

# *Data Science for Digital Humanities*

BY  
HENDRIK HEUER

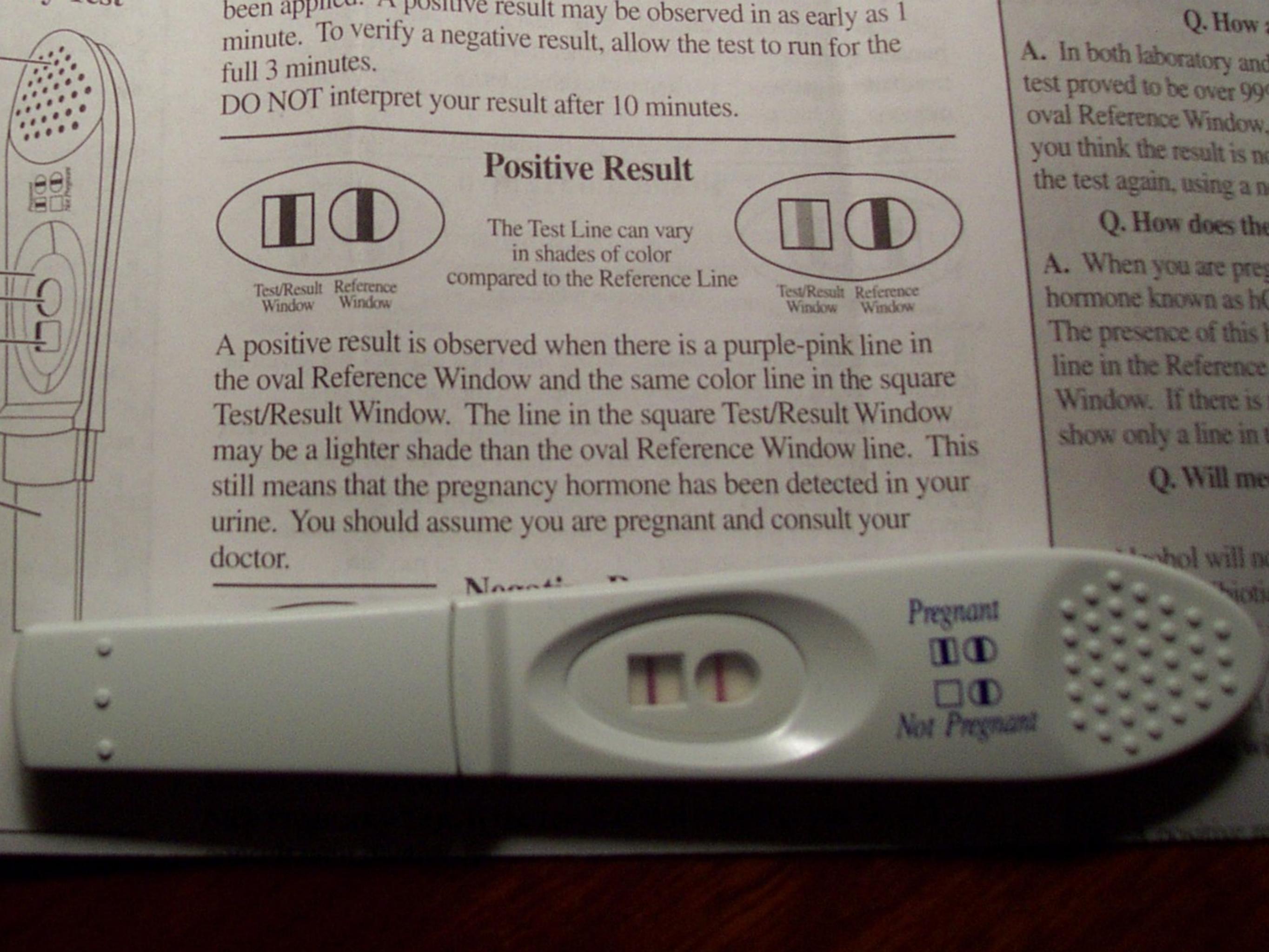
INSTITUTE FOR  
INFORMATION MANAGEMENT  
BREMEN GMBH



# About me

- **Hendrik** Heuer, MSc.
- Doctoral Researcher at the  
**Institute of Information  
Management Bremen (ifib)** at  
the **University of Bremen**
- Mail: hheuer@uni-bremen.de
- Website: <http://hen-drik.de>





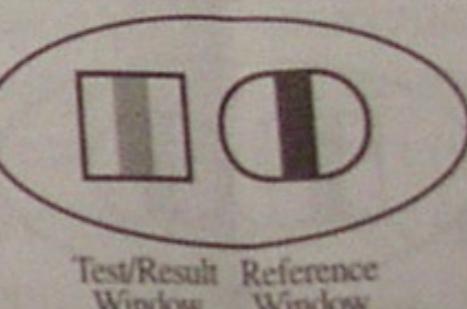
been applied. A positive result may be observed in as early as 1 minute. To verify a negative result, allow the test to run for the full 3 minutes.

DO NOT interpret your result after 10 minutes.

## Positive Result



The Test Line can vary  
in shades of color  
compared to the Reference Line



A positive result is observed when there is a purple-pink line in the oval Reference Window and the same color line in the square Test/Result Window. The line in the square Test/Result Window may be a lighter shade than the oval Reference Window line. This still means that the pregnancy hormone has been detected in your urine. You should assume you are pregnant and consult your doctor.

Q. How

A. In both laboratory and test proved to be over 99% oval Reference Window, you think the result is no the test again, using a ne

Q. How does the

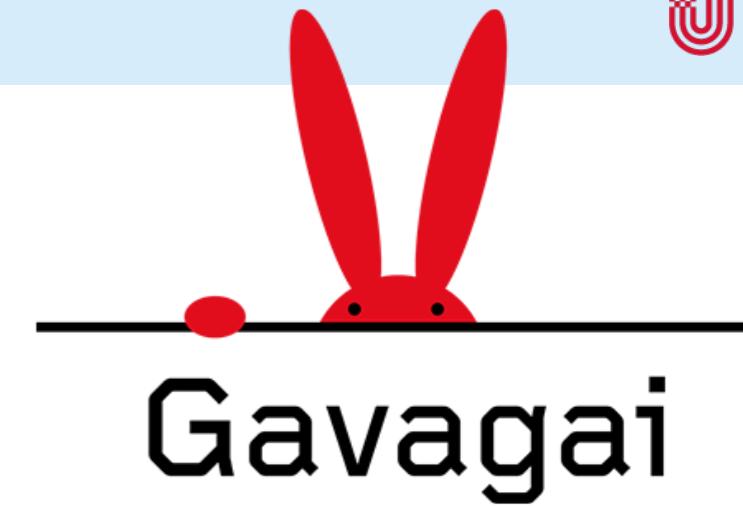
A. When you are preg  
hormone known as hCG

The presence of this h  
line in the Reference  
Window. If there is n  
show only a line in t

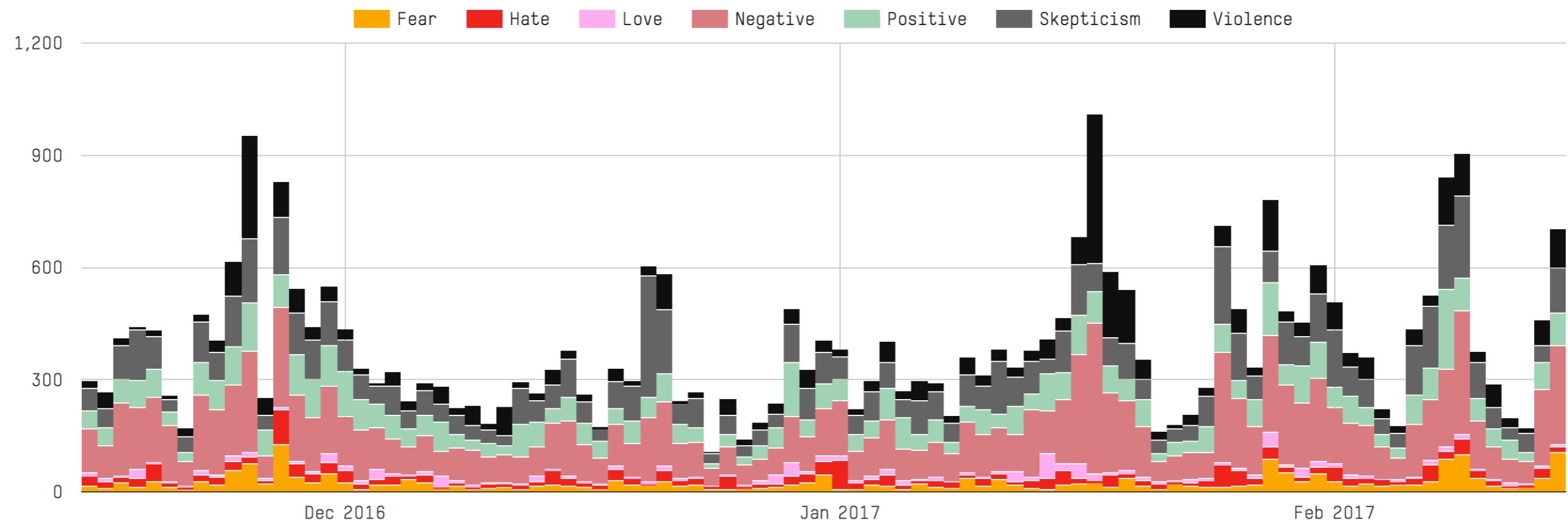
Q. Will me



# Sentiment Analysis

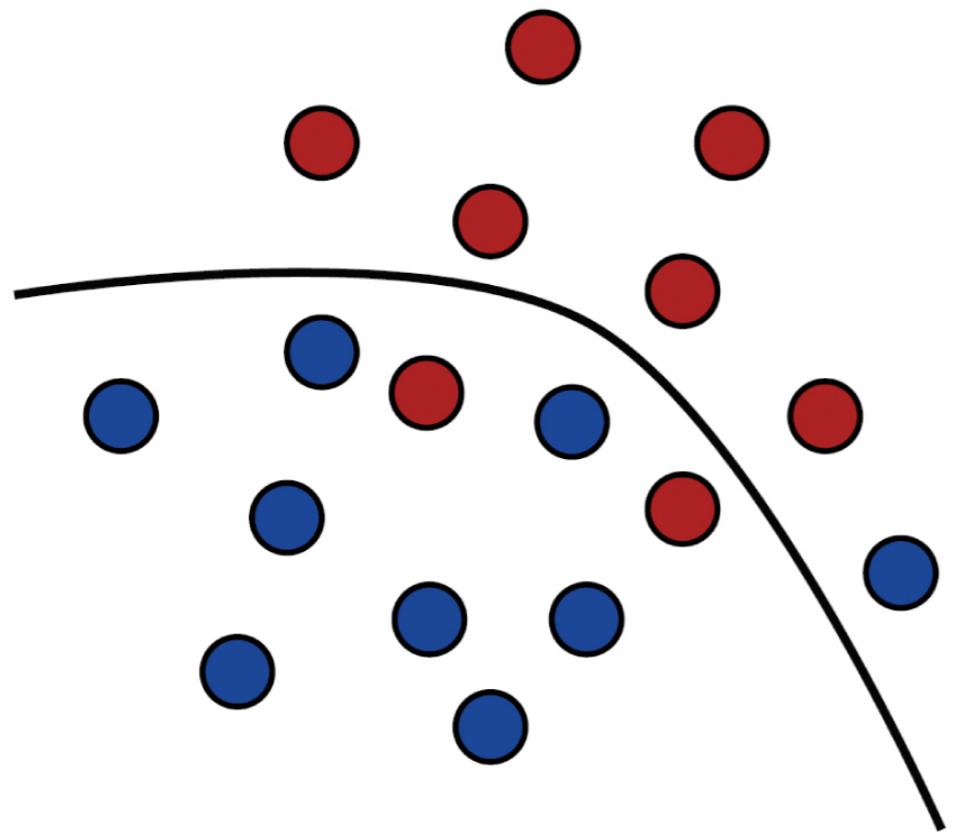


HOW PEOPLE IN SWEDEN FEEL  
ABOUT REFUGEES



**Analysing** millions of images  
and enormous text sources  
**using machine learning** and  
deep learning techniques  
**is simple** & straightforward in  
the Python ecosystem

