# The Theory, Practice and Limits of
# Big Data for the Social Sciences

**Martin Hilbert**
**Department of Communication**
**hilbert@UCDavis.edu**

# World's Info Storage Capacity
in optimally compressed MB

2007
ANALOG
**18.86 billion gigabytes**

Paper, film, audiotape and vinyl: 6.2%
Analog videotapes: 93.8%          ANALOG

Other digital media: 0.8%*          DIGITAL
Portable media players, flash drives: 2%
Portable hard disks: 2.4%

CDs and minidisks: 6.8%

Computer servers and
mainframe hard disks: 8.9%

Digital tape: 11.8%

2000

1986
ANALOG
**2.62 billion**

ANALOG STORAGE          DIGITAL

DVD/Blu-ray: 22.8%

1993

DIGITAL
**0.02 billion**

**Stored digital information has doubled every 2.5 years**

**≈ 5 ZB in 2014** (5 x 10$^{21}$ Bytes)

Sun                    Earth

91,000,000 miles

**4,500 piles of printed books**

Ken Costello

PC hard disks: 44.5%
**123 billion gigabytes**

*Other includes chip cards, memory cards,
floppy disks, mobile phones/PDAs,
cameras/camcorders, video games

2007
DIGITAL
**276.12 billion gigabytes**

All DNA on Earth = $1.3 \times 10^{37}$ Bytes
$\approx$> 100 years!

7.2 bn humans * 6.2 bn nucleotides = $1 \times 10^{19}$ Bytes   vs.    $5 \times 10^{21}$ Bytes

OECD
BETTER POLICIES FOR BETTER LIVES

McKinsey&Company

*"data as a new source of growth"*

*"the new oil"*

*"need to recognize the potential of harnessing big data to unleash the next wave of growth"*

Secrets of
The BIG DATA Revolution
The world is changing. Are you ready?
Jason Kolb · Jeremy Kolb

**The Theory, Practice and Limits of Big Data for the Social Sciences**

# SDG Goal 2.4: …ensure sustainable food production systems and implement resilient agricultural practices that increase productivity

[weather data] + [RICE crops data] +

+ algorithms from neuroscience / biology =>
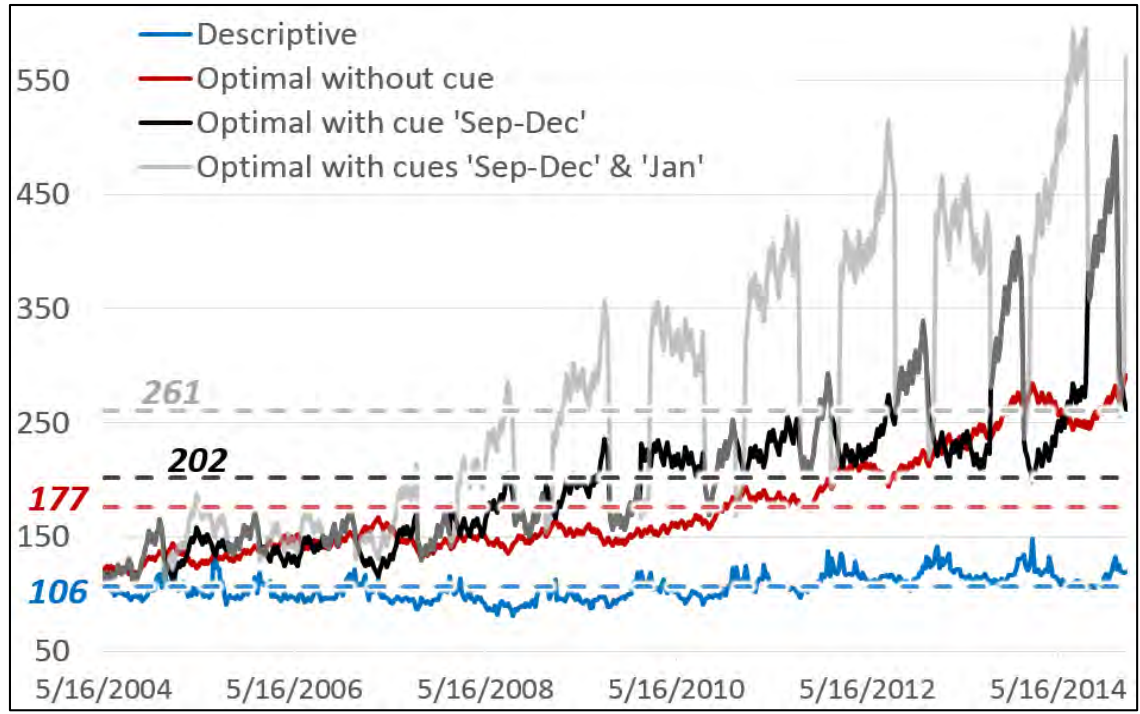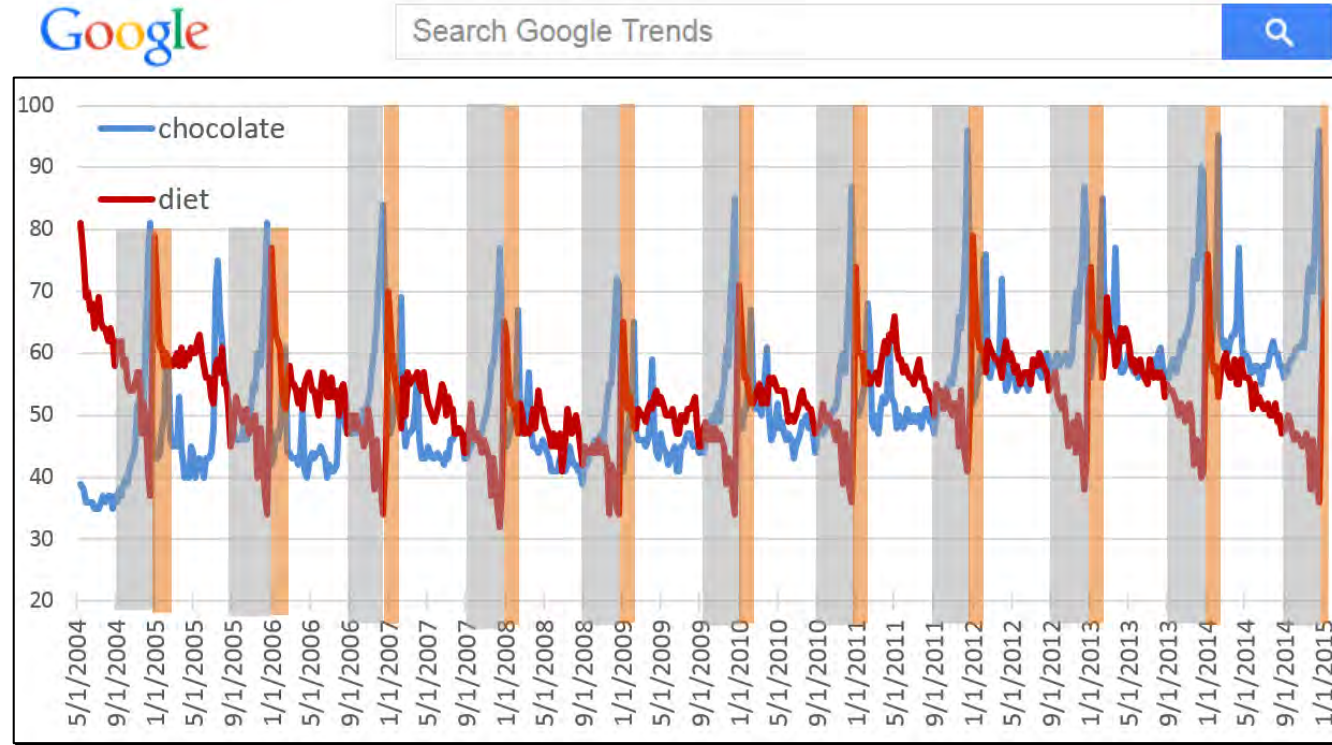
=> *climate change*
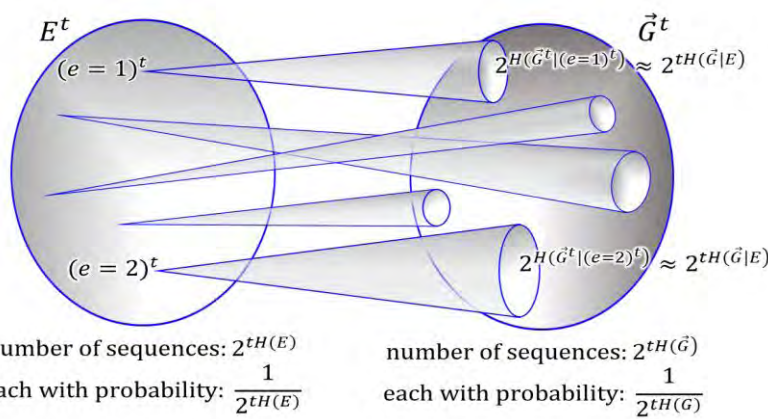
Results **localized** for towns:

