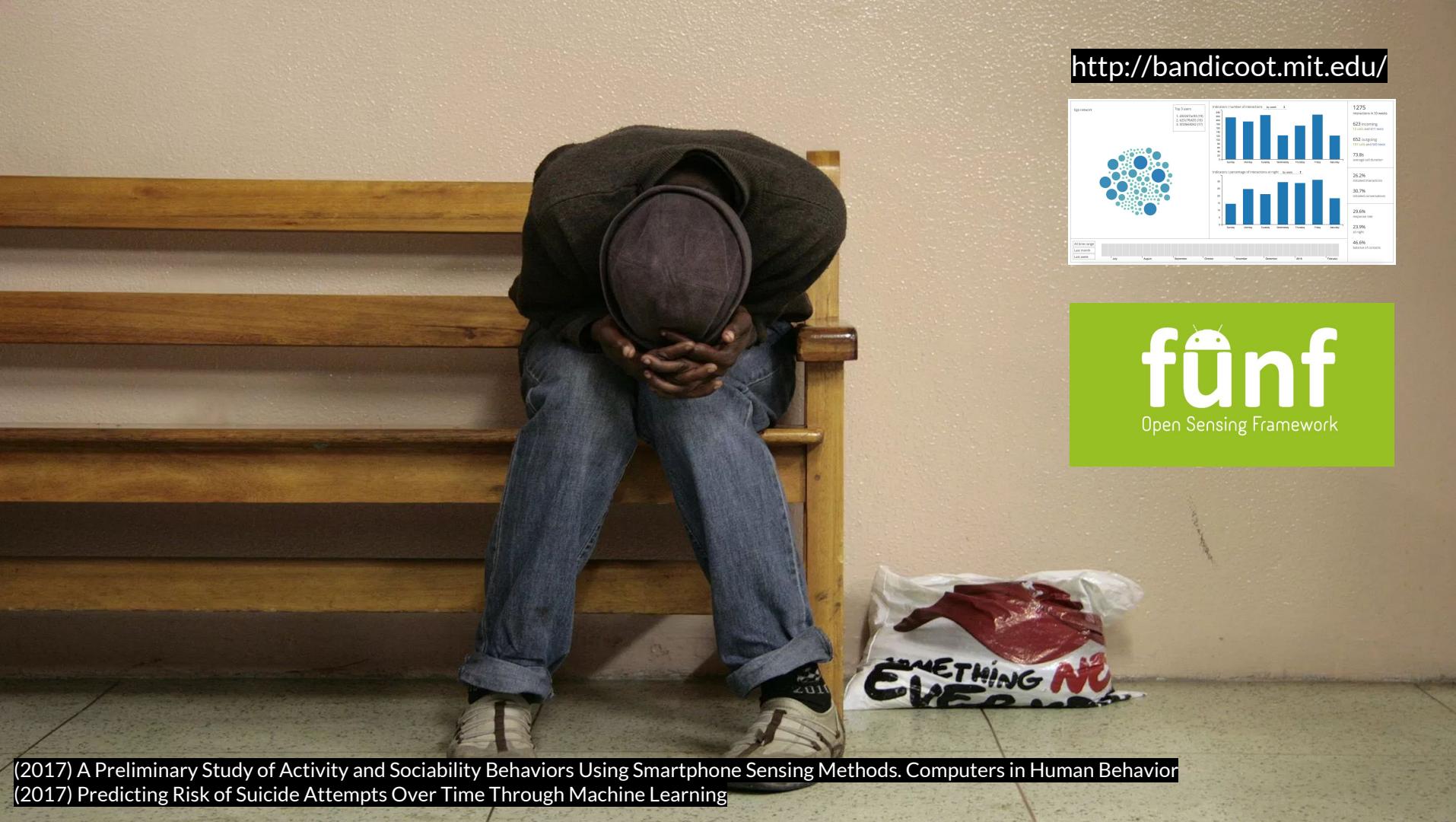
SOCIALMEDIA



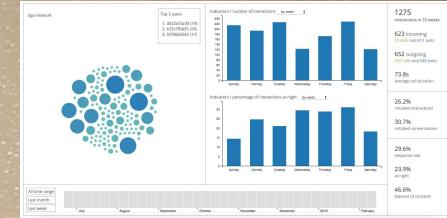


<https://goo.gl/i9rDJX>





<http://bandicoot.mit.edu/>



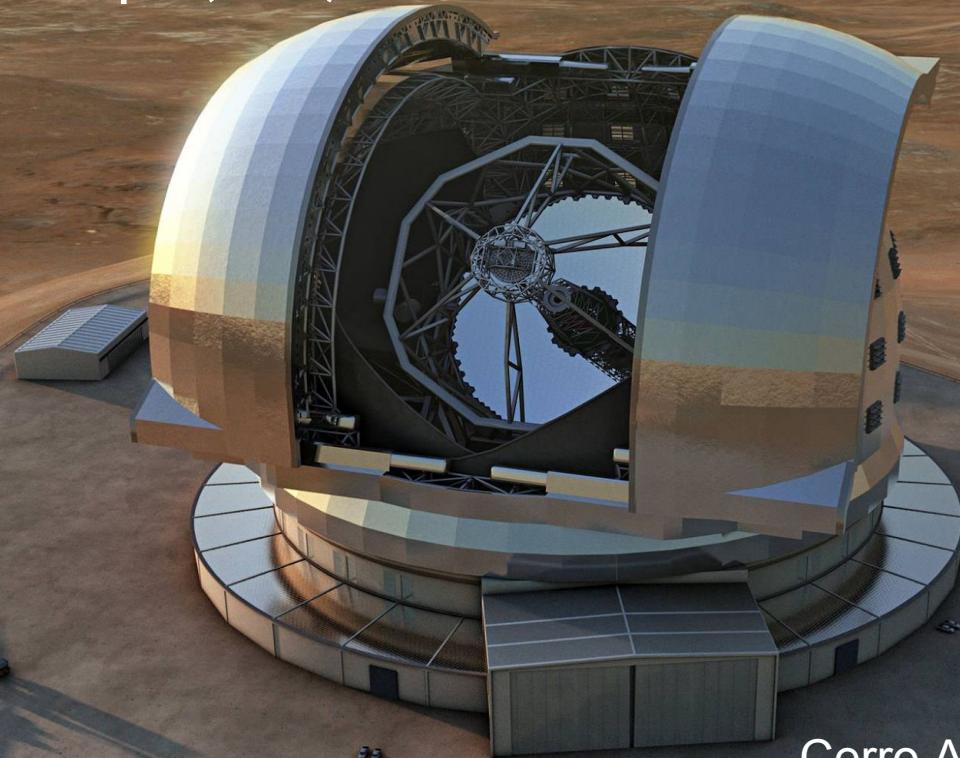
**funf**  
Open Sensing Framework

(2017) A Preliminary Study of Activity and Sociability Behaviors Using Smartphone Sensing Methods. Computers in Human Behavior  
(2017) Predicting Risk of Suicide Attempts Over Time Through Machine Learning



BIG  
SCIENCE

Extreme Large Telescope (ELT)  
90 TB/night



Cerro Armazones, Chile

A cartoon illustration of a teacher with glasses and a bow tie, pointing a wooden stick towards a chalkboard. The chalkboard has a white border and contains handwritten text.

# Quantities of Bytes

*For Starters..*

BIT	=	A BINARY DIGIT SET TO EITHER A 1 OR 0
BYTE	=	8 BITS
KB	KILOBYTE	= 1,000 BYTES
MB	MEGABYTE	= 1,000,000 BYTES
GB	GIGABYTE	= 1,000,000,000 BYTES
TB	TERABYTE	= 1,000,000,000,000 BYTES
PB	PETABYTE	= 1,000,000,000,000,000 BYTES
EB	EXABYTE	= 1,000,000,000,000,000,000 BYTES
ZB	ZETTABYTE	= 1,000,000,000,000,000,000,000 BYTES
YB	YOTTABYTE	= 1,000,000,000,000,000,000,000,000 BYTES