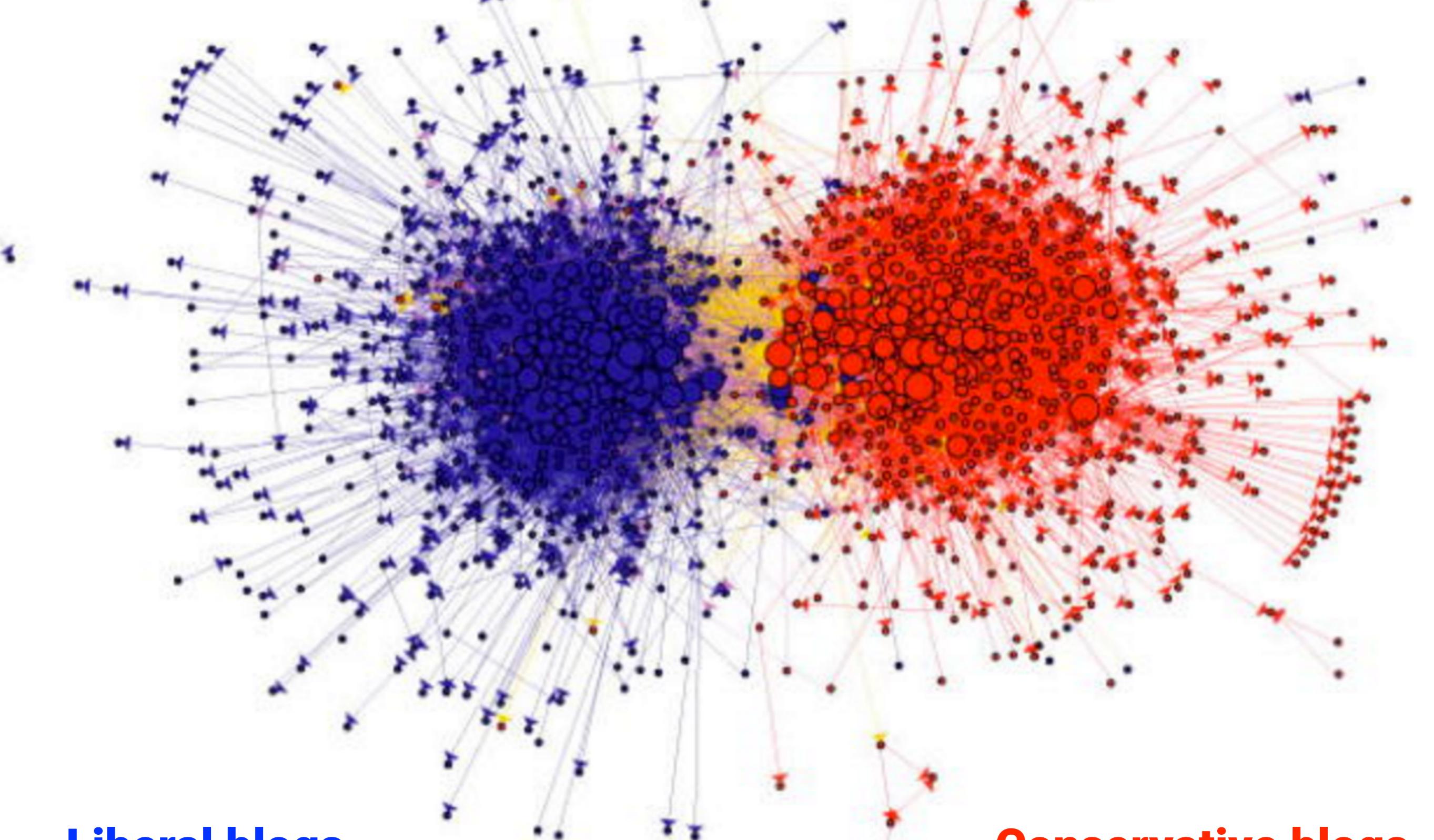
# "Data Science is **statistics** on a Mac."



**"A data scientist is someone who is better at statistics than any software engineer and better at software engineering than any statistician."**

# Data Science

- introduced by **Cleveland** in 2001
- enlargement of **statistics**
- enabling analysts to "**learn from data**"
- "**Sexiest job of the 21st century**"
  - **Davenport & Patil, 2012**  
**Harvard Business Review**



**Liberal blogs**

**Conservative blogs**

**liberal => conservative**

**liberal <= conservative**

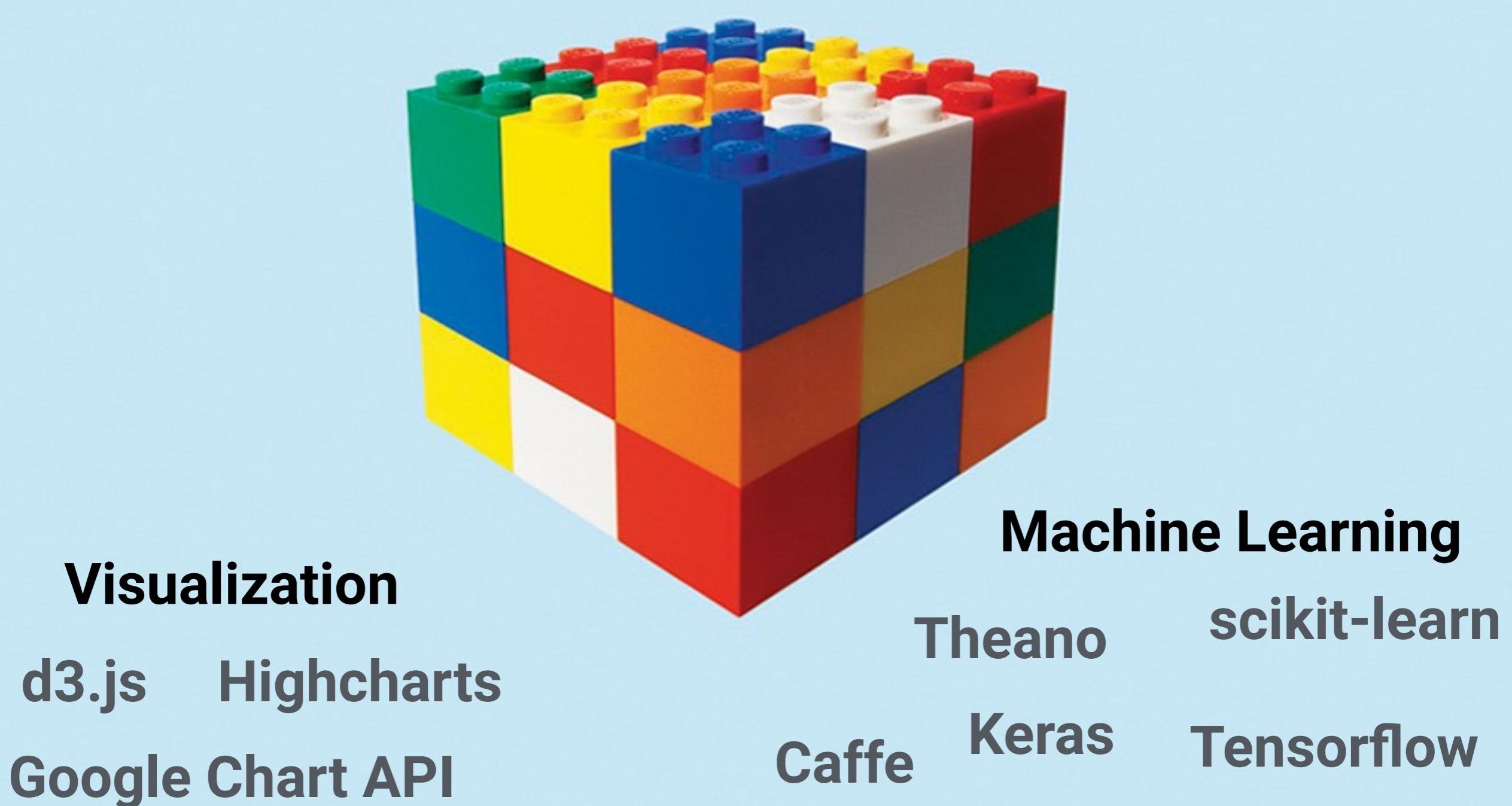
**2004**

*informatics = mathematics + x*

3 4 2 1 9 5 6 2 1 8  
8 9 1 2 5 0 0 6 6 4  
6 7 0 1 6 3 6 3 7 0  
3 7 7 9 4 6 6 1 8 2  
2 9 3 4 3 9 8 7 2 5  
1 5 9 8 3 6 5 7 2 3  
9 3 1 9 1 5 8 0 8 4  
5 6 2 6 8 5 8 8 9 9  
3 7 7 0 9 4 8 5 4 3  
7 9 6 4 7 0 6 9 2 3

