

# DATA SCIENCE IN THE REAL WORLD *CHALLENGES & OPPORTUNITIES*

Olivier Verscheure, PhD  
Swiss Data Science Center  
EPFL & ETH Zurich

# About me

Academia



**ETH** zürich



1999

2016

IBM Research  
CS Data Management, MH Mathematics  
Discovery & Data Mining, Networking &  
ers, Privacy, Programming Languages, Mate  
Language Processing, Computer Arch  
ages, Computational Biology, Relational

New York



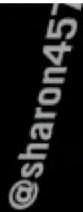


# What Do You See?



O'Connell Bridge / D'ollier St. Dublin City CCTV  
8 Apr 2013 18:31:50 GMT Daylight Time

100





# Data is the New Oil

**The  
Economist**

MAY 6TH–12TH 2017

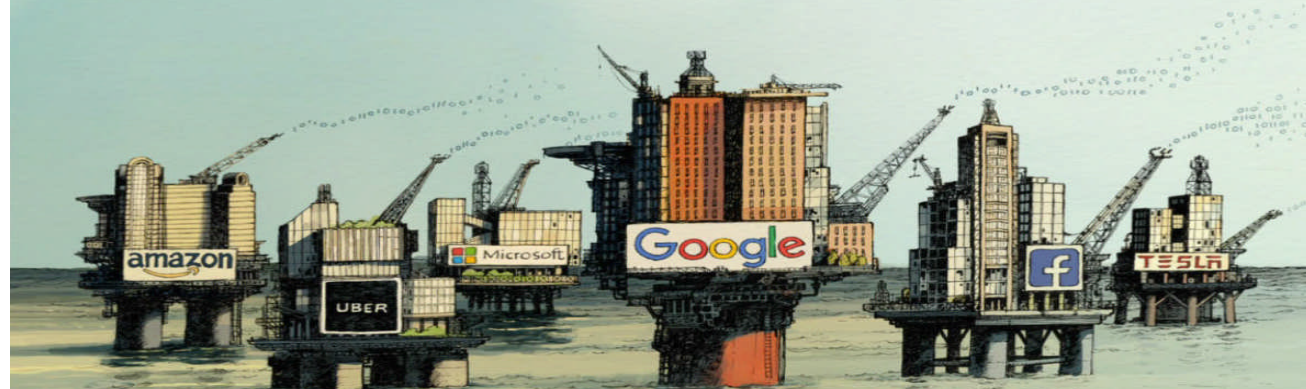
**Crunch time in France**

**Ten years on: banking after the crisis**

**South Korea's unfinished revolution**

**Biology, but without the cells**

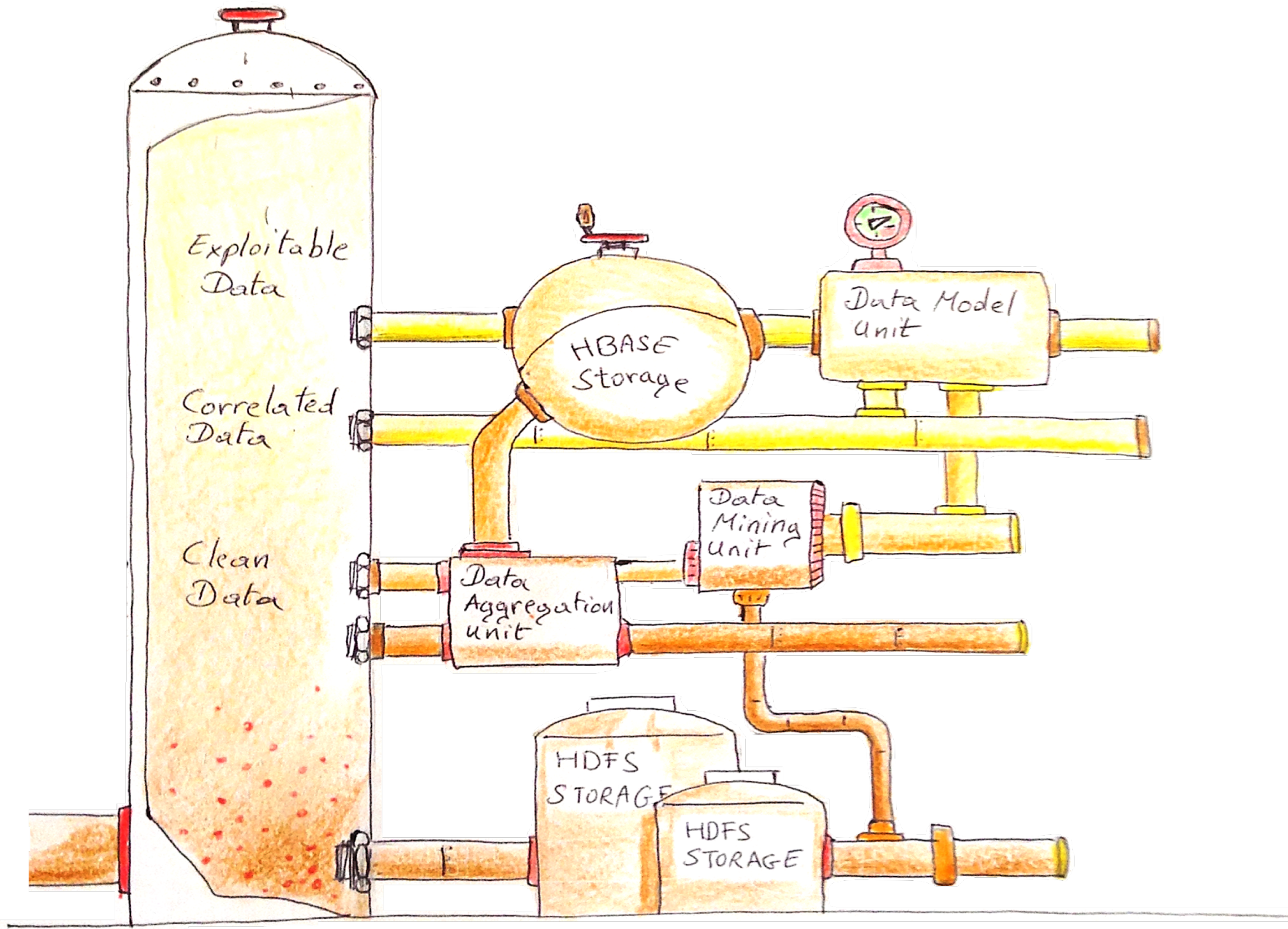
## **The world's most valuable resource**