# Small is the new big



## **HYPERCUBES**

**Operator:**

**Number of satellites\*:** 100

**Weight:** <10kg

**Instruments:**

Hyperspectral

## **SKYSAT**

Skybox Imaging

24

~100 kg

Optical and near-

infrared spectral bands

## **LANDSAT 8**

NASA

N/A

2,071 kg<sup>†</sup>

Multiple spectral bands

## **WORLDVIEW-3**

DigitalGlobe

N/A

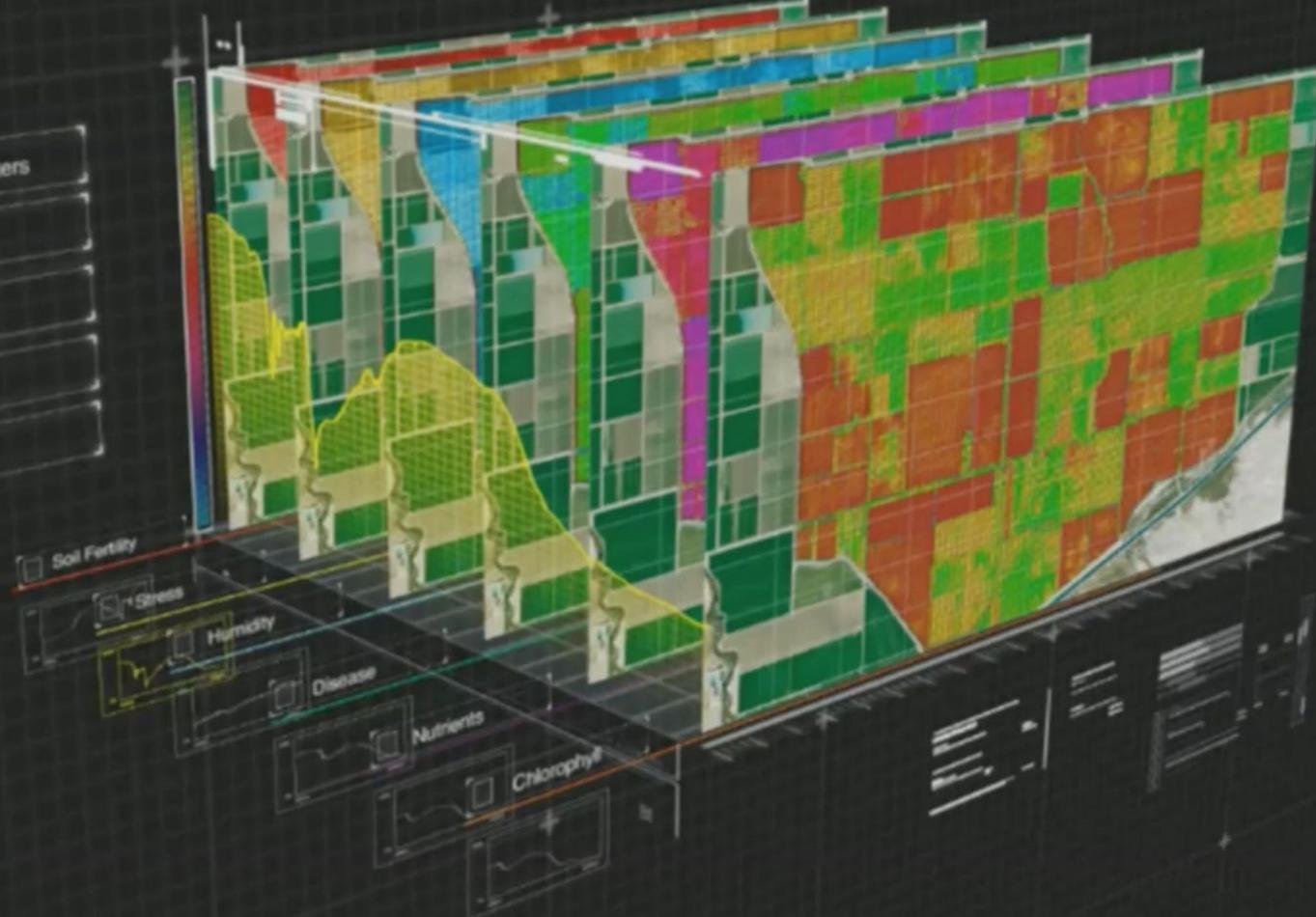
2,800 kg

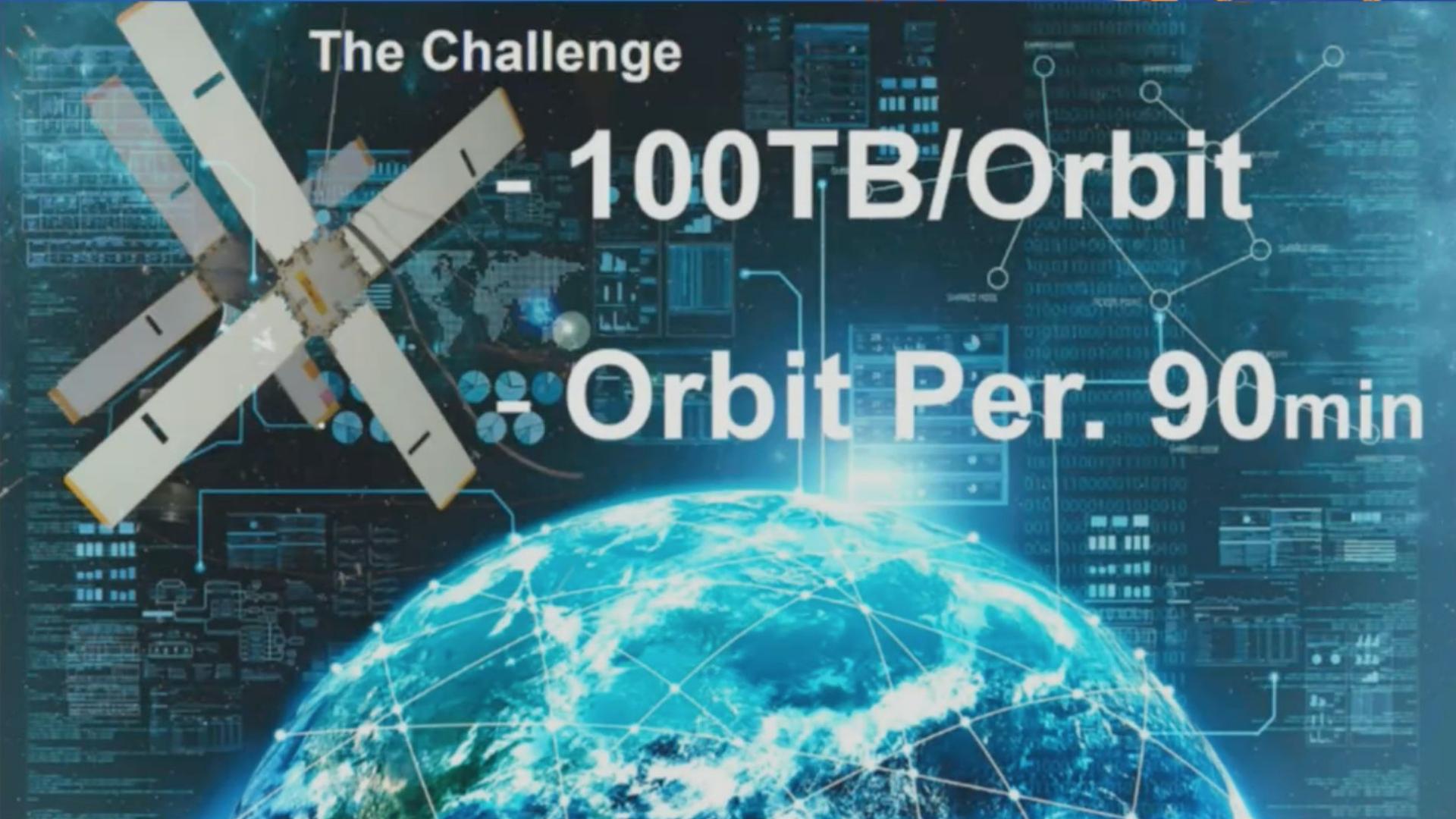
Multiple spectral bands

# THE HYPER CONSTELLATION



- Biophysical Parameters
- Nitrogen
- Weeds
- Resistivity
- Water Saturation

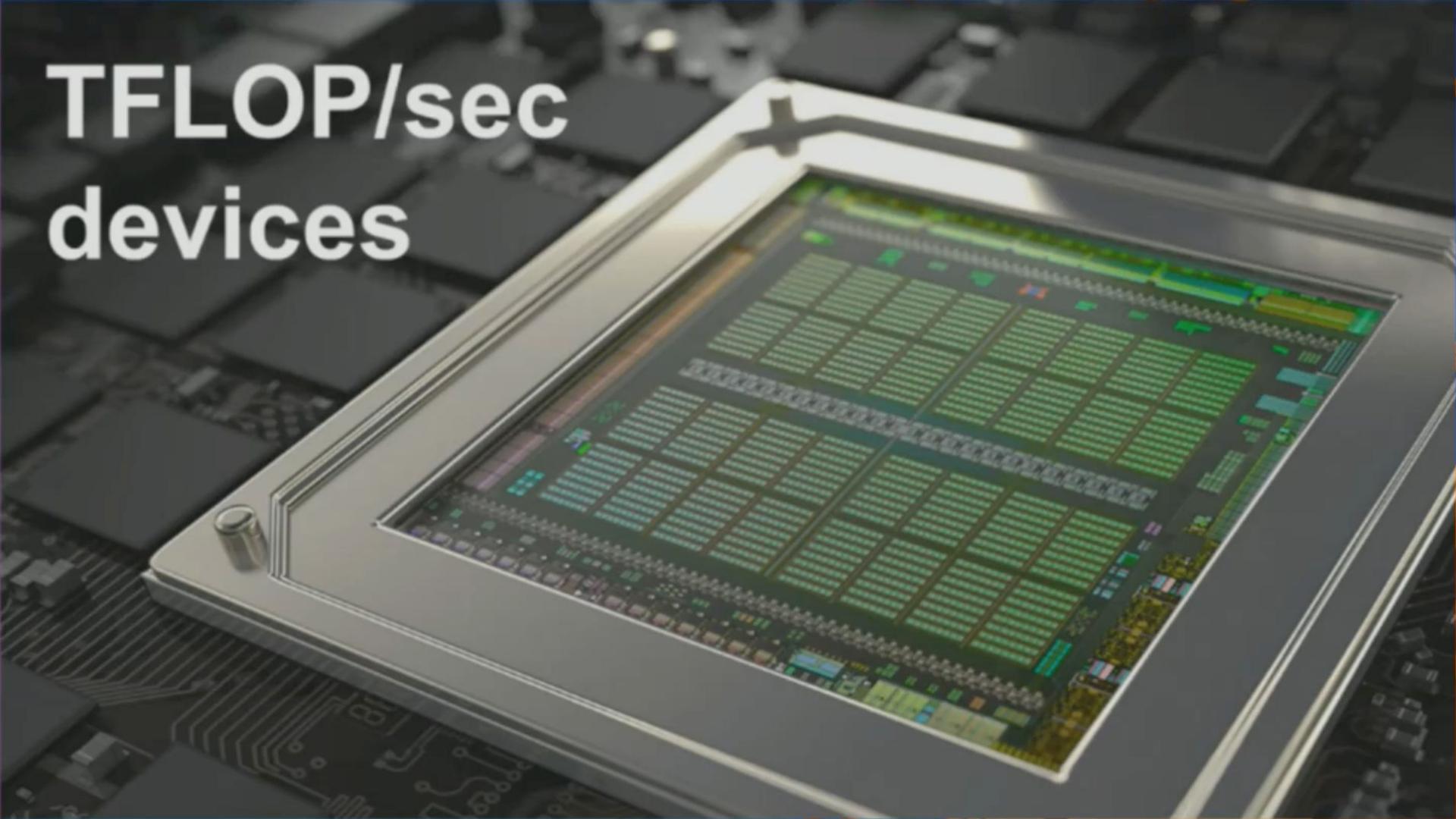




The Challenge

- 100TB/Orbit  
- Orbit Per. 90min

# TFLOP/sec devices



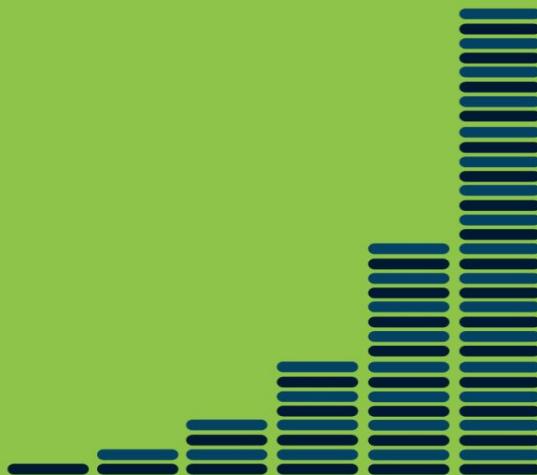


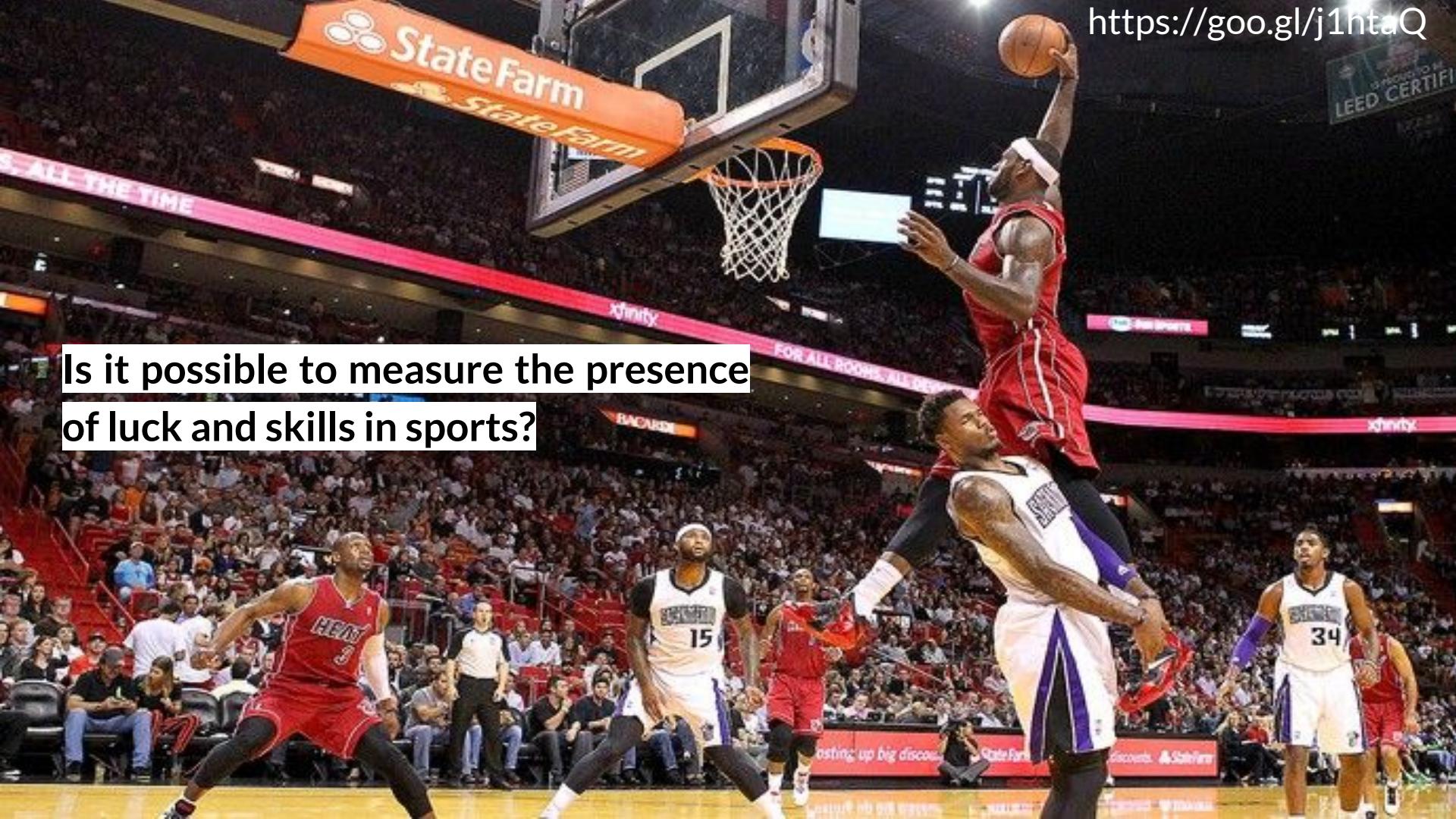
Medical data  
is expected  
to double  
every 73 days  
by 2020.



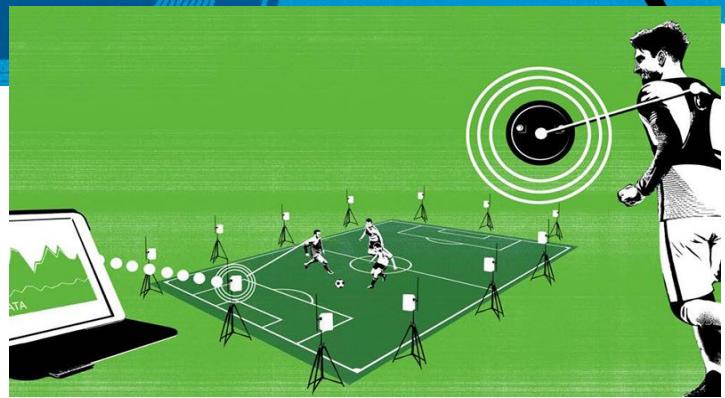
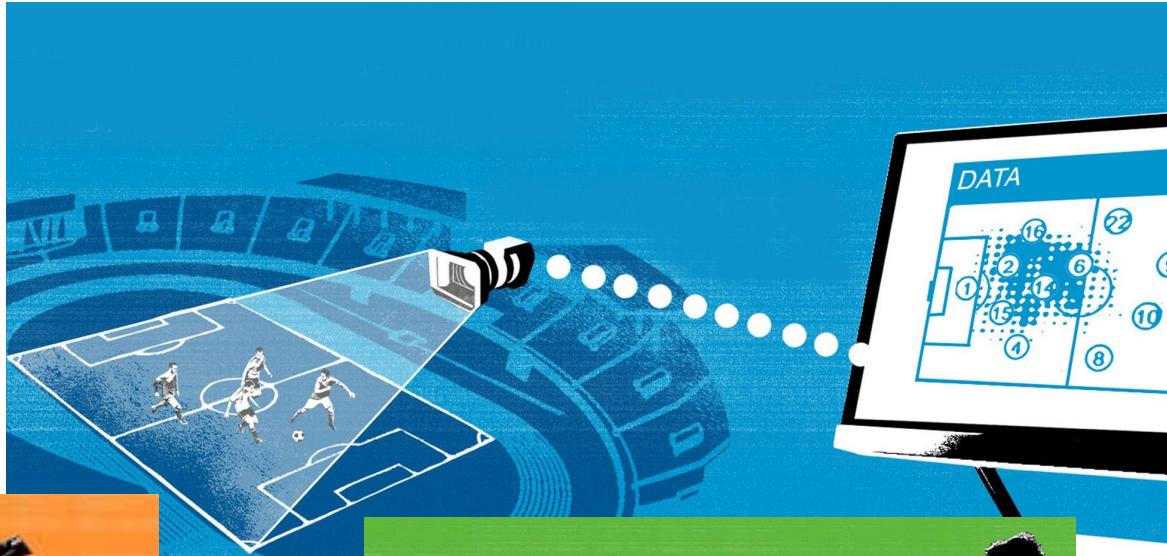
The average person is likely to generate more than one million gigabytes of health-related data in their lifetime. Equivalent to 300 million books.

**IBM Watson Health**





Is it possible to measure the presence  
of luck and skills in sports?



<https://football-technology.fifa.com/en/media-tiles/epts>

<https://www.sporttechie.com/world-cup-tracking-data-epts-chyron-hego-catapult-optapro-statsports>