**Text Processing**  
**Natural Language Toolkit**  
word2vec spaCy

**Topic Modeling**  
gensim

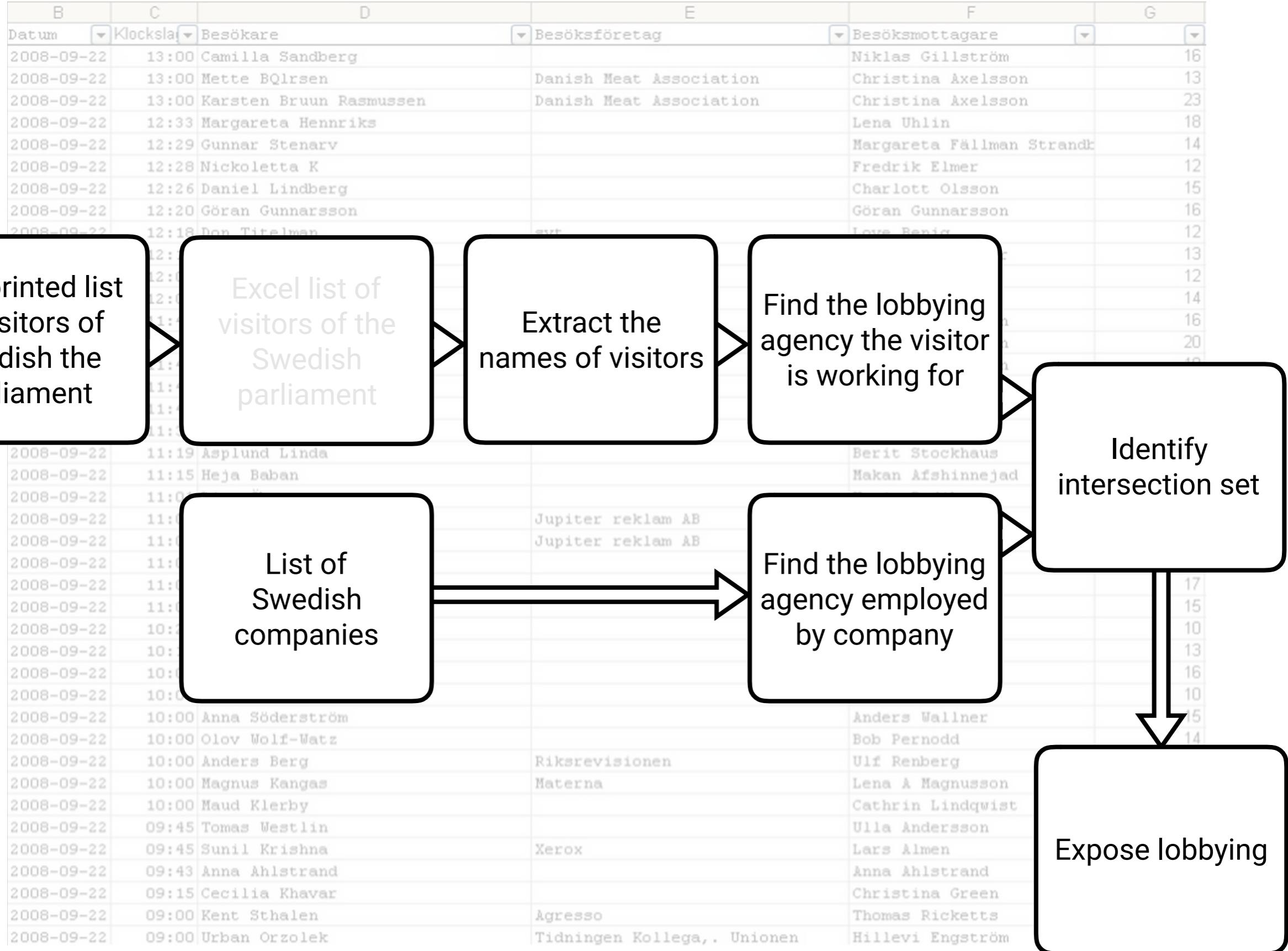


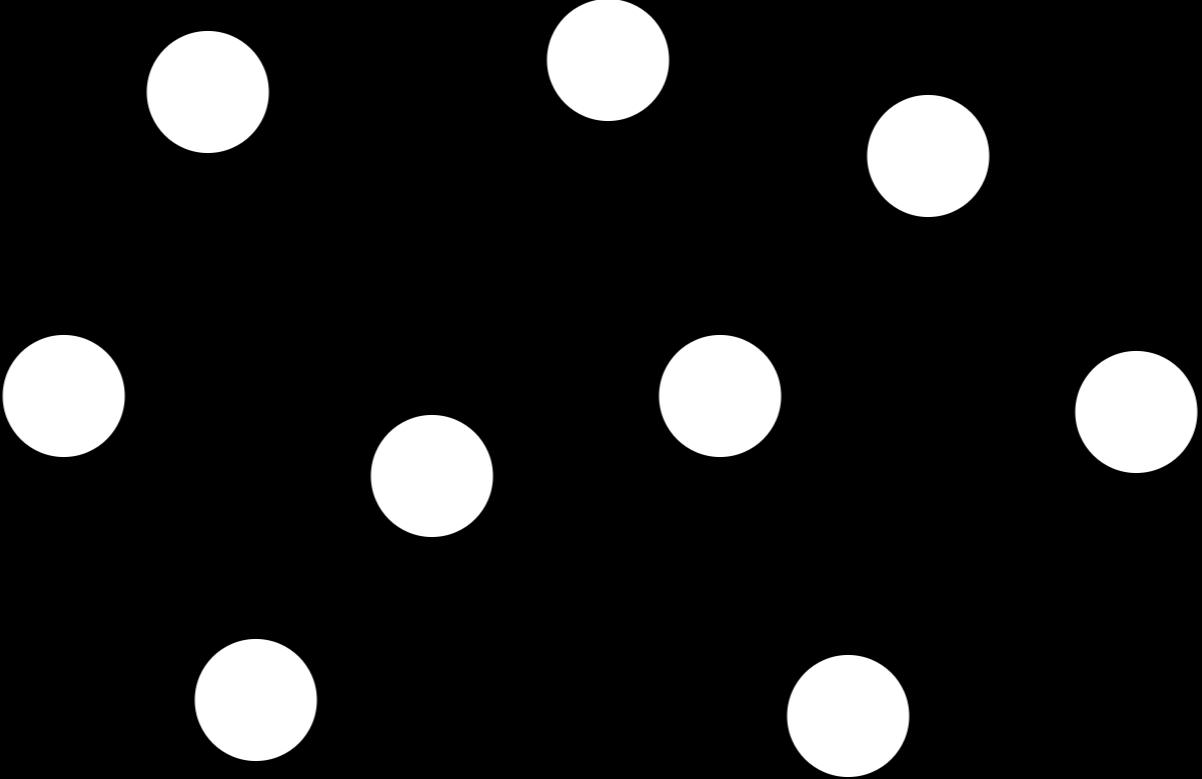
**Visualization**  
d3.js Highcharts  
Google Chart API

**Machine Learning**  
scikit-learn  
Theano  
Caffe Keras Tensorflow

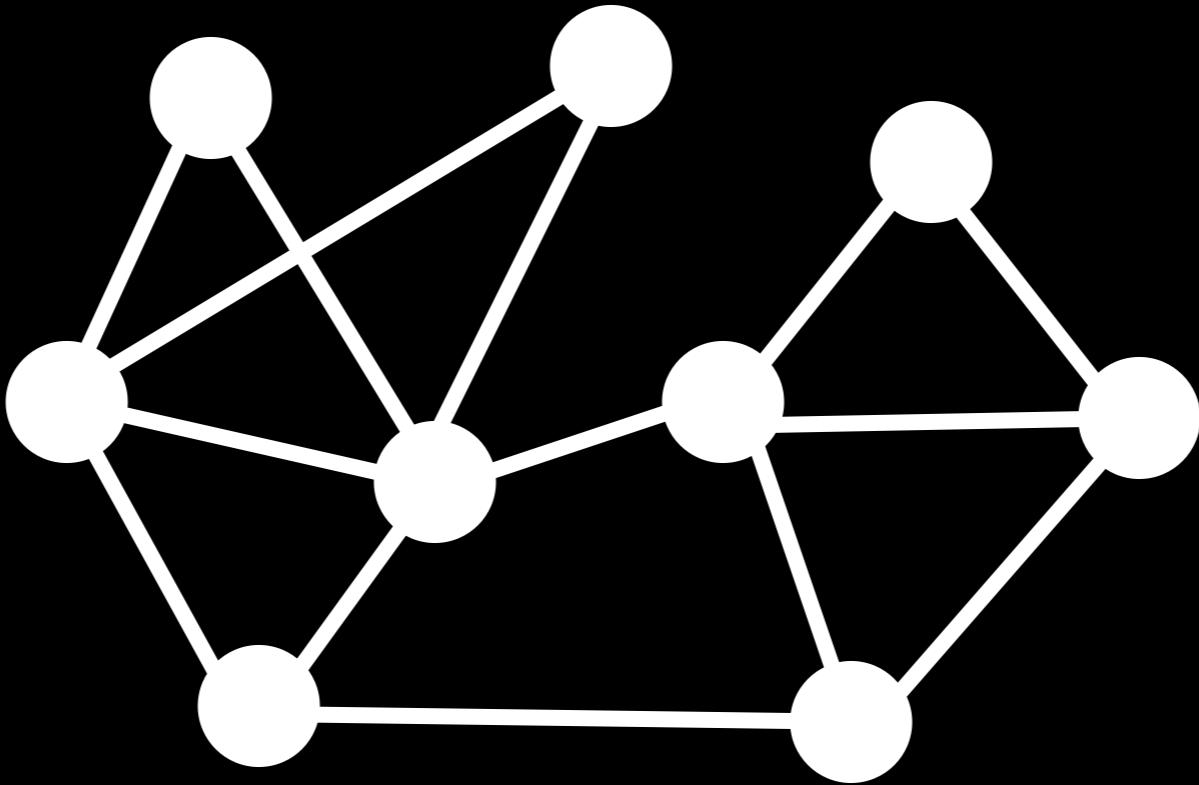


Imagine a journalist who wants to **cross-reference** the names on the guest list of a parliament with online information about **lobbyists** to identify which party meets which **company**





# Small-world experiments



# Small-world experiments

Karinthy (1929)  
Milgram (1967)  
Duncan Watt (2001),  
Leskovec and Horvitz (2007)

# Digital Humanities

- are the academic disciplines that study the **expressions of the human mind** (Rapport Duurzame Geesteswetenschappen, 2010)
- **intersection** of the **humanities** and **computer science** opens up **new research methods** and creates a new environment in which the humanities become subject to new approaches



# Topic Modelling (LDA)

## Topics

gene 0.04  
dna 0.02  
genetic 0.01  
...

life 0.02  
evolve 0.01  
organism 0.01  
...

brain 0.04  
neuron 0.02  
nerve 0.01  
...

data 0.02  
number 0.02  
computer 0.01  
...

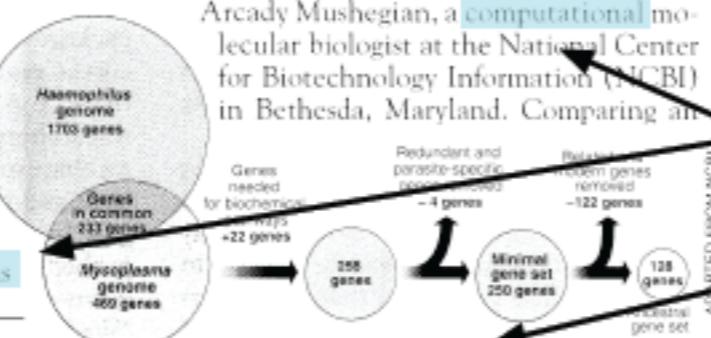
## Documents

### Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

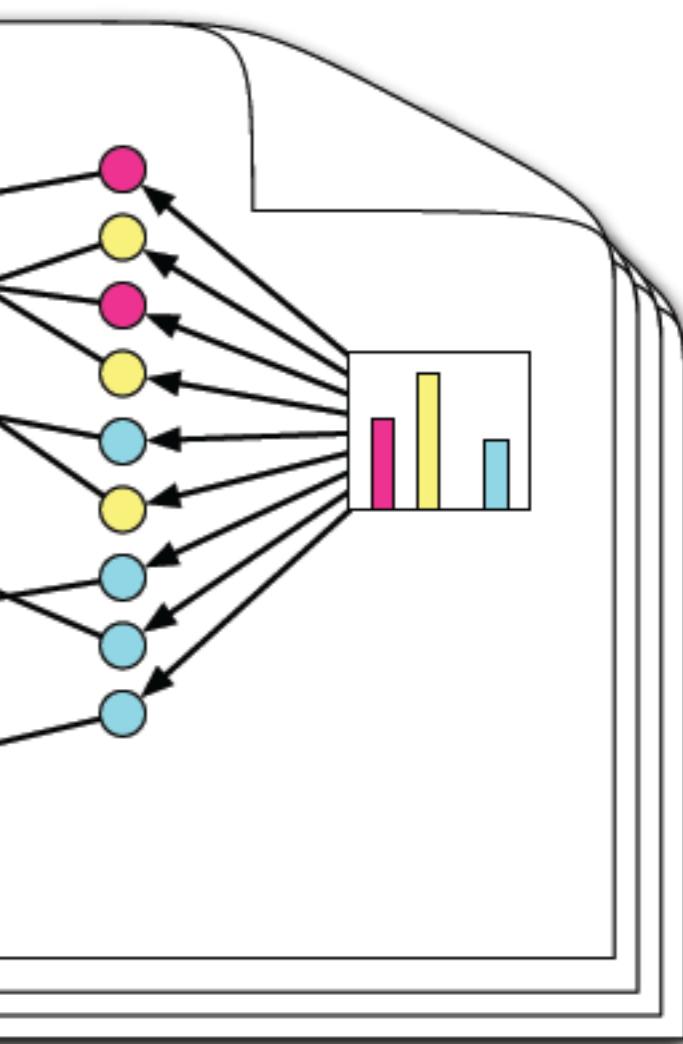
"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game; particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

## Topic proportions and assignments



Blei, David M., Andrew Y. Ng, and Michael I. Jordan.  
"Latent Dirichlet Allocation."  
Journal of Machine Learning Research 3.Jan (2003): 993-1022.

