

Challenges in Making the Hidden Visible

*Based on material in 05-389
Interactive Data Science
data.cmu.edu for more info
t/th 9-10:30, Spring*

Jennifer Mankoff and many other collaborators



Information can change how we act in the world [e.g., CHI'09, ICWSM'09]



stepgreen^{beta} enrich your life.

Report Actions Share Account Help About Admin Logged in as: **jmankoff**

Show time graphs Each square is a morewood gardens east tower resident.

Your savings	This week saved	Overall saved	Other person savings	This week saved	Overall saved
Turn off work screen saver	127 lbs	23260		18 lbs	1514 lbs
Use CFLs					
Teach proper hand washing.					
Plant a tree					
Insulate water heater					
Programmable thermostat					
Turn off home screen saver					
Lower water heater					
Use sleep mode at home			6 lbs		
				Use sleep mode at work	

Your saving displayed on left. The space contains the s... Information anonymous p... residing in mo... gardens east... Exit dorm

Information For and About People

As a prosthetic, changes what we can do in the world

As a motivator, changes how we (or our machines) behave

Shared with others, may

- Change the balance of power [UbiComp '09, '10]

- Build new kinds of action and knowledge

Beyond individuals, may support policy, politics, economics [DIS In Submission]



Making Data Actionable

Data

- Collecting and Interpreting



Information

- Sense



Knowledge

- Visualize
Act
Adapt



Value

Collecting the right data

What is the problem? What data will solve the problem? How can we get that data?

Techniques needed: Careful analysis & thought

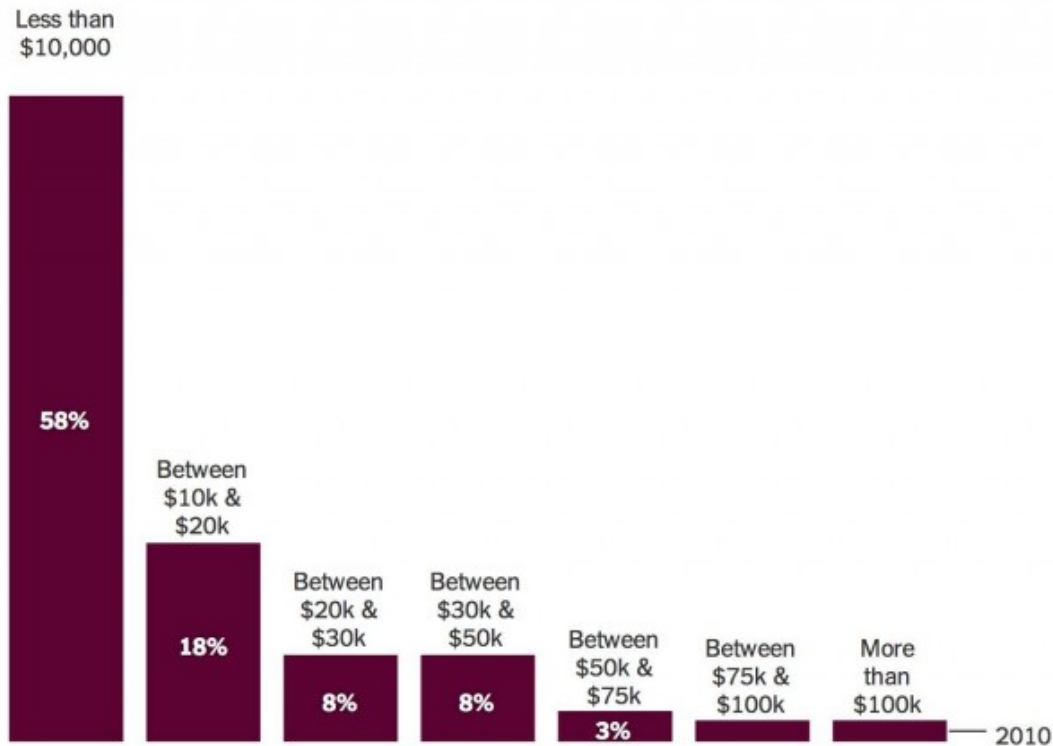
Tools: simulation & prototyping



Can you trust your data?

Large Amounts of Student Debt Are Not Common

Only 7 percent of young-adult households with student debt have more than \$50,000 in such debt.



