

# Curses and Blessings of an (almost) Data-Complete Science: Big Data and the Social Sciences

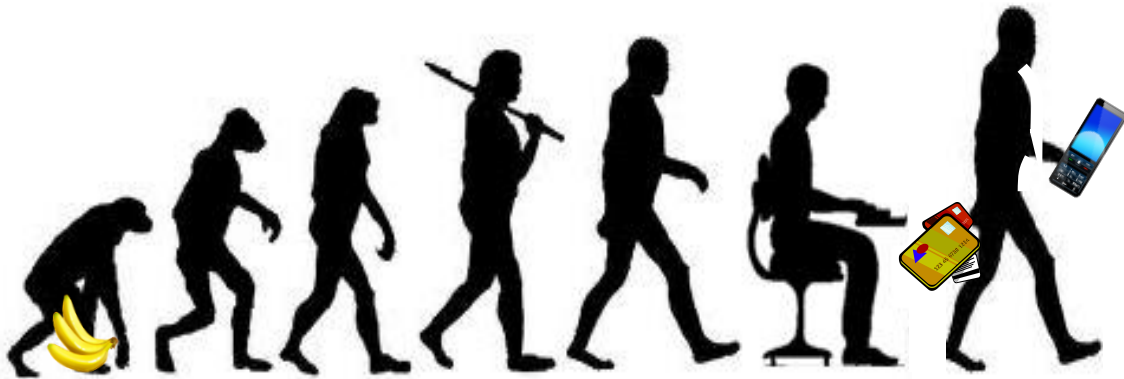
**UCDAVIS**  
UNIVERSITY OF CALIFORNIA

Martin Hilbert  
Department of Communication  
[hilbert@ucdavis.edu](mailto:hilbert@ucdavis.edu)

## Abstract:

Big Data has turned the social sciences from a traditionally data-poor science into arguably the most data-complete science to date – and this basically “overnight”. With over 99% of all of human kinds’ technologically mediated information in digital format, and a mobile penetration of 98% worldwide, the digitalization of human interaction produces an impressively detailed digital footprint of everything that’s relevant for the social sciences. Each and every digital communication inevitably leaves a trace that can be analyzed to better understand and influence social conduct. This renders many traditional survey and data collection and production processes obsolete. While creating unprecedented opportunities for private actors and lots of low hanging fruits for academic research, it also creates challenges that call for a profound paradigm shift in our relation to data.

# Computational Social Science



Volume 493 Number 7432 pp271-446

17 January 2013



“The **biggest stumbling block... is obtaining the data** to parameterize and validate...  
...using automated cameras and image recognition... motion-activated cameras... continuous plankton recorders towed beneath ships...”

A hyena surveys a flock of flamingos in South Africa.

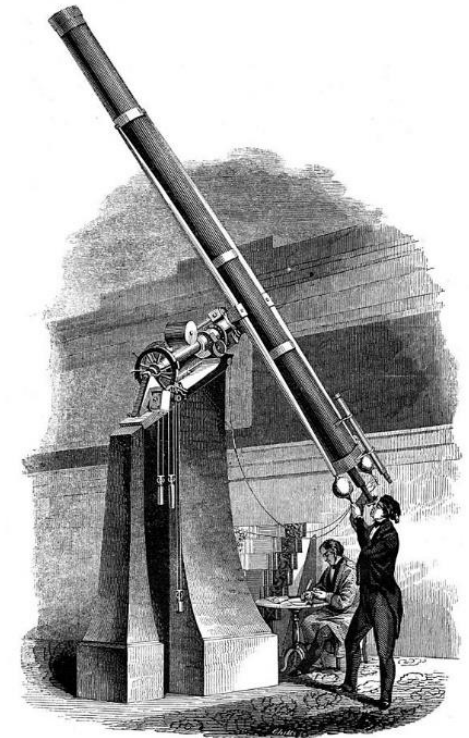
FAQ's

What, in a nutshell, is 'The Madingley Model'?

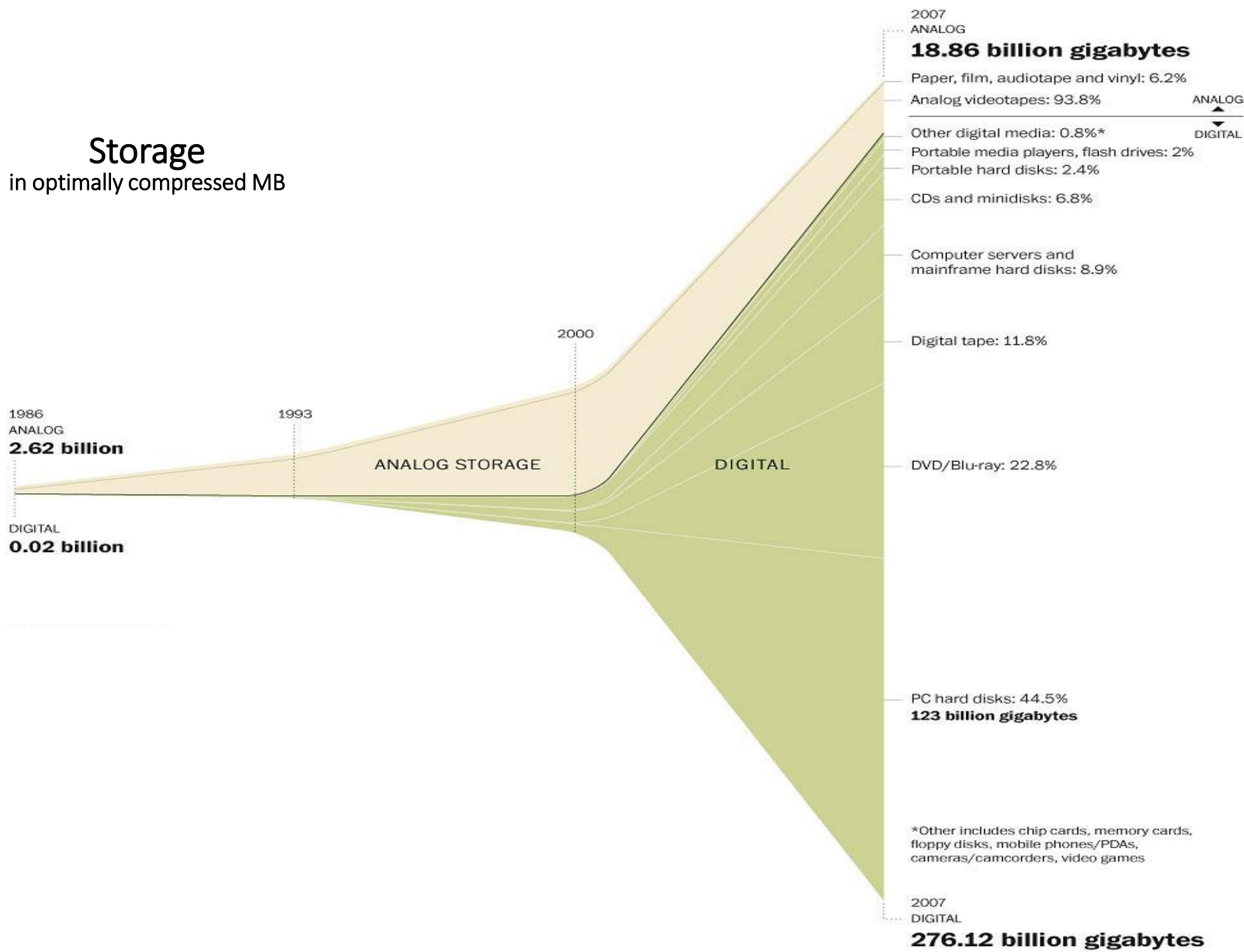
A huge computer simulation of all life on Earth.