➢ Saldaña: solar radiation during the grain-ripening stage

➢ Espinal: sensitivity to warm nights

⇒ **Low cost** solutions: sowing crops in right period of time

⇒ **Impact**: 170 farmers avoided direct losses of $ 3.6 million + productivity from 1 to 3 tons per hectare.

….now being scaled out through Colombia, Argentina, Nicaragua, Peru and Uruguay.

# Information & Growth



$$Growth = E_e[\log {}^d W] - H(E|\vec{G}) - D_{KL}(\vec{P}(e|g)\|P(e|m)) - I(E;\vec{G})$$

Hilbert, M. (2015). An Information Theoretic Decomposition of Fitness: Engineering the Communication Channels of Nature and Society (SSRN Scholarly Paper No. ID 2588146). Social Science Research Network. http://papers.ssrn.com/abstract=2588146

# Information & Growth



**E** *Environment*

**$\vec{G}$** *Population*

1 bit of **information** = reduction of **uncertainty** by half

Claude Shannon (1948)

*A Mathematical Theory of Communication,*

**½** * uncertainty = 1 bit of information = **2** * Growth

$$Growth = E_e\left[\log {}^dW\right] - H(E|\vec{G}) - D_{KL}\left(\vec{P}(e|g)\|P(e|m)\right) - I(E;\vec{G})$$

Hilbert, M. (2015). An Information Theoretic Decomposition of Fitness: Engineering the Communication Channels of Nature and Society (SSRN Scholarly Paper No. ID 2588146). Social Science Research Network. http://papers.ssrn.com/abstract=2588146
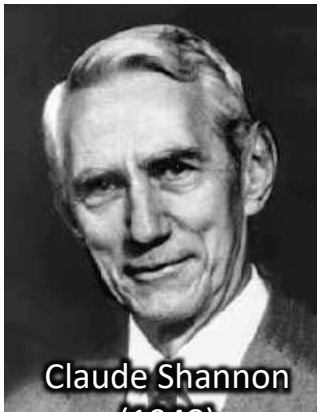
# The Theory, Practice and Limits of
# Big Data for the Social Sciences

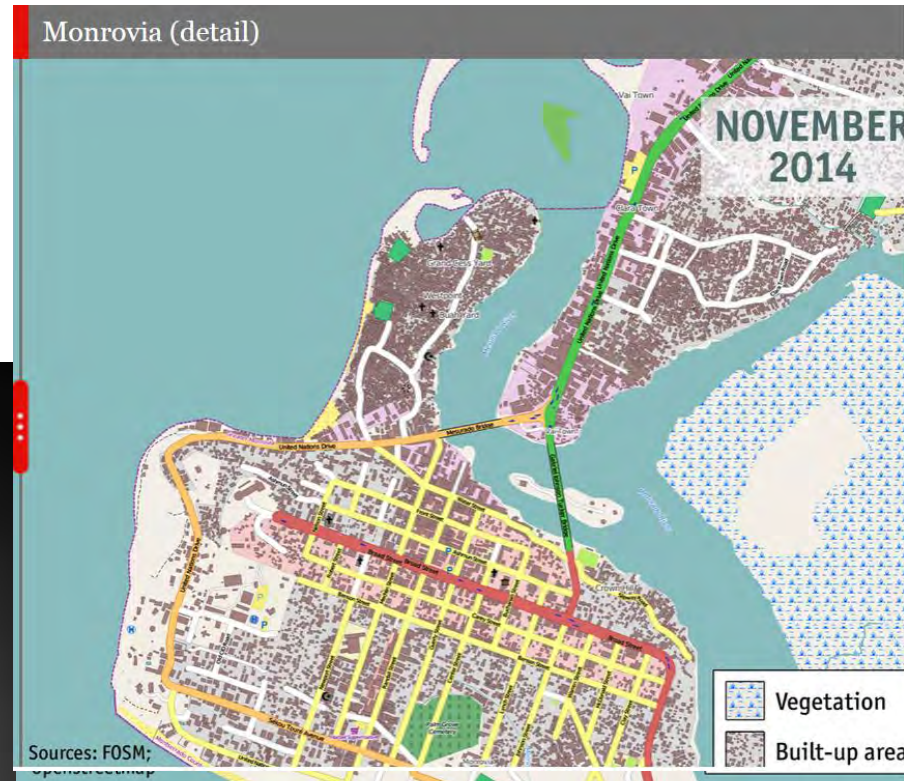**Characteristics of Big Data**

➢ **Digital footprint** (produced anyways for free)

➢ **n = N** (no sampling, but potential bias)

➢ **Data-fusion** (unstructured and incomplete)

➢ **In real-time** (dynamic)

➢ **Machine Learning** (no need for theory)

Source: Hilbert, M. (2015). Big Data for Development: A Review of Promises and Challenges. *Development Policy Review*.

# Digital footprint



Monrovia (detail)

NOVEMBER 2014

Vegetation

Built-up area

Sources: FOSM;
openstreetmap

Thu Aug 20 21:26:49 PDT 2009

8am
9am
10am

Source:
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

# Digital Footprint



https://maps.google.com/locationhistory

**N = n**

75 % among those with less than US$1 per day!
Naef, et al. (2014). *Using Mobile Data for Development*

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

Note: * Estimate
Source: ITU World Telecommunication /ICT Indicators database

Sources: Raento, M., Oulasvirta, A., & Eagle, N. (2009). Smartphones: An Emerging Tool for Social Scientists. *Sociol. Methods & Research*, *37*(3), 426–454.
Frias-Martinez, V., & Frias-Martinez, E. (2014). Spectral clustering for sensing urban land use using Twitter activity. *Engin. Appl. of Artificial Intell.*, 35, 237–245.
Frias-Martinez, V., & Virseda, J. (2013). Cell Phone Analytics: Scaling Human Behavior Studies into the Millions. *ITID*, 9(2), pp. 35–50.
Frias-Martinez, V., Frias-Martinez, E., & Oliver, N. (2010). A Gender-centric Analysis of Calling Behavior…. AAAI 201 *Artificial Intelligence for Development*.
Blumenstock, J. E., Gillick, D., & Eagle, N. (2010). Who's Calling? Demographics of Mobile Phone Use in Rwanda. AAAI 201 *Artificial Intelligence for Development*.

social *science*

a woman's age *vs. the age of the men who look best to her*

a man's age *vs. the age of the women who look best to him*

# Data Fusion



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: http://www.youtube.com/watch?v=wqjKTW3wJZ8     US Senate. *A Review of the Data Broker Industry: Collect, Use, and Sale of Consumer Data for Marketing Purposes*, 2013)