# Latent Dirichlet Allocation

- I ate a banana and spinach smoothie for breakfast
- I like to eat broccoli and bananas.
- Chinchillas and kittens are cute.
- My sister adopted a kitten yesterday.
- Look at this cute hamster munching on a piece of broccoli.

# Latent Dirichlet Allocation

- I ate a banana and spinach smoothie for breakfast
- I like to eat broccoli and bananas.
- Chinchillas and kittens are cute.
- My sister adopted a kitten yesterday.
- Look at this cute hamster munching on a piece of broccoli.
- Topic A: food
- Topic B: cute animals

# Latent Dirichlet Allocation

- I ate a banana and spinach smoothie for breakfast
- I like to eat broccoli and bananas.
- Chinchillas and kittens are cute.
- My sister adopted a kitten yesterday.
- Look at this cute hamster munching on a piece of broccoli.
- Topic A: food
- Topic B: cute animals
- Sentences 1 and 2: 100% Topic A
- Sentences 3 and 4: 100% Topic B
- Sentence 5: 60% Topic A, 40% Topic B

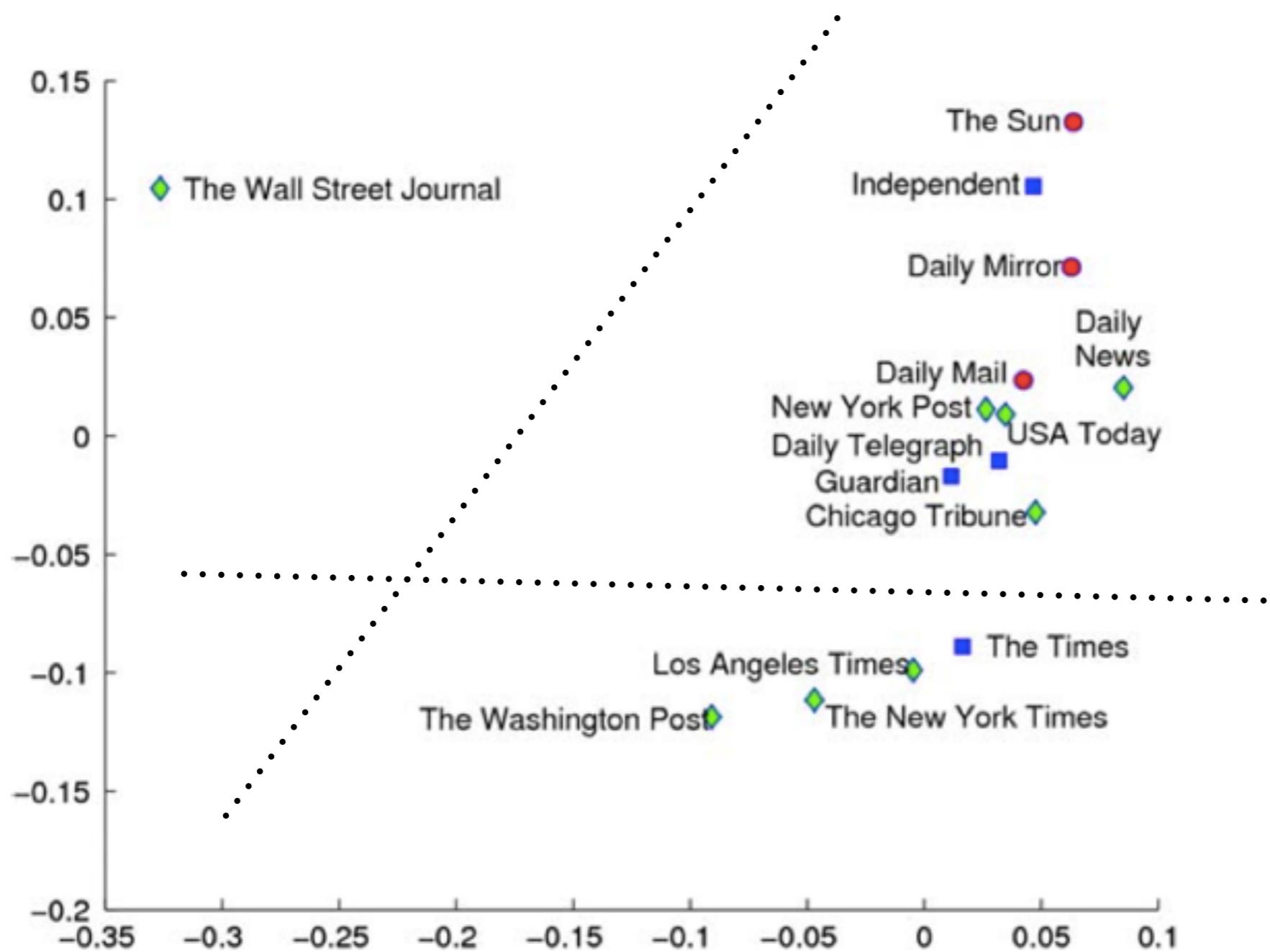
# Latent Dirichlet Allocation

- I ate a banana and spinach smoothie for breakfast
- I like to eat broccoli and bananas.
- Chinchillas and kittens are cute.
- My sister adopted a kitten yesterday.
- Look at this cute hamster munching on a piece of broccoli.
- Topic A: 30% broccoli, 15% bananas, 10% breakfast, 10% munching
- Topic B: 20% chinchillas, 20% kittens, 20% cute, 15% hamster

# Massive-scale automated analysis of news-content

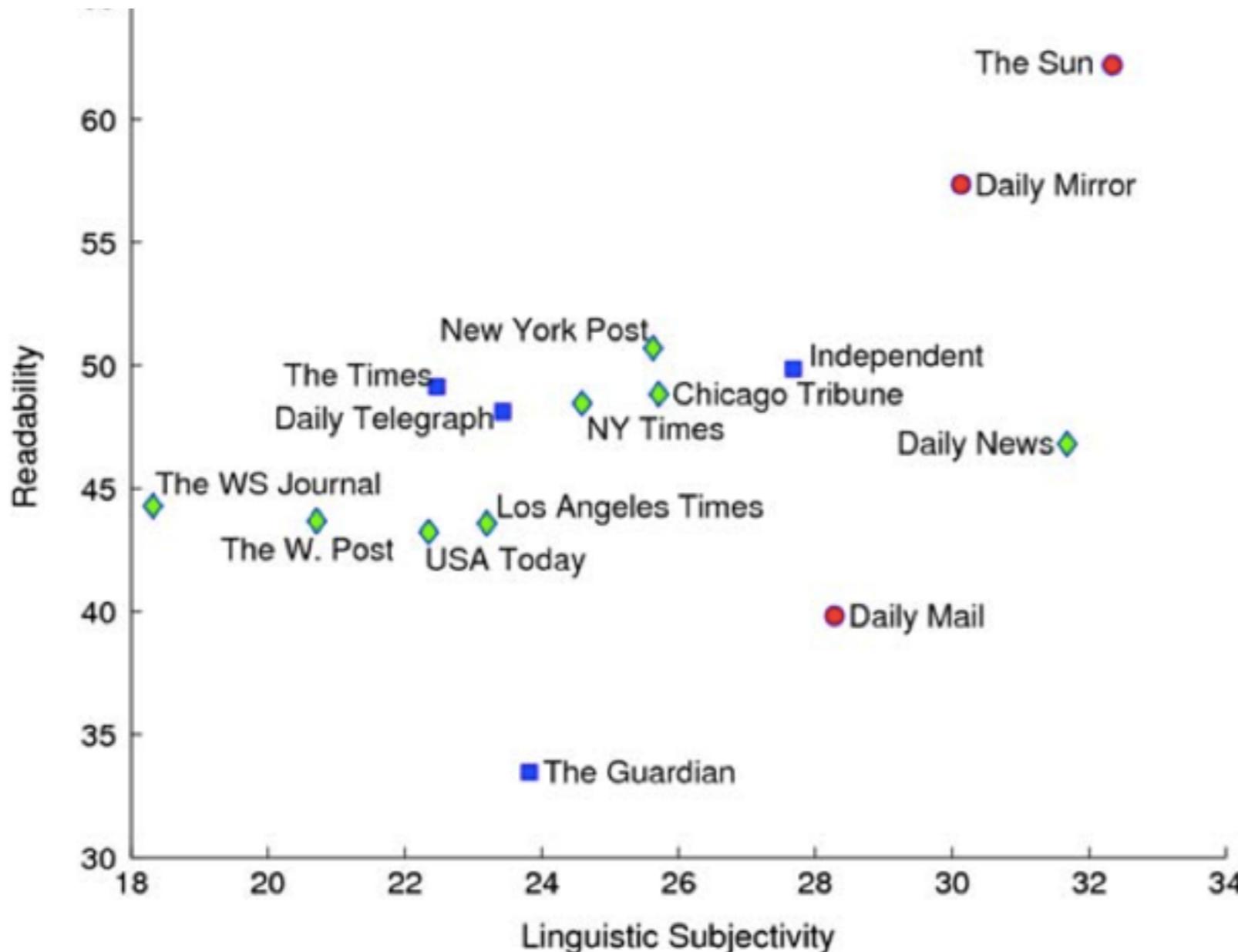
- 2.5 million articles from 498 different English-language news outlets
  - Reuters
  - New York Times Corpus
- automatically annotated into 15 topic areas
- the topics were compared in regards to readability, linguistic subjectivity and gender imbalances