## Time to model all life on Earth

To help transform our understanding of the biosphere, ecologists — like climate scientists — should simulate whole ecosystems, argue Drew Purves and colleagues.

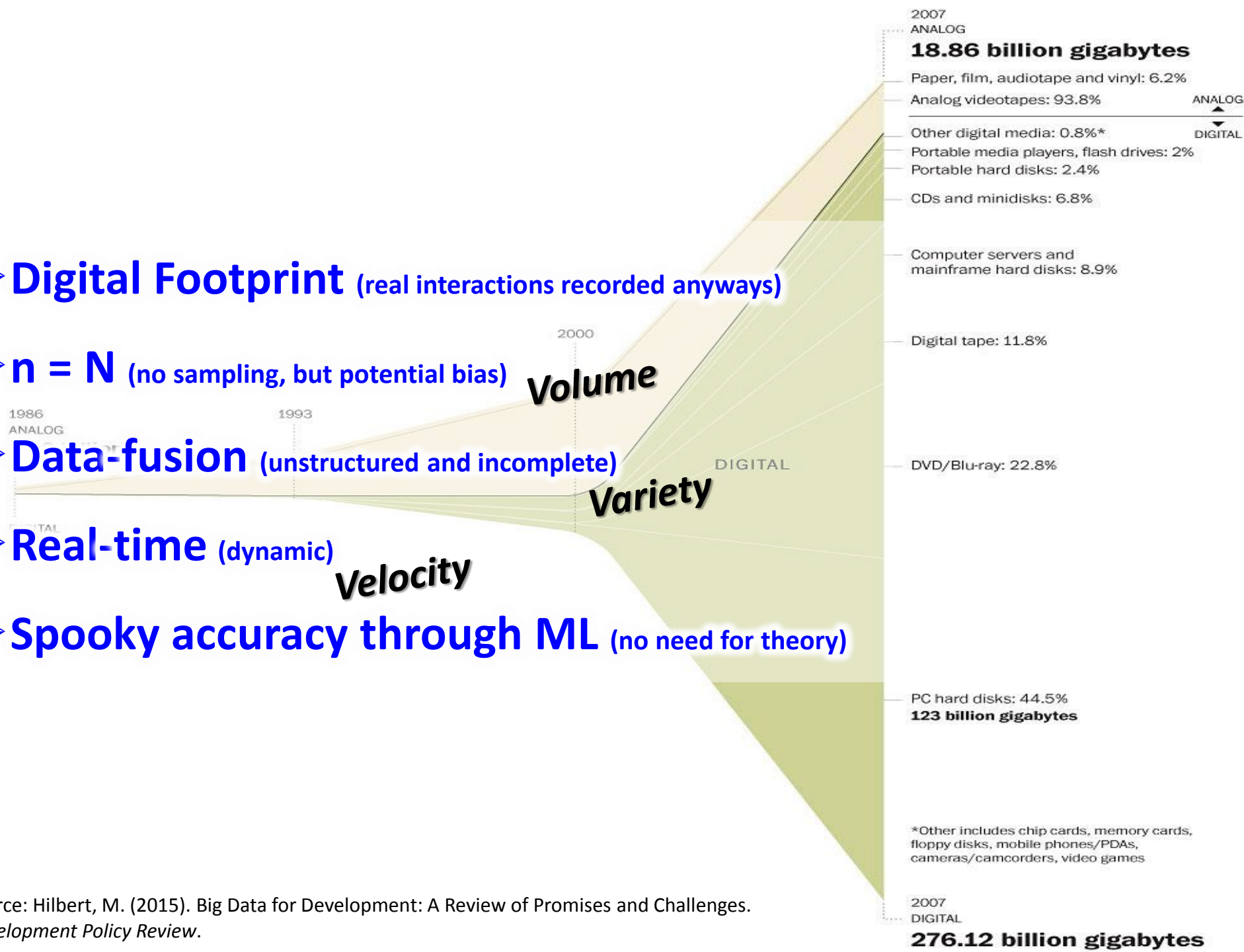


# Storage in optimally compressed MB

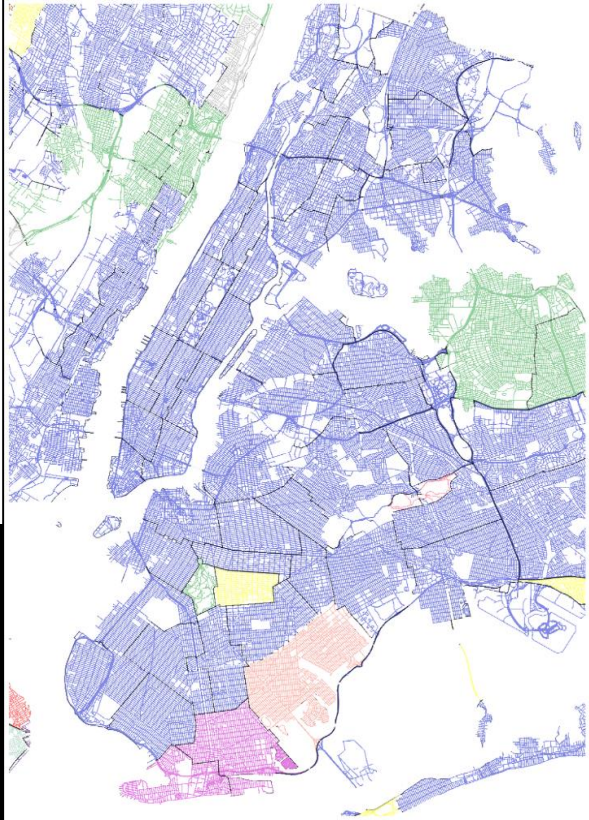


Source: animation by Washington Post, based on Hilbert, M., P. López (04/2011). The world's technological capacity to store, communicate and compute information. Science, 332, 6025, 60-65 [www.martinhilbert.net/WorldInfoCapacity.html](http://www.martinhilbert.net/WorldInfoCapacity.html)