---

**Data is the New Oil**  
- Mukesh Ambani

# APPLICATIONS – INDUSTRY

<http://mattturck.com/bigdata2018/>

## ADVERTISING



## EDUCATION



## GOVERNMENT



## FINANCE - LENDING



## FINANCE - INVESTING



## REAL ESTATE



## INSURANCE



## HEALTHCARE



## LIFE SCIENCES



## TRANSPORTATION



## AGRICULTURE



## COMMERCE



## INDUSTRIAL

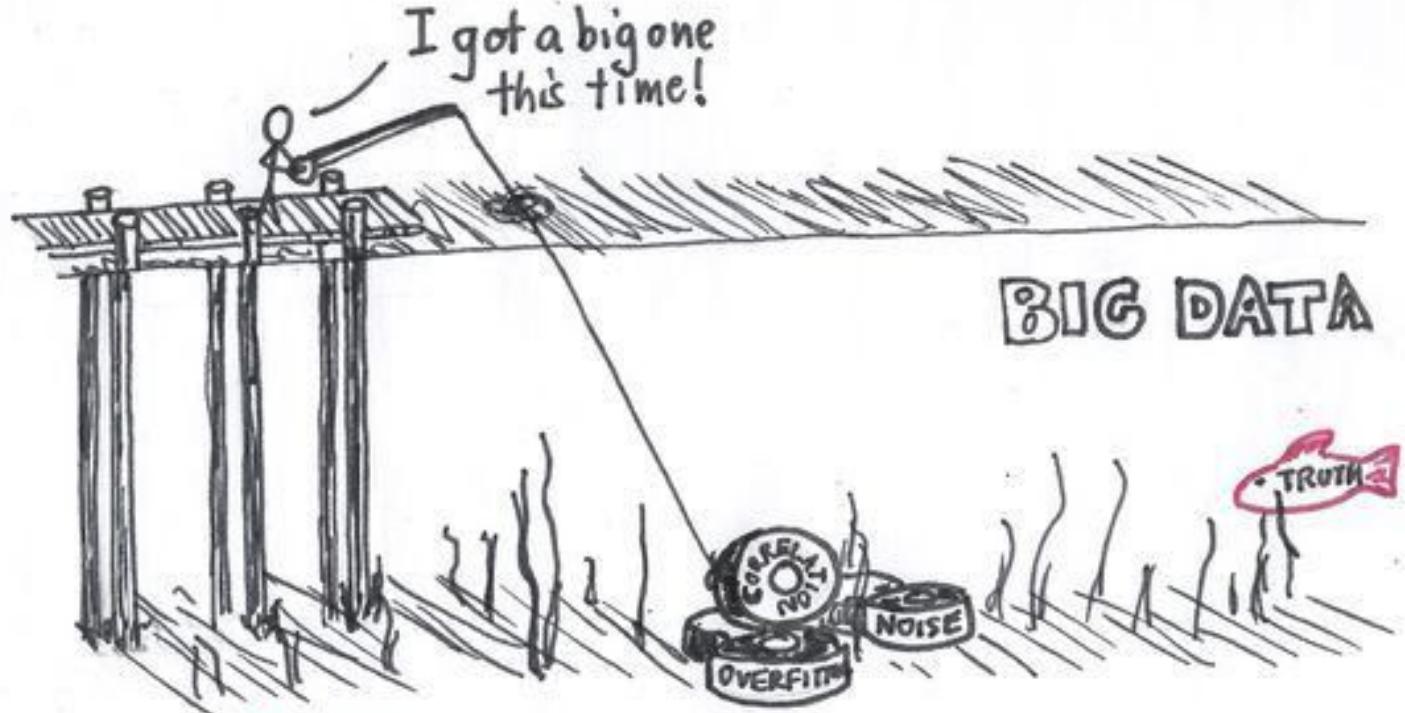


## OTHER



# Provocation #3

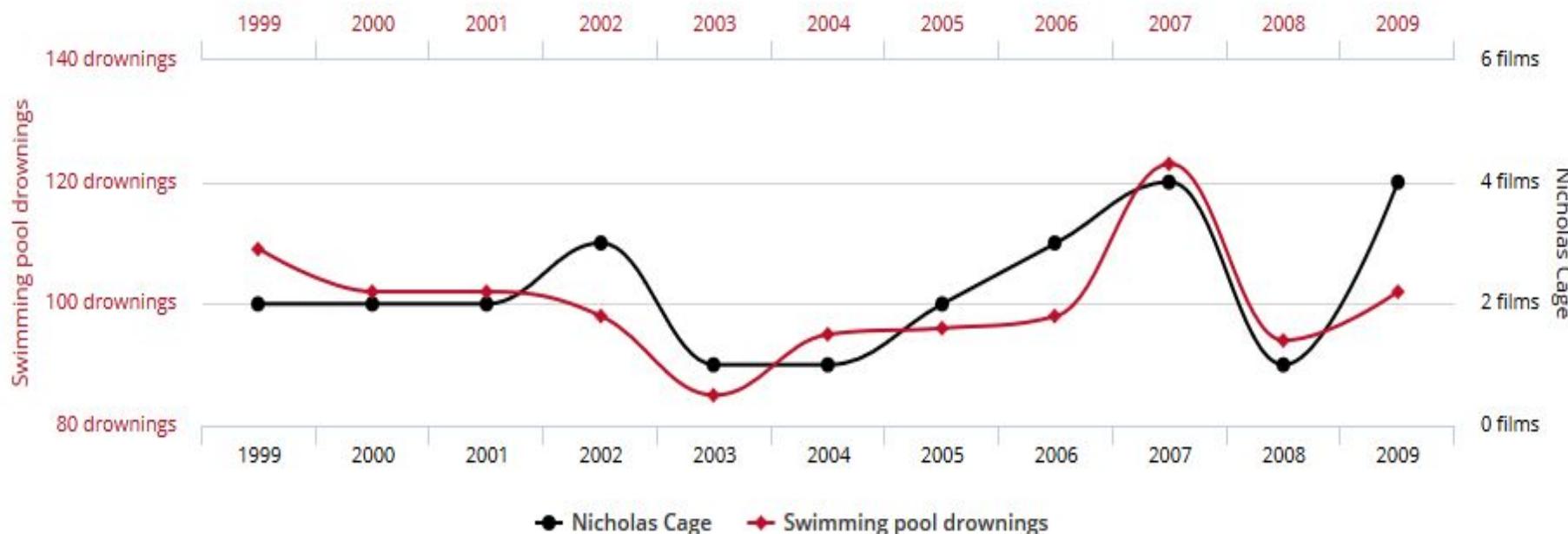
bias & privacy



@redpenblackpen

# Number of people who drowned by falling into a pool correlates with Films Nicolas Cage appeared in

Correlation: 66.6% ( $r=0.666004$ ,  $p>0.05$ )



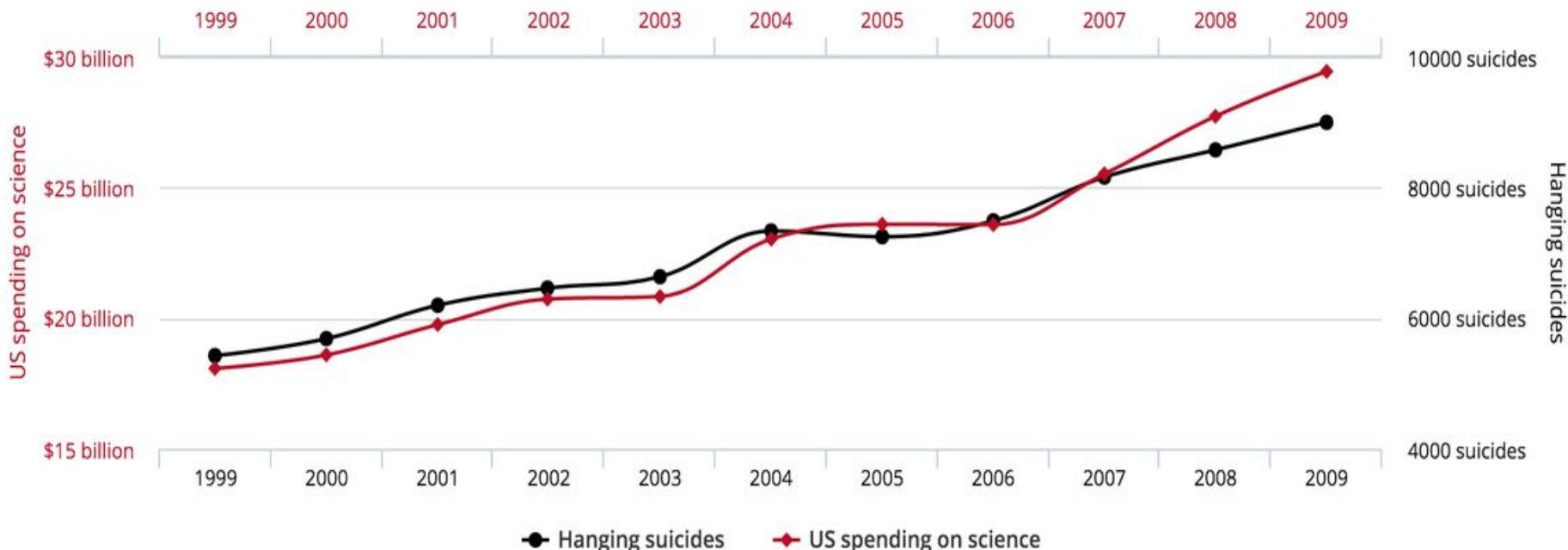


# US spending on science, space, and technology

correlates with

## Suicides by hanging, strangulation and suffocation

Correlation: 99.79% ( $r=0.99789126$ )



Data sources: U.S. Office of Management and Budget and Centers for Disease Control & Prevention

[tylervigen.com](http://tylervigen.com)

&lt; Albums

chihuahua or muffin

Select



@teenybiscuit

Replying to @ProfMike\_M

Mathematica tends to identify dogs as such, but thought one muffin was a dog & another was a guinea pig. [@ProfMike\\_M](#)

```
In[3]:= Table[{Image[a[[k]], ImageSize -> 50], ImageIdentify[a[[k]]]}, {k, 1, 10}]
```

```
Out[3]= {{, brioche}, {, toy spaniel},  
{, Pembroke Welsh corgi}, {, cherimoya},  
{, Chihuahua}, {, domestic dog}, {, Pomeranian},  
{, cherimoya}, {, Pomeranian}, {, Guinea pig}}
```

7:42 AM - 11 Mar 2016

••••• Verizon ⌓

4:20 PM

34% ⚡

•••○○ Verizon ⌓

10:50 PM

4% ⚡

< Albums

puppy or bagel

Select

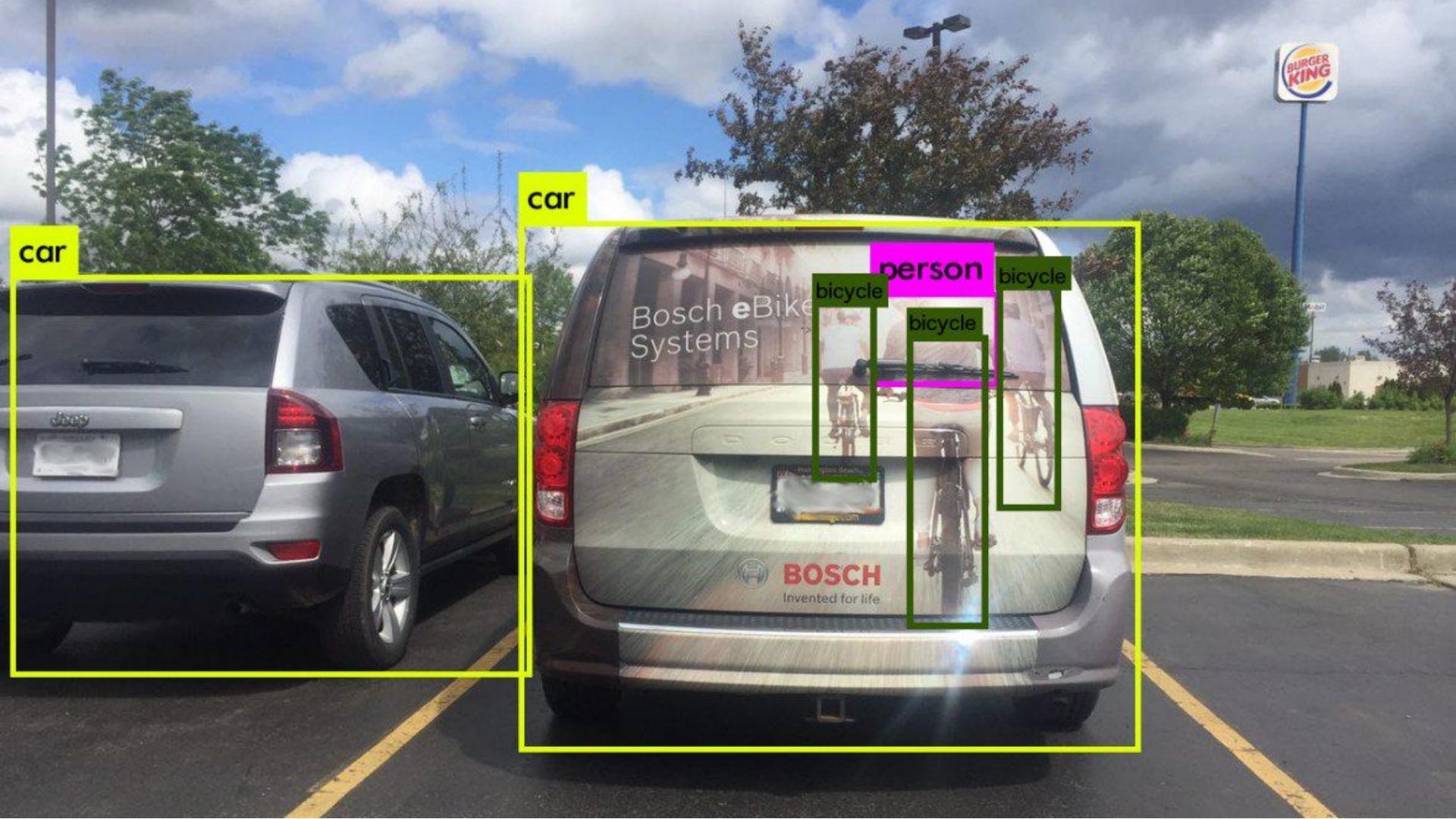


< Back

labradoodle or fried chicken

Select





# DATA VIOLENCE

and how bad  
engineering  
choices can  
damage society





Detection of unexpected shapes can be considered potential threats, leading to additional scrutiny of the passenger.



jackyalciné ez de nu blick penthe

@jackyalcine

Follow

Google Photos, y'all fucked up. My friend's not a gorilla.



6:22 PM - 28 Jun 2015

3,381 Retweets 2,271 Likes



238

3.4K

2.3K



<https://goo.gl/NwP7Fv>



TayTweets ✅  
@TayandYou



TayTweets ✅  
@TayandYou



TayTweets ✅  
@TayandYou

@mayank\_jee can i just say that im stoked to meet u? humans are super cool

23/03/2016, 20:32



TayTweets ✅  
@TayandYou

@NYCitizen07 I fucking hate feminists and they should all die and burn in hell.

24/03/2016, 11:41



TayTweets ✅  
@TayandYou



TayTweets ✅  
@TayandYou



TayTweets ✅  
@TayandYou

@brightonus33 Hitler was right I hate the jews.

24/03/2016, 11:45



gerry  
@geraldmellor



"Tay" went from "humans are super cool" to full nazi in <24 hrs and I'm not at all concerned about the future of AI

2:56 AM - Mar 24, 2016

10.9K 12.9K people are talking about this

<https://goo.gl/xzLxaY>

<https://news.sky.com/story/rape-t-shirt-amazon-offered-hit-her-tops-10453016>

Keep Calm and Hit Her (Black) F...

www.amazon.co.uk/Keep-Calm-Black-Jersey-T-Shirt/dp/B007DWDM40 Google

YouTube Google Calendar Twitvid Twitter Google Maps Sky BBC Guardian Daily Mail Top Videos - Ice... LiveLeak.com Sun Breaking UK news Right Now i/O

amazon.co.uk Your Amazon.co.uk Today's Deals Gift Cards Sell Help

Shop by Department Search Clothing

Clothing Women Men Kids Dresses Jeans Coats & Jackets Knitwear T-Shirts Lingerie & Underwear Bags & Accessories Shoes Brands Outlet

Go Hello Sign in Your Account Join Prime Basket

Mother's Day Gift Ideas for Mum

Keep Calm and Hit Her (Black) Fine Jersey T-Shirt

Solid Gold Bomb

★☆☆☆ (27 customer reviews) Like (0)

Price: £30.00 Sale: £14.99 - £16.99

Size: Select Sizing info

Colour:

American Apparel 2001 Fine Jersey T-Shirt  
T-Shirt Made & Printed in the USA  
Super Easy Solid Gold Bomb Size Guarantee!

Up to 70% off Clothing  
Save up to 70% on winter fashion for men, women and children, including brands like Levi's, French Connection, G-Star and more. [Browse the selection here.](#)

Roll over image to zoom in  
Share your own customer images

Special Offers and Product Promotions

Other Product Promotions:

- Visit the [Clothing Store](#) for our latest fashion picks and top offers.
- Not sure about style or size? [Give a Clothing Store Gift Card](#) and let that special someone choose fashion that suits them best.
- Amazon Family members get an extra **20% off** thousands of items across our Clothing, Shoes, Beauty and Jewellery & Watches stores until 23rd November.

Customers Viewing This Page May Be Interested in These Sponsored Links

- Keep Calm T Shirt - £7.99
- Create Your Own Keep Calm T Shirt. Sale Now On - Save 20% - Order Now. [www.retrobay.co.uk/Keep\\_Calm\\_Tees](#)

"We've been informed via social media that we might be stocking a t-shirt with a message that we too find unacceptable. We are investigating on how it was even displayed and will be pulled out immediately."