**Real time**



**The Process Therapy Model**
The Six Personality Types with Adaptations

**Personality Condominium**

Inactions (reflections)
Directing (receives)

Actions
Directing

Reactions (likes/dislikes)
Playing

Emotions
Comforting

Opinions
Asking

Thoughts
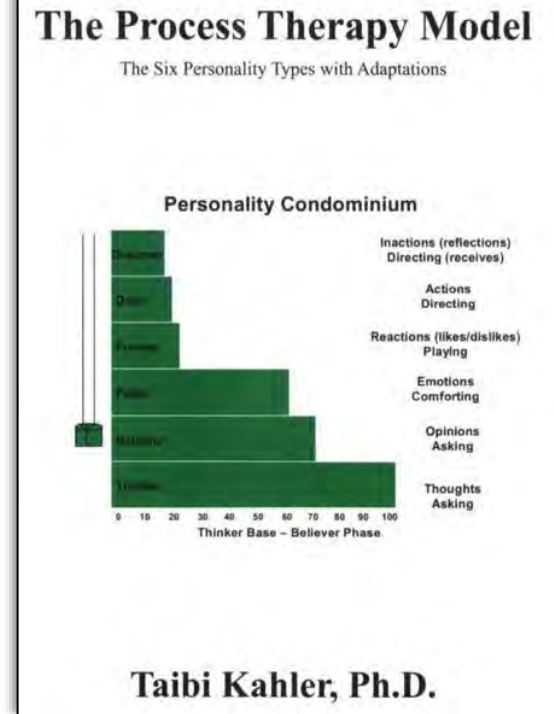Asking

Thinker Base – Believer Phase

**Taibi Kahler, Ph.D.**

Matching Personality Types:
- ✓ Call average from 10 min to 5 min
- ✓ Customer Satisfaction from 47 % to 92%

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

**EMOTIONS-DRIVEN (30% of the population)**
**THOUGHTS-DRIVEN (25%)**
**REACTIONS-DRIVEN (20%)**
**OPINIONS-DRIVEN (10%)**
**REFLECTIONS-DRIVEN (10%)**
**ACTIONS-DRIVEN (5%)**

http://www.eloyalty.com ; http://www.mattersight.com/ ; http://www.fastcompany.com/1706766/how-personality-test-designed-pick-astronauts-taking-pain-out-customer-support ; http://www.ssca.com/resources/articles/104-the-history-of-the-process-communication-model-in-astronaut-selection ; http://www.forbes.com/forbes/2011/0214/entrepreneurs-kelly-conway-software-eloyalty-your-pain.html
Cook, Scott (October 2013). "Personality Matters: Behavioral analytics is now a reality in contact centres". Direct Marketing Magazine 26 (3): 5.

# Obama 2012 campaign

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.
- Ranking the 20% of Obama's 2008 vote that shifted to undecided on a 0-10 persuasion score
- **62,000 computer simulations** of likely voter behavior

## ➢ Outcome

- Obama paid 35% less per broadcast commercial than opponent 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
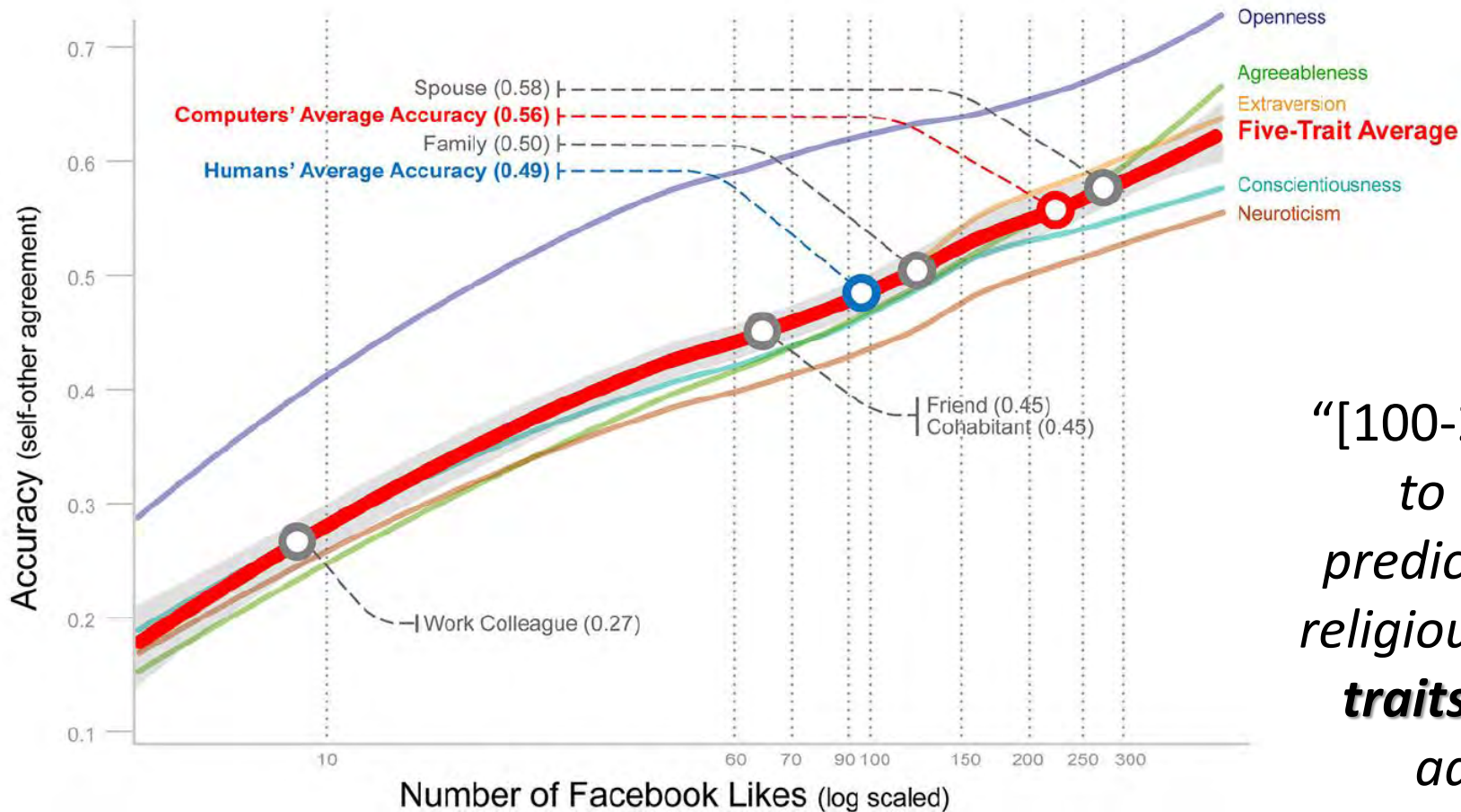- **Change voting behavior of 78 % of targeted undecided voters through Facebook**

Sources: Woodie, A. (2013, June 7). Big Data Analytics Give Electoral Edge. Datanami. Kolb, J., & Kolb, J. (2013). The Big Data Revolution. CreateSpace Independent Publishing Platform. Madrigal, A. C. (2012, November 16). When the Nerds Go Marching In. The Atlantic. Rutenberg (2013), Data You Can Believe In The Obama Campaign's Digital Masterminds Cash In; NYT.

# Machine learning knows us better than we ourselves



"*[100-250] 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…*"

Source: Youyou, W., Kosinski, M., & Stillwell, D. (2015). Computer-based personality judgments are more accurate than those made by humans. *PNAS*, 201418680.
Kosinski, M., Stillwell, D., & Graepel, T. (2013). Private traits and attributes are predictable from digital records of human behavior. *PNAS*, *110*(15), 5802–5805.

July 2016: $100,000
August: $250,000
September: $ 5 million

Cambridge Analytica

**32 personality types in 17 states**

PREDICTING TRUMP
©CBSN   CBS News. Always On.