- **Digital Footprint** (real interactions recorded anyways)
- **n = N** (no sampling, but potential bias)
- **Data-fusion** (unstructured and incomplete)
- **Real-time** (dynamic)
- **Spooky accuracy through ML** (no need for theory)



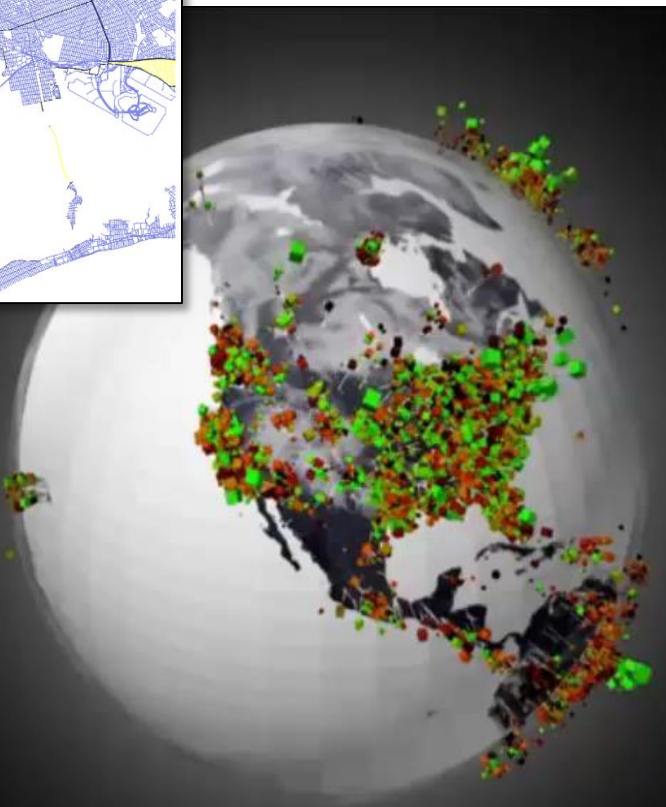
Source: animation by Washington Post, based on Hilbert, M., P. López (04/2011). The world's technological capacity to store, communicate and compute information. Science, 332, 6025, 60-65 [www.martinhilbert.net/WorldInfoCapacity.html](http://www.martinhilbert.net/WorldInfoCapacity.html)



TWITTER: 2<sup>nd</sup> language

Blue - Spanish,  
 Light Green - Korean,  
 Fuchsia - Russian,  
 Red - Portuguese,  
 Yellow - Japanese,  
 Pink - Dutch,  
 Grey - Danish,  
 Coral - Indonesian.

# Digital Footprint



Monrovia (detail)



TED-Ed. (2013). Visualizing the world's Twitter data - Jer Thorp. <http://www.youtube.com>

The Economist. (2014, November 15). Off the map. *The Economist*. <http://www.economist.com>

Mocanu, et al.(2013). The Twitter of Babel: Mapping World Languages through Microblogging Platforms. *PLoS ONE*, 8(4), e61981.

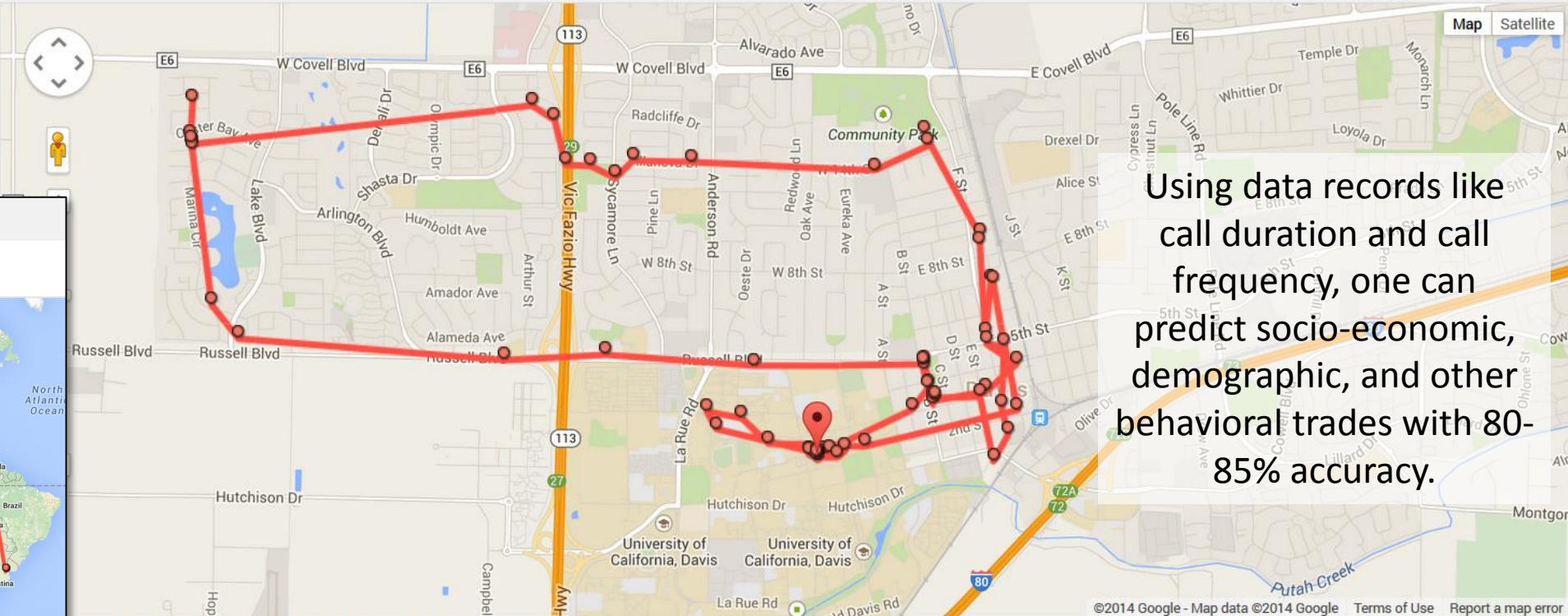
# N = n

## <https://maps.google.com/locationhistory>



### Location history

November 2014								
«	Sun	Mon	Tue	Wed	Thu	Fri	Sat	»
	26	27	28	29	30	31	1	
	2	3	4	5	6	7	8	
	9	10	11	12	13	14	15	
	16	17	18	19	20	21	22	
	23	24	25	26	27	28	29	
	30	1	2	3	4	5	6	



Using data records like call duration and call frequency, one can predict socio-economic, demographic, and other behavioral traits with 80-85% accuracy.