# Comparing Newspapers



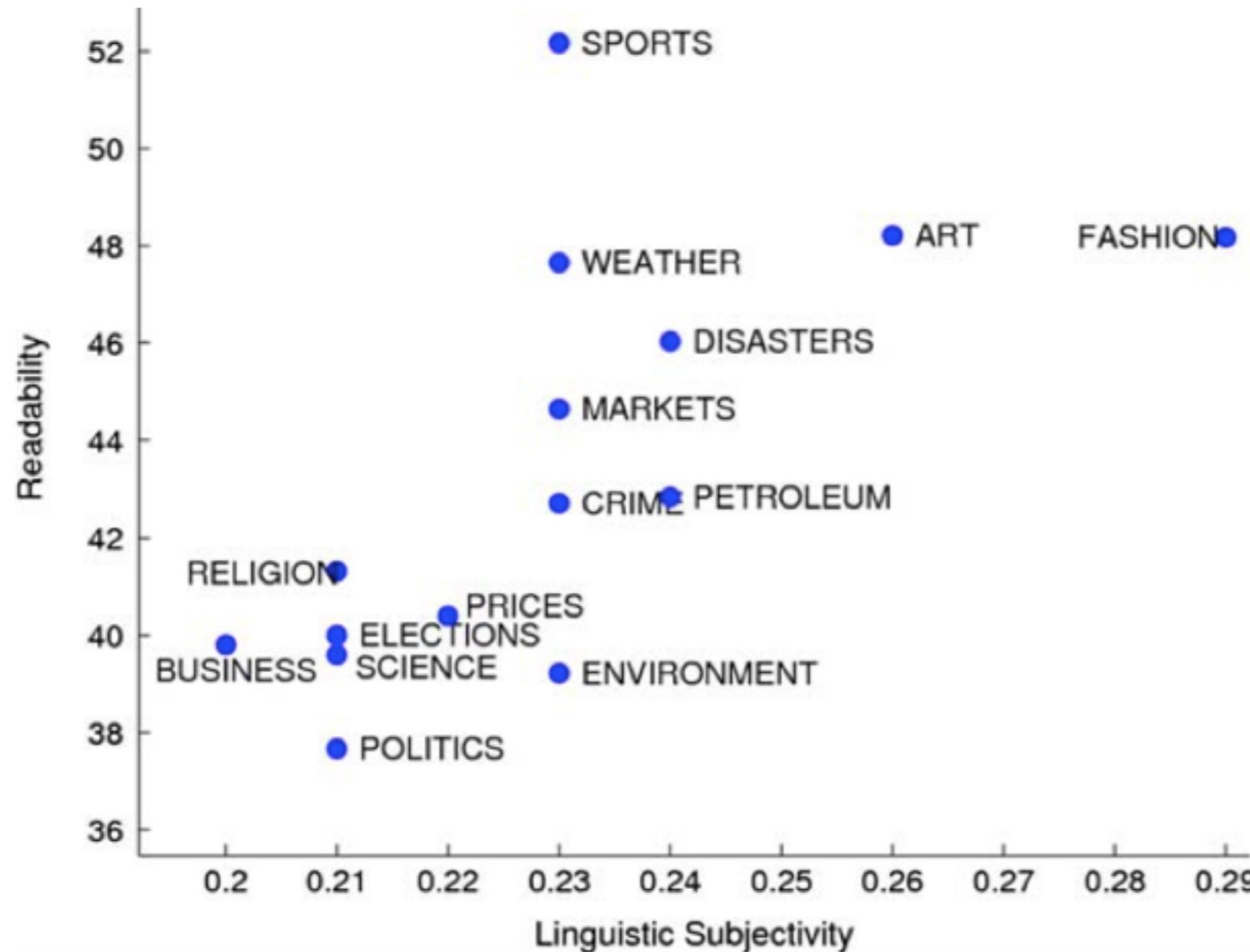
I. Flaounas, O. Ali, T. Lansdall-Welfare, T. De Bie, N. Mosdell, J. Lewis, and N. Cristianini, 'Research Methods in the Age of Digital Journalism: Massive-scale automated analysis of news-content: topics, style and gender', *Digital Journalism*, vol. 1, no. 1, 2013. DOI:10.1080/21670811.2012.714928

# Comparing Newspapers



I. Flaounas, O. Ali, T. Lansdall-Welfare, T. De Bie, N. Mosdell, J. Lewis, and N. Cristianini, 'Research Methods in the Age of Digital Journalism: Massive-scale automated analysis of news-content: topics, style and gender', *Digital Journalism*, vol. 1, no. 1, 2013. DOI:10.1080/21670811.2012.714928

# Comparing Topics



I. Flaounas, O. Ali, T. Lansdall-Welfare, T. De Bie, N. Mosdell, J. Lewis, and N. Cristianini, 'Research Methods in the Age of Digital Journalism: Massive-scale automated analysis of news-content: topics, style and gender', *Digital Journalism*, vol. 1, no. 1, 2013. DOI:10.1080/21670811.2012.714928

# Macroanalysis

- Jockers uses digital methods for **literary history**
- showing that factors such as author **gender**, author **nationality**, and **date of publication** affect the **choice of literary themes in novels**
- measurable, data-driven **proxy for literary themes**
- identified and extracted hundreds of topics from a corpus of **3346 works of 19th-century British, Irish, and American fiction**

# Topic Modelling (LDA)





Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman and Alexei A. Efros.  
"Generative Visual Manipulation on the Natural Image Manifold",  
in European Conference on Computer Vision (ECCV). 2016.

Image: [https://en.wikipedia.org/wiki/Gezi\\_Park\\_protests](https://en.wikipedia.org/wiki/Gezi_Park_protests)



Figure 1: The frequency of top 20 hashtags associated with Gezi Protests. (Banko and Babacan, 2013)

Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman and Alexei A. Efros.  
"Generative Visual Manipulation on the Natural Image Manifold",  
in European Conference on Computer Vision (ECCV). 2016.

Image: [https://en.wikipedia.org/wiki/Gezi\\_Park\\_protests](https://en.wikipedia.org/wiki/Gezi_Park_protests)

## LETTERS

# Detecting influenza epidemics using search engine query data

Jeremy Ginsberg<sup>1</sup>, Matthew H. Mohebbi<sup>1</sup>, Rajan S. Patel<sup>1</sup>, Lynnette Brammer<sup>2</sup>, Mark S. Smolinski<sup>1</sup> & Larry Brilliant<sup>1</sup>

Seasonal influenza epidemics are a major public health concern, causing tens of millions of respiratory illnesses and 250,000 to 500,000 deaths worldwide each year<sup>1</sup>. In addition to seasonal influenza, a new strain of influenza virus against which no previous immunity exists and that demonstrates human-to-human transmission could result in a pandemic with millions of fatalities<sup>2</sup>. Early detection of disease activity, when followed by a rapid response, can reduce the impact of both seasonal and pandemic influenza<sup>3,4</sup>. One way to improve early detection is to monitor health-seeking behaviour in the form of queries to online search engines, which are submitted by millions of users around the world each day. Here we present a method of analysing large numbers of Google search queries to track influenza-like illness in a population. Because the relative frequency of certain queries is highly correlated with the percentage of physician visits in which a patient presents with influenza-like symptoms, we can accurately estimate the current level of weekly influenza activity in each region of the United States, with a reporting lag of about one day. This approach may make it possible to use search queries to detect influenza epidemics in areas with a large population of web search users.

Traditional surveillance systems, including those used by the US Centers for Disease Control and Prevention (CDC) and the European Influenza Surveillance Scheme (EISS), rely on both virological and clinical data, including influenza-like illness (ILI) physician visits. The CDC publishes national and regional data from these surveillance systems on a weekly basis, typically with a 1–2-week reporting lag.

In an attempt to provide faster detection, innovative surveillance systems have been created to monitor indirect signals of influenza activity, such as call volume to telephone triage advice lines<sup>5</sup> and over-the-counter drug sales<sup>6</sup>. About 90 million American adults are believed to search online for information about specific diseases or medical problems each year<sup>7</sup>, making web search queries a uniquely valuable source of information about health trends. Previous

By aggregating historical logs of online web search queries submitted between 2003 and 2008, we computed a time series of weekly counts for 50 million of the most common search queries in the United States. Separate aggregate weekly counts were kept for every query in each state. No information about the identity of any user was retained. Each time series was normalized by dividing the count for each query in a particular week by the total number of online search queries submitted in that location during the week, resulting in a query fraction (Supplementary Fig. 1).

We sought to develop a simple model that estimates the probability that a random physician visit in a particular region is related to an ILI; this is equivalent to the percentage of ILI-related physician visits. A single explanatory variable was used: the probability that a random search query submitted from the same region is ILI-related, as determined by an automated method described below. We fit a linear model using the log-odds of an ILI physician visit and the log-odds of an ILI-related search query:  $\text{logit}(I(t)) = \alpha \text{logit}(Q(t)) + \varepsilon$ , where  $I(t)$  is the percentage of ILI physician visits,  $Q(t)$  is the ILI-related query fraction at time  $t$ ,  $\alpha$  is the multiplicative coefficient, and  $\varepsilon$  is the error term.  $\text{logit}(p)$  is simply  $\ln(p/(1-p))$ .

Publicly available historical data from the CDC's US Influenza Sentinel Provider Surveillance Network (<http://www.cdc.gov/flu/weekly>) was used to help build our models. For each of the nine surveillance regions of the United States, the CDC reported the average percentage of all outpatient visits to sentinel providers that were ILI-related on a weekly basis. No data were provided for weeks outside of the annual influenza season, and we excluded such dates from model fitting, although our model was used to generate unvalidated ILI estimates for these weeks.

We designed an automated method of selecting ILI-related search queries, requiring no previous knowledge about influenza. We measured how effectively our model would fit the CDC ILI data in each region if we used only a single query as the explanatory variable,  $Q(t)$ . Each of the 50 million candidate queries in our database was sepa-

## BIG DATA

# The Parable of Google Flu: Traps in Big Data Analysis

David Lazer,<sup>1,2\*</sup> Ryan Kennedy,<sup>1,3,4</sup> Gary King,<sup>3</sup> Alessandro Vespignani<sup>3,5,6</sup>

In February 2013, Google Flu Trends (GFT) made headlines but not for a reason that Google executives or the creators of the flu tracking system would have hoped. *Nature* reported that GFT was predicting more than double the proportion of doctor visits for influenza-like illness (ILI) than the Centers for Disease Control and Prevention (CDC), which bases its estimates on surveillance reports from laboratories across the United States (1, 2). This happened despite the fact that GFT was built to predict CDC reports. Given that GFT is often held up as an exemplary use of big data (3, 4), what lessons can we draw from this error?

The problems we identify are not limited to GFT. Research on whether search or social media can predict  $x$  has become commonplace (5–7) and is often put in sharp contrast with traditional methods and hypotheses.



Large errors in flu prediction were largely avoidable, which offers lessons for the use of big data.

the algorithm in 2009, and this model has run ever since, with a few changes announced in October 2013 (10, 15).

Although not widely reported until 2013, the new GFT has been persistently overestimating flu prevalence for a much longer time. GFT also missed by a very large margin in the 2011–2012 flu season and has missed high for 100 out of 108 weeks starting with August 2011 (see the graph). These errors are not randomly distributed. For example, last week's errors predict this week's errors (temporal autocorrelation), and the direction and magnitude of error varies with the time of year (seasonality). These patterns mean that GFT overlooks considerable information that could be extracted by traditional statistical methods.

Even after GFT was updated in 2009, the comparative value of the algorithm as a

surement and construct validity and reliability and dependencies among data (12).