**Data and the new rules  
of competition**

*The Economist, May 2017*

# Like Oil, Data Must be Refined



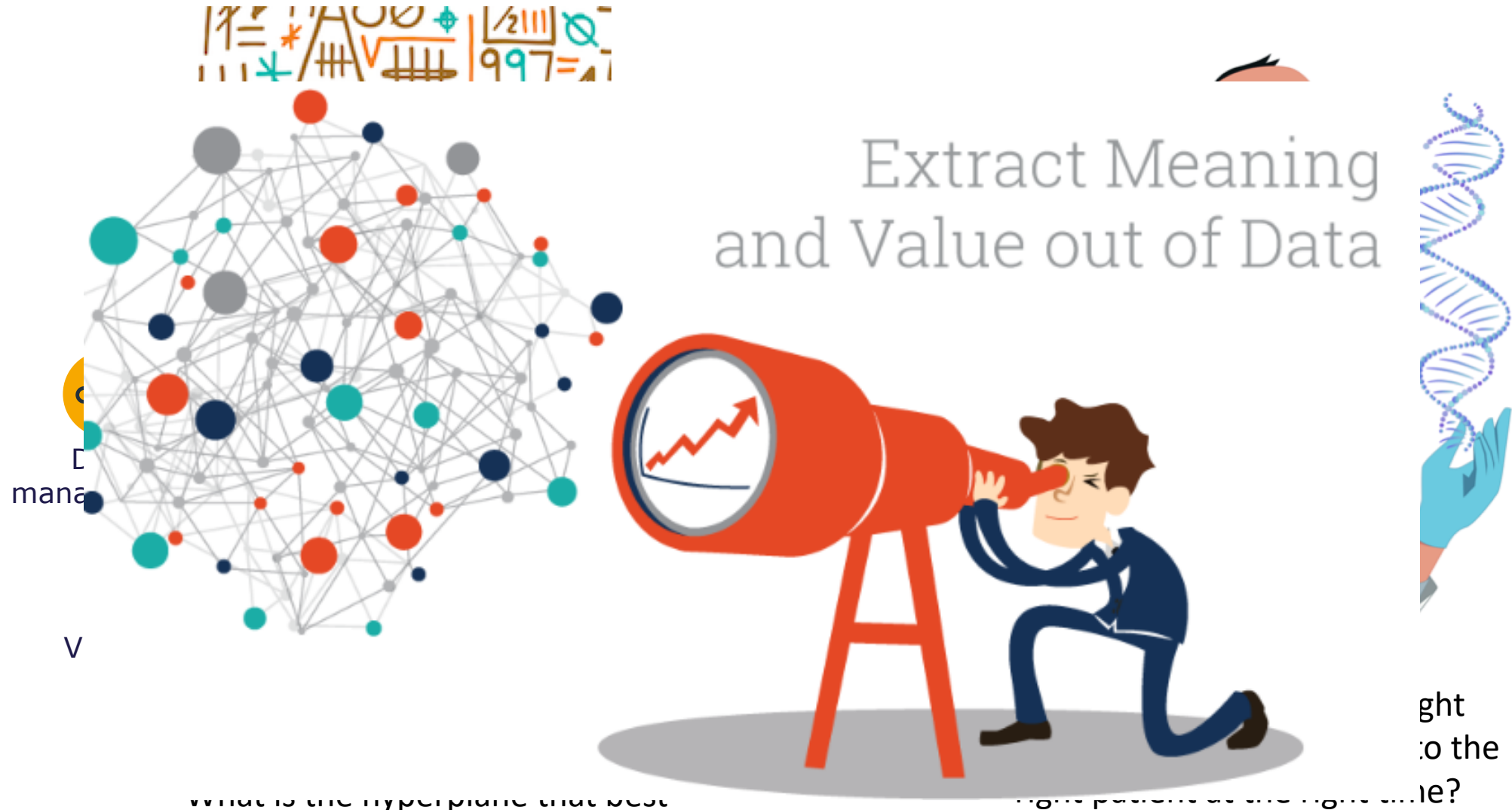
# Big Data, Bad Data



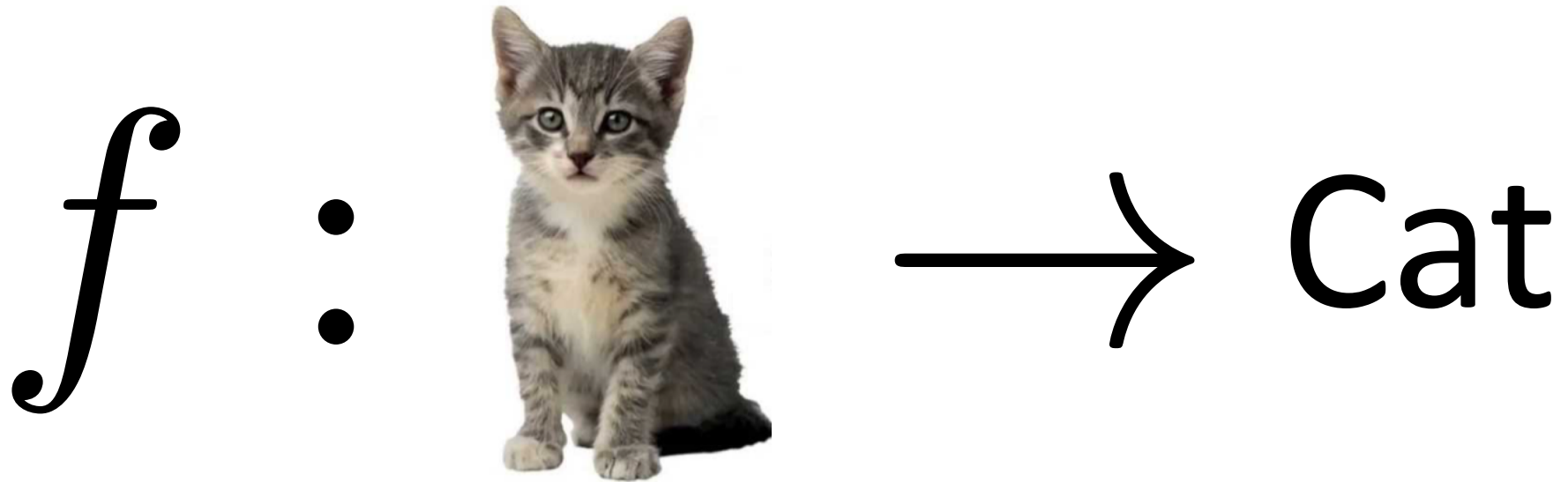
BY LOWE FOR THE SUN-SENTINEL, FLO



# Data Science, a Fragmented Ecosystem

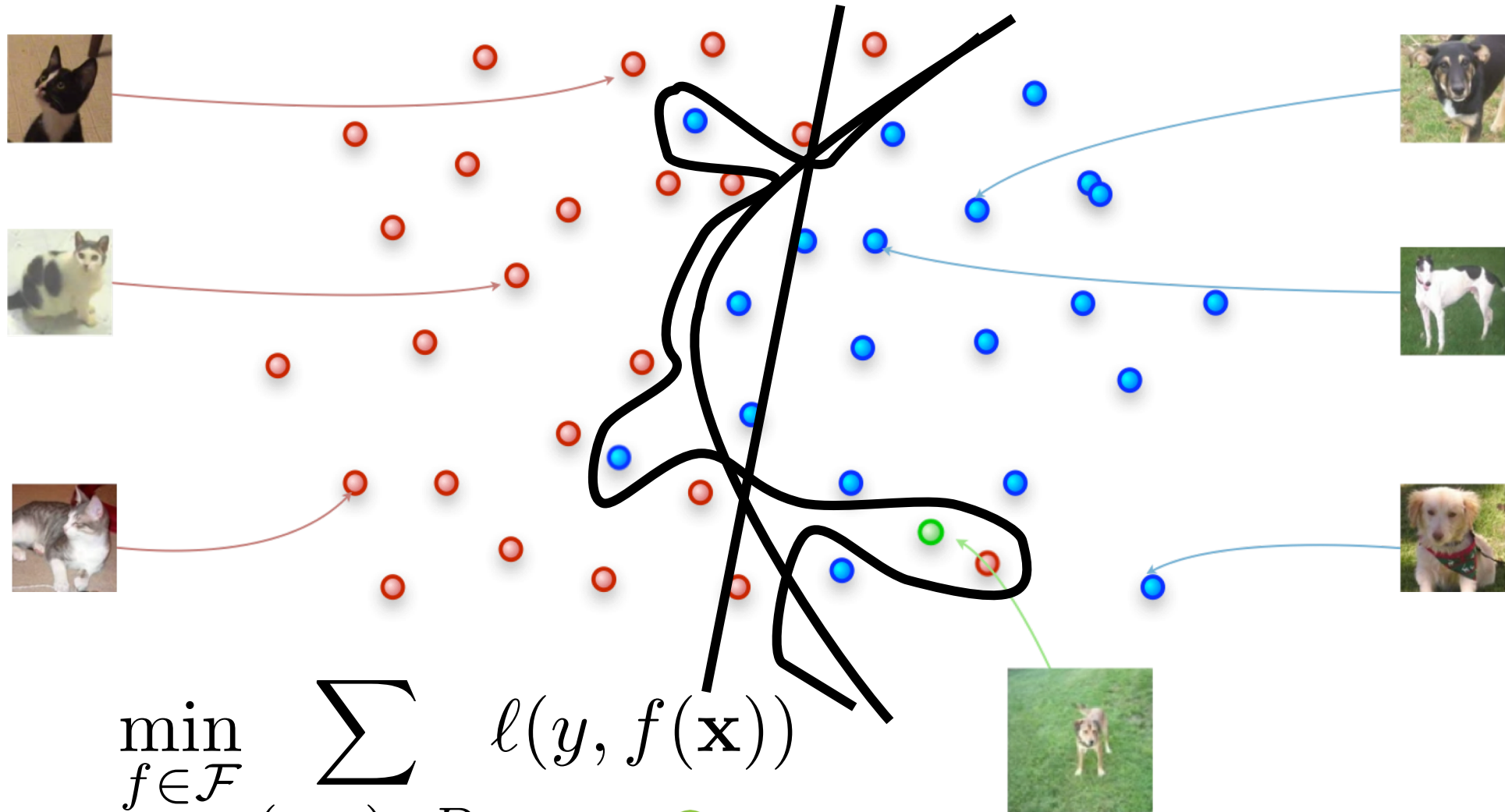


# Machine Learning in a Nutshell





# Machine Learning (ML) in a Nutshell



$$\min_{f \in \mathcal{F}} \sum_{(\mathbf{x}, y) \in D} \ell(y, f(\mathbf{x}))$$

# Some Key Players



# Recent advances in ML: Deep Learning

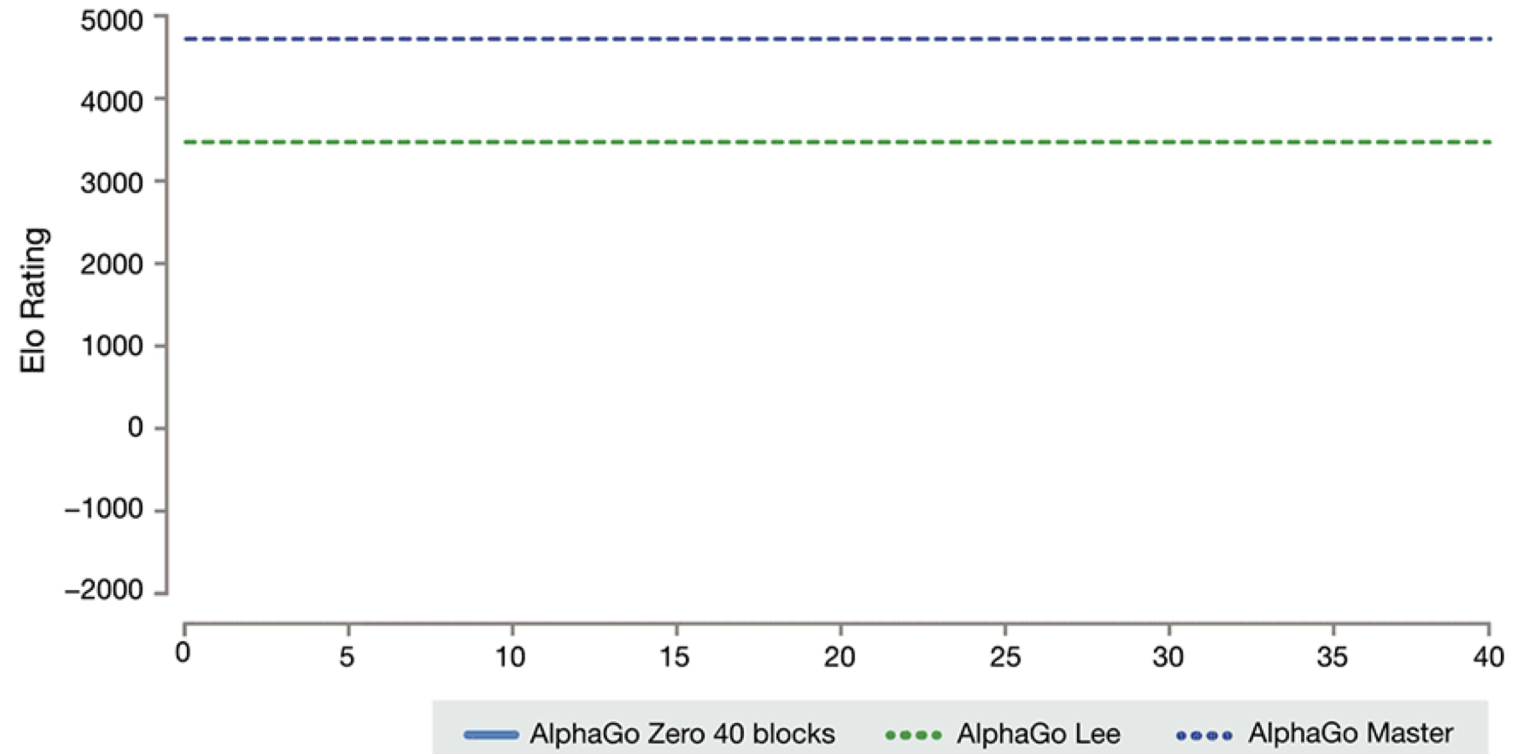
## ImageNet Challenge

IMAGENET

- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.



# Recent advances in ML: Deep Reinforcement Learning



# This Success Relies On...

1. Large dataset of labelled data
2. Good quality data
3. Enough computing power
4. Clear and measurable objectives



# An Unexpected Outcome

**Title:** Universal adversarial perturbations

**Authors:** [Moosavi-Dezfooli, Seyed-Mohsen](#); [Fawzi, Alhussein](#); [Fawzi, Omar](#); [Frossard, Pascal](#)

**Publication:** eprint arXiv:1610.08401

**Publication Date:** 10/2016



- It's an Indian elephant!
- At least after adding a universal noise to the image
- Deep learning models do not mimic brain activity