Location history

November 2014								
«	Sun	Mon	Tue	Wed	Thu	Fri	Sat	»
	26	27	28	29	30	31	1	
	2	3	4	5	6	7	8	
	9	10	11	12	13	14	15	
	16	17	18	19	20	21	22	
	23	24	25	26	27	28	29	
	30	1	2	3	4	5	6	

Show: 1 Day

November 25, 2014

- Show timestamps
- Export to KML
- Delete history from this day
- Delete all history

Some points have been hidden from view. Show All Points Learn More

Distance from starting location (farthest distance: 5,300 miles)  
Move mouse over graph to show location on map

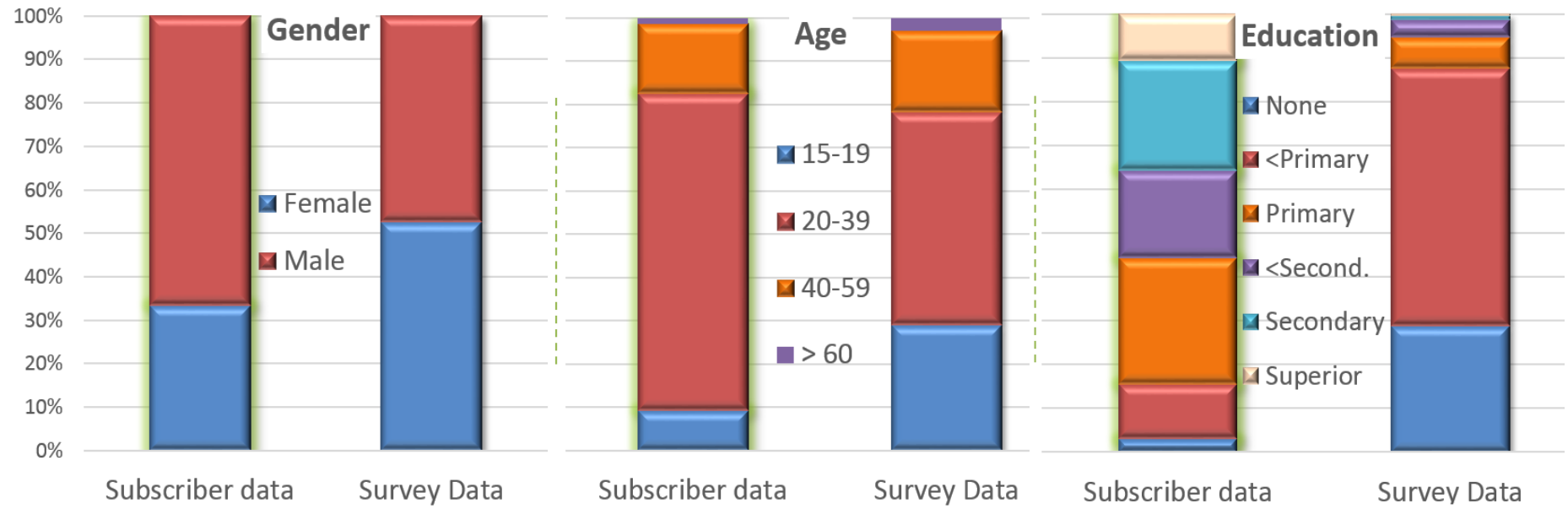


“...human mobility traces are highly unique. ...four spatio-temporal points are enough to uniquely identify 95% of the individuals.”

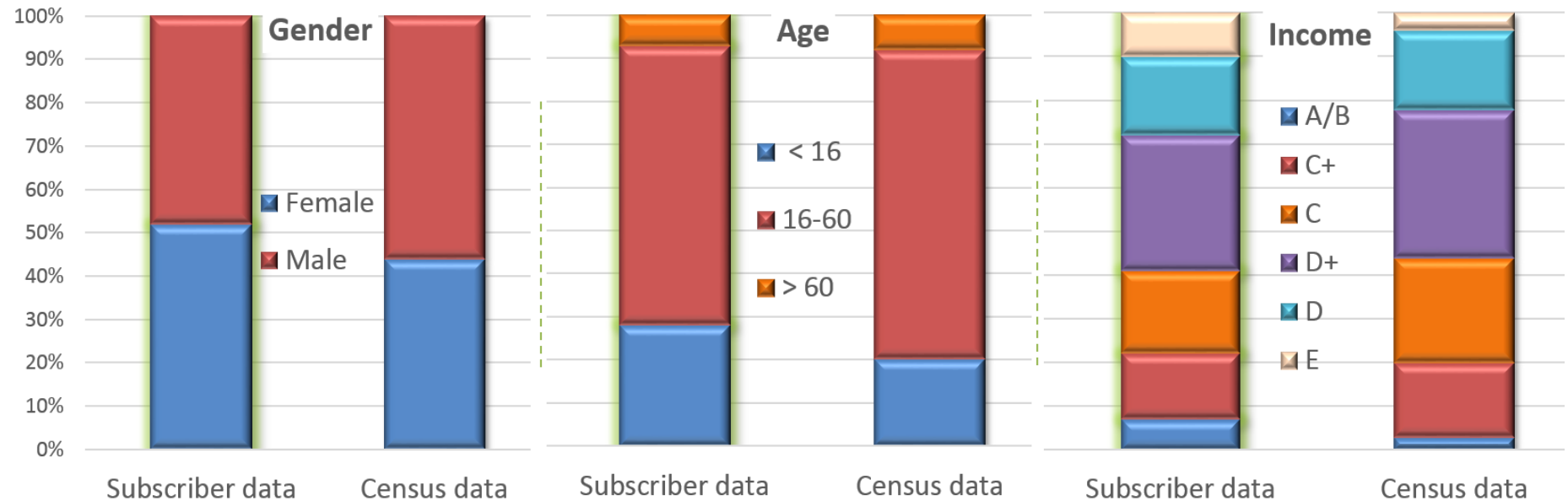
Sources: Raento, et al. (2009). Smartphones: *Sociol. Methods & Research*, 37(3), 426–454. Frias-Martinez, et al. (2014). Spectral clustering for sensing urban land use using Twitter activity. *Engin. Appl. of Artificial Intell.*, 35, 237–245. Frias-Martinez, et al. (2013). Cell Phone Analytics: *ITID*, 9(2), pp. 35–50. Frias-Martinez, et al. (2010). A Gender-centric Analysis of Calling Behavior.... *AAAI 201 Artificial Intelligence for Development*. Blumenstock et al. (2010). Who’s Calling? *AAAI 201 Artificial Intelligence for Development*. De Montjoye, et al. (2013). Unique in the Crowd: *Scientific Reports*, 3

**N = n ?**

*(a) Rwanda 2005/09:  
mobile phone penetration of  
2-20%*



*(b) LatAm economy 2009/10:  
mobile phone penetration of  
60-80%*



Source: (a) Blumenstock and Eagle (2012); (b) Frias-Martinez and Virseda (2013).