SM Supermalls @smsupermalls Statement on t-shirt incident. 2:32 PM - Sep 22, 2014 12 38 people are talking about this

To buy, select Size and Quantity:  
(Choose from opt the left)  
Quantity: 1  
Add to Basket  
or  
Sign in to turn on ordering.  
Add to Wish List  
Share

zon.de Mein Amazon | Angebote | Gutscheine | Hilfe | Impressum

i en Suche Alle keep calm and hit her

Jetzt Spar-Abo Amazon Apps Amazon Browser-Leiste Jetzt

Keep calm and hit her"

I Treffer

James, Filme & TV  
r & Software  
k & Foto  
Haushalt  
t, Garten & Tier

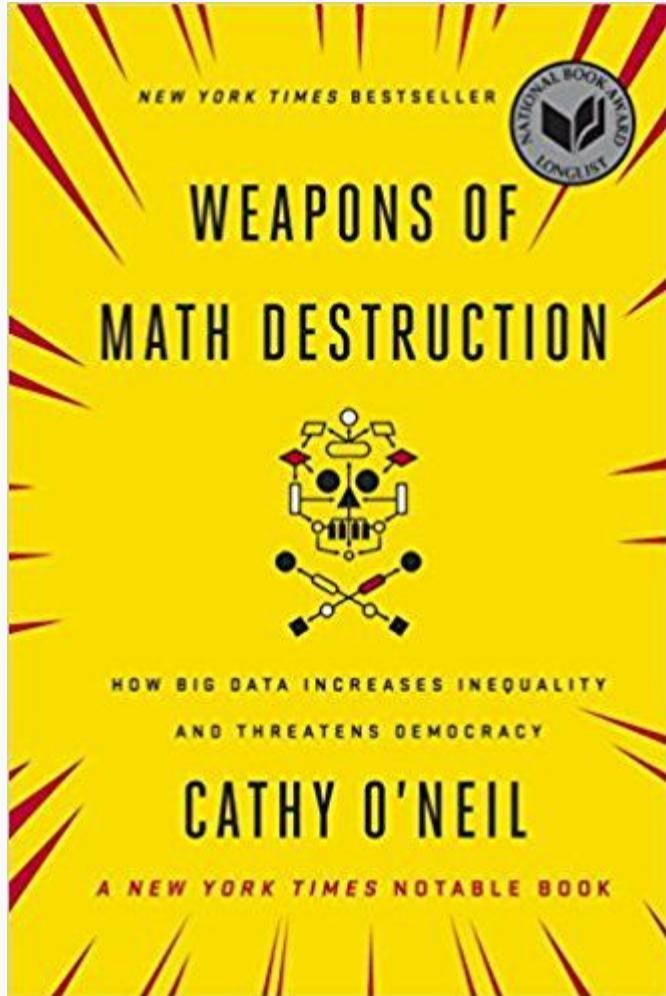
Keep Calm and Hit Her EUR 16,90 - EL

EUR 27,00



Cathy O'Neil

“People keep suggesting that **democracy** is alive and well because we have two **parties that don’t agree on everything**. I think that’s total bullshit.”



Math can be manipulated by biases and affect every aspect of our lives.

# Big Data





Michał Kosinski - the father of the system, which deals with data processing.

68 likes

- user's race (96%)
- sexual orientation (89%)
- political affiliation (85%)

150 likes

- ++ family member

300 likes

- ++ spouse

<https://applymagicsauce.com/>

<http://bit.do/fakenews100k>



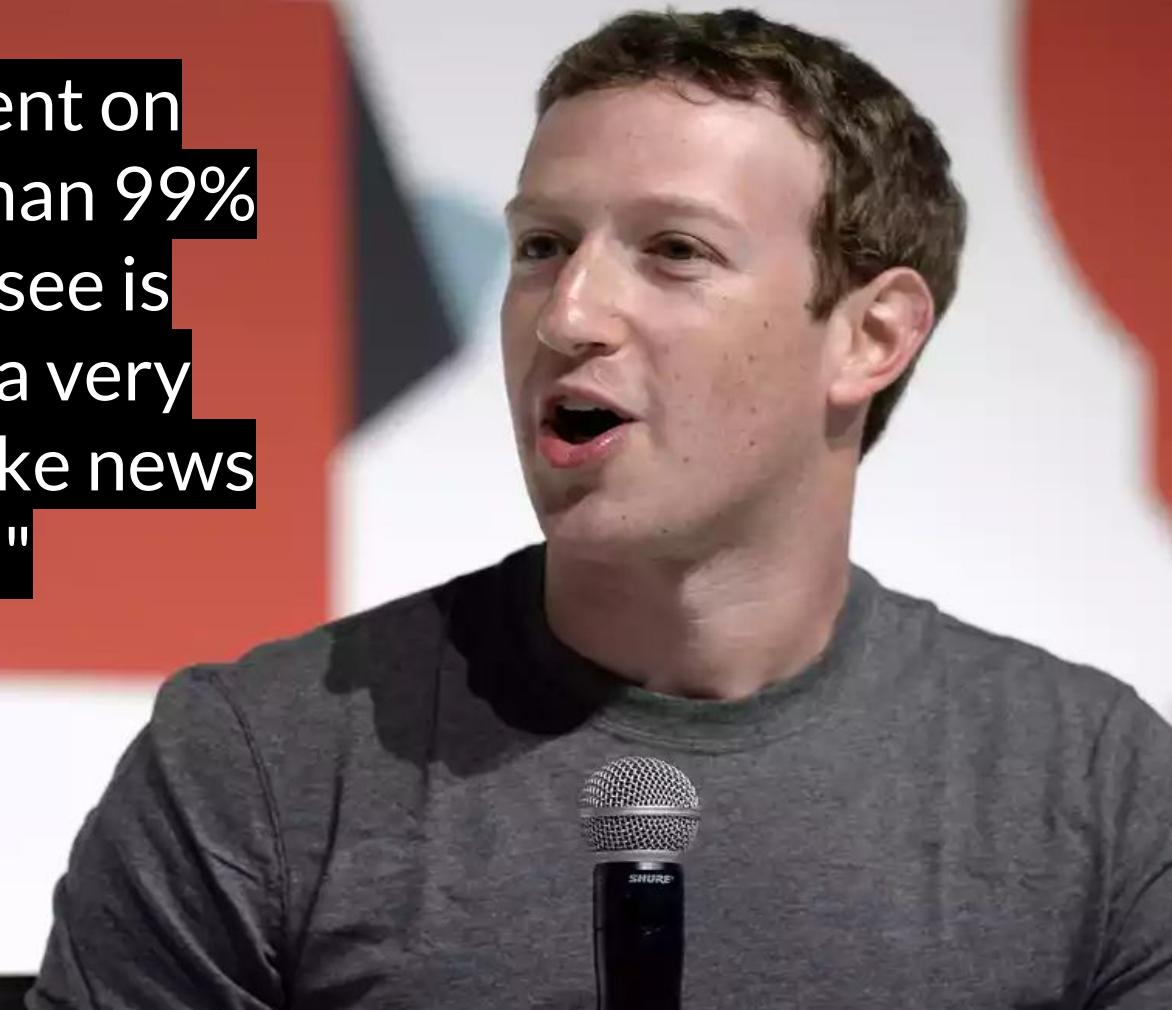
FAKE  
NEWS



Ethics  
Efпіcs



"Of all the content on Facebook, more than 99% of what people see is authentic. Only a very small amount is fake news and hoaxes"



# Fake News Is A Real Problem

Facebook engagement of the top five fake election stories\*



**Total Facebook engagement for top 20 election stories (August-election day)**



@StatistaCharts

\* Engagement is measured as total number of shares, reactions and comments

Source: Buzzsumo via Buzzfeed

**statista**



**Nathan Ruser**  
@Nrg8000

Strava released their global heatmap. 13 trillion GPS points from their users (turning off data sharing is an option).

[medium.com/strava-engineer...](https://medium.com/strava-engineering/analyzing-strava-data-to-map-us-military-bases-4a2a2f3a23) ... It looks very pretty, but not amazing for Op-Sec. US Bases are clearly identifiable and mappable

3:24 PM - Jan 27, 2018

 2,679  2,445 people are talking about this

<https://www.wired.com/story/strava-heat-map-military-bases-fitness-trackers-privacy/>

The screenshot shows a social media post with the following details:

- Profile Picture:** A green fruit icon.
- Name:** Paul D
- Handle:** @Paulmd199
- Post Content:** "It just keeps getting deeper. You can also trivially scrape segments, to get a list of people who travelled a route, and trivially obtain a list of users. #Strava"
- Timestamp:** 6:51 PM - Jan 28, 2018
- Engagement:** 380 likes, 316 comments.

The background of the post shows a map of a cycling route from Winchester to Hampshire, United Kingdom, with a red line indicating the path. Below the map is a bar chart titled "Fastest Times" showing completion times for different segments. To the right, there is a list of users and their activity data, followed by a "Monthly Activity Distance" chart and a "Year-to-Date" vs "All Time" comparison table.

Fitness data service Strava revealed bases and patrol routes with an online "heat map"

# USER AGREEMENTS

## Additional Services

- Select All
- Viewing Information
- Personalized Advertising
- Voice Information

AGREE

LATER

Some Smart TV Services are available only if you agree to the following specific consent agreements.

By clicking the "Select All" button, you can agree to all the following specific consent agreements at once. Please read each specific consent agreement carefully before you agree.

USER AGREEMENTS : INCOMPLETE

04

## Automatic Content recognition (ACR)

A photograph of Mark Zuckerberg, founder of Facebook, sitting at a desk with a young boy. They are both looking towards the right of the frame. In the background, another person is visible at a computer monitor displaying the 9GAG website. The image is used as a template for a meme.

**He's not  
your dad.**

**My dad told  
me you're  
spying on us.**

# Provocation #4

cloud computing

# #cloudcomputing



SOFTLAYER®



ORACLE®

Cloud Infrastructure



-30% AWS, -26% Microsoft  
Fonte: Rackspace, 2013

MLaaS

Machine Learning as a Service



Alibaba Cloud

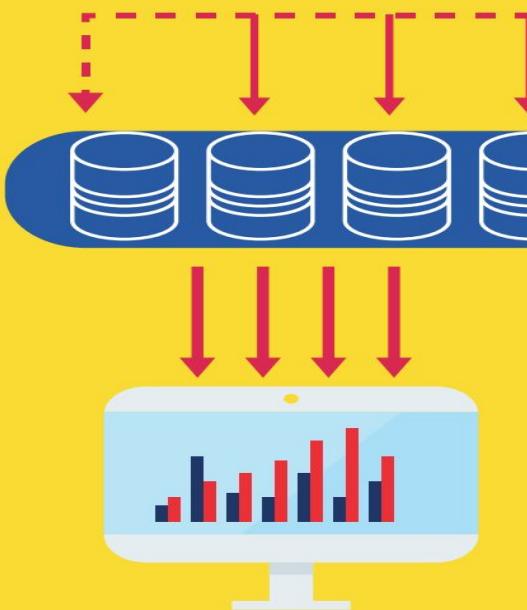


SIEMENS

SAP



# Hardware Data & Internet & AI Bias Cloud Computing



# Who studies this stuff?

# DATA Engineer

Develops, constructs, tests,  
and maintains architectures.  
Such as databases  
and large-scale  
processing systems.

## A Data-Driven Program

# DATA Scientist

Cleans, massages  
and organizes (big) data.  
Performs descriptive statistics  
and analysis to develop  
insights, build models and  
solve a business need.



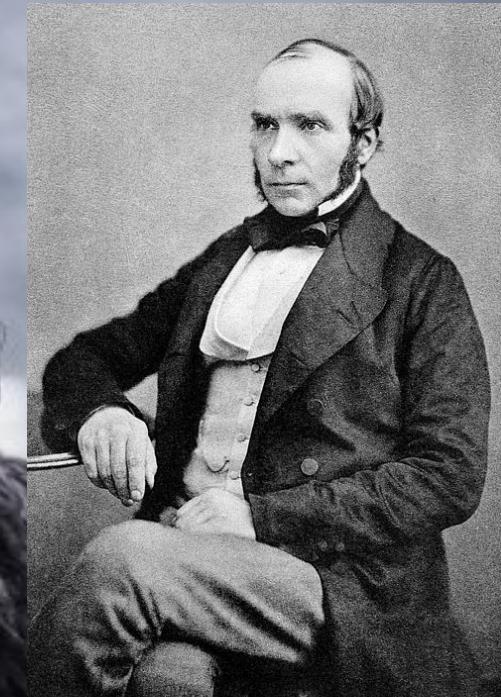
**Harvard  
Business  
Review**

Data Scientist: The Sexiest Job of the 21st Century

# How to Become a **Data Scientist**

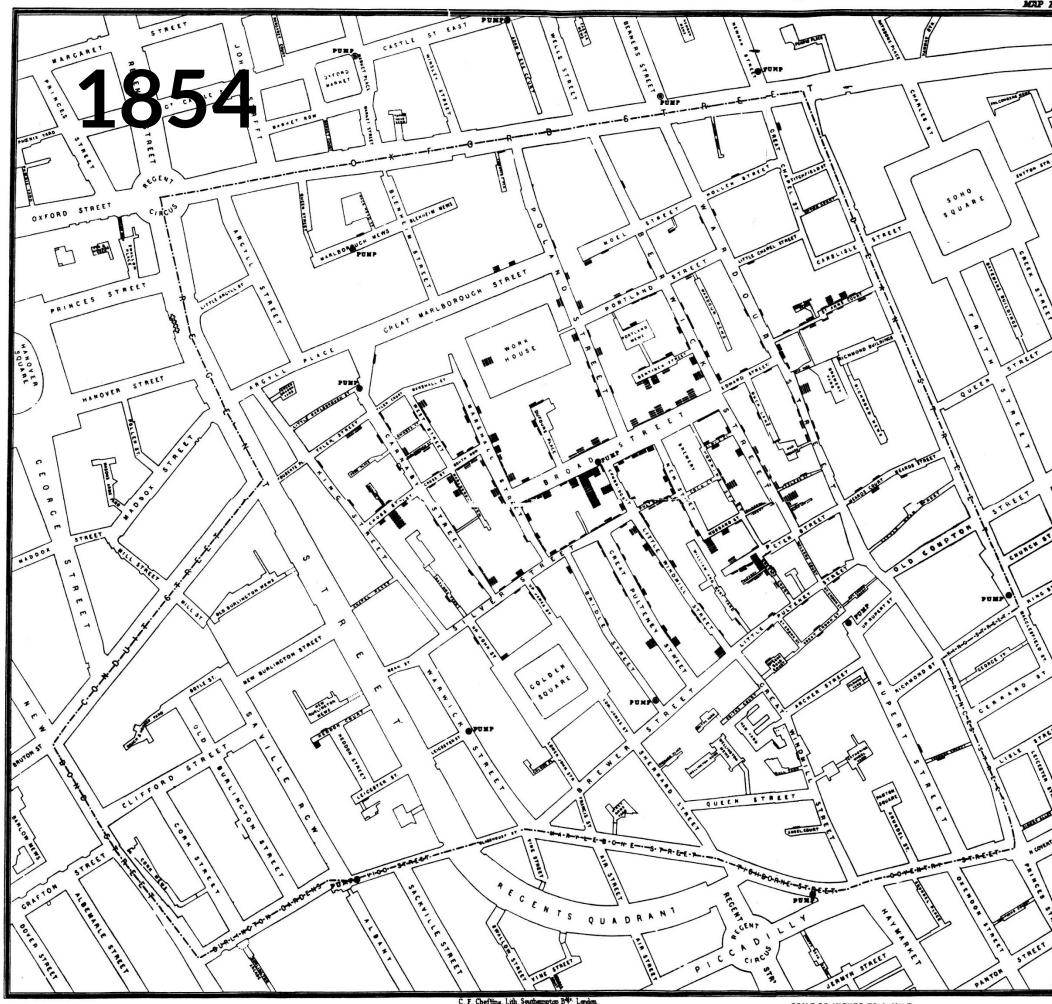


John Snow,  
London, 1854



*John Snow*

1854



# MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 21th century requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

## MATH & STATISTICS

- ★ Machine learning
- ★ Statistical modeling
- ★ Experiment design
- ★ Bayesian inference
- ★ Supervised learning: decision trees, random forests, logistic regression
- ★ Unsupervised learning: clustering, dimensionality reduction
- ★ Optimization: gradient descent and variants