*This is not a sock*

# A Disturbing Outcome



Turning a **STOP** sign into a **45 mph speed limit**



# Structured vs Unstructured Data

Tissue Cask Cell ID	Total Cells	Tissue Cask Cell Density X Posit	Process Re Distance In Category R	Distance In Nucleus An Nucleus An Nucleus Min Nucleus Min Nucleus PC Nucleus PD1 (FITC) Count (Normalized Counts, Total Weighting)	Nucleus PC Nucleus PC Nucleus PC Nucleus PC															
86-1A_image_1.m3	No tissue	751	4	INFA	05	NA	101	0.09%	0.66	7.11	17.91	23.99	0.56%	0.812	1.138	0.108	0.021	0.787		
86-1A_image_1.m3	Tumor	2	878	57	INFA	1	NA	0.01%	0.72	0.01%	2.16	0.01%	0.01%	1.289	1.199	0.019	75.433	0		
86-1A_image_1.m3	Tumor	3	910	4	INFA	1	INFA	0.1	0.01%	0.9	7.04	1.3	1.85	0.126	0.42	0.038	0.094	25.642	0	
86-1A_image_1.m3	Tumor	4	1002	6	INFA	1	INFA	96	0.01%	0.66	10.02	14.02	1.4	0.281	0.48	0.711	0.063	0.474	0	
86-1A_image_1.m3	Stroma	5	1214	8	INFA	02	INFA	0.08	0.08	0.65	1.75	4.08	0.067	0.65	1.026	0.07	46.499	0		
86-1A_image_1.m3	Stroma	6	1306	5	INFA	13	INFA	67	0.02%	0.73	9.05	10.97	1.21	0.28	0.484	0.08	0.095	32.437	0	
86-1A_image_1.m3	No tissue	7	175	6	INFA	13	INFA	66	0.02%	0.68	9.03	11.07	1.22	0.989	1.267	1.599	0.141	83.995	0.688	
86-1A_image_1.m3	No tissue	8	789	107	INFA	1	INFA	2	0.76	1	12	1.91	1.2	0.301	0.641	1.309	0.186	68.641	0	
86-1A_image_1.m3	Tumor	9	1056	5	INFA	1	INFA	11	0.01%	0.75	9	14	1.87	0.327	0.541	0.79	0.083	62.226	0	
86-1A_image_1.m3	Tumor	10	484	7	INFA	1	INFA	88	0.03%	0.78	8.69	15.5	1.52	0.726	0.15	1.684	0.216	98.887	0	
86-1A_image_1.m3	Tumor	11	980	1	INFA	1	INFA	0.73	0.02%	0.9	12.69	1.49	2.40	0.093	0.088	0.043	60.804	0		
86-1A_image_1.m3	Stroma	12	373	9	INFA	0	INFA	65	0.02%	0.71	8.16	11.29	1.38	0.282	0.44	0.611	0.073	28.808	0	
86-1A_image_1.m3	No tissue	13	1262	3	INFA	1	INFA	1	0.43	0.28	27.28	0.43	1.51	0.132	0.143	0.132	14.443	0		
86-1A_image_1.m3	No tissue	14	295	9	INFA	04	INFA	66	0.02%	0.72	6.08	10.9	1.35	1.087	1.464	1.882	0.199	96.306	0.847	
86-1A_image_1.m3	No tissue	15	774	8	INFA	1	INFA	2	1.31	0.02%	0.66	12.01	14.03	1.17	0.554	0.27	2.651	447	158.170	
86-1A_image_1.m3	No tissue	16	1118	10	INFA	1	INFA	3	0.45	0.23	25.1	0.45	1.23	0.21	0.36	0.36	0.36	171.263	0	
86-1A_image_1.m3	Tumor	17	636	11	INFA	1	INFA	161	0.02%	0.64	11.77	20.61	1.75	0.304	0.529	0.809	0.094	85.12	0	
86-1A_image_1.m3	Tumor	18	1017	9	INFA	1	INFA	241	0.03%	0.6	17.99	18.03	1	0.287	0.035	0.812	0.069	129.007	0	
86-1A_image_1.m3	Tumor	19	1035	207	INFA	0.54	INFA	0.54	13.78	0.54	13.78	0.54	1.41	0.31	0.645	0.056	0.391	191.28	0	
86-1A_image_1.m3	No tissue	20	1241	10	INFA	57	INFA	11	6.02%	0.48	9.41	12.51	1.33	0.42	0.73	1.257	0.141	48.66	0	
86-1A_image_1.m3	Tumor	21	865	1	INFA	6	INFA	6	0.01%	0.62	10.86	16.6	2.89	0.027	0.391	0.622	0.067	78.634	0	
86-1A_image_1.m3	Tumor	22	1402	14	INFA	288	INFA	1	0.01%	0.37	14.29	3.73	2.22	0.217	0.099	0.804	0.091	140.025	0	
86-1A_image_1.m3	No tissue	23	1269	13	INFA	1	INFA	1	105	0.03%	0.79	11.03	12.03	1.09	0.275	0.581	0.85	1.24	0.31	0.31
86-1A_image_1.m3	No tissue	24	910	14	INFA	130	INFA	0.03	0.16	0.63	0.03	0.03	0.03	0.03	0.03	0.03	0.03	68.78	0	
86-1A_image_1.m3	Stroma	25	783	14	INFA	10	7	70	0.03%	0.55	8.45	12.38	1.46	0.383	0.754	1.509	0.262	52.812	0	
86-1A_image_1.m3	No tissue	26	972	17	INFA	1	INFA	169	0.02%	0.57	11.59	20.18	1.74	0.291	0.525	0.793	0.1	88.741	0	
86-1A_image_1.m3	Tumor	27	1005	19	INFA	119	0.01%	0.51	12.34	13.97	1.06	0.33	0.607	0.769	0.884	0.525	61.251	0		
86-1A_image_1.m3	Tumor	28	1073	13	INFA	1	7	271	0.03%	0.53	14.58	25.85	1.77	0.219	0.749	0.777	0.095	128.886	0	
86-1A_image_1.m3	Stroma	29	514	18	INFA	0.02	INFA	0.02	179	0.02%	0.59	10.69	21.81	2.04	0.263	0.627	0.947	1.19	112.241	0.126
86-1A_image_1.m3	Tumor	30	820	178	INFA	0.06	INFA	0.06	15.01	0.06	15.01	0.06	1.07	0.305	0.591	0.942	0.095	66.261	0	
86-1A_image_1.m3	Tumor	31	484	18	INFA	489	0.06%	0.45	23.3	30.32	1.3	0.315	0.198	0.38	0.146	0.292	30	0	0	
86-1A_image_1.m3	Tumor	32	494	18	INFA	107	INFA	0.15	0.75	1.18	0.15	0.15	0.15	0.15	0.15	0.15	0.15	61.832	0.063	
86-1A_image_1.m3	Tumor	33	861	18	INFA	218	0.03%	0.5	15.64	21.95	1.35	0.303	0.643	0.672	0.943	1.03	140.025	0		
86-1A_image_1.m3	Stroma	34	484	18	INFA	0	INFA	0.05	10.04%	0.6	10.21	15.17	14.9	0.485	0.775	1.5	0.17	81.383	0.885	
86-1A_image_1.m3	Tumor	35	1360	17	INFA	4	INFA	4	11	0.03%	0.73	11.63	13.99	1.12	0.365	0.733	0.391	63.673	0	
86-1A_image_1.m3	No tissue	36	1380	22	INFA	57	INFA	162	0.04%	0.34	9.37	28.07	2.99	0.312	0.645	1.072	0.199	104.444	0	
86-1A_image_1.m3	No tissue	37	1241	23	INFA	1	INFA	205	0.05%	0.56	14.01	20.03	2.43	0.737	1.696	0.767	0.947	347.887	0	
86-1A_image_1.m3	No tissue	38	408	93	INFA	0.01	INFA	0.01	11.1	0.04%	0.91	11.03	10.47	0.17	0.768	1.024	0.136	71.217	0	
86-1A_image_1.m3	Tumor	39	935	25	INFA	1	INFA	124	0.02%	0.65	12.38	14.3	1.16	0.271	0.747	0.727	0.087	57.176	0	
86-1A_image_1.m3	Tumor	40	892	24	INFA	280	0.04%	0.55	15.03	20.01	1.75	0.327	0.611	0.526	0.925	1.026	1.29	171.102	0	
86-1A_image_1.m3	No tissue	41	1256	18	INFA	57	7	7	0.02%	0.68	91	7.12	10.69	0.68	0.511	1.16	0.107	120.619	0	
86-1A_image_1.m3	Stroma	42	360	26	INFA	9	19	122	0.05%	0.74	10.4	15.32	14.7	0.41	0.003	0.959	0.107	73.507	0.034	
86-1A_image_1.m3	Stroma	43	2492	8	INFA	8	123	16	0.02%	0.57	16.56	14.62	14.62	0.462	0.34	1.12	0.116	67.647	0	
86-1A_image_1.m3	Tumor	44	632	26	INFA	1	INFA	128	0.02%	0.53	11.89	15.02	13.33	0.392	0.71	0.708	0.511	101.163	0	
86-1A_image_1.m3	Tumor	45	1150	25	INFA	1	INFA	2	84	0.01%	0.81	9.03	11.61	1.22	1.166	0.339	0.491	263.988	0	
86-1A_image_1.m3	Tumor	46	733	22	INFA	1	INFA	194	0.06%	0.46	20.38	1.37	28.88	0.482	0.481	0.283	0.042	59.143	0	
86-1A_image_1.m3	Tumor	47	867	26	INFA	1	INFA	59	0.01%	0.46	7.63	13.38	1.9	0.367	0.571	0.837	0.088	33.715	0	
86-1A_image_1.m3	Tumor	48	1012	27	INFA	1	INFA	192	0.02%	0.45	15.22	20.22	1.33	0.305	0.501	0.901	0.094	96.208	0	
86-1A_image_1.m3	Tumor	49	286	57	INFA	1	INFA	0.91	0.02%	0.39	0.02	0.39	0.02	0.39	0.39	0.02	0.39	0.02	0.39	0