# Real-time

## U.S. Bureau of Labor Statistics

- > 100s staff visiting > 90 cities
- 80,000 prices
- Cost: US\$ 250 million/year

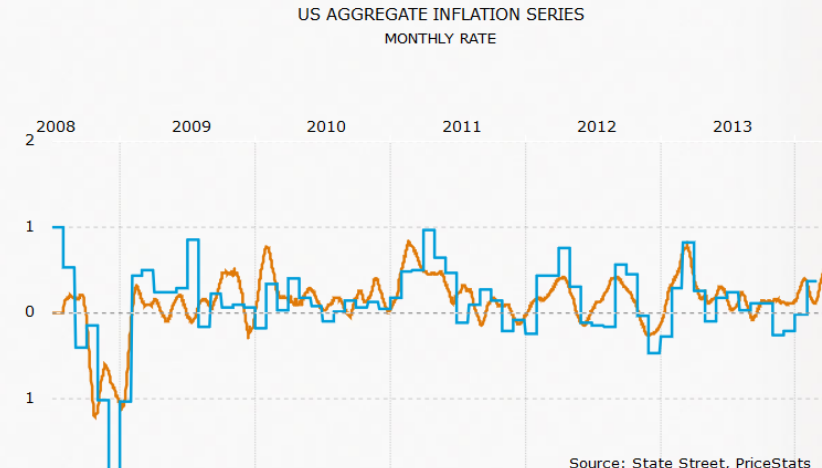
## PriceStats

- 17 staff
- 300 online retailers; > 70 countries
- daily 5,000,000 prices

### US INFLATION SERIES

PriceStats estimates aggregate inflation in the US using online prices. The objective of this series is to anticipate major changes in US inflation trends, but not to forecast monthly CPI announcements. At any point in time, our index can be different from the CPI. Our data anticipates changes in inflation trends not only because we observe prices sooner, but also because online prices tend to react to shocks more quickly.

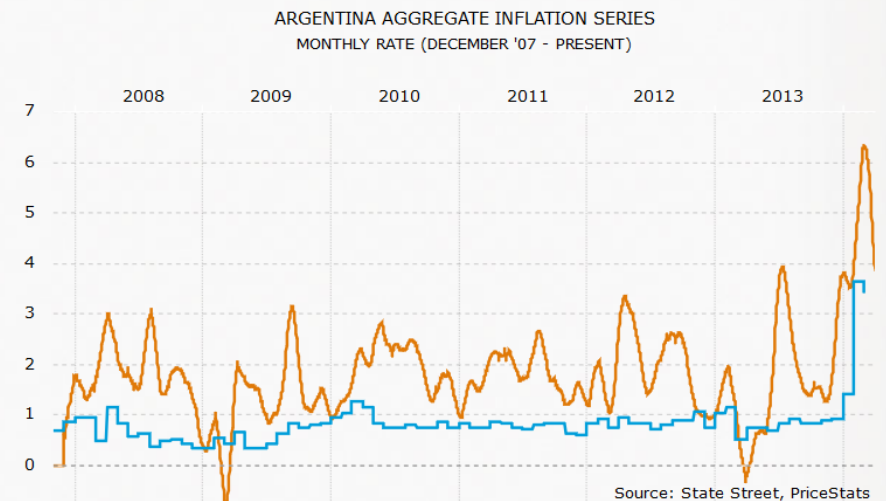
Official CPI  
PriceStats Index



### ARGENTINA INFLATION SERIES

The State Street PriceStats Argentina Index is now published in *The Economist* on a weekly basis. The publication chose PriceStats' statistics as an alternative version to official figures for Argentina. Unlike in other countries, our Argentina series show a significant departure from official numbers. Our 2013 annual inflation rate is ~23%; official figures account for ~11%.  
[Read the article](#)

Official CPI  
PriceStats Index



The Economist | World politics | Business & finance | Economics | Science & technology | Culture

## Official statistics Don't lie to me, Argentina

Why we are removing a figure from our indicators page  
Feb 25th 2012 | From the print edition



Sources:  
<http://www.pricestats.com/about-us/meet-the-team> ;  
[www.economist.com/node/21548242](http://www.economist.com/node/21548242) ;  
<http://www.inflacionverdadera.com>

# “Big” doesn’t need to know Why

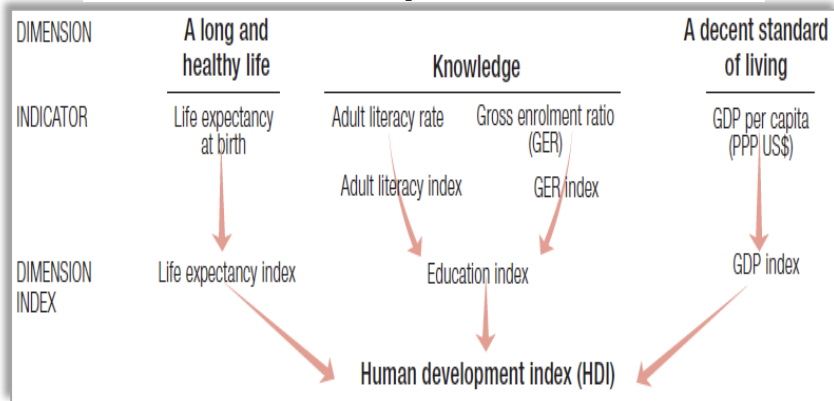
Choose a representation that can use **unsupervised learning on unlabeled data**, which is so much more plentiful than labeled data.

Translate

The screenshot shows the Google Translate interface. At the top, there are language selection buttons for English, German, Esperanto, and Detect language. A dropdown menu is open showing Esperanto, Swahili, and English. A blue 'Translate' button is visible. The input text is 'Ĉu vi povas kompreni tion?' and the output text is 'Unaweza kuelewa kwamba?'. There are also icons for voice input, star, list, and edit.