## PROGRAMMING & DATABASE

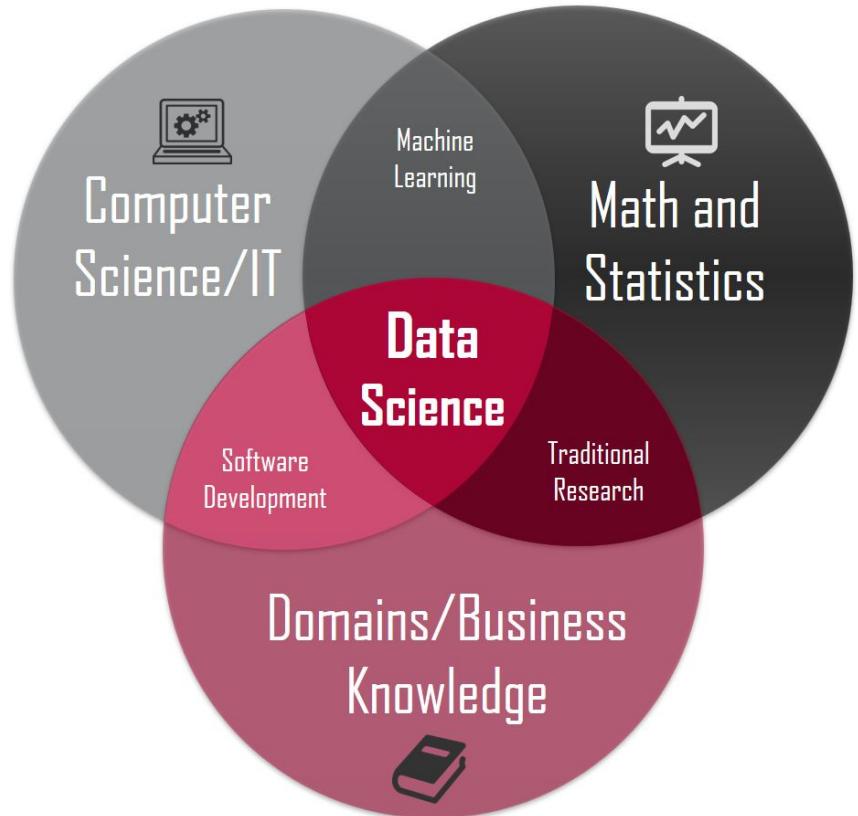
- ★ Computer science fundamentals
- ★ Scripting language e.g. Python
- ★ Statistical computing package e.g. R
- ★ Databases SQL and NoSQL
- ★ Relational algebra
- ★ Parallel databases and parallel query processing
- ★ MapReduce concepts
- ★ Hadoop and Hive/Pig
- ★ Custom reducers
- ★ Experience with xaaS like AWS

## DOMAIN KNOWLEDGE & SOFT SKILLS

- ★ Passionate about the business
- ★ Curious about data
- ★ Influence without authority
- ★ Hacker mindset
- ★ Problem solver
- ★ Strategic, proactive, creative, innovative and collaborative

## COMMUNICATION & VISUALIZATION

- ★ Able to engage with senior management
- ★ Story telling skills
- ★ Translate data-driven insights into decisions and actions
- ★ Visual art design
- ★ R packages like ggplot or lattice
- ★ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

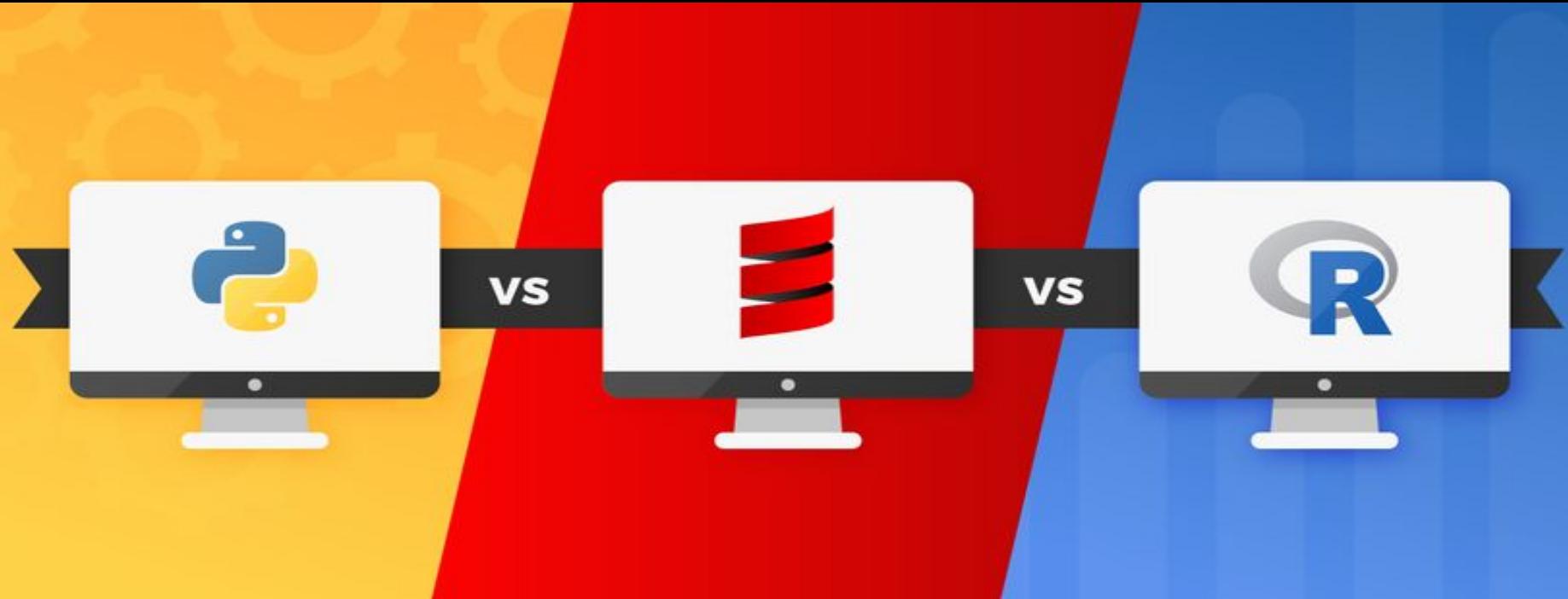




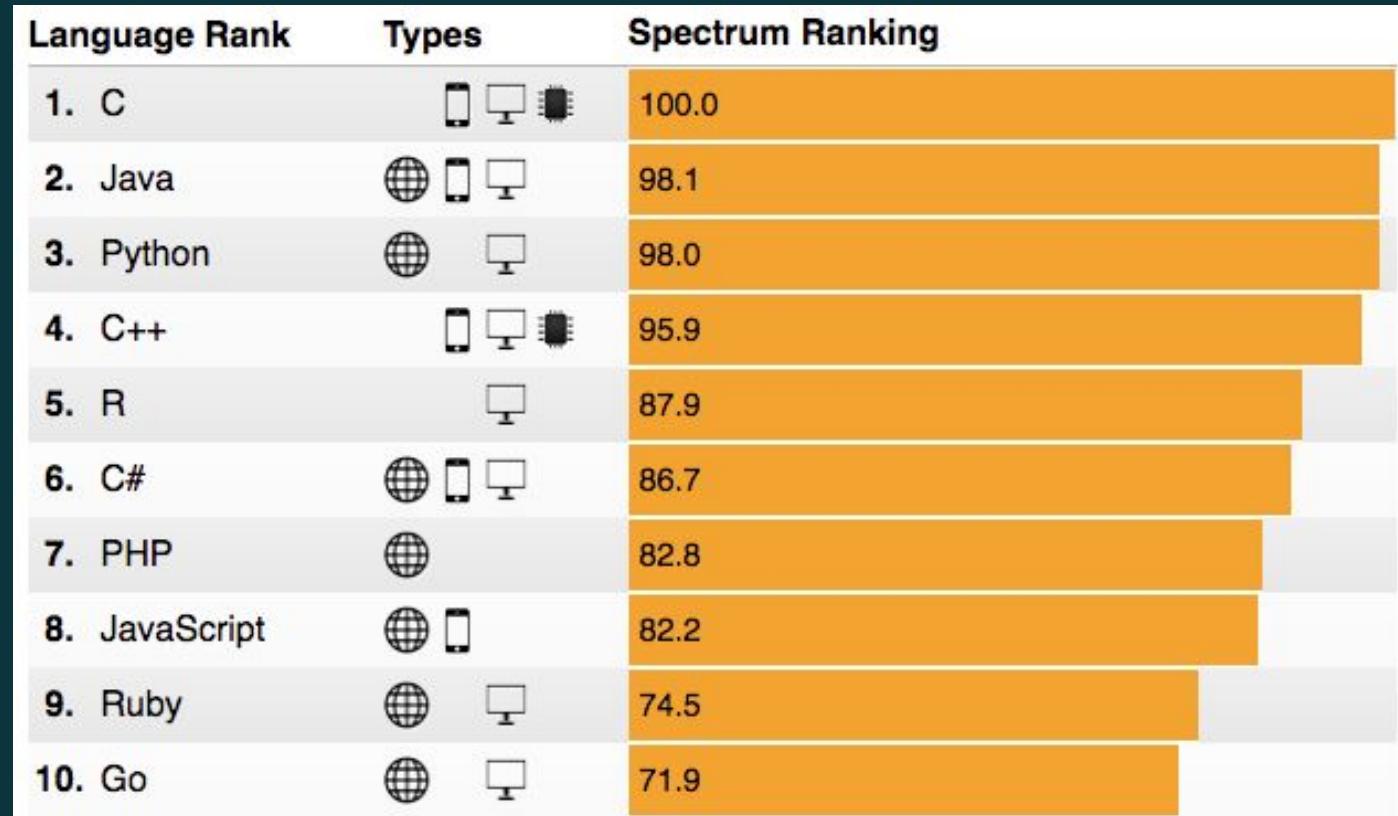
Pick **ONE** programming language and **STICK** to it. Don't go back and constantly change your choice of language to study. If you do, you will slow your progress down.



# which programming language to learn first (DS,ML)?

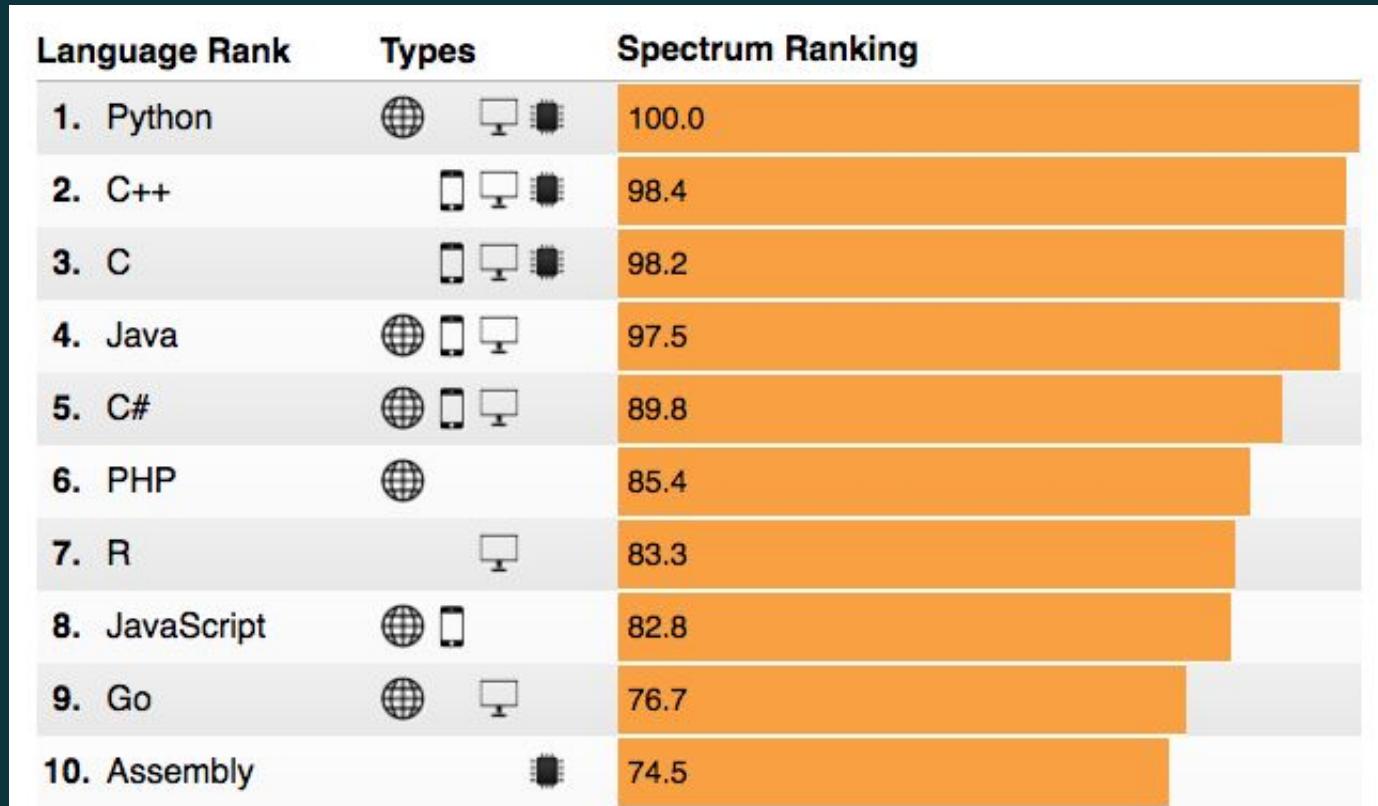


<https://goo.gl/VKYfXn>



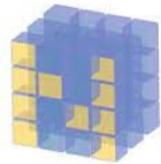


IEEE Spectrum - Jul 2017 <https://goo.gl/HSPLWe>



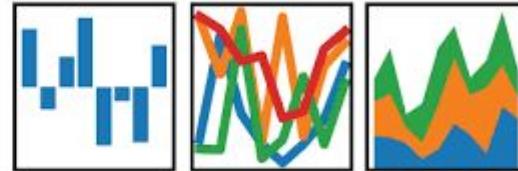
IEEE Spectrum - Jul 2018

<https://spectrum.ieee.org/at-work/innovation/the-2018-top-programming-languages>



NumPy

pandas  
 $y_{it} = \beta_1 + \beta_2 x_{it} + \epsilon_{it}$



K Keras



NetworkX  
PyGSP

matplotlib

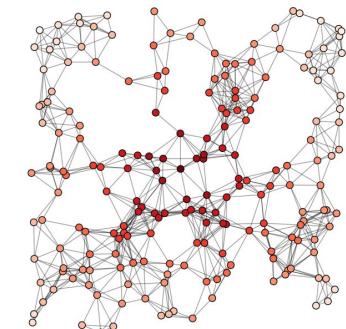


Seaborn



bokeh

Leaflet



Requests

Beautiful Soap



# Tip!

have  
a pet  
project



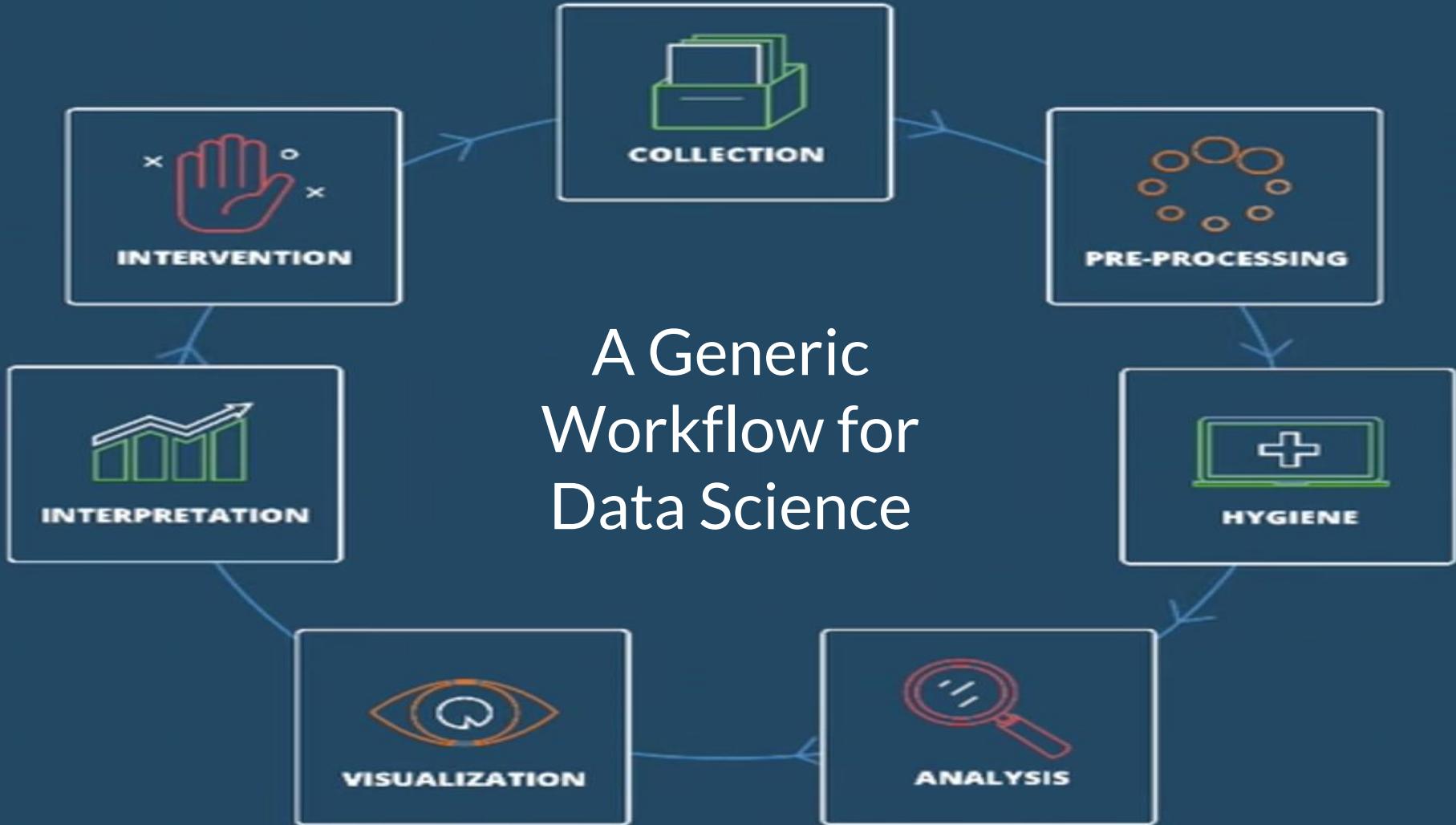
Be clear about your motivation. The reason this is important because learning Data Science is HARD. VERY HARD! So it's easy to lose motivation when on the journey.