## Structured



## Unstructured



## Semi-structured

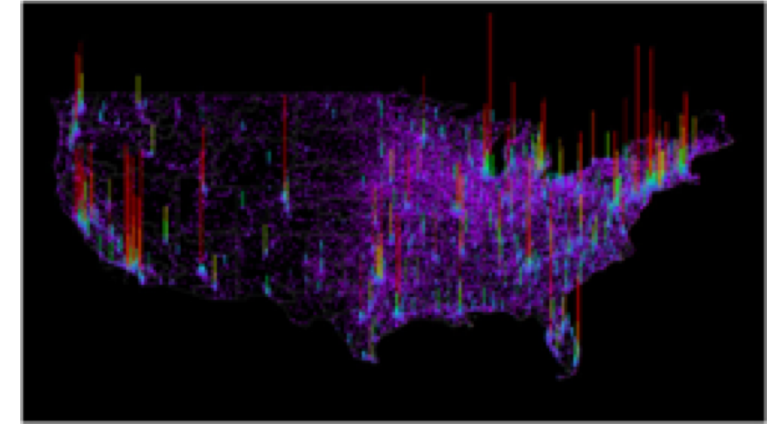
# The Structured World



Health care



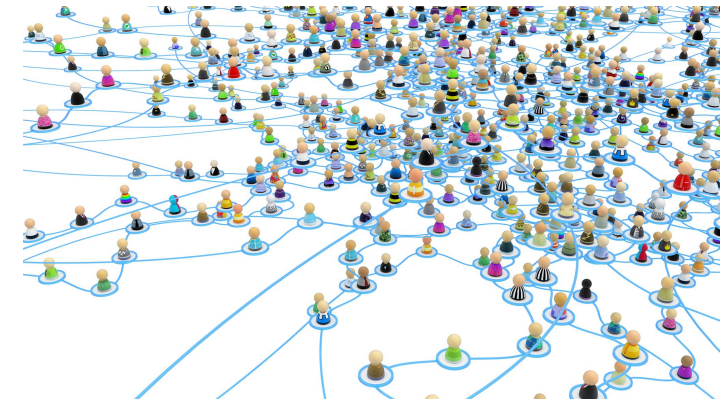
Predictive maintenance



Energy networks



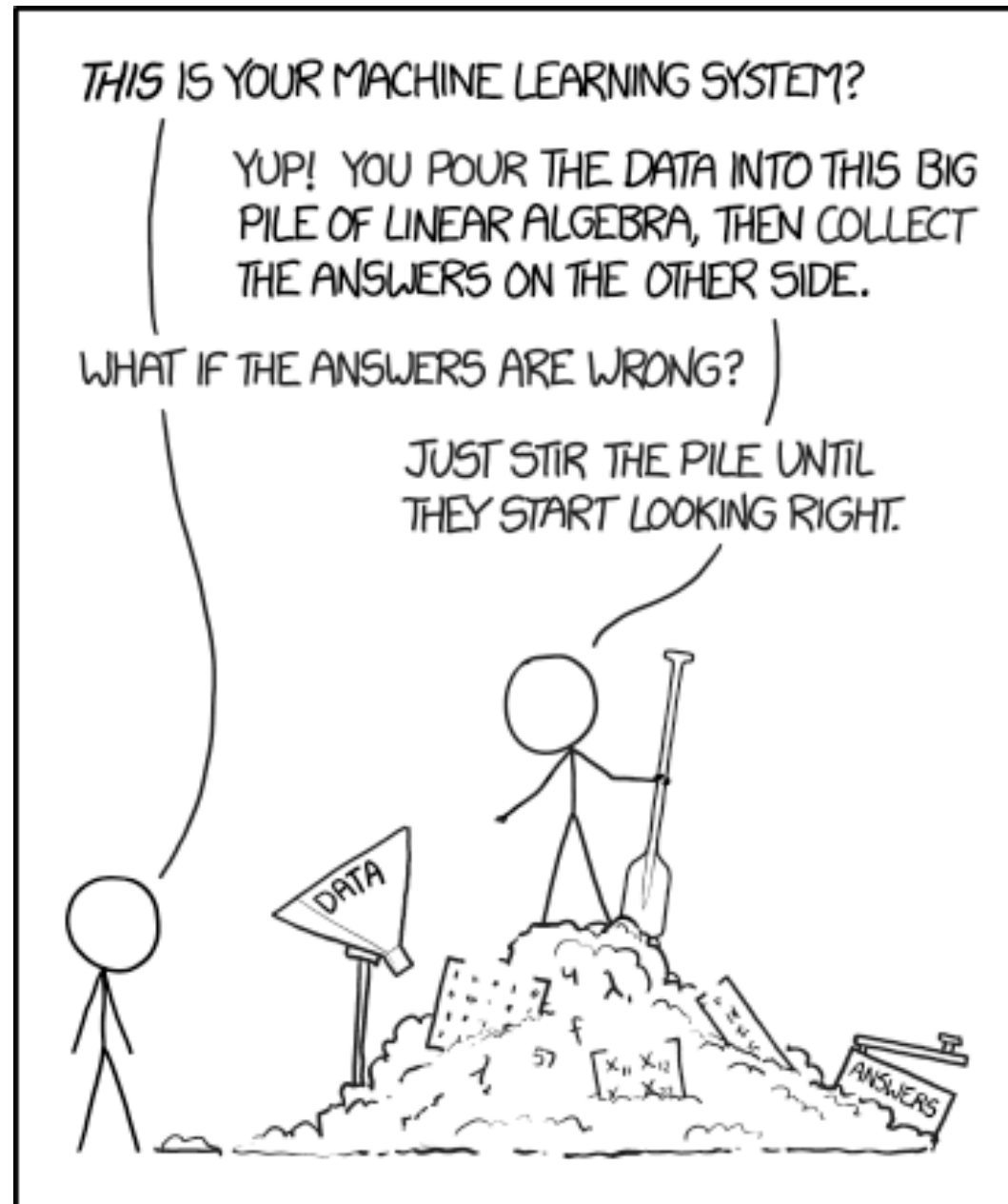
Financial time series



Social networks



# A Sobering View of Data Science



# Obstacles to a Wider Adoption in a Structured World

1. ~~Large dataset of labelled data~~ -> Labelling is expensive

2. ~~Good quality data~~ -> Data is usually missing/Increased uncertainty