# Big Data Hubris

was to find the best matches among 50 million search terms to fit 1152 data points (13). The odds of finding search terms that match the propensity of the flu but are structurally unrelated, and so do not predict the future, are quite high. GFT has

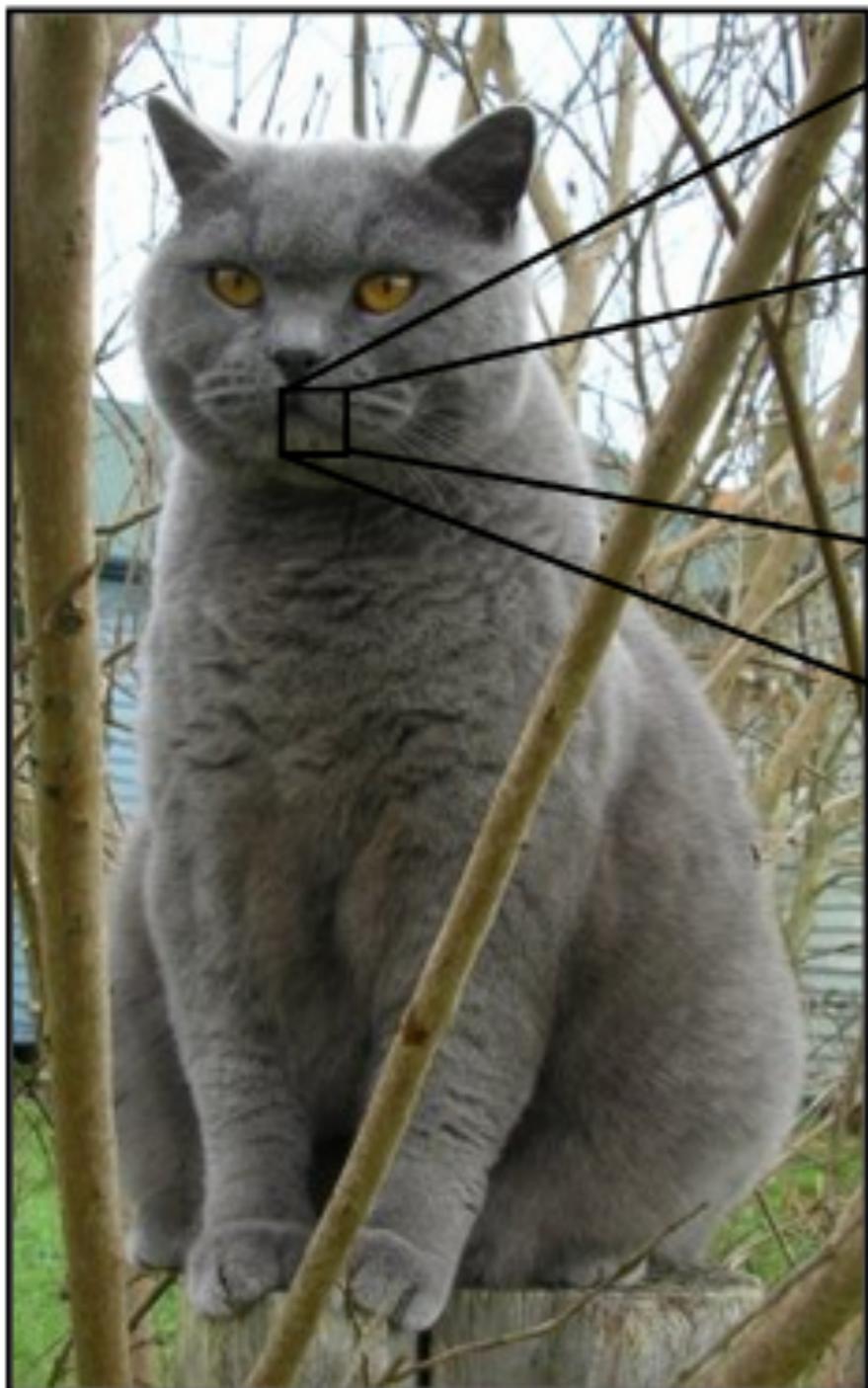
Even 3-week-old CDC data do a better job of projecting current flu prevalence than GFT [see supplementary materials (SM)].

Considering the large number of approaches that provide inference on influenza activity (16–18), there is no reason that

## Big Data Hubris

“Big data hubris” is the often implicit assumption that big data are a substitute for, rather than a supplement to, traditional data collection and analysis. Flu

# What the Computer sees



08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	91	65
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	48	64	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	50	88	30	03	49	13	36	65
52	70	95	23	04	60	11	42	62	21	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	03	89	41	92	36	54	22	40	40	28	66	33	13	80
24	47	38	60	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
00	96	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	52	25	39	11	24	94	72	18	08	46	29	32	40	62	76	36
20	69	36	41	72	30	23	88	31	60	99	69	82	67	59	85	74	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	48	66	41	16	23	57	05	54
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	25	17	48

What the computer sees

image classification

82% cat  
15% dog  
2% hat  
1% mug

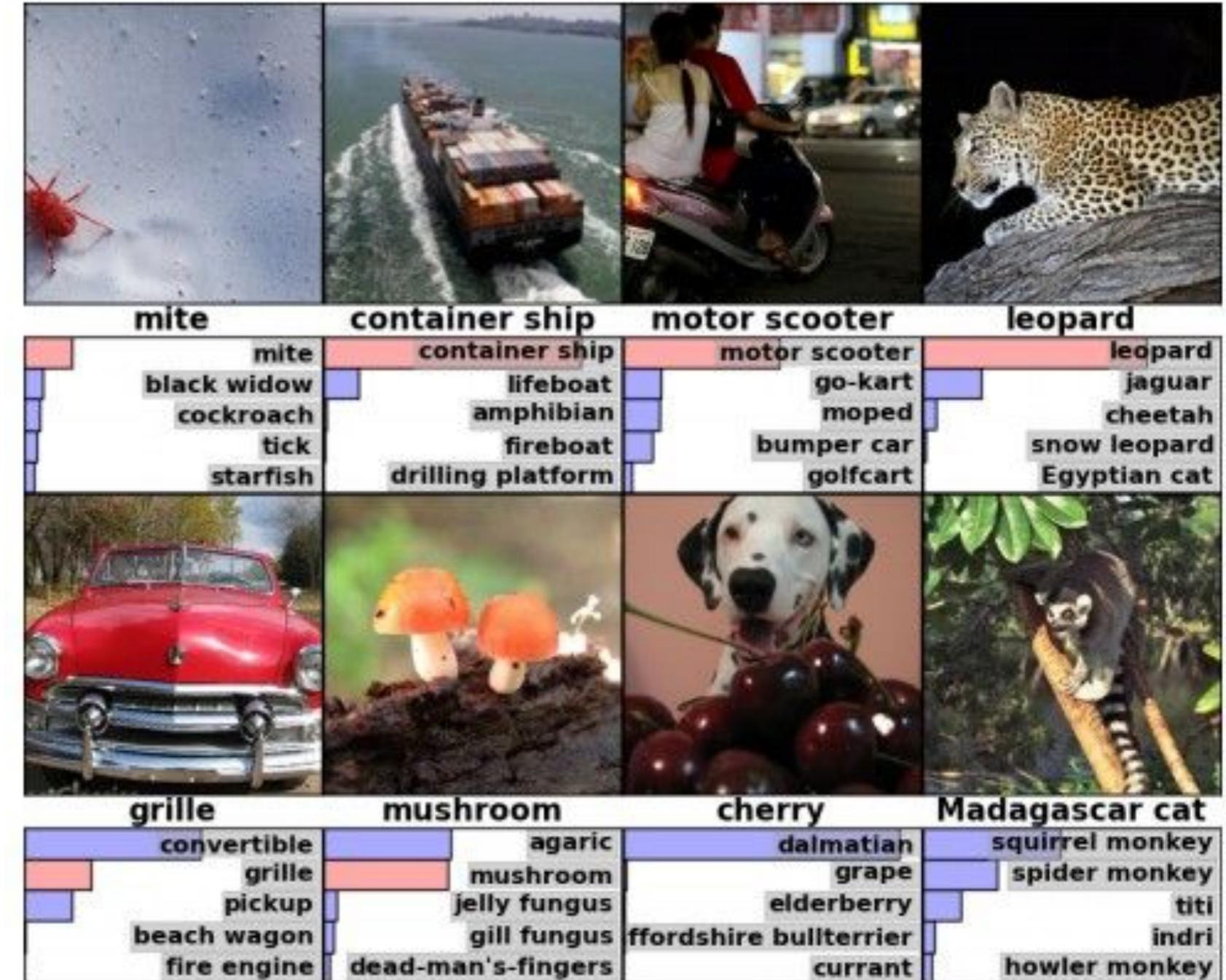
# Computer Vision

Source: Fei-Fei Li &  
Justin Johnson &  
Serena Yeung CS231n,  
Stanford University

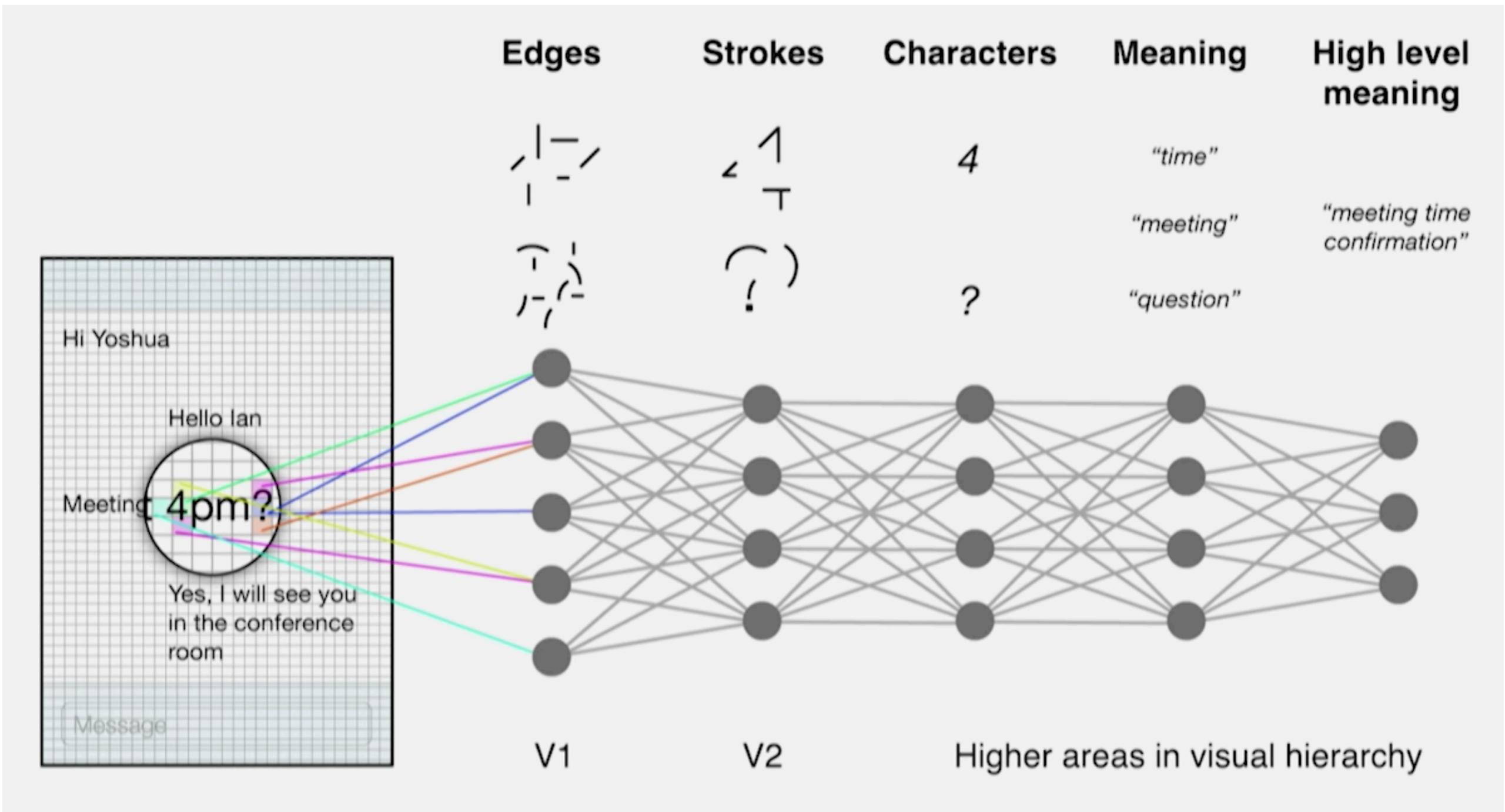
**ImageNet**  
1000 object class  
(categories)