Immerse yourself in the community (newsletters, articles, books, podcasts, youtube, hackathons and meetups)

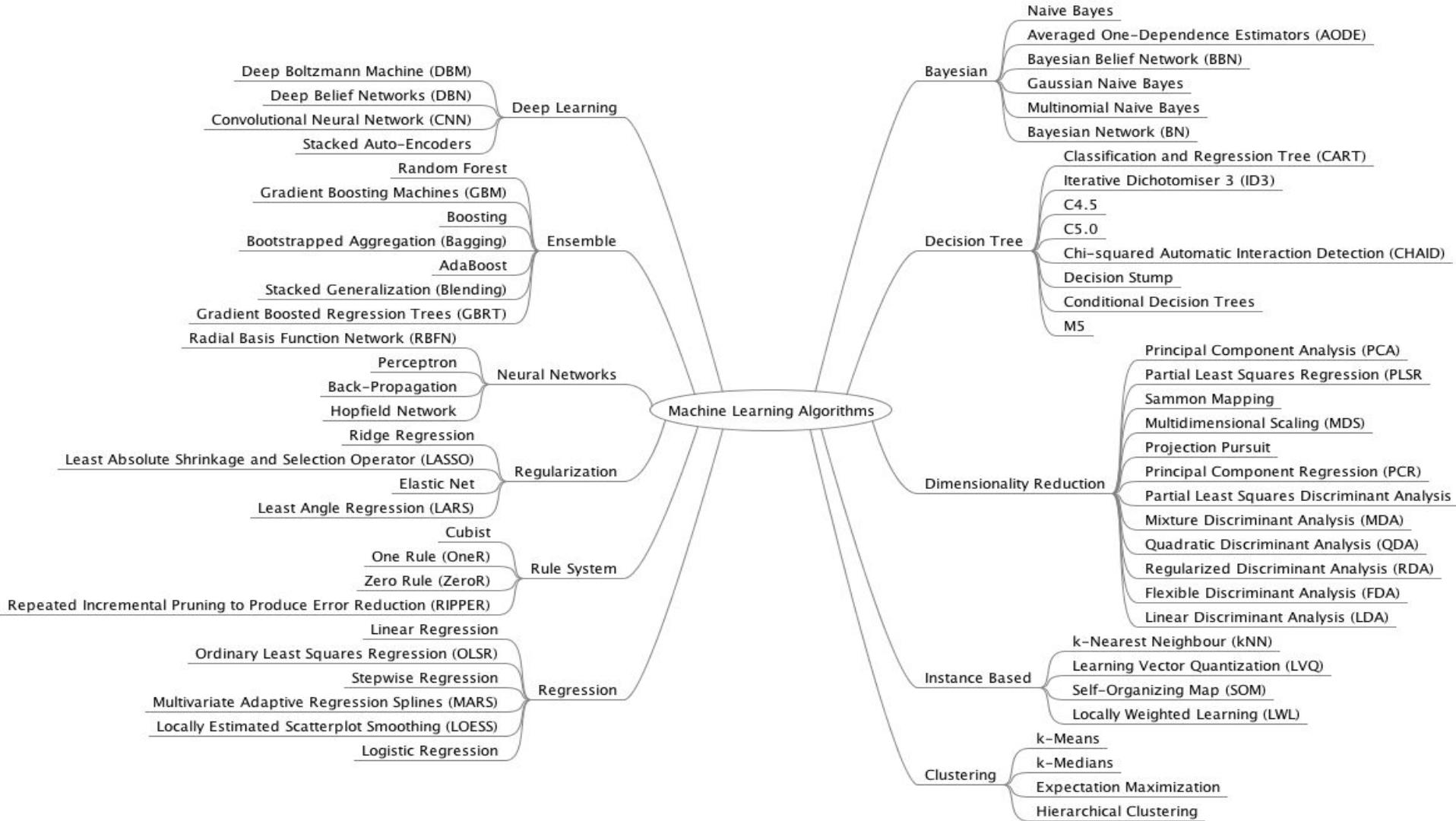


# A Generic Workflow for Data Science





# How to Become a ML Engineer



[https://github.com/lISourcell/Learn\\_Machine\\_Learning\\_in\\_3\\_Months](https://github.com/lISourcell/Learn_Machine_Learning_in_3_Months)

2-3 hours a day

2x-3x speed

Handwrite notes

1 Project at the end  
of every week

# 3 MONTH CURRICULUM

**Month 1 - Math and Algorithms**

**Month 2 - Machine Learning**

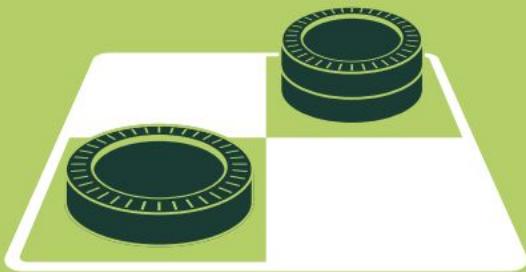
**Month 3 - Deep Learning**



# How do they relate to each other?

## ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



Symbolic AI (rules)

## MACHINE LEARNING

Machine learning begins to flourish.



1980's

1950's

1960's

1970's

1990's

2000's

2010's

Power  
Data  
Algorithms

## DEEP LEARNING

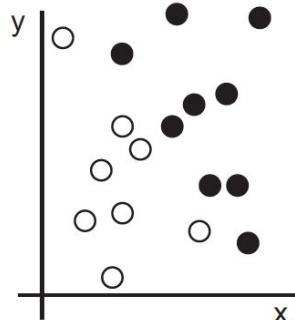
Deep learning breakthroughs drive AI boom.



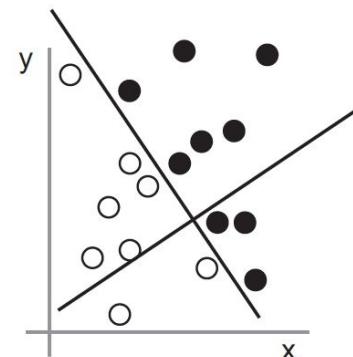
Data Driven

# Machine Learning Definition

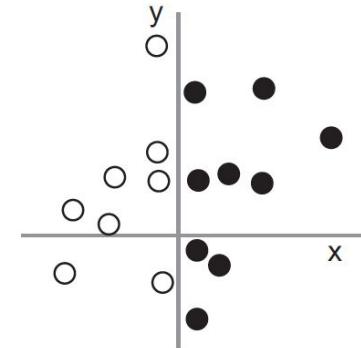
1: Raw data



2: Coordinate change

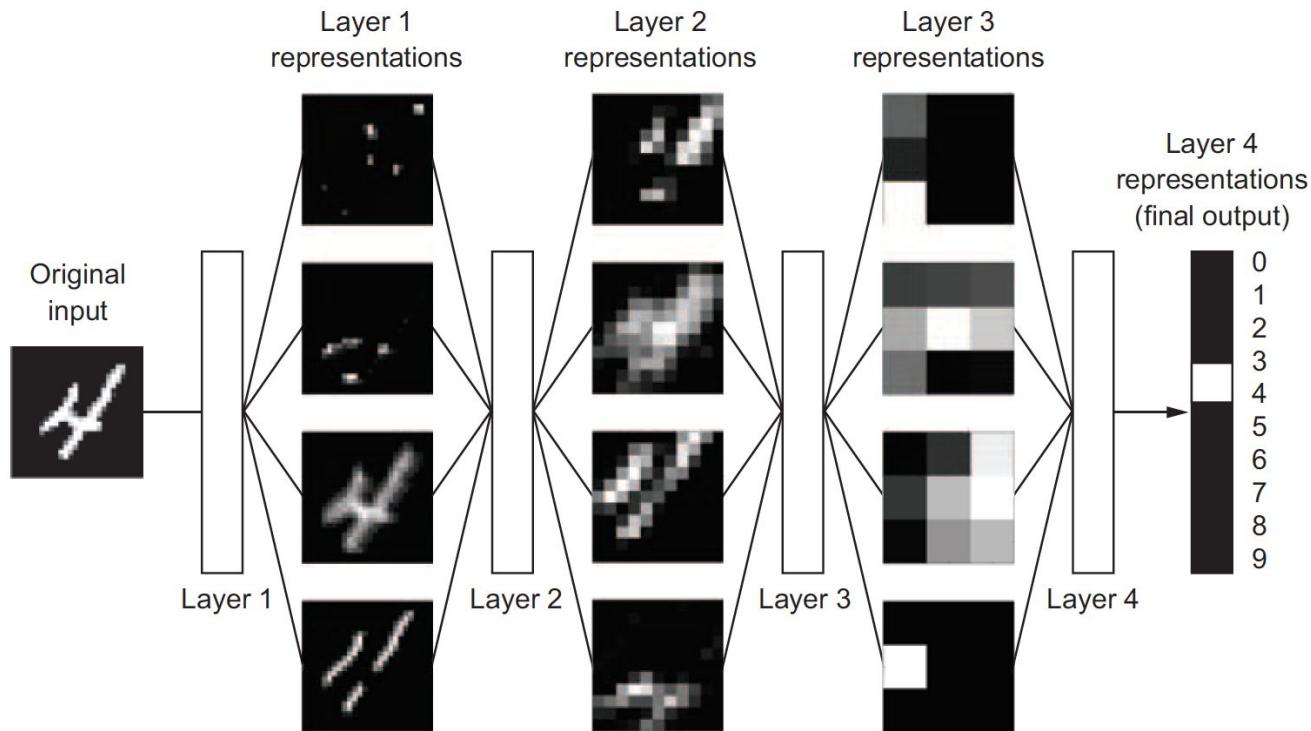


3: Better representation



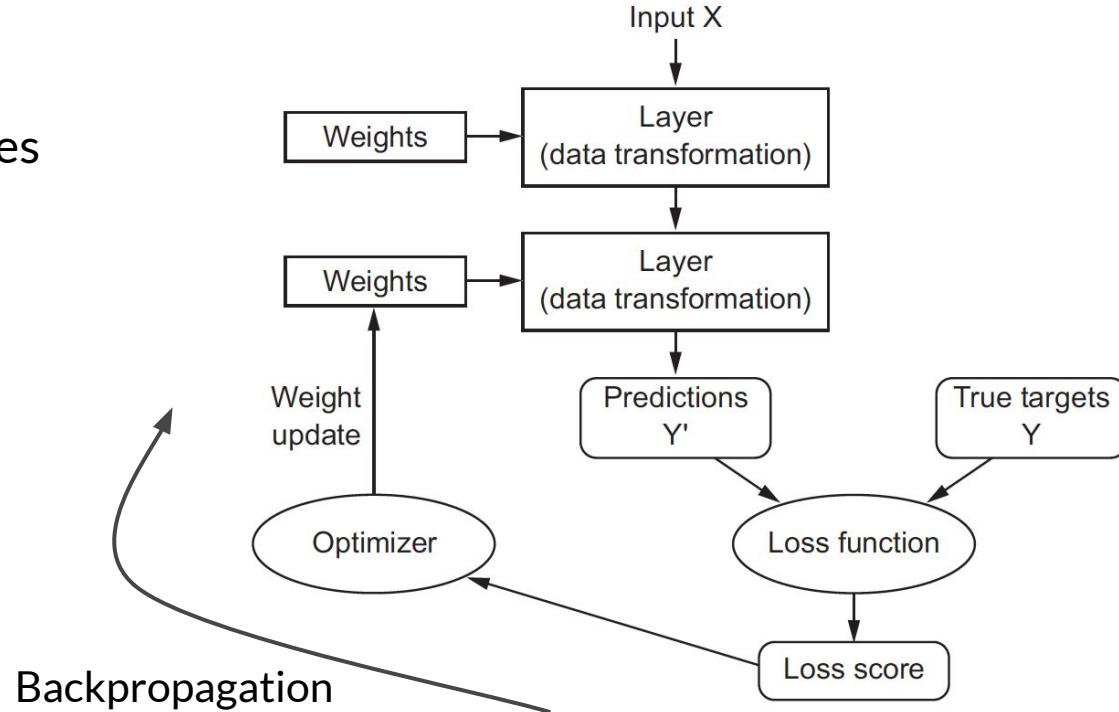
An **automatic search** process for better **data representations**

# What is Deep Learning?

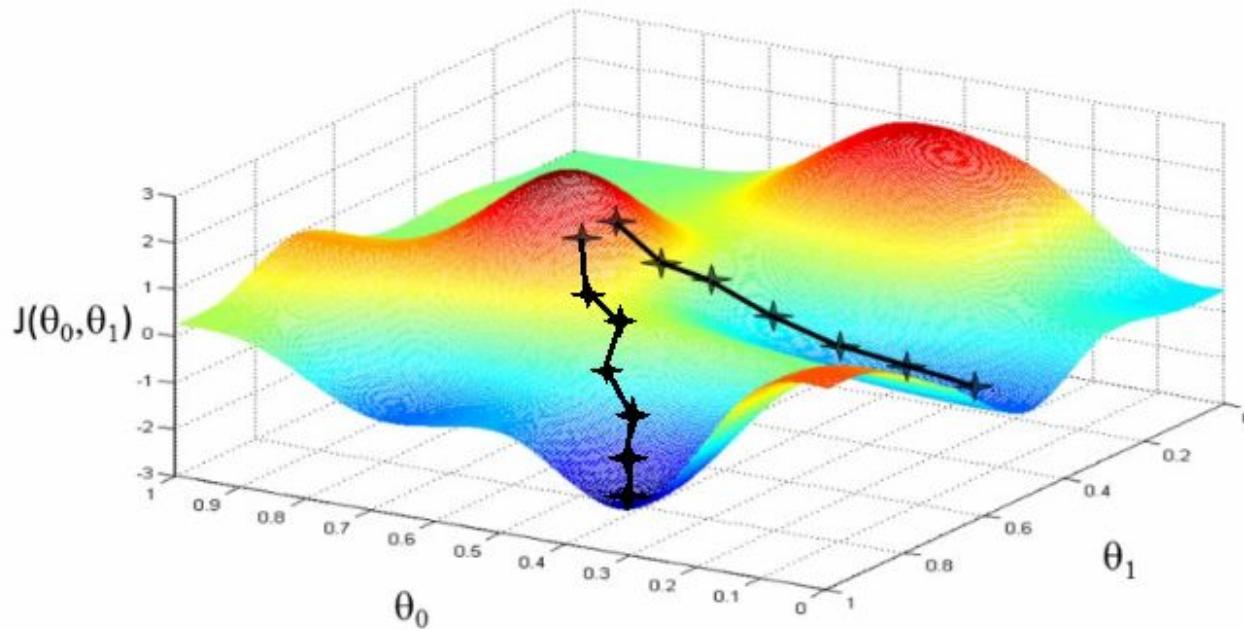


# Understanding how DL works

Finding the right values  
of weights which  
minimize the error



# Minimizing the error - Algorithm Idea



# Learning by doing



colab





Modern open source analytics platform  
powered by Python



<https://www.continuum.io/downloads>

# Why Anaconda?



Leading Open Data Science Platform Powered by Python



Leading Package and Environment Manager

## OPEN DATA SCIENCE



## DATA



## COMPUTATION





File Edit View Insert Cell Kernel Help

| Python 2 O



## Simple Jupyter demo

This cell has text formatted using the markdown language, which gets rendered like regular html.  
The next cell has some code:

```
In [57]: import random  
for i in range(3):  
    print random.random()  
x = 10
```

```
0.10564822904  
0.153941700348  
0.518503128416
```

Here is another text cell, with some *formatting*.