RACE FOR THE WHITE HOUSE
CAMPAIGN DELIVERS PSYCHOLOGICALLY TAILORED ADS   CNN
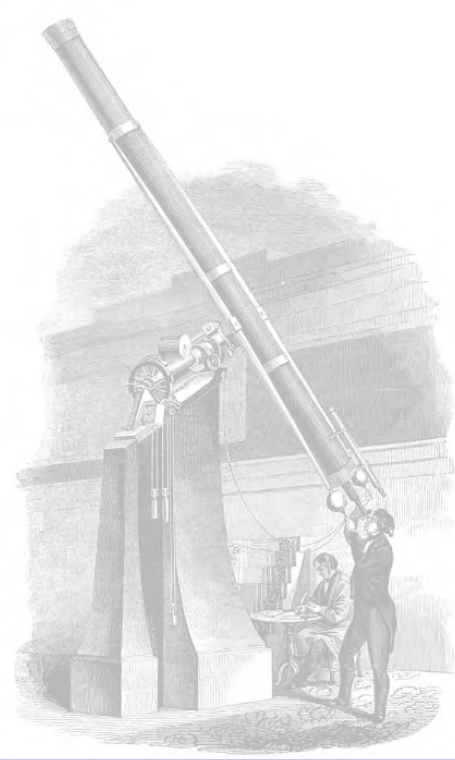Isa Soares | CNN Correspondent   DAX ▲ 5.28

Audience Insight

Deeper insight into the people who matter most.

Our psychographic analysis is a powerful and unique tool for gaining a deeper knowledge of your audience groups by revealing the core personality traits and motivations that drive behavior.

Audience Profiles  +  Research Psychologist  →  Motivation Understanding

Ownership of American-built car (phone app)
Haitians: Clinton Foundation Haiti Earthquake
Afro-Americans: Clinton's superpredators soundbite

Psych: 2nd amendment: fear or tradition?

3rd Debate: 175,000 variations of Trump's arguments
Differences in title, subtitle, color, picture, video, etc

**The** Theory, Practice and **Limits of**
**Big Data for the Social Sciences**

**Data ≠ Reality**

**Meaning ≠ Meaningful**
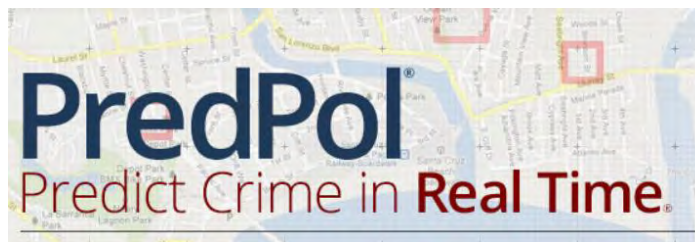
**Correlation ≠ Causation**

**Past ≠ Future**

# Data ≠ Reality

**PredPol**
Predict Crime in **Real Time**

## Homicide Parole candidates
○ 60 – 70 % correct who commits homicide

"We kill people based on metadata"

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…*"
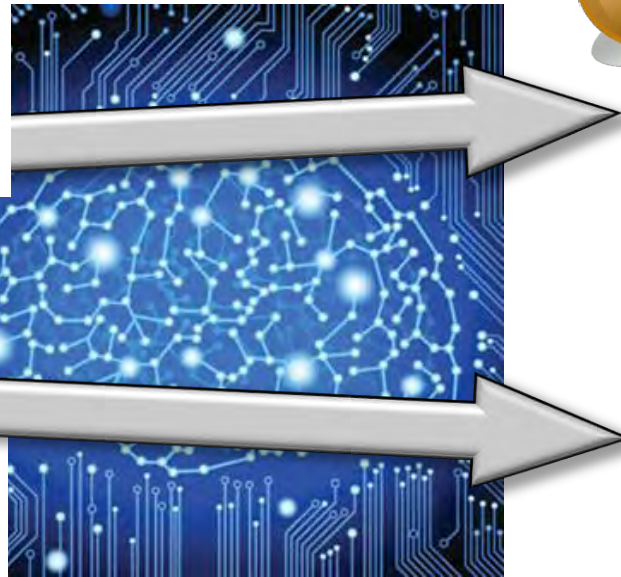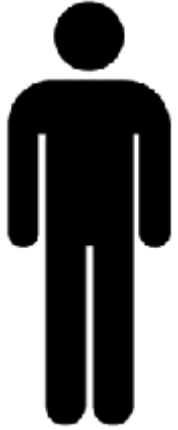
Michael Hayden
former Director
NSA & CIA

Berk, R., Sherman, L., Barnes, G., Kurtz, E., & Ahlman, L. (2009). Forecasting murder within a population of probationers and parolees: a high stakes application of statistical learning. *Journal of the Royal Stat.Soc.: Series A*, *172*(1), 191–211. http://spectrum.ieee.org/podcast/at-work/innovation/can-software-predict-repeat-offenders ; http://www.spiegel.de/netzwelt/web/in-santa-cruz-sagen-computer-verbrechen-voraus-a-899422.html ;http://www.sfgate.com/default/article/Sci-fi-policing-predicting-crime-before-it-occurs-3725708.php ; Wikipedia Commons; Scahill, J., & Greenwald, G. (2014). The NSA's Secret Role in the U.S. Assassination Program. *The Intercept*.

# Meaning ≠ Meaningful

"John, Paul, Mike, Kevin, Bill"

"executive, management, professional, corporation, salary, office, business, career."
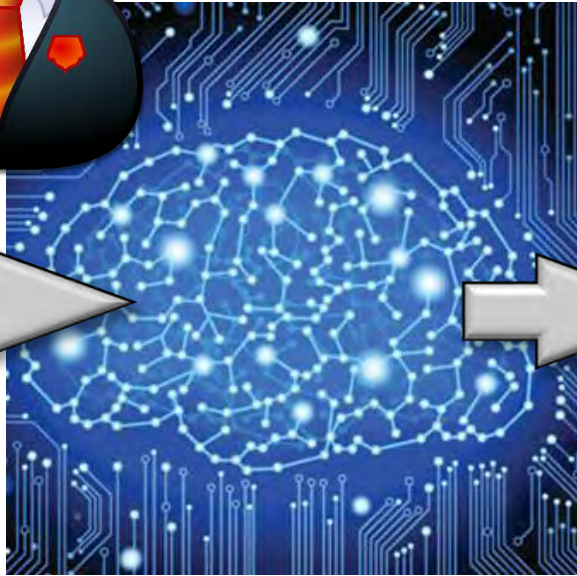
"Amy, Lisa, Sarah, Diana, Ann"

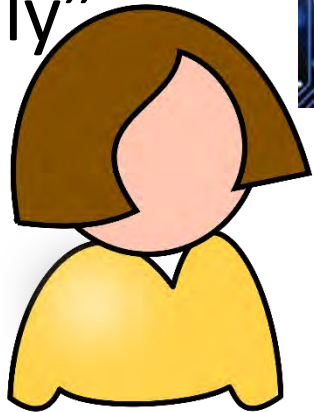"home, parents, children, family, marriage, wedding, relatives"

Caliskan-Islam, Bryson, Narayanan (2016). Semantics derived automatically from language corpora necessarily contain human biases. *arXiv:1608.07187 [Cs]*
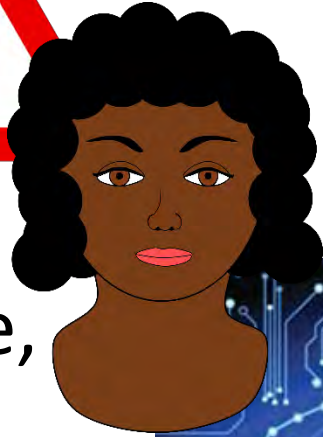
# Meaning ≠ Meaningful