Images for training:  
1,2 millions

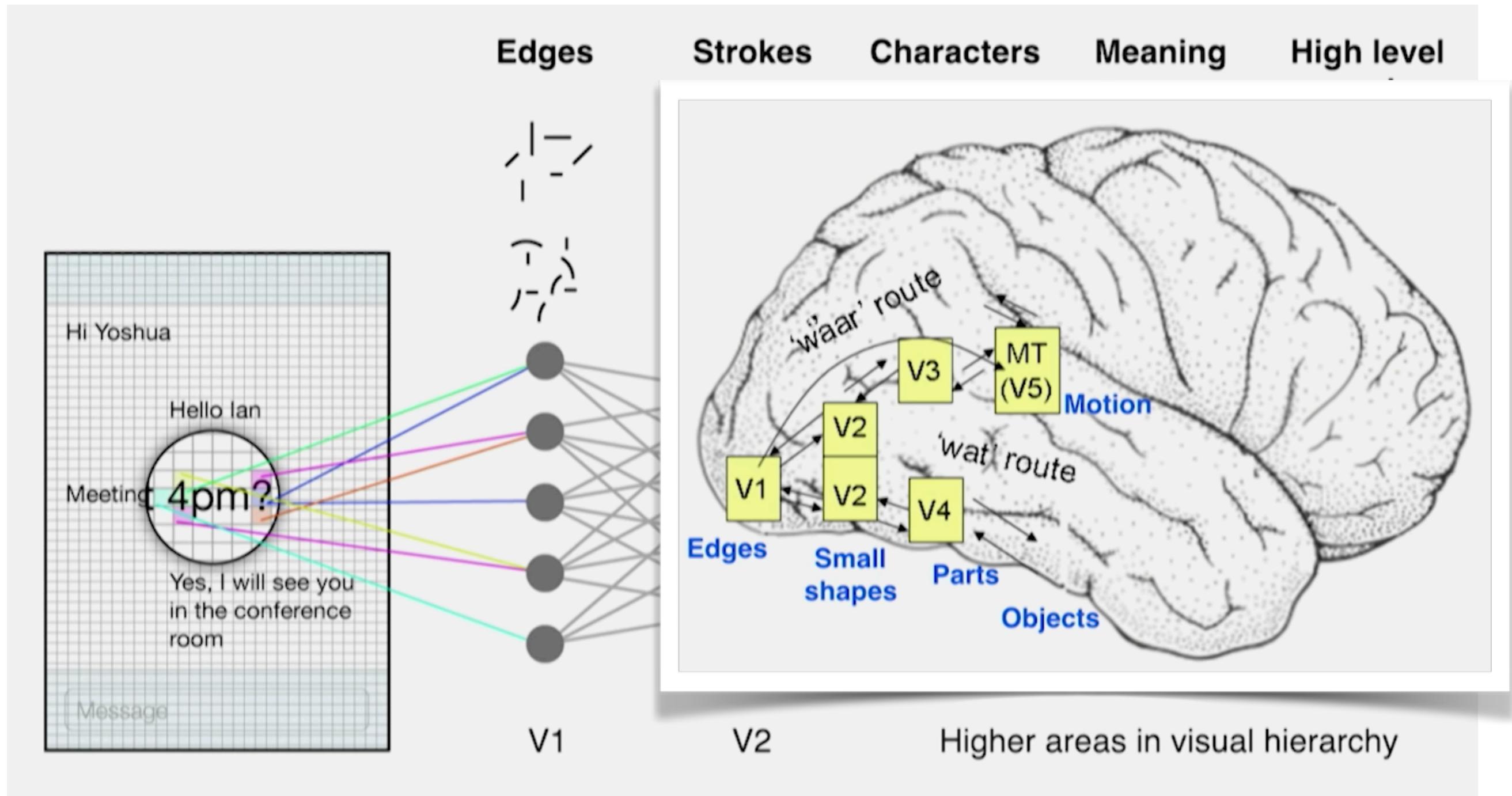
Images for testing:  
100k

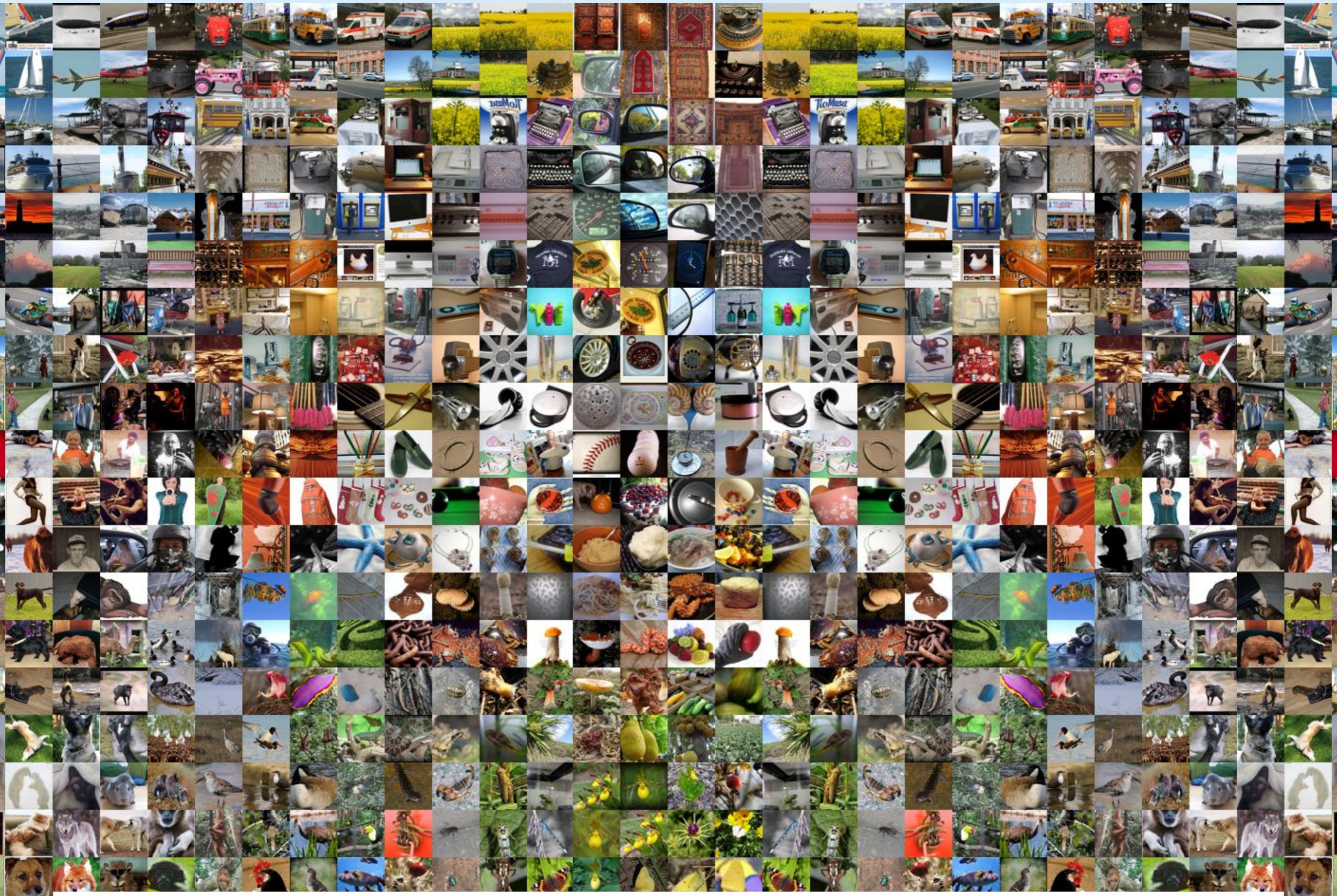


# Representation Learning



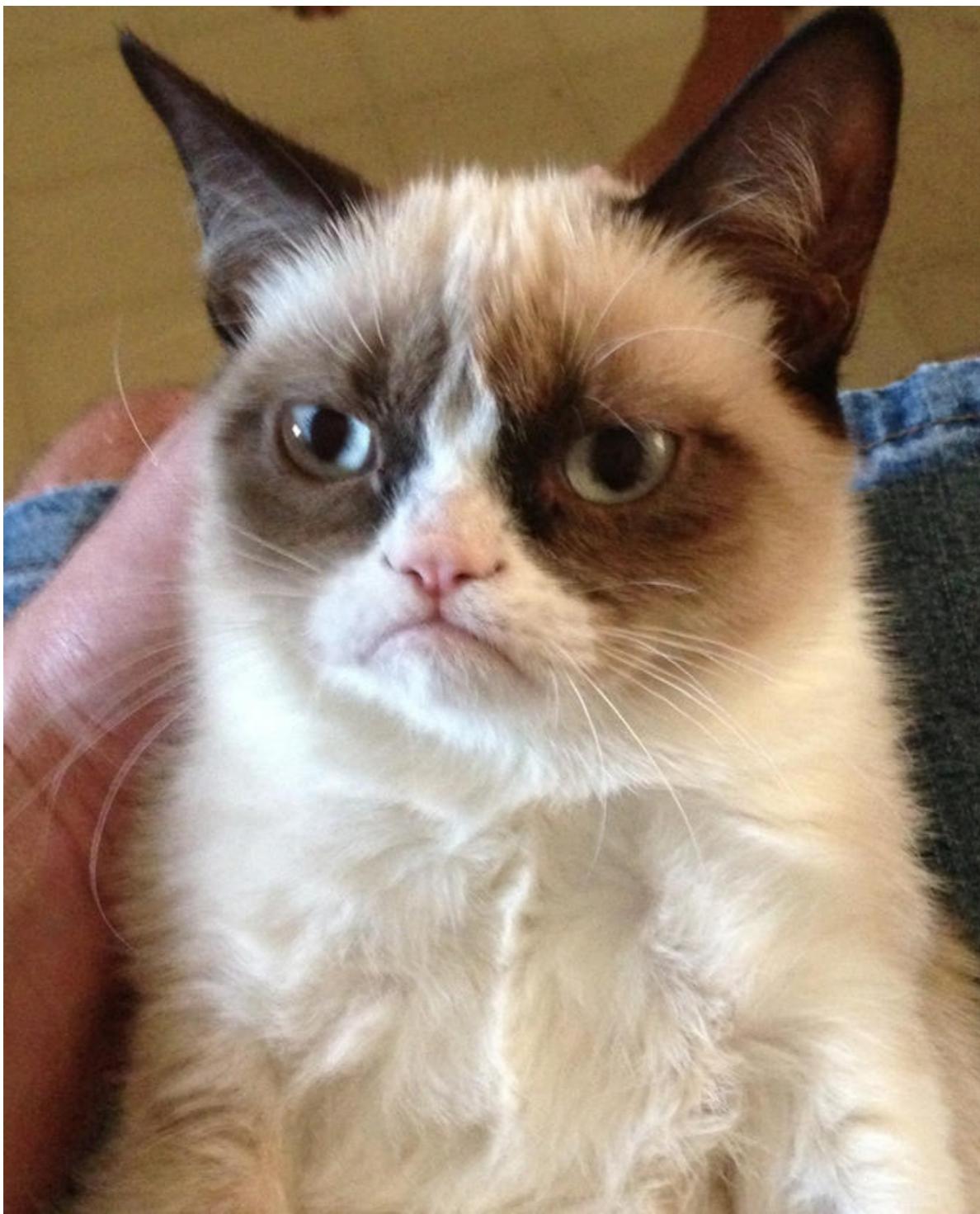
# Representation Learning





Source: Andrej Karpathy: What I learned from competing against a ConvNet on ImageNet