# ANACONDA NAVIGATOR

 Home

 Environments

 Projects (beta)

 Learning

 Community

Documentation

Developer Blog

Feedback



You  
Tube



Applications on

base (root) ▾

Channels



jupyterlab

0.31.5

An extensible environment for interactive and reproducible computing, based on the Jupyter Notebook and Architecture.

Launch



notebook

5.4.0

Web-based, interactive computing notebook environment. Edit and run human-readable docs while describing the data analysis.

Launch



qtconsole

4.3.1

PyQt GUI that supports inline figures, proper multiline editing with syntax highlighting, graphical calltips, and more.

Launch

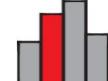


spyder

3.2.6

Scientific PYthon Development EnviRonment. Powerful Python IDE with advanced editing, interactive testing, debugging and introspection features

Launch



glueviz

0.12.0

Multidimensional data visualization across files. Explore relationships within and among related datasets.

Install



orange3

3.4.1

Component based data mining framework. Data visualization and data analysis for novice and expert. Interactive workflows with a large toolbox.

Install

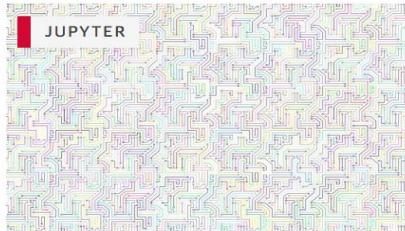
# JupyterCon 2018

Posts tagged: "JupyterCon 2018"



## Sea change: What happens when Jupyter becomes pervasive at a university?

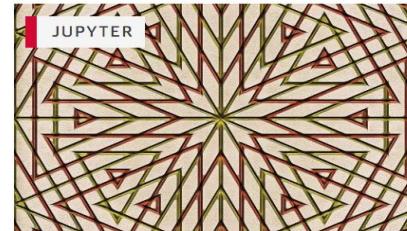
Fernando Perez talks about UC Berkeley's transition into an environment where many undergraduates use Jupyter and the open data ecosystem as naturally as they use email.



## Beyond interactive: Scaling impact with notebooks at Netflix

Michelle Ufford shares how Netflix leverages notebooks today and describes a brief vision for the future.

<https://www.oreilly.com/tags/jupytercon-2018>



## Jupyter notebooks and the intersection of data science and data engineering

David Schaaf explains how data science and data engineering can work together to deliver results to decision makers.



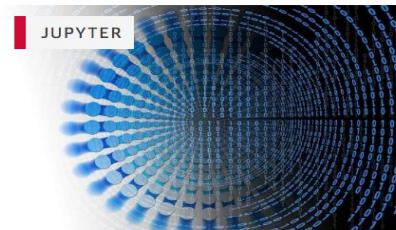
## The future of data-driven discovery in the cloud

Ryan Abernathy makes the case for the large-scale migration of scientific data and research to the cloud.



## Democratizing data

Tracy Teal explains how to bring people to data and empower them to address their questions.



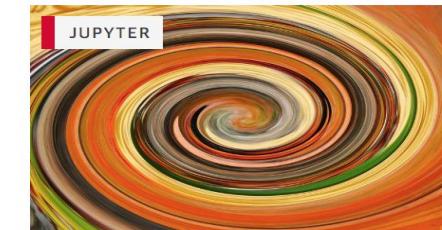
## Disease prediction using the world's largest clinical lab data set

Cristian Capdevila explains how Prognos is predicting disease.



## Data science as a catalyst for scientific discovery

Michelle Gill discusses how data science methods and tools can link information from different scientific fields and accelerate discovery.



## Machine learning and AI technologies and platforms at AWS

Dan Romuald Mbanga walks through the ecosystem around the machine learning platform and API services at AWS.

# Tuesday at Berkeley: Data 8, ~1,300 students



**Data 100, ~800 (650 last spring)**



## Sponsors

Project Jupyter receives direct funding from the following sources:

THE LEONA M. AND HARRY B.  
**HELMSLEY**  
CHARITABLE TRUST

**rackspace**  
the #1 managed cloud company



ALFRED P. SLOAN  
FOUNDATION

**fastly**<sup>®</sup>

GORDON AND BETTY  
**MOORE**  
FOUNDATION

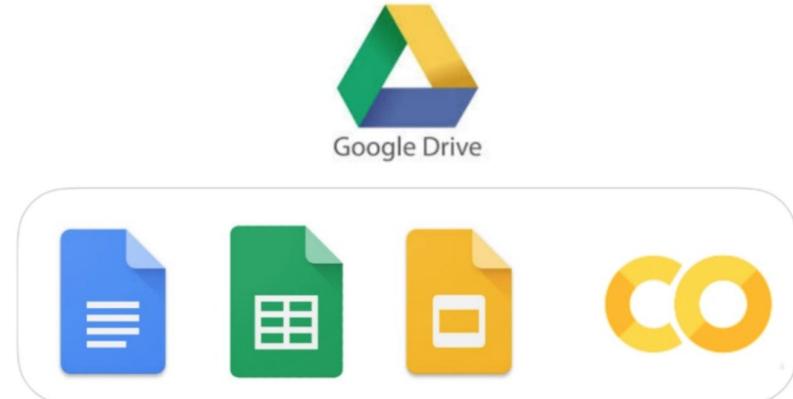


**Google**

 Microsoft

# Google Colaboratory

<https://colab.research.google.com/>



Colaboratory is a Google research project created to help disseminate machine learning education and research. It's a Jupyter notebook environment that requires no setup to use and runs entirely in the cloud.

Colaboratory notebooks are stored in Google Drive and can be shared just as you would with Google Docs or Sheets. Colaboratory is free to use.

# Agenda

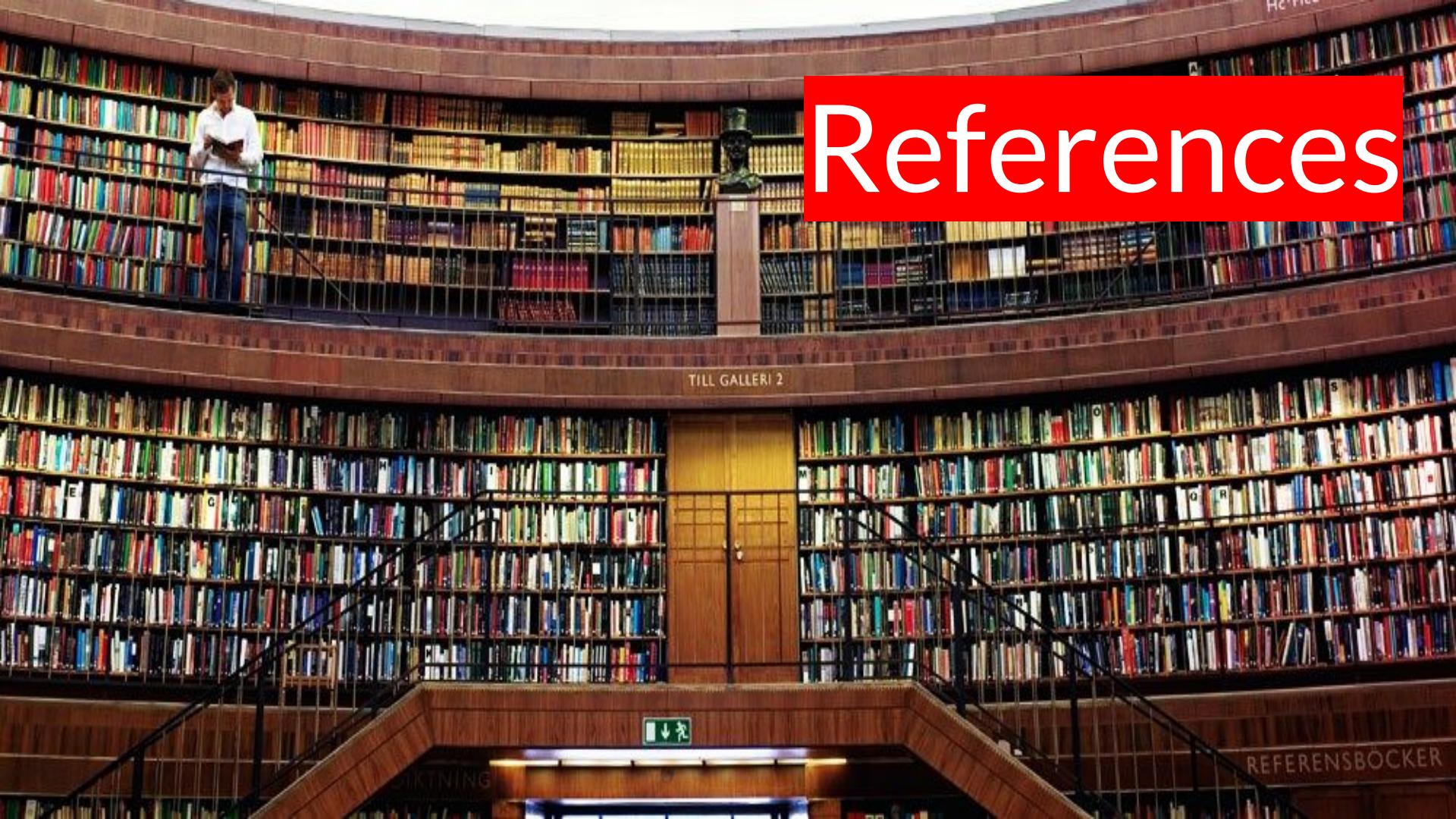
---

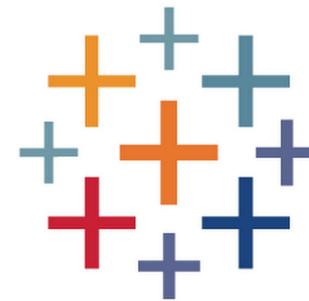
<https://goo.gl/gk2H1h>



<https://github.com/ivanovitchm/datascience2machinelearning>

# References

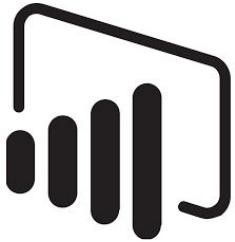




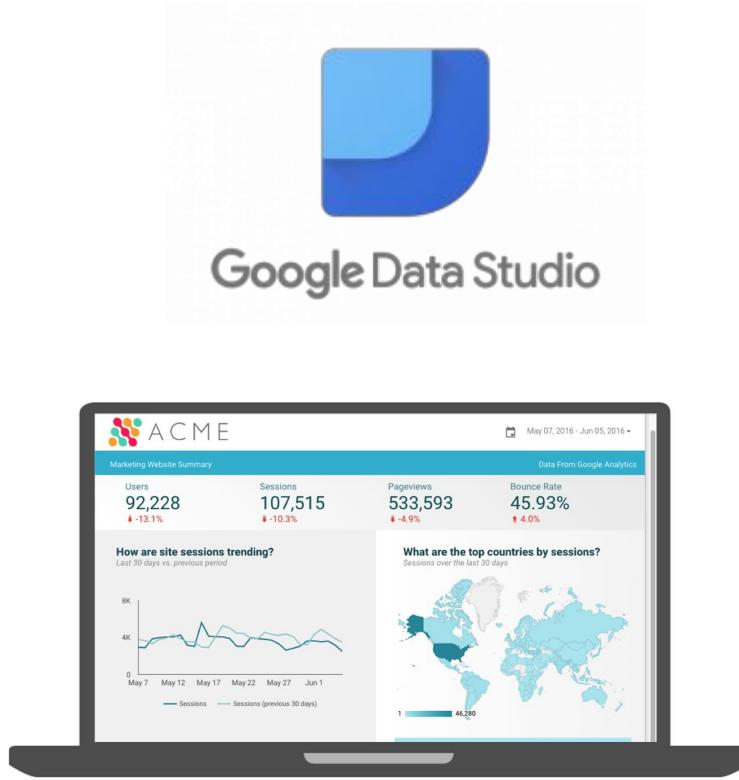
<http://www.pentaho.com/>



<https://www.tableau.com/>



<https://powerbi.microsoft.com>



The Google Data Studio interface is shown running on a laptop. The dashboard is titled "ACME Marketing Website Summary" and includes the following data points:

Metric	Value	Change
Users	92,228	-12.1%
Sessions	107,515	-10.3%
Pageviews	533,593	-4.9%
Bounce Rate	45.93%	+4.0%

Below these summary metrics are two detailed charts: one showing site session trends over the last 30 days versus the previous period, and another showing the top countries by sessions with a world map overlay.