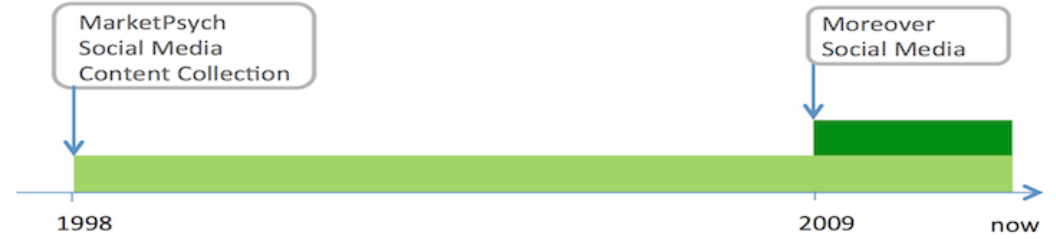
Afrikaans	Cebuano	Finnish	Hmong	Korean	Nepali	Somali	Welsh
Albanian	Chinese (Simplified)	French	Hungarian	Lao	Norwegian	Spanish	Yiddish
Arabic	Chinese (Traditional)	Galician	Icelandic	Latin	Persian	Swahili	Yoruba
Armenian	Croatian	Georgian	Igbo	Latvian	Polish	Swedish	Zulu
Azerbaijani	Czech	German	Indonesian	Lithuanian	Portuguese	Tamil	
Basque	Danish	Greek	Irish	Macedonian	Punjabi	Telugu	
Belarusian	Dutch	Gujarati	Italian	Malay	Romanian	Thai	
Bengali	English	Haitian Creole	Japanese	Maltese	Russian	Turkish	
Bosnian	Esperanto	Hausa	Javanese	Maori	Serbian	Ukrainian	
Bulgarian	Estonian	Hebrew	Kannada	Marathi	Slovak	Urdu	
Catalan	Filipino	Hindi	Khmer	Mongolian	Slovenian	Vietnamese	

# Human Development Index HDI

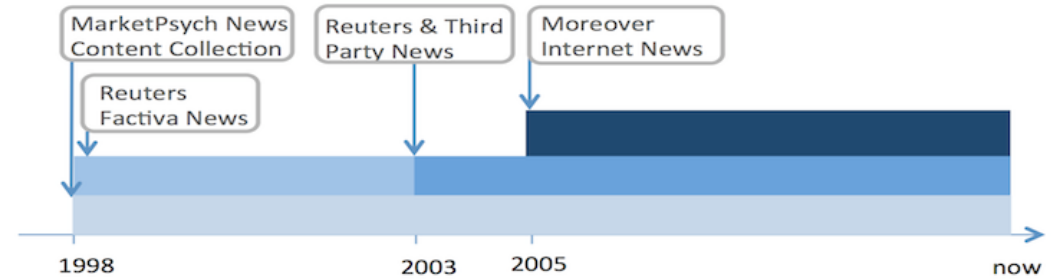


# Data-fusion

## SOCIAL MEDIA SOURCES

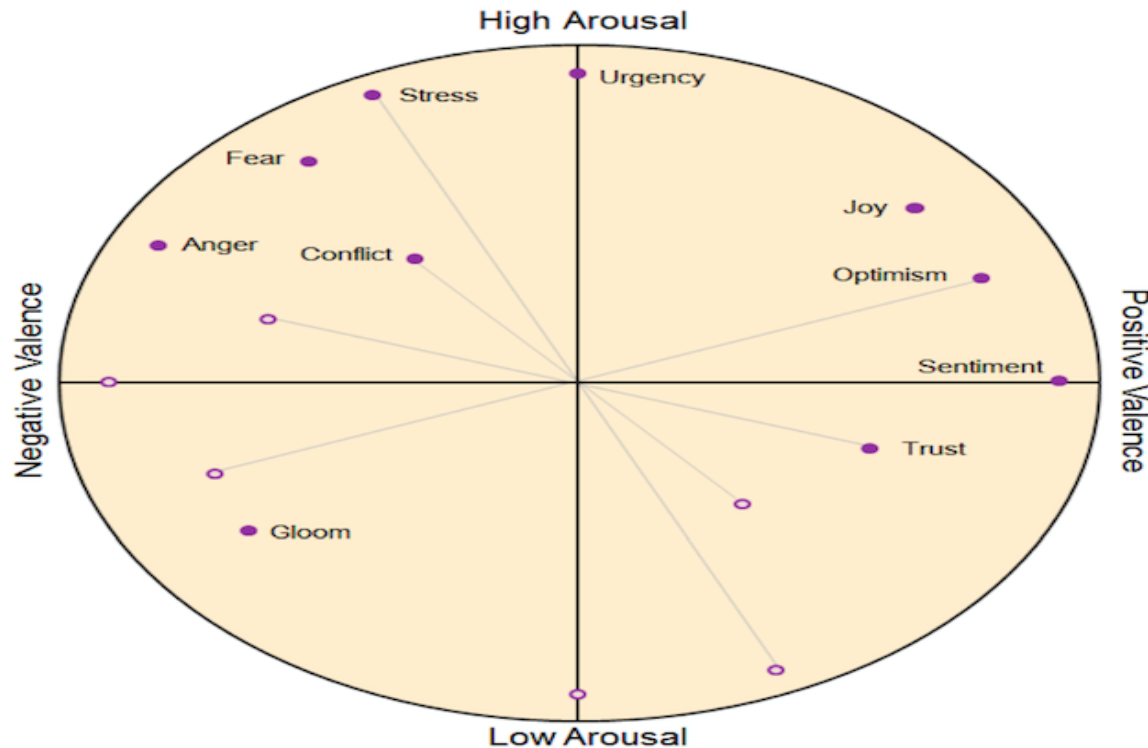


## NEWS MEDIA SOURCES



MarketPsych Data

**Thomson Reuters MarketPsych Indices (TRMI)**  
 18,864 separate indices, 119 countries, updated each minute (!)



# Network Data fusion

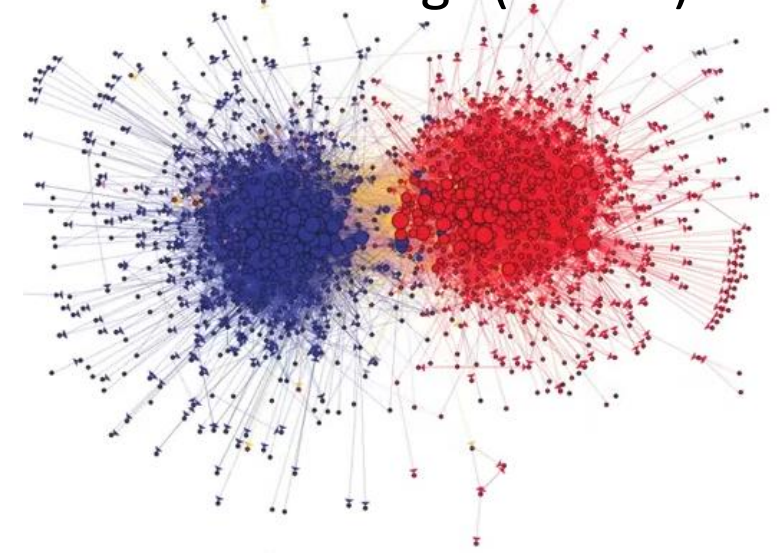
Traditional database of attributes

	Gender	Location	Income	Educat.
Jorge	M	Urban	700	Tertiary
Maria	F	Urban	500	Second.
Juan	M	Rural	300	Primary
Magda	F	Rural	200	---

Network database of links

	Jorge	Maria	Juan	Magda
Jorge	Self	↗	---	↘
Maria	↖	Self	---	↗
Juan	---	---	Self	---
Magda	---	↖	↖	Self

Political blogs (online)



Without any information about a Facebook user beyond a list of his friends, one can accurately predict his sexual orientation.