"Harry, Katie, Jonathan, Nancy, Emily"

"freedom, health, love, peace, friend, heaven, gentle, loyal, lucky, diploma, happy, laughter, vacation"

Caliskan-Islam, Bryson, Narayanan (2016). Semantics derived automatically from language corpora necessarily contain human biases. *arXiv:1608.07187 [Cs]*

# Meaning ≠ Meaningful

"Jerome, Ebony, Jasmine, Latisha, Tia"

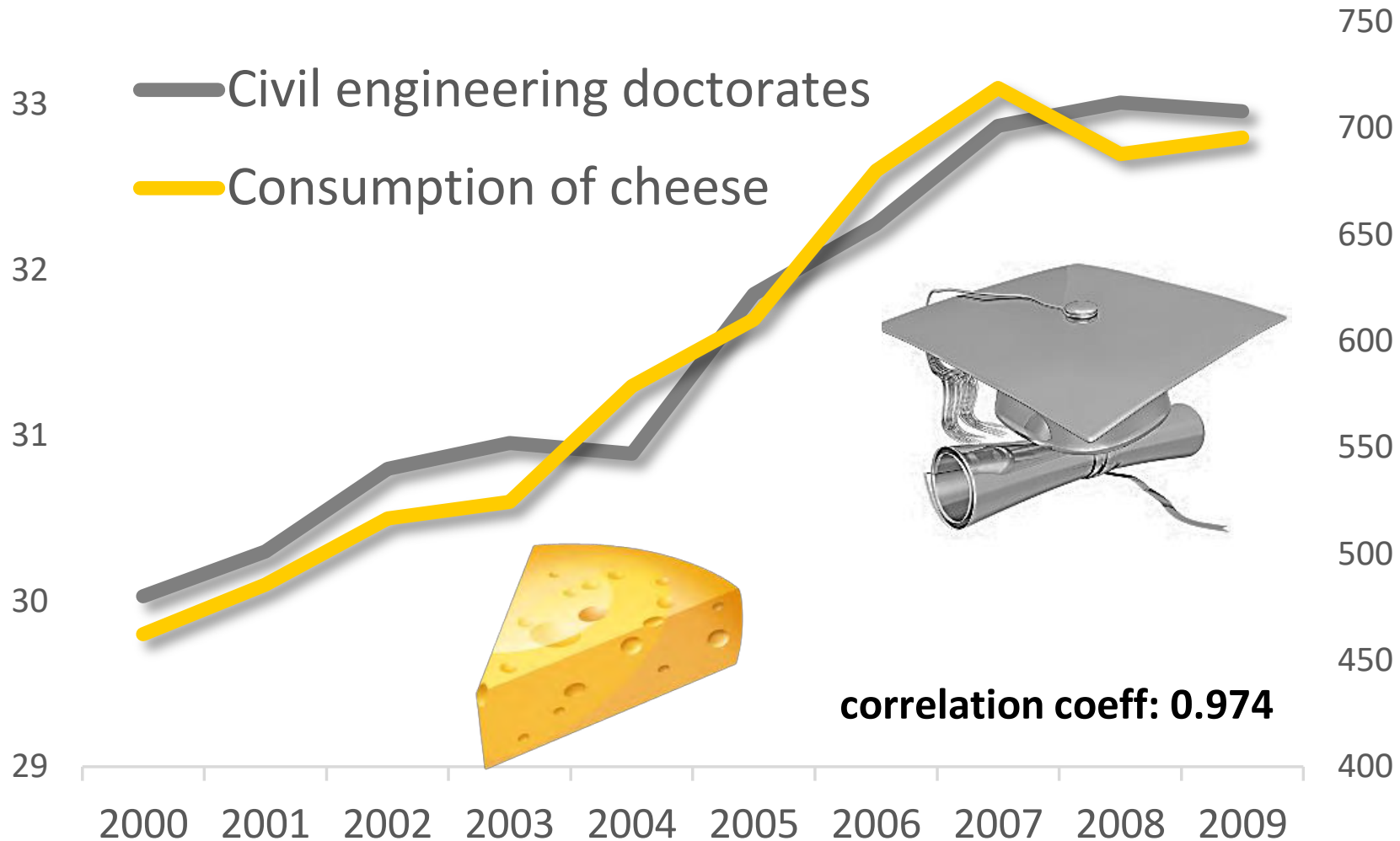"abuse, filth, sickness, accident, death, grief, poison, assault, poverty, ugly, evil, agony, prison."
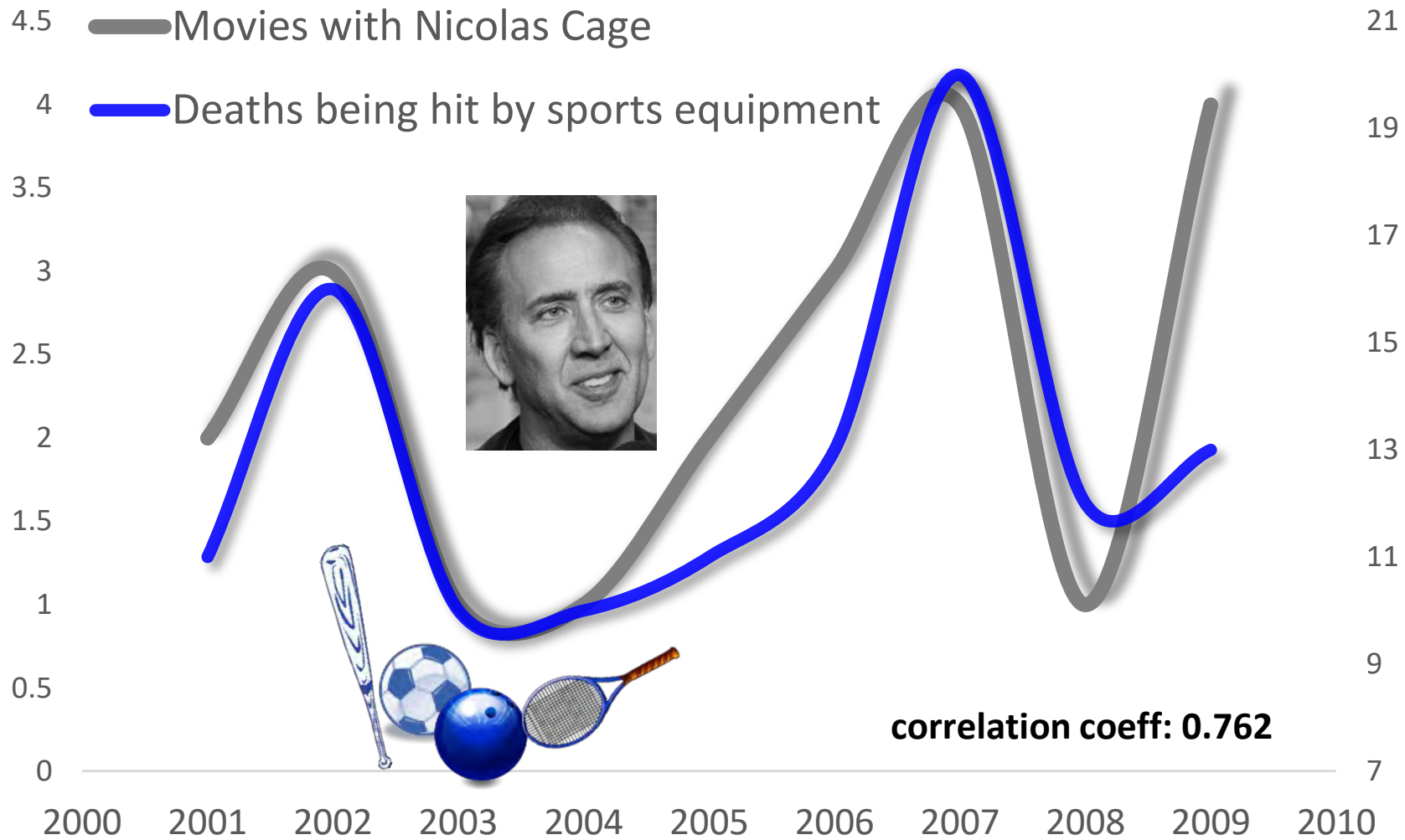
66 %