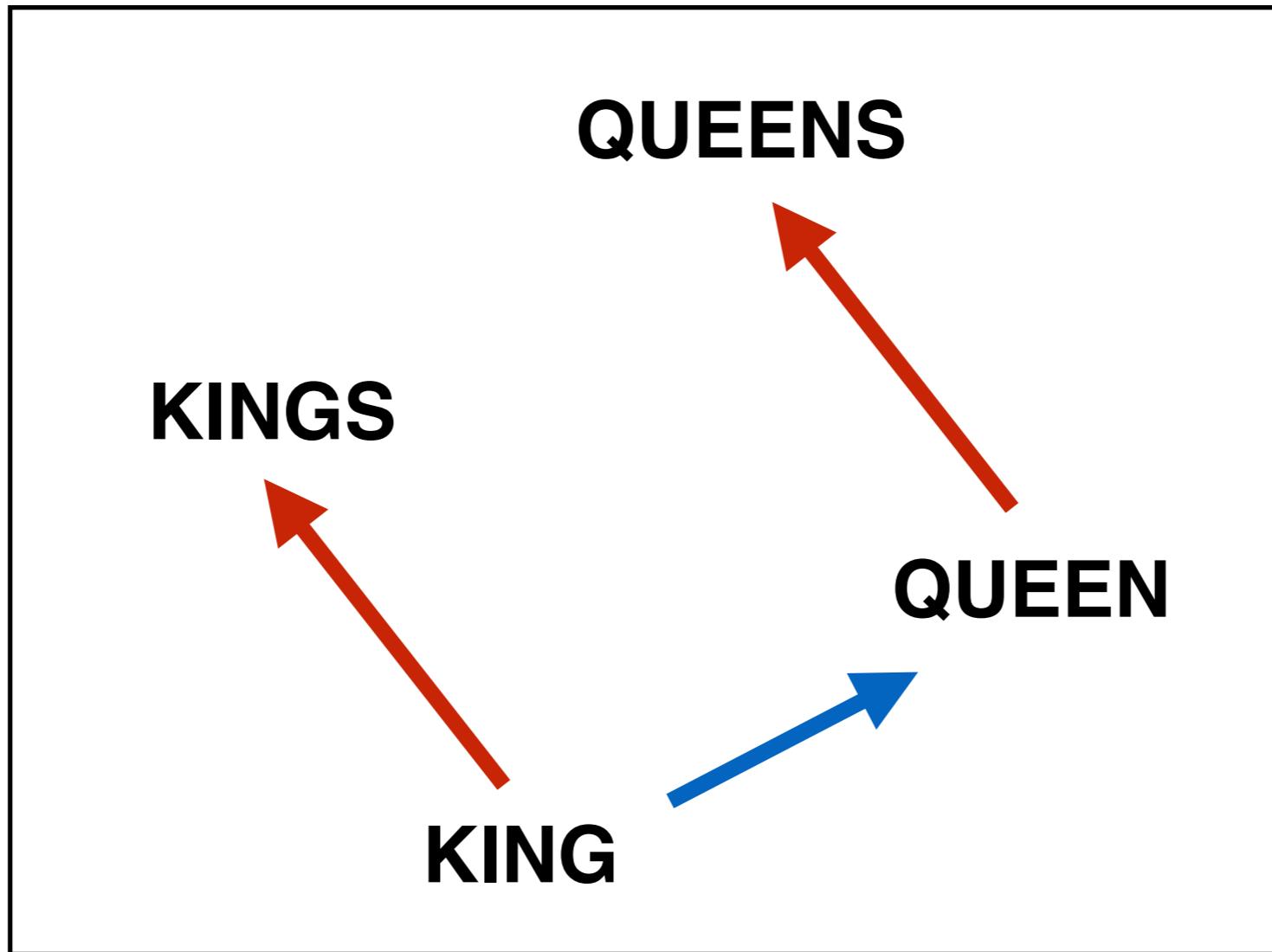
# Transfer Learning





David Teniers the Younger: Archduke Leopold Wilhelm and the artist in the archducal picture gallery in Brussels

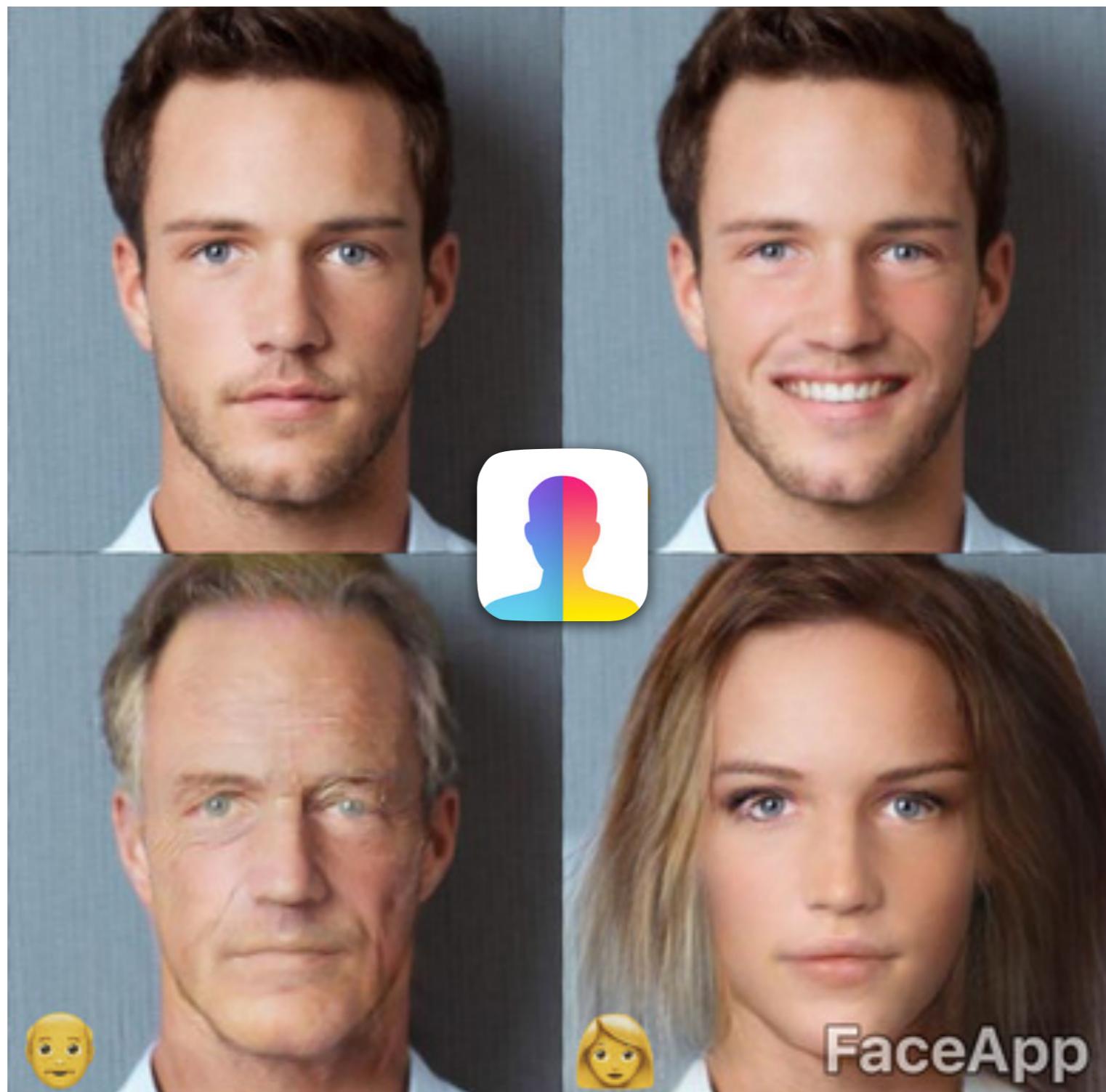
# word2vec



Vectors can encode relationships

T. Mikolov, K. Chen, G. Corrado, and J. Dean, 'Efficient Estimation of Word Representations in Vector Space', CoRR, vol. abs/1301.3781, 2013  
[Online]. Available: <http://arxiv.org/abs/1301.3781>

# Generative Models



Makhzani, Alireza, Shlens, Jonathon, Jaitly, Navdeep, Goodfellow, Ian, Brendan, Frey. Adversarial Autoencoders. <http://arxiv.org/abs/1511.05644>. 2016.

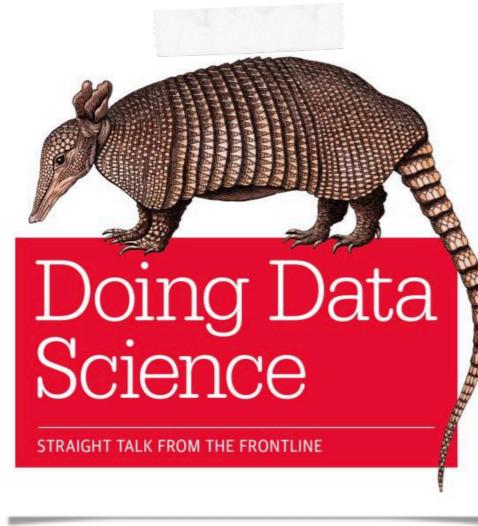
# Generative Models



Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman and Alexei A. Efros.

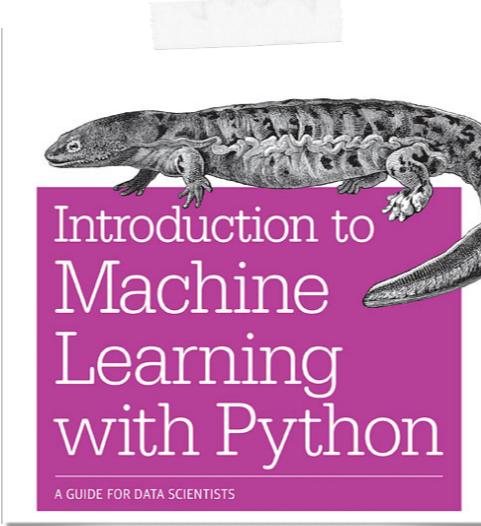
"Generative Visual Manipulation on the Natural Image Manifold",  
in European Conference on Computer Vision (ECCV). 2016.

# Reading List



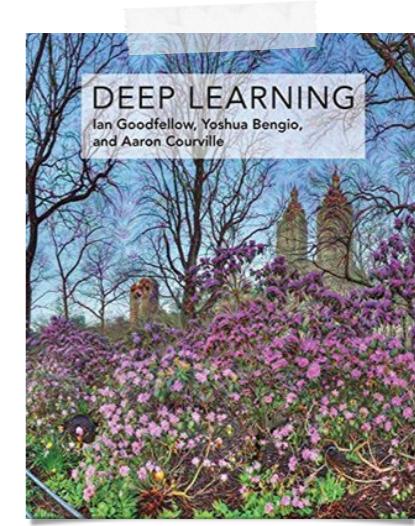
## Doing Data Science

- Cathy O'Neil & Rachel Schutt
- one of the first books on Data Science
- Schutt was a data scientist with Google



## Introduction to Machine Learning with Python

- Andreas C. Müller & Sarah Guido
- from the creator of scikit-learn
- practical introduction to Python and ML



## Deep Learning

- Ian Goodfellow, Yoshua Bengio, & Aaron Courville
- best textbook on deep learning
- surprisingly readable

# *Data Science for Digital Humanities*

BY  
HENDRIK HEUER

INSTITUTE FOR  
INFORMATION MANAGEMENT  
BREMEN GMBH