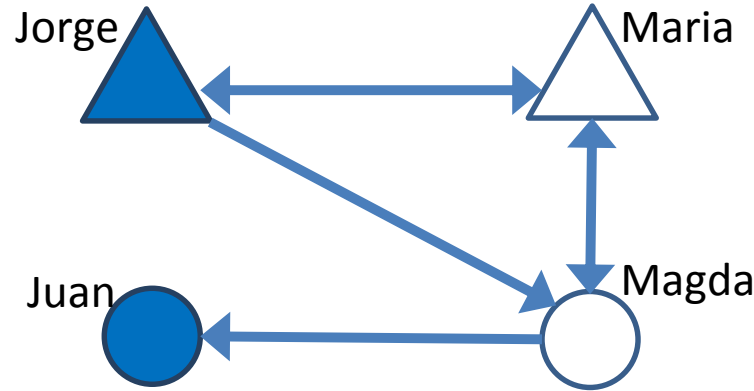
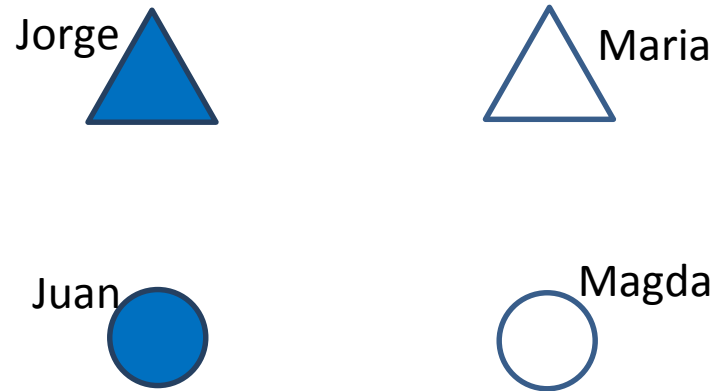
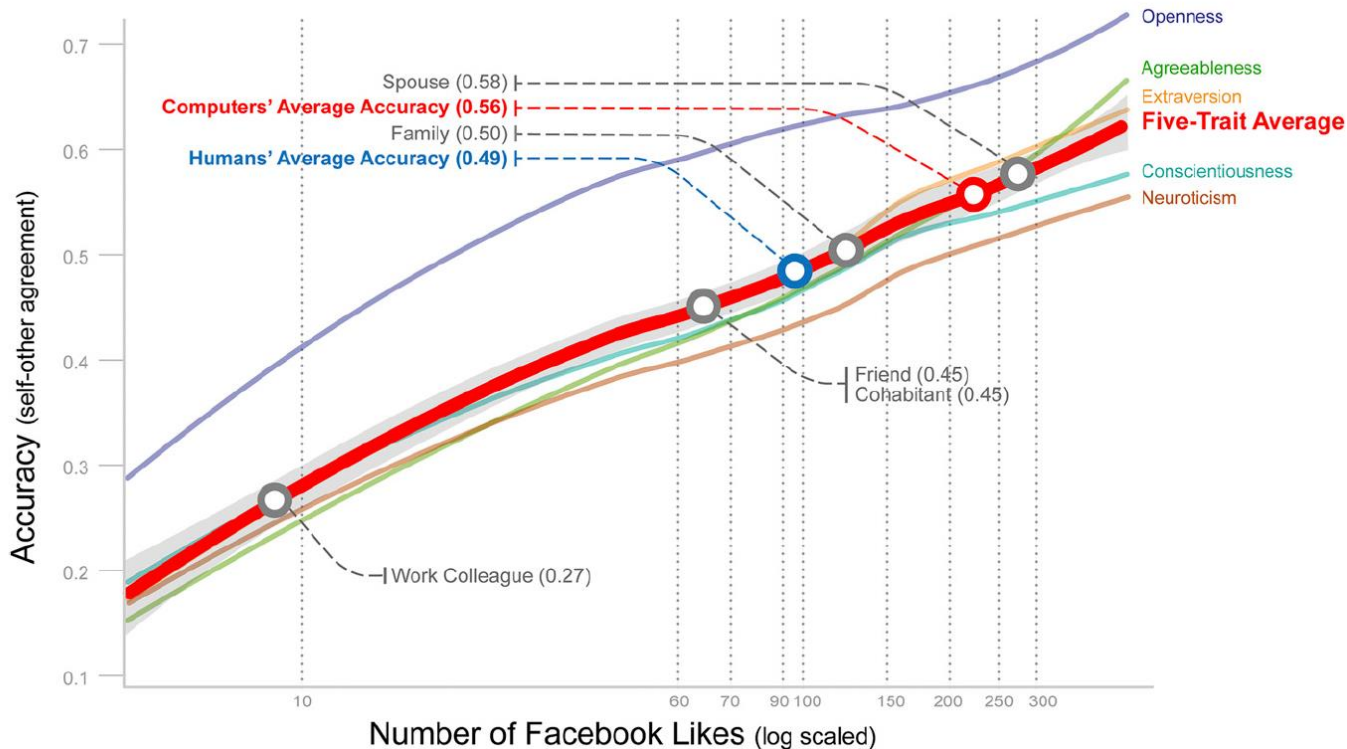


Table 6: Subjects presented in Table 5 with private profiles accurately classified as gay male.