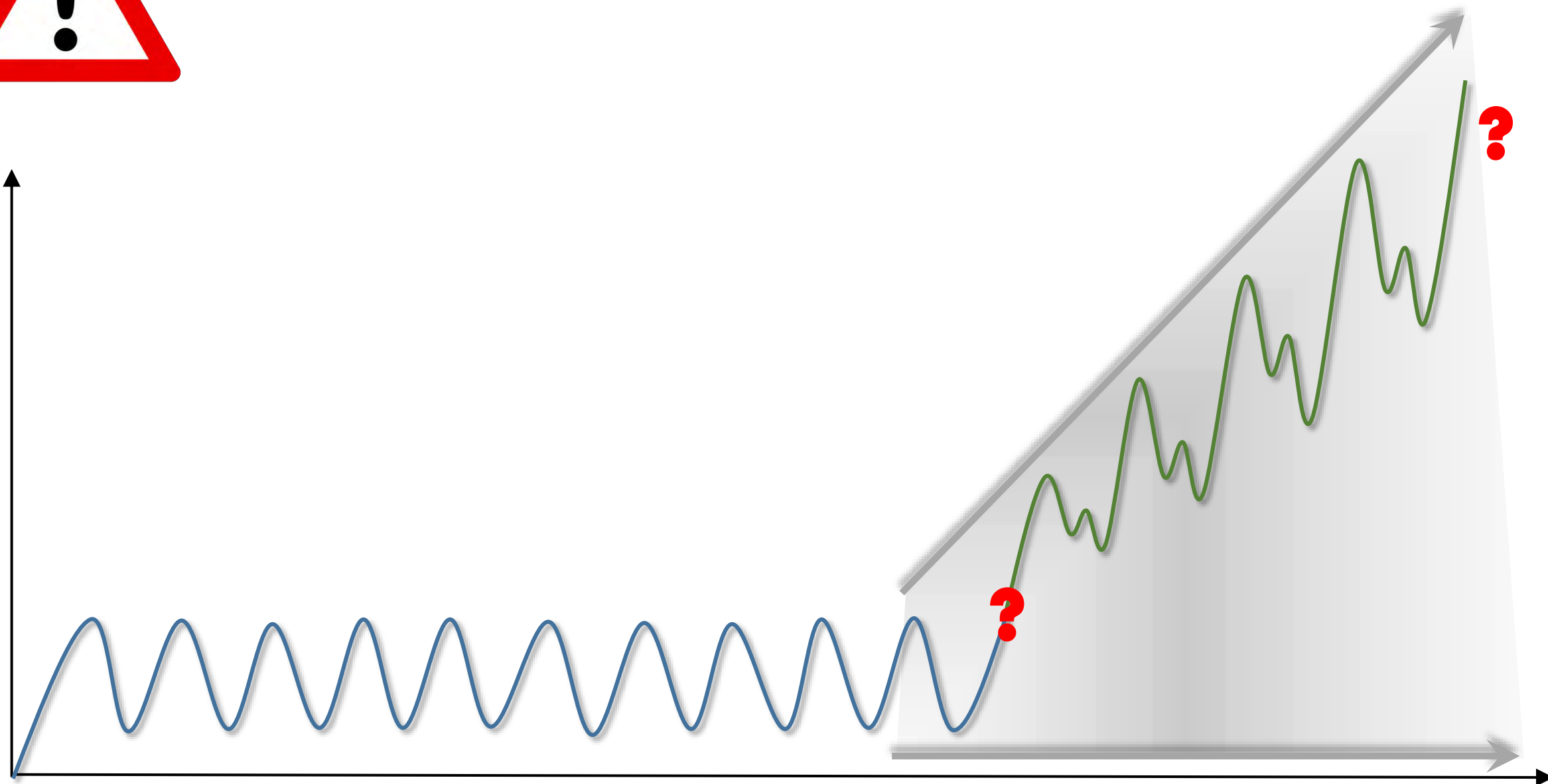Caliskan-Islam, Bryson, Narayanan (2016). Semantics derived automatically from language corpora necessarily contain human biases. *arXiv:1608.07187 [Cs]*

# Correlation ≠ Causation



Movies with Nicolas Cage

Deaths being hit by sports equipment

correlation coeff: 0.762

Past ≠ Future

Sources: Bohemia Interactive Simulations, http://youtu.be/G9P9bUTCdpA ; TRANSIMS: http://www.youtube.com/watch?v=mN7kq0ITAys ; Epstein, http://www.youtube.com/watch?v=wZZJClGtVkw



# Computational Social Science

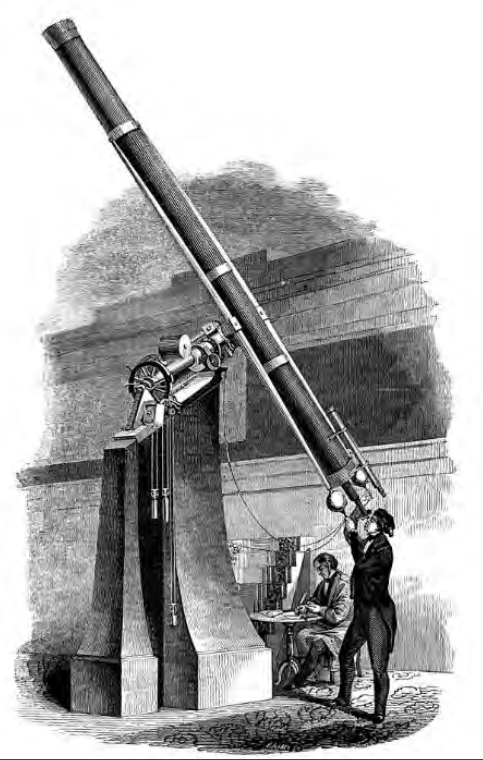ECONOMETRIC POLICY EVALUATION: A CRITIQUE

Robert E. Lucas, Jr.

"...any change in policy will systematically alter the structure of econometric models"
**(1976)**

# Main References:

> Hilbert (2016). Big Data for Development: A Review of Promises and Challenges. *Development Policy Review*, 34(1), 135–174. https://doi.org/10.1111/dpr.12142

> Hilbert, M. (2015). ICT4ICTD: Computational Social Science for Digital Development. 48th (HICSS) (pp. 2145–2157). *IEEE Computer Society*. https://doi.org/10.1109/HICSS.2015.258

> Gillings, Hilbert & Kemp (2016). Information in the Biosphere: Biological and Digital Worlds. *Trends in Ecology & Evolution*, *31(3), 180–189* www.martinhilbert.net/information-in-the-biosphere/

> Hilbert & López (2011). The world's technological capacity to store, communicate and compute information. *Science*, 332, 6025, 60-65 www.martinhilbert.net/WorldInfoCapacity.html

> Hilbert (2016). The bad news is that the digital access divide is here to stay: Domestically installed bandwidths among 172 countries for 1986–2014. *Telecommunications Policy*. www.martinhilbert.net/the-bad-news-is-that-the-digital-access-divide-is-here-to-stay