3. ~~Clear and measurable objectives~~ -> Knowledge discovery/causality

4. Lack of interpretability/lack of trust

# Upcoming Challenges for Structured Data

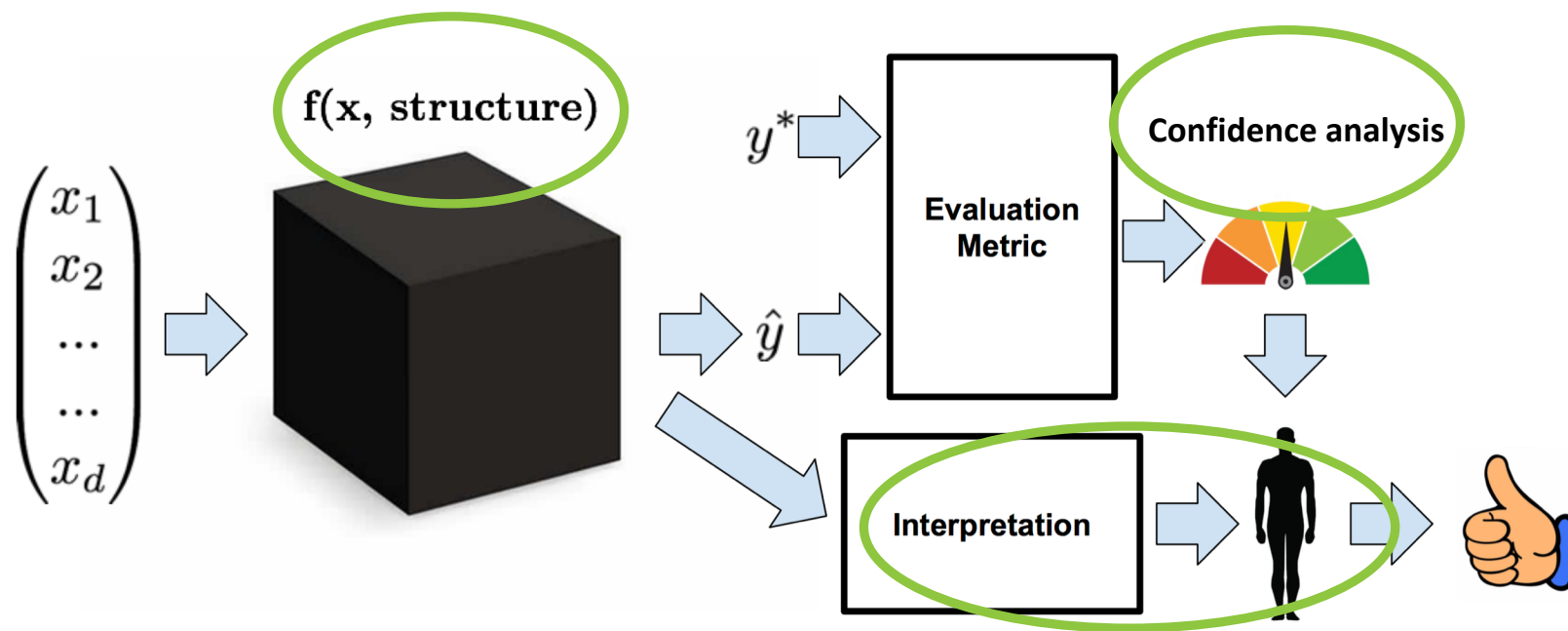


Figure adapted from Z. Lipton

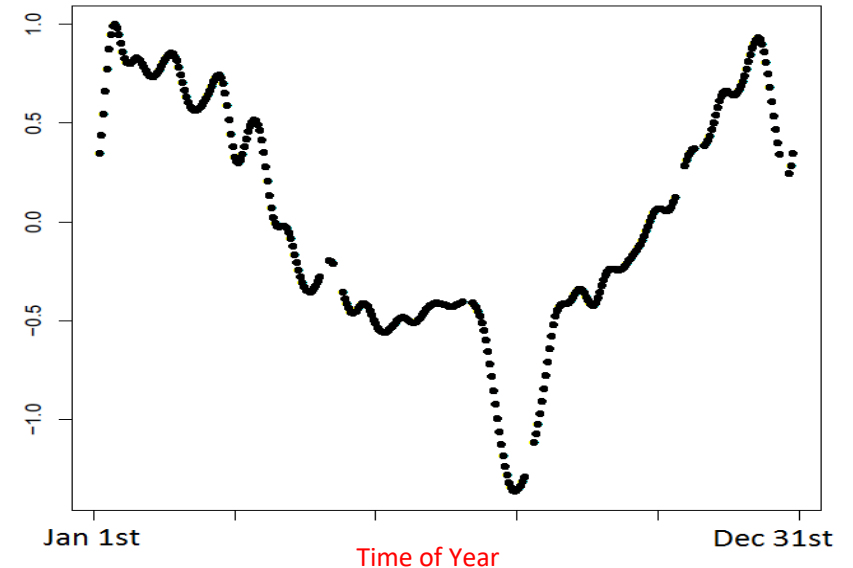
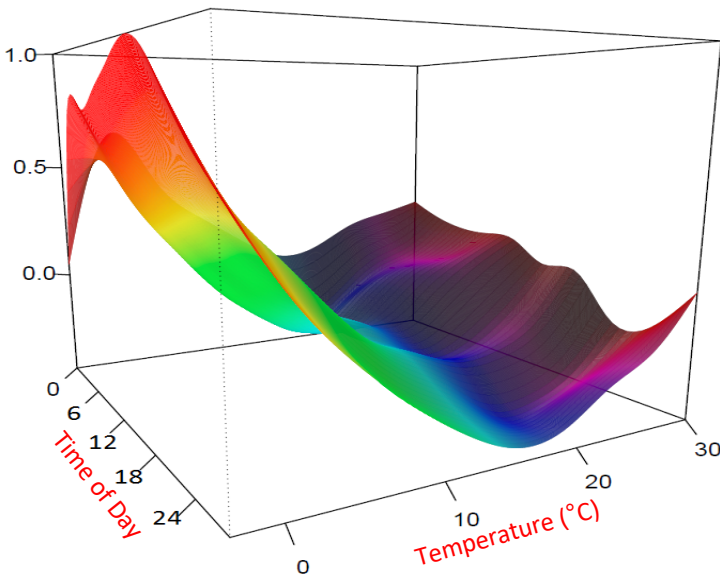
1. Incorporating structure knowledge in the model for data curation
2. Dealing with uncertainties
3. Promoting causality/interpretability
4. Focusing on unsupervised/semi-supervised learning