Name	Profile privacy setting	Reported sex	Reported orientation	Percentage gay friends	Classified as gay
A	Private	Unknown	Unknown	13.21%	True
H	Private	Unknown	Unknown	4.56%	True
I	Private	Unknown	Unknown	4.19%	True
K	Private	Unknown	Unknown	3.80%	True
P	Private	Unknown	Unknown	2.86%	True
R	Private	Unknown	Unknown	2.65%	True

# Big Network Data:

## Digital Footprint + N=n + Data-fusion + Real-time + ML



*“Facebook Likes, can be used to automatically and accurately predict...: sexual orientation, ethnicity, religious and political views, personality traits, intelligence, happiness, use of addictive substances, parental separation, age, and gender...”*

*“(i) computer predictions based on ... Facebook Likes are more accurate ( $r = 0.56$ ) than those made by the participants' Facebook friends ( $r = 0.49$ ); (iii) computer personality judgments have higher external validity when predicting life outcomes such as substance use, political attitudes, and physical health; for some outcomes, they even outperform the self-rated personality scores...”*

*"This call might be recorded for quality and training purposes."*

## Social transparency, homophily & polarization



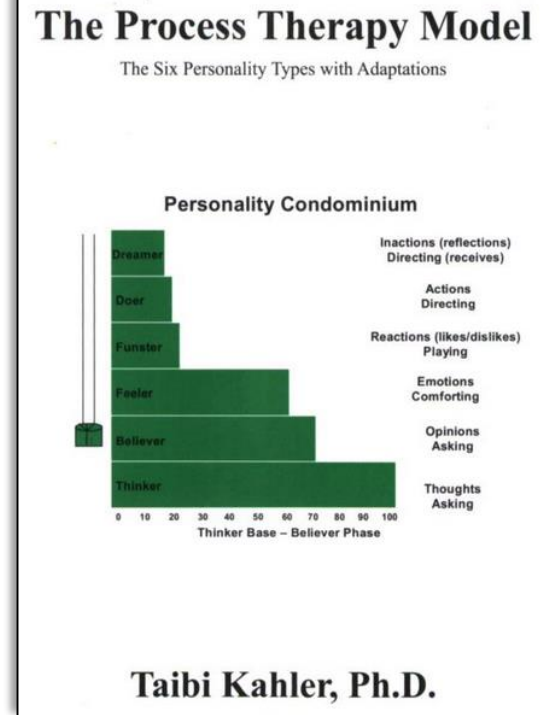
### Meaningful business impact

Mattersight, a leader in enterprise analytics focused on customer-employee interactions, offers analysis showing that a behavioral *mismatch* between a customer and employee, as opposed to a favorable behavioral *connection*, has a significant impact on business outcomes across industries, even in highly specialized contact center functions:

Business Outcome	Impact of Behavioral Connection
Sales Conversion Rate	85% to 230%
Customer Retention/Attrition Rate	25% to 50%
Customer Service Cost/Efficiency	35% to 45%
Student Enrollment Rate (Private Sector Education)	200% to 700%
Debt Cure Rate (Collections Organization/Function)	Over 1400 Basis Points

### Matching Personality Types:

- ✓ Call average from 10 min to 5 min
- ✓ Customer Satisfaction from 47 % to 92%



**EMOTIONS-DRIVEN (30% of the population)**

**THOUGHTS-DRIVEN (25%)**

**REACTIONS-DRIVEN (20%)**

**OPINIONS-DRIVEN (10%)**

**REFLECTIONS-DRIVEN (10%)**

**ACTIONS-DRIVEN (5%)**

# Obama 2012 campaign

YOLO: MEET THE OBAMA CAMPAIGN'S CHIEF TECHNOLOGY OFFICER



The President hugging Harper Reed as shown on his Instagram feed.

## ➤ Data

- **US\$1 billion investment; core group of 40 engineers**  
(from Twitter, Google, Facebook, Craigslist, stem cell, professional poker players...)
- **Project Narwhal: 16 million unique voter profiles:**  
email sign-ups, zip codes, profession, voter registrations, volunteering & donation record, Tweets, Facebook postings and network ties, TV Watching behavior through 20 million set-top boxes, etc.
- **62,000 computer simulations** of likely voter behavior

## ➤ Outcome

- Identified the 20% of Obama's 2008 vote that shifted into the undecided column, ranking them on a 0-10 persuasion score
- Obama paid 35% less per broadcast commercial than Romney  
(40,000 more spots on the air, spending \$90 million less!)
- Present tailor made campaign promises (agreeable adds; etc)
- Guide volunteers in phone and door-to-door campaigns
- Email donation requests, raising \$181 million/month
- Predict States voting outcome at an accuracy of 0.5 percent
- **Change voting behavior of 78 % of targeted undecided voters through Facebook**

# Big Data as commodity



## Consumers' financial vulnerability:

- *"Social Influencer"*
- *"Rural and Barely Making It"*
- *"Ethnic Second-City Strugglers"*
- *"Retiring on Empty: Singles"*
- *"Tough Start: Young Single Parents"*
- *"Credit Crunched: City Families"*
- *"Transitory lifestyles: military personnel"*
- *"Elderly Opportunity Seekers: elderly looking for ways to make money"*
- *"Oldies but Goodies: gullible, want to believe their luck can change"*

Source: US Senate. *A Review of the Data Broker Industry: Collect, Use, and Sale of Consumer Data for Marketing Purposes*, 2013)



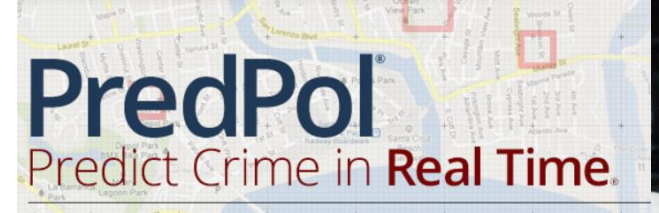
Dynamic Insights

Smart Steps





# The long-begun end of “free will”?



## Amazon to ship things before you've even thought of buying them?

Amazon files a patent for "anticipatory package shipping." The idea is that the company will know from your buying patterns what you're likely to want next.



by Chris Matyszczyk | January 19, 2014 1:07 PM PST



Get email alerts

## Predictive Policing LADP & SantaCruz

- Data on crimes, weather, buses, parks...
- Models from Earthquake aftershock
- Predictions to 500<sup>2</sup> feet / 50<sup>2</sup> m
  - ✓ Crimes down 13 %; burglaries 11 %; car theft 8 % (while other districts went up during same period)

## Homicide Parole candidates

- dataset > 60,000 crimes
- with some 300 predictors (nature of crime, *age*, repetition)
  - ✓ 60 – 70 % correct who commit homicide



## Pre-punishment vs. Free Will?

- already insurance premiums per age or gender (*punishment*)



"Your recent Amazon purchases, Tweet score and location history makes you 23.5% welcome here."

# Proxies are just proxies

Fruit prices to detect violence in Jalalabad Afghanistan



...a “big breakthrough” (DARPA’s Director Regina Dugan)  
that impressed a group of four-star generals...

JSOC drone operator: “It’s of course assumed that the phone belongs to a human being who is nefarious and considered an ‘unlawful enemy combatant.’ This is where it gets very shady...”



# Computational Social Science

## I. Big Data Blessings

- Produced anyways
- $n = N$  (Volume)
- Data-fusion (Variety)
- In real-time (Velocity)
- Spooky accuracy through ML

## II. Big Data Curses:

- Social Transparency & Polarization
- Data as commodity
- Confusing statistical variables (proxies) with reality