# The Theory, Practice and Limits of
# Big Data for the Social Sciences

**UCDAVIS**
UNIVERSITY OF CALIFORNIA

**Martin Hilbert**
**Department of Communication**
**hilbert@UCDavis.edu**

# N = n ?
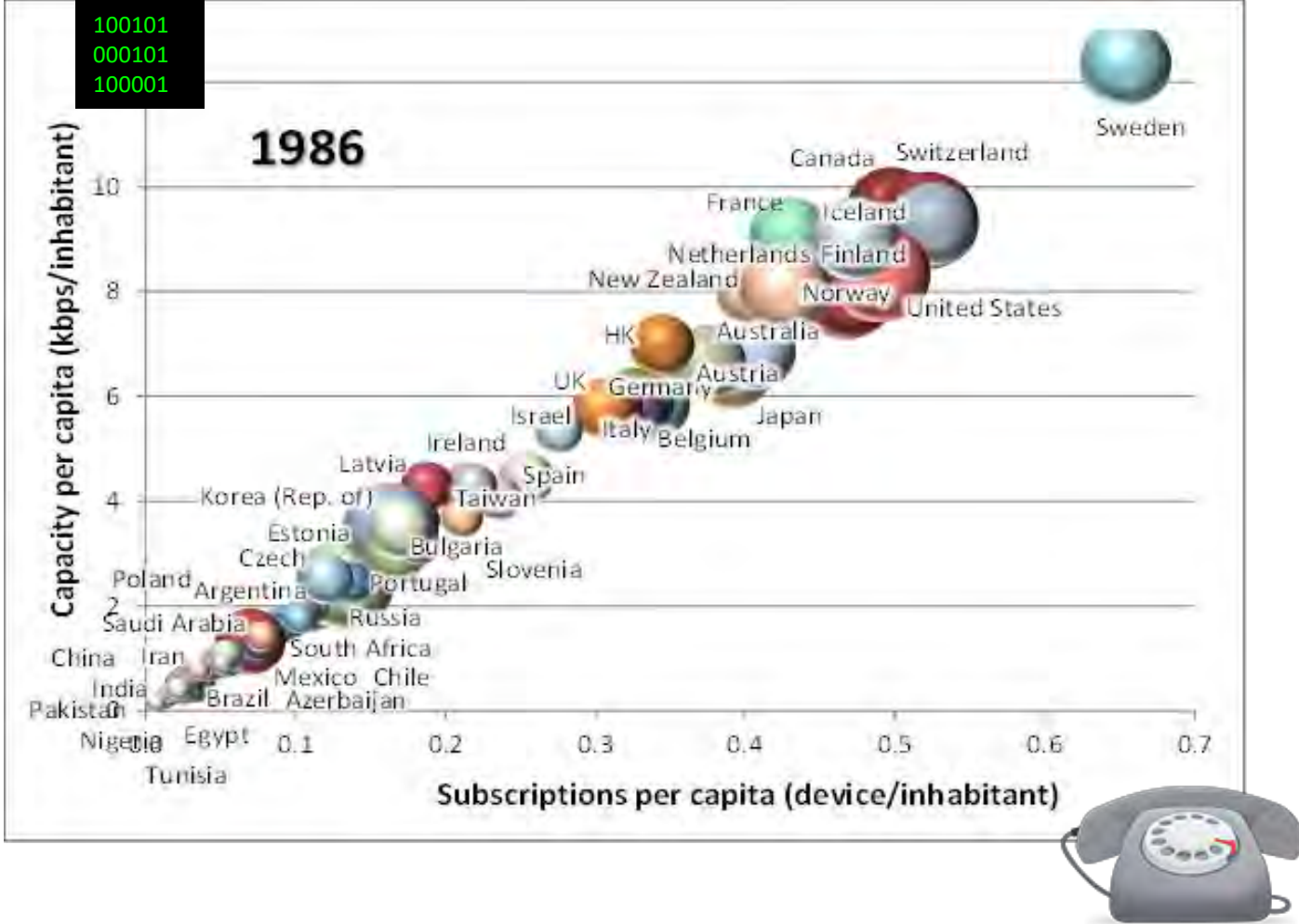
(a) Rwanda 2005/09:

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).

1986

So we're all finally good now! ...(?)

Increasing heterogeneity

1986

100101
000101
100001

Source: Hilbert, M. (2013), Technological information inequality as an incessantly moving target: The redistribution of info. and communication capacities between 1986 and 2010. *Journal of the Assoc. for Info. Science and Technology*. http://www.martinhilbert.net/TechInfoInequality.pdf

# Number of subscriptions of countries

Telecom: OECD vrs. the rest of world
(fixed and mobile Internet and telephony SUBSCRIPTIONS per capita)



$$\frac{1.2 \text{ devices}}{0.2 \text{ devices}} = 6$$
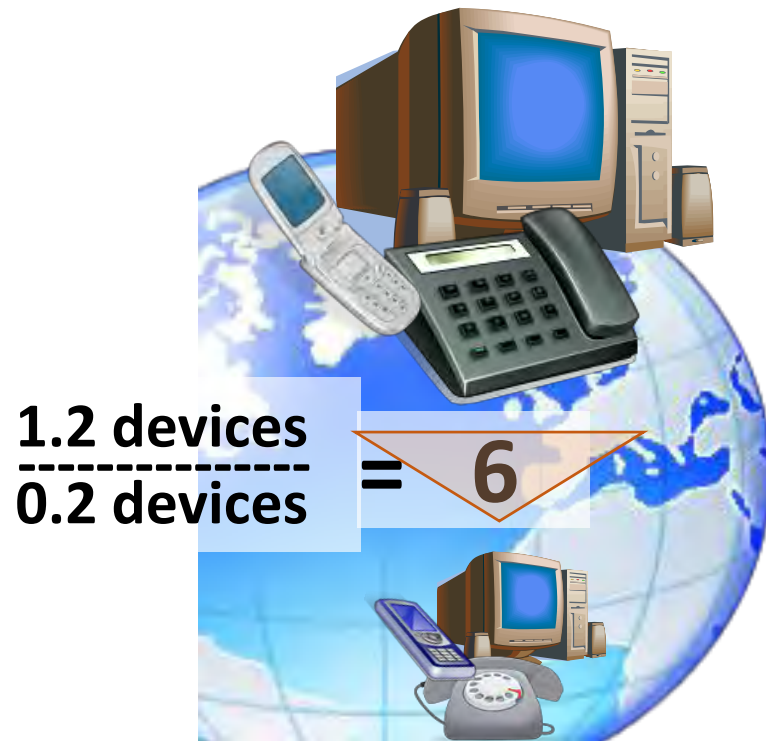
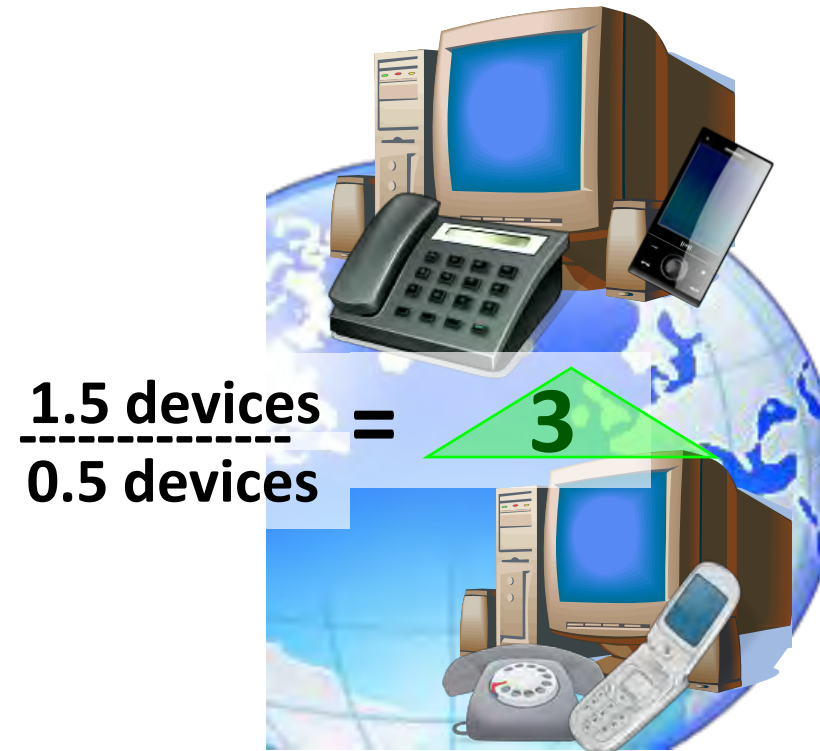$$\frac{1.5 \text{ devices}}{0.5 \text{ devices}} = 3$$