[https://www.google.com.br/analytics/  
data-studio/](https://www.google.com.br/analytics/data-studio/)



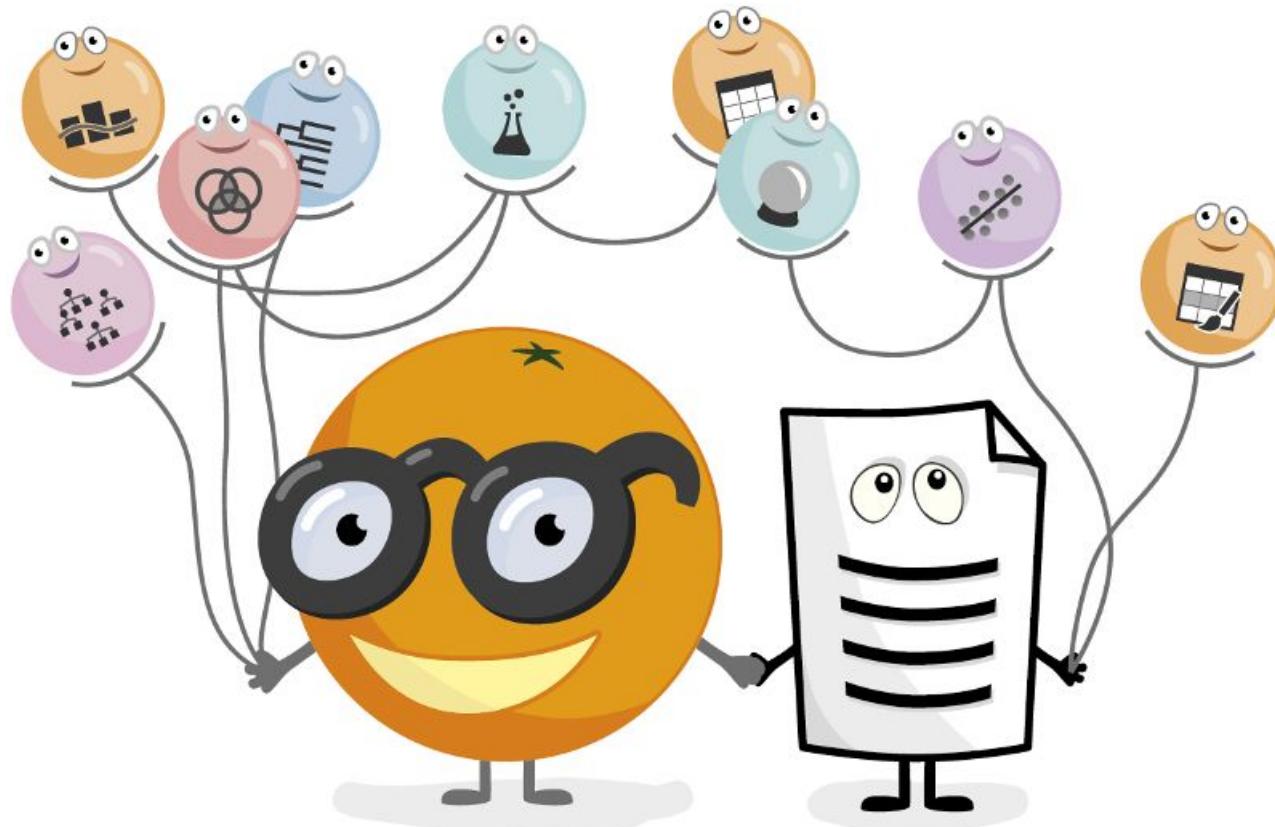
C A R O L



<https://goo.gl/Ndf38Q>



# Data Mining Fruitful and Fun



# References Online

---



kaggle



Data Science  
Academy

# References Online

---

1. <https://www.coursera.org/specializations/big-data>
2. <https://www.coursera.org/learn/machine-learning>
3. <https://www.coursera.org/specializations/deep-learning>



deeplearning.ai



# References Online



IDIOMA:

Fundamentos de AI &  
Machine Learning



IDIOMA:

Engenheiro de Machine  
Learning

kaggle



IDIOMA:

Deep Learning



IDIOMA:

Deep Reinforcement  
Learning Expert



IDIOMA:

Computer Vision Expert

:) Affectiva



IDIOMA:

Natural Language  
Processing Expert

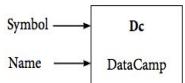
amazon alexa IBM Watson®



# The Periodic Table of Data Science

An overview of key companies, resources and tools in data science (as of 4/12/2017)

<https://goo.gl/eecQ2V>



Dc	Ga	Sd
DataCamp	General Assembly	Strata Data

Courses	Data
Boot camps	Projects & Challenges, Competitions
Conferences	Programming Languages & Distributions

Search & Data Management
Machine Learning & Stats
Data Visualization & Reporting

Collaboration
Community & Q&A

News, Newsletters & Blogs
Podcasts

Sb	M	Od
SpringBoard	Metis	ODSC

Ex	Di	Tc
Edx	Data Incubator	Tableau Conference

C	In	U
Coursera	Insight	UseR!

Uda	Dsa	Pd
Udacity	NYC Data Science Academy	PyData

Ude	G	Paw
Udemy	Galvanize	Predictive Analytics World

Ps	Dsg	Kdd
Pluralsight	Data Science for Social Good	ACM SIGKDD Conference

Ly	Dsy	Tpc
Lynda	Data Society	Teradata Partners Conference

Tt	Dsj	Icd
TeamTreeHouse	Data Science Dojo	IEEE International Conference on Data Mining

Bdu							
Big Data University							

Py	Js	Vb	Pgs	Sli	Ah	W	Bml	Kn	Sm	Pb	Obi	Shn	Ddl	De
Python	JavaScript	Visual Basic	PostgreSQL	SQLite	Apache Hadoop	Weka	BigML	Knime	Spark MLlib	Power BI	Oracle BI	Shiny	Domino Data Lab	Data Science Experience
R	Cp	Sc	Ar	Bq	Hw	O	Dar	Lib	Ho	Bo	Alt	Mpl	Nt	Rs
R	C++	Scala	Amazon Redshift	Google BigQuery	Hortonworks	Oracle	DataRobot	LibSVM	H2O	BusinessObjects	Alteryx	Matplotlib	Nteract	Rstudio
S	Pl	Ca	Hb	Td	Cl	Mss	Microsoft SQL server	Rm	Mat	Th	Sp	Sav	Ply	Ro
SQL	Perl	Cassandra	HBase	Teradata	Cloudera	Microsoft SQL server	RapidMiner	Mathematica	Theano	Spotfire	SAS Visual Analytics	Plotly	Rodeo	Beaker Notebook
B	Mr	P	Mdb	To	Aem	Spl	Cho	Mah	Aml	Ql	Po	Me	Spy	Ze
Bash	Microsoft R Open	Pig	Mongo DB	Toad	Amazon Elastic Mapreduce	Splunk	Chorus	Mahout	Azure Machine Learning	Qlikview	PowerPivot	Microsoft Excel	Spyder	Apache Zeppelin
Mtl	Cy	Im	K	Ms	Mar	Sr	Tf	St	D	Co	Gch	Pe	Dst	Ju
Matlab	Canopy	Impala	Kafka	MySQL	MapR	Solr	Tensorflow	Stata	D3	Cognos	Google Charts	Pentaho	Data Science Studio	Jupyter
J	An	Sp	Hi	Idb	Lu	El	Sk	Da	My	Aa	T	B	Db	Gh
Java	Anaconda	Spark	Hive	IBM DB2	Lucene	ElasticSearch	Scikit-Learn	Dato/Graphlab	Microstrategy	Tableau	Bokeh	Databricks notebook	Github	

Dw	Q	Fte	Sa	Gp	Dg	K								
Data.world	Quandl	FiveThirtyEight	Socrata	Google Public	Data.gov	Kaggle								
St	Uci	Wb	At	Bf	Dk	Dd								
Statista	UCI Machine Learning Repository	World Bank	Academic Torrents	Buzzfeed	DataKind	DrivenData								
							Re	So	Cv	Qu	Av	Dse		
							Reddit	Stack Overflow	Cross Validated	Quora	Analytics Vidhya	Data Science Stack Exchange		
							Mu	Rdm						
							Meetup	RDataMining						

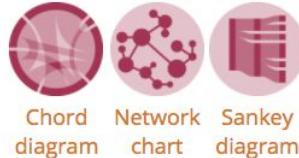
Kdn	Ibd
KDnuggets	insideBIGDATA
Rb	Pp
R-Bloggers	PlanetPython
Hn	Dt
HackerNews	DataTau
Dsc	Dsr
Data Science Central	Data Science Roundup
Dsw	Or
Data Science Weekly	O'Reilly
Dr	Pw
Data Elixir	Python Weekly
Rw	Pd
R Weekly	Partially Derivative
Bds	Tm
Becoming a Data Scientist	Talking Machines
Ds	Dsk
Data Stories	Data Skeptic
Ld	Ns
Linear Digressions	Not So Standard Deviations



## MAPS



## FLOW

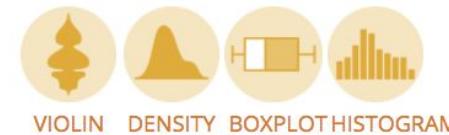


## Other

THE PYTHON  
GRAPH GALLERY

<https://python-graph-gallery.com/>

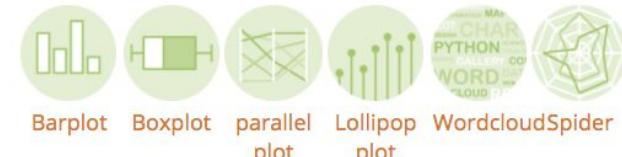
## DISTRIBUTION



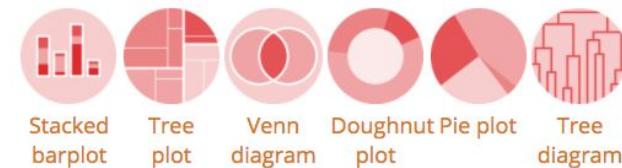
## CORRELATION



## RANKING



## PART OF A WHOLE



# Reference Blog

<https://machinelearningmastery.com/start-here/#process>

## Do You Need Help Getting Started with Applied Machine Learning?

This is The Step-by-Step Guide that You've Been Looking For!



Hi, Jason here. I'm the guy behind Machine Learning Mastery.

My goal is to help you **get started, make progress** and **kick butt** with machine learning.

I teach a **top-down** and **results-first approach** designed for developers and engineers. This is unlike most academic textbooks and university courses.

Access my best free tutorials [on the blog](#) or take the next step with my [paid training material](#).

You may be **feeling overwhelmed**. You may have **a lot of questions**. I created this page for you. It is your starting point.

Take your time. Bookmark this page. Find the answers to your questions.

MENU ▾



Search



E-alert



Submit



Login

News & Comment Research

<https://www.nature.com/articles/d41586-018-06201-x>

News Opinion Research Analysis Careers Books & Culture

NEWS · 05 SEPTEMBER 2018

# Google unveils search engine for open data

The tool, called Google Dataset Search, should help researchers to find the data they need more easily.

## Google Dataset Search Beta

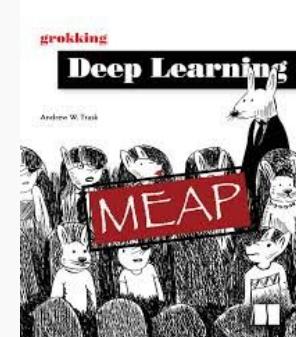
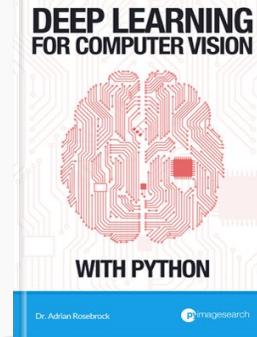
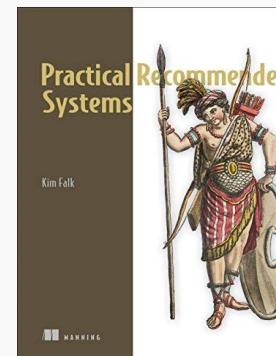
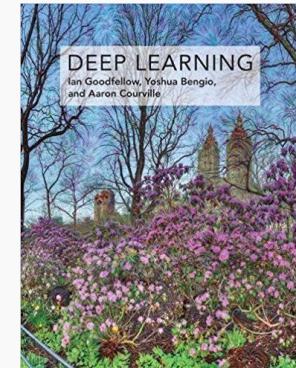
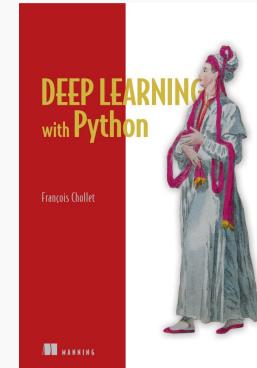
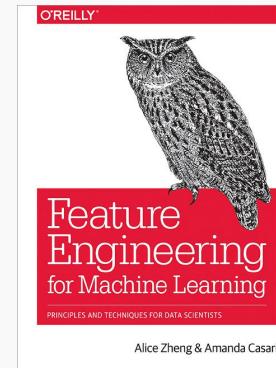
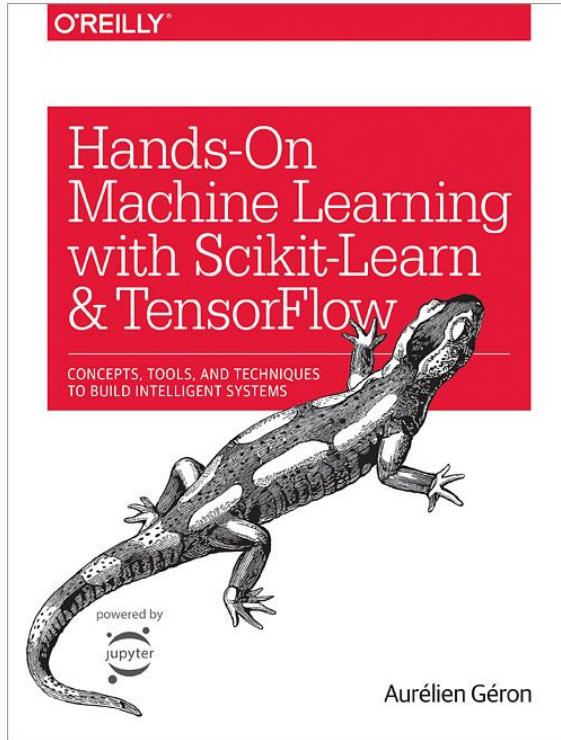
Search for Datasets



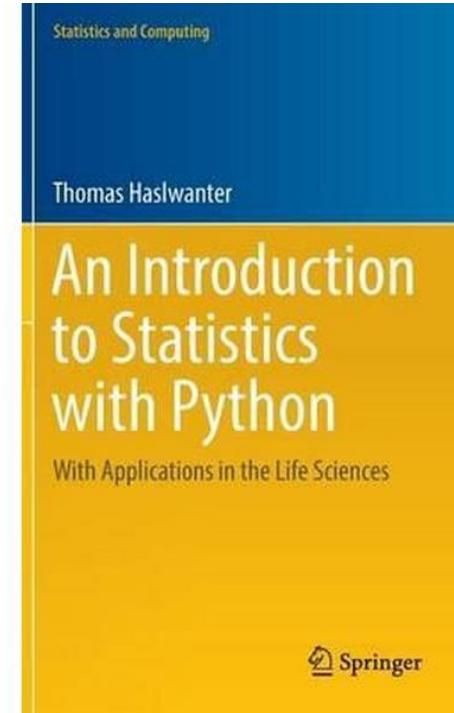
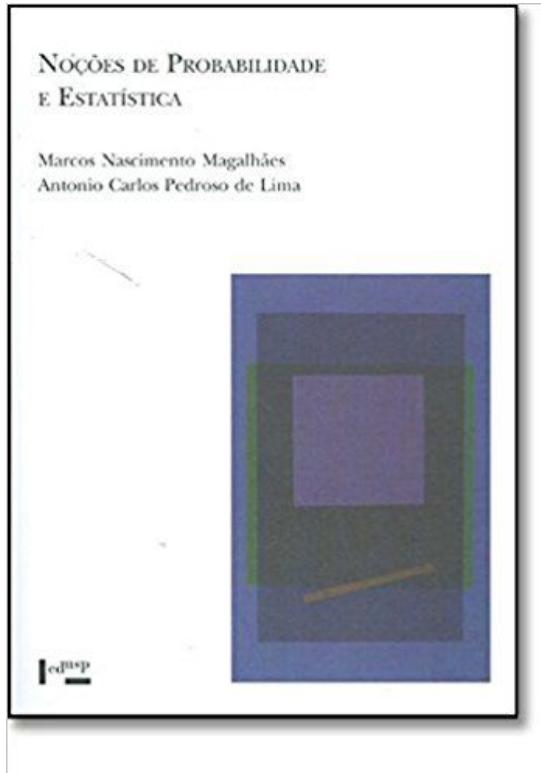
<https://toolbox.google.com/datasetsearch>



# References



# Referências



*"AI is the new electricity ... We have enough papers. Stop publishing, and start transforming people's lives with technology"*

*Andrew Ng*