# Interpretable Machine Learning – Use Case

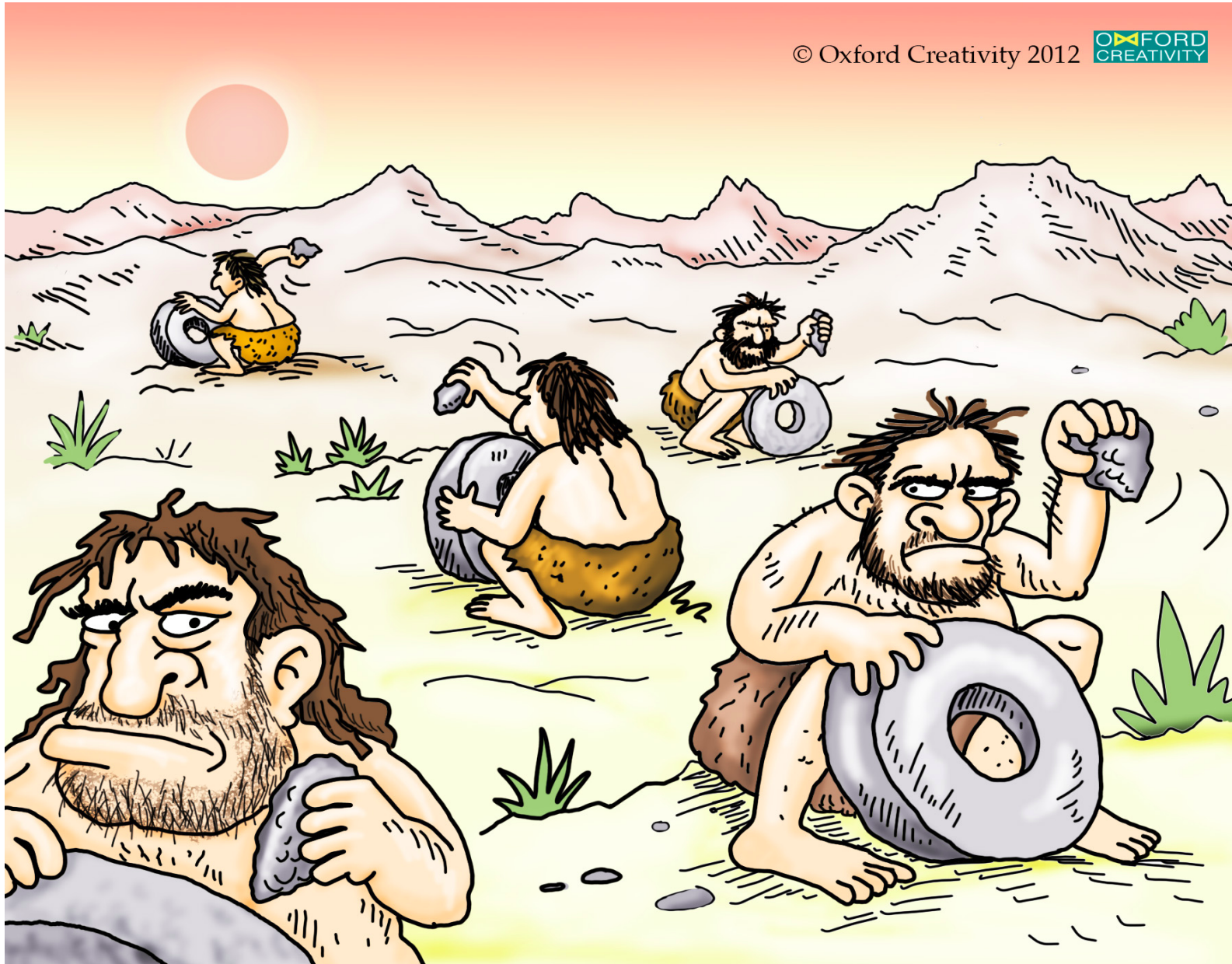
- Forecasting demand in electricity (France)

$$\begin{aligned}
 y_k = & \beta^{\text{Intercept}} + \overset{\text{Trend}}{f^{\text{Trend}}(k)} + \overset{\text{Lag load}}{f^{\text{LagLoad}}(y_{k-48})} + \overset{\text{Day-type specific daily pattern}}{\sum_{l=1}^6 \mathbf{1}(x_k^{\text{DayType}} = l)(\beta_l^{\text{DayType}} + f_l^{\text{TimeOfDay}}(x_k))} \\
 & + f^{\text{CloudCover}}(x_k) + \underbrace{f^{\text{Temperature/TimeOfDay}}(x_k)}_{\text{Lag temperature (accounting for thermal inertia)}} + f^{\text{LagTemperature}}(x_{k-48}) \\
 & - \underbrace{f^{\text{TimeOfYear}}(x_k)}_{\text{Lag temperature (accounting for thermal inertia)}} + x_k^{\text{LoadDecrease}} f^{\text{LoadDecrease}}(x_k) + \epsilon_k.
 \end{aligned}$$

Transfer functions learned from data:



# Other Roadblocks in Data Science Ventures



credit: oxford creativity, <https://www.triz.co.uk/>

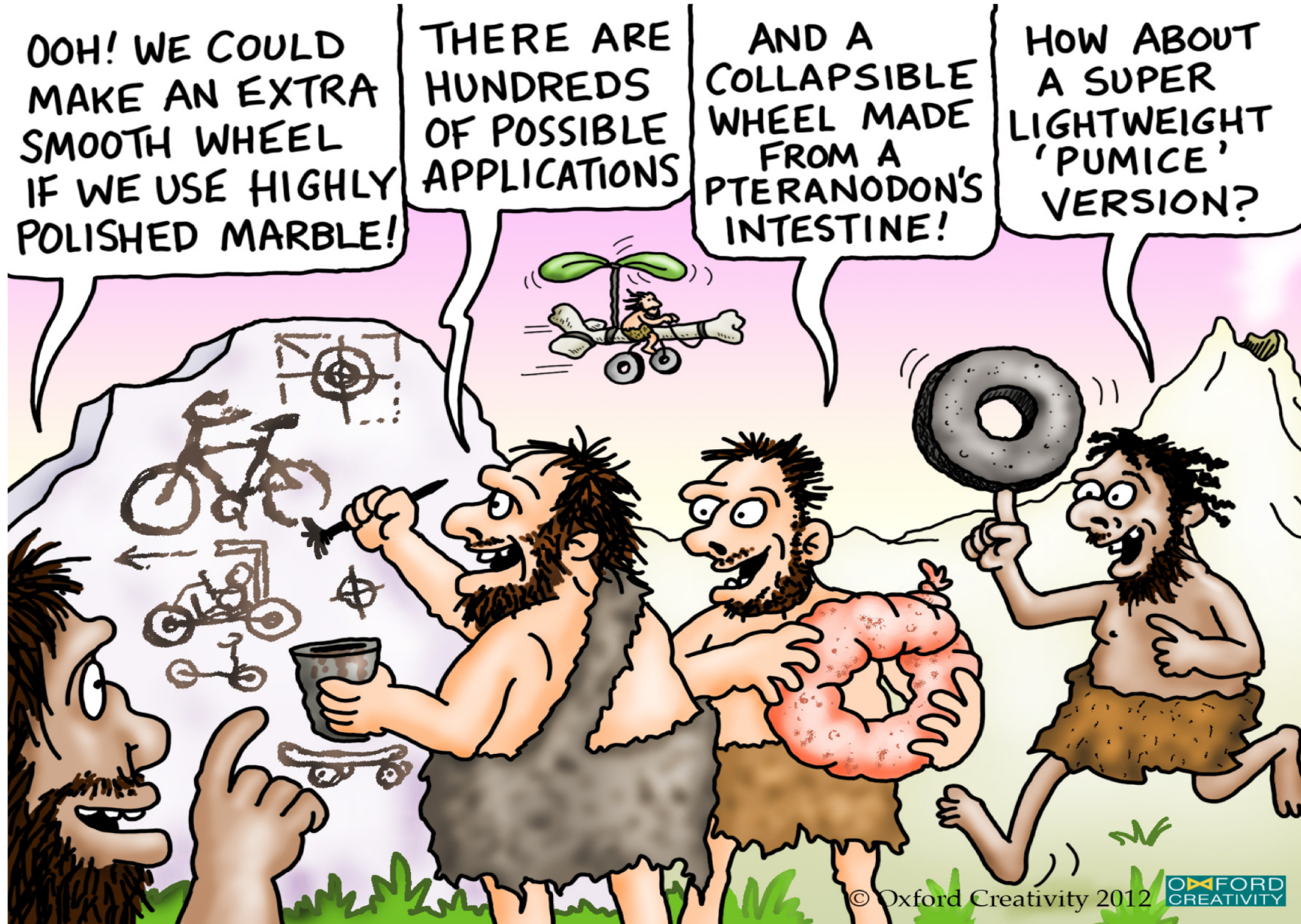


# Facilitate communication to foster innovation





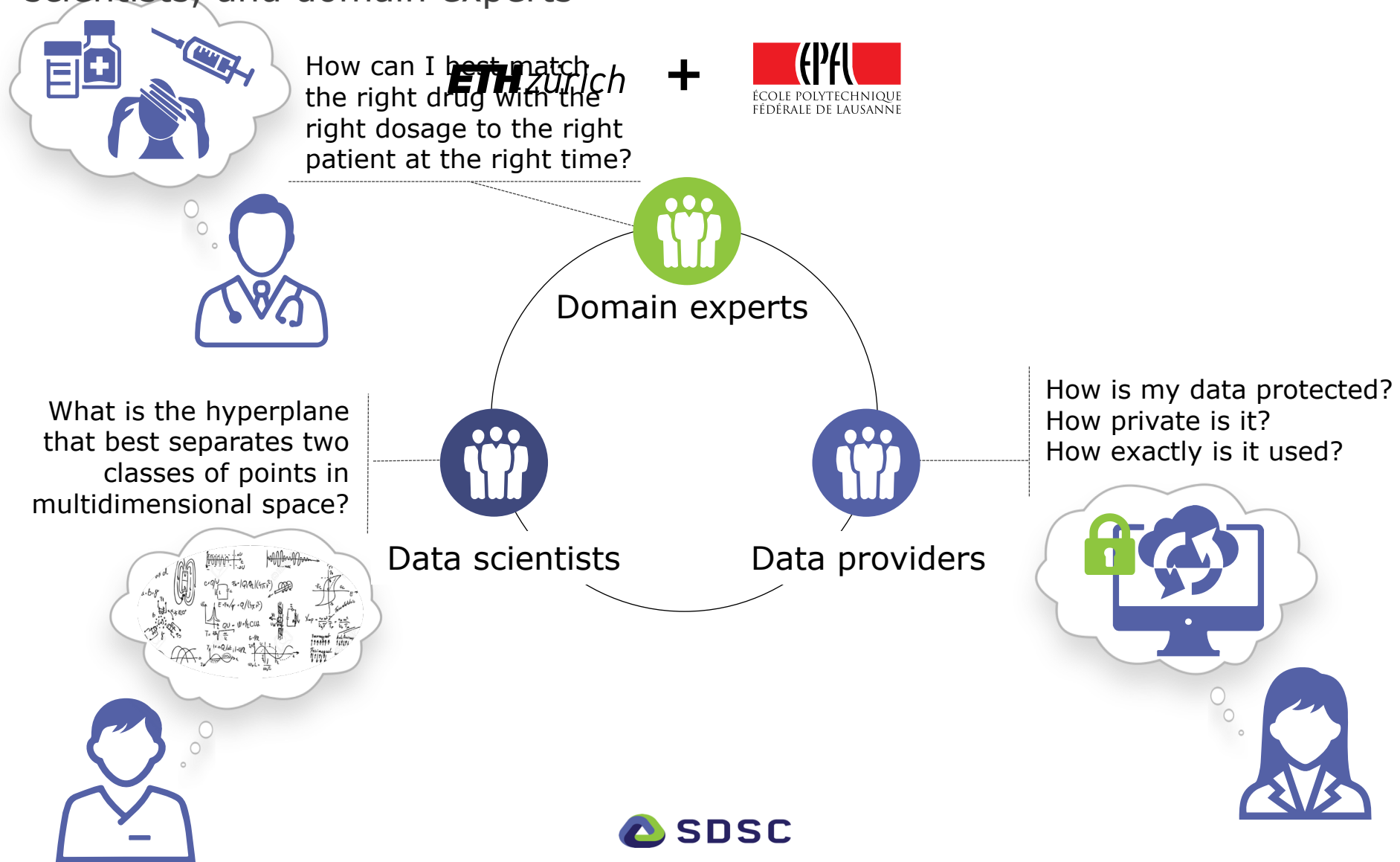
# Foster multidisciplinary collaborations



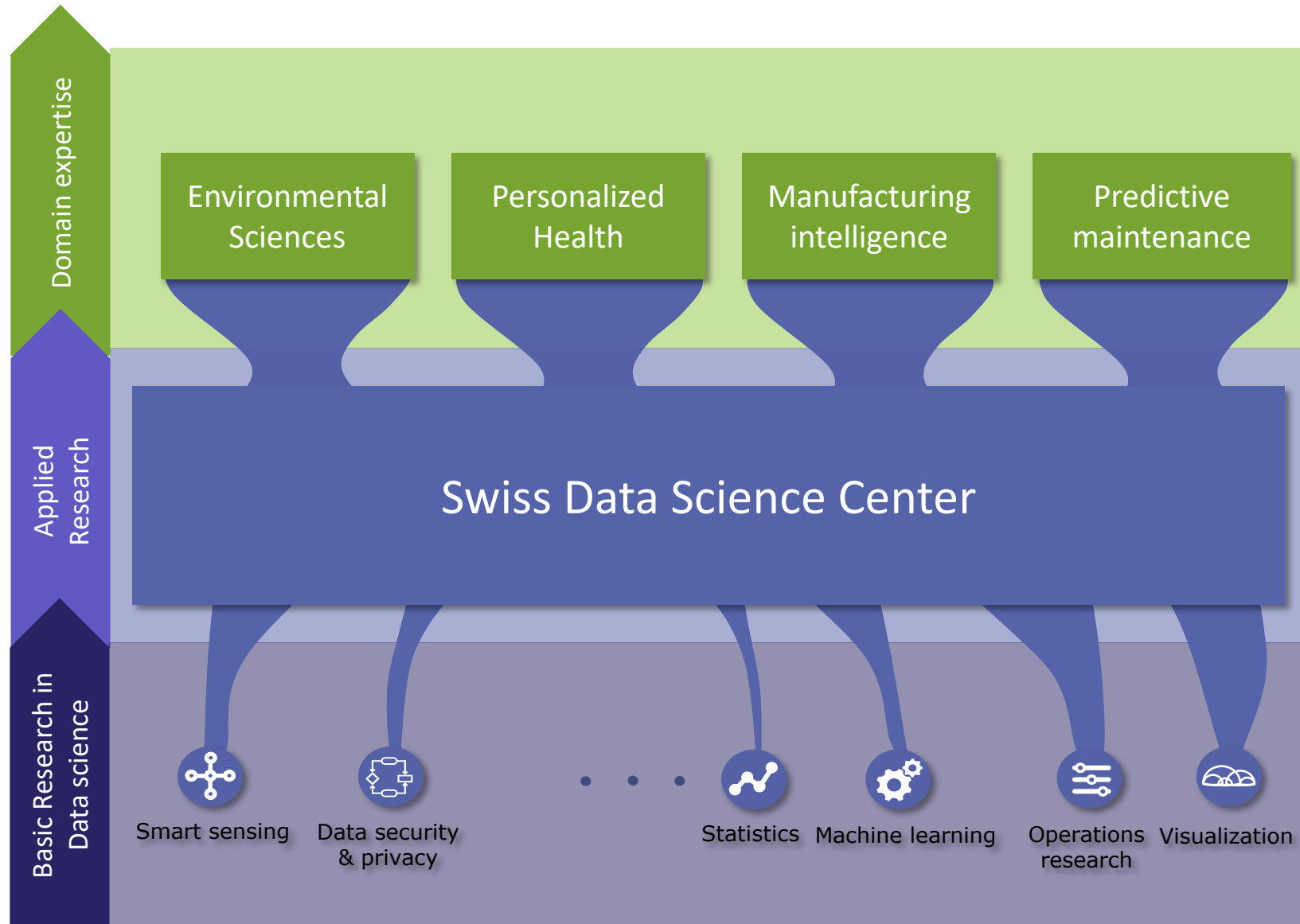
credit: oxford creativity, <https://www.triz.co.uk/>

# Swiss Data Science Center (SDSC)

Multi-disciplinary team of 40 full-time computer and data scientists, and domain experts



# Key Actor in a Complex Ecosystem





IN GOD WE TRUST.



ALL OTHERS MUST BRING DATA.

– W. EDWARDS DEMING, STATISTICIAN, PROFESSOR, AUTHOR