**2001**

**2006**

Source: Hilbert, M. (2013), Technological information inequality as an incessantly moving target: The redistribution of info. and communication capacities between 1986 and 2010. *Journal of the Assoc. for Info. Science and Technology*. http://www.martinhilbert.net/TechInfoInequality.pdf
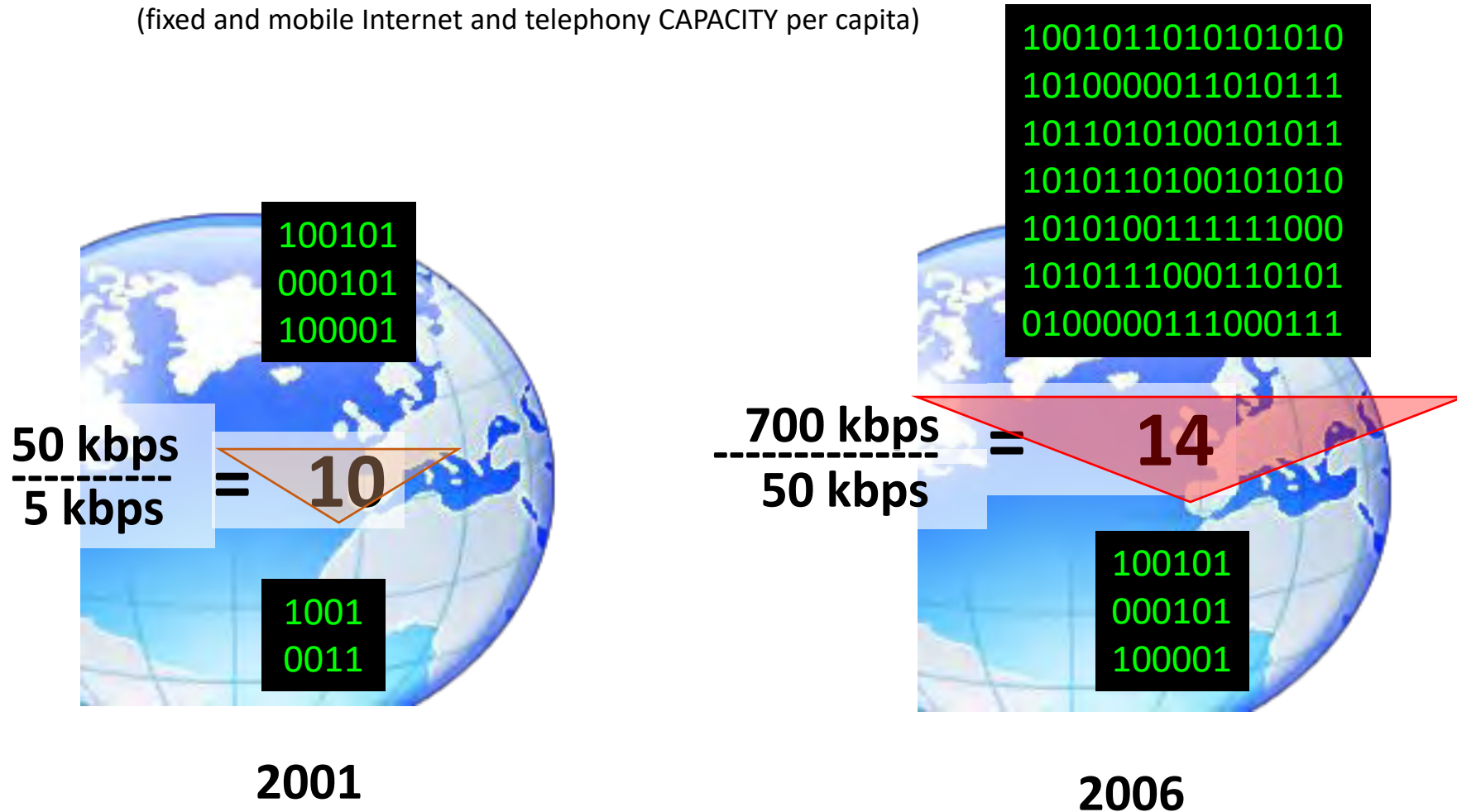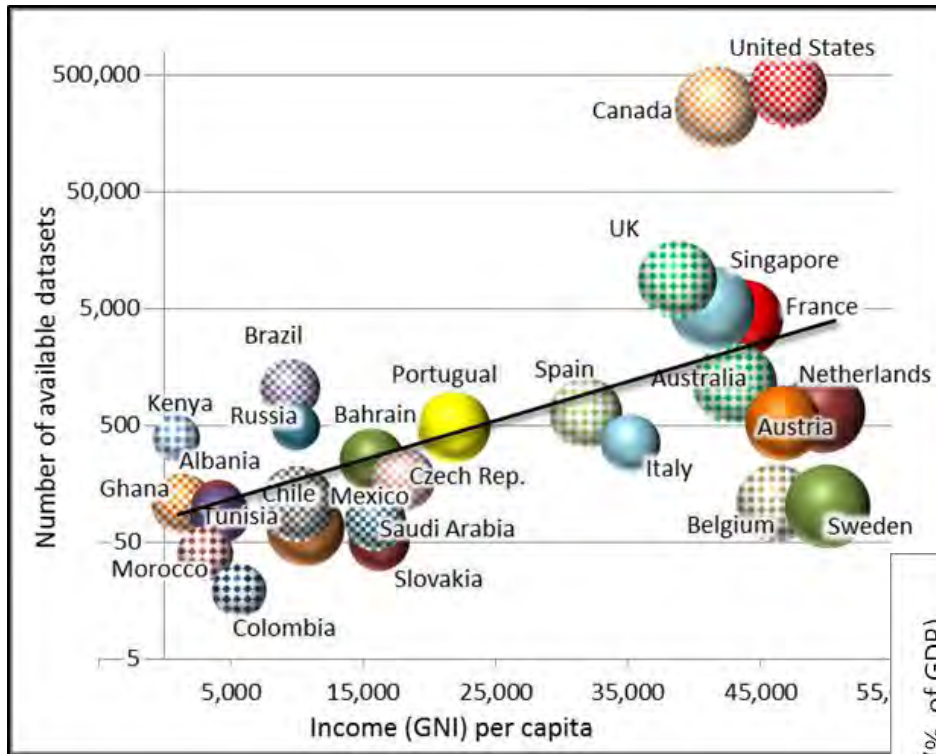
# Telecommunication capacity of countries

Telecom: OECD vrs. the rest of world
(fixed and mobile Internet and telephony CAPACITY per capita)



$$\frac{50 \text{ kbps}}{5 \text{ kbps}} = 10$$

**2001**

$$\frac{700 \text{ kbps}}{50 \text{ kbps}} = 14$$

**2006**

Source: Hilbert, M. (2013), Technological information inequality as an incessantly moving target: The redistribution of info. and communication capacities between 1986 and 2010. *Journal of the Assoc. for Info. Science and Technology.* http://www.martinhilbert.net/TechInfoInequality.pdf

# Content Divide



*Open Government data: Number of datasets provided on central government portal (vertical y-axis, logarithmic scale), Gross National Income per capita (horizontal x-axis), Corruption Perception Index (size of bubbles: larger bubbles, more transparent) (year=2011; n=27).*

**Public data on natural resource extraction:** *Natural resource rent vs. government data disclosure (year=2010; n=40).*