Source: Elizabeth Akers and Matthew Chingos, Brookings Institution

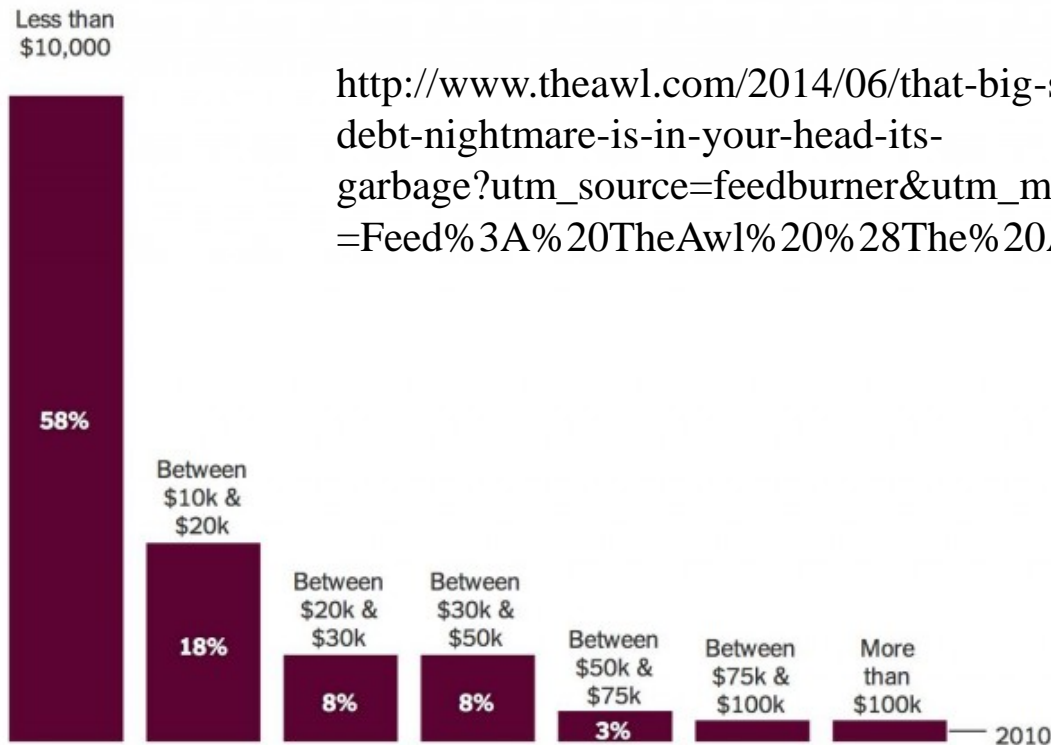
2010 data; based on households with people between 20 to 40 years old with at least some education debt



Can you trust your data?

Large Amounts of Student Debt Are Not Common

Only 7 percent of young-adult households with student debt have more than \$50,000 in such debt.



http://www.theawl.com/2014/06/that-big-study-about-how-the-student-debt-nightmare-is-in-your-head-its-garbage?utm_source=feedburner&utm_medium=feed&utm_campaign=Feed%3A%20TheAwl%20%28The%20Awl%29

Source: Elizabeth Akers and Matthew Chingos, Brookings Institution

2010 data; based on households with people between 20 to 40 years old with at least some education debt



Can you trust your data?

Just because you have a lot of data, does not mean that it is *good* data

The plural of “anecdote” is not “data”.



Sources of Error

Sampling errors

Random Errors due to the sample forming only part of the population

Systematic Bias in sampling

Bias During Data Collection

Demand Characteristics

Illusory Superiority

Data entry / processing errors

Data is generated accurately but errors introduced during recording or processing



What Makes a Good Sample?

Representative of the population

(Along dimensions that matter to the question being asked)

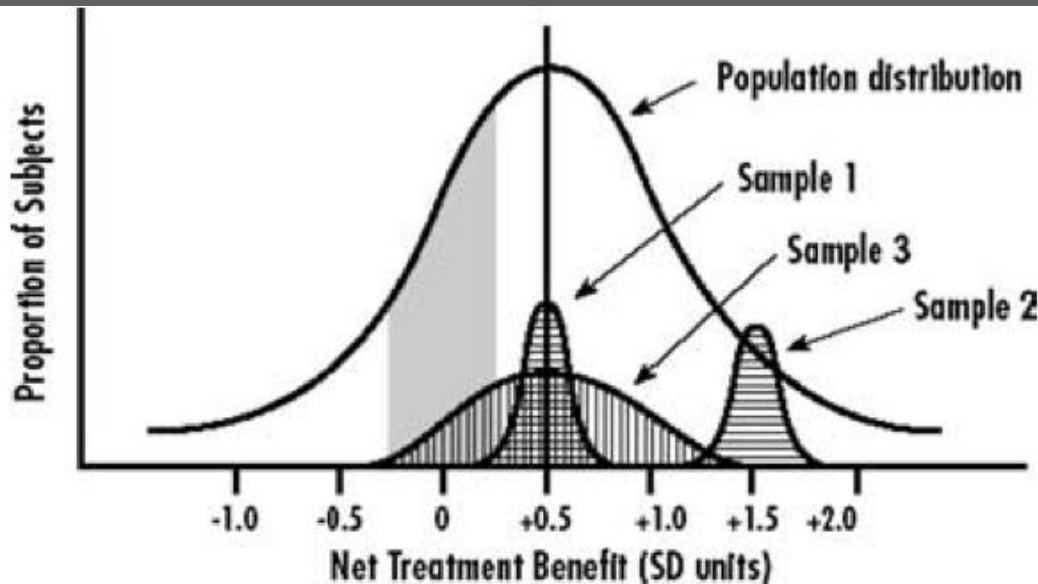


Figure 1: Kravitz *et al*, (2004) Evidence-based medicine, heterogeneity of treatment effects, and the trouble with averages. *The Milbank Quarterly* **82**(4):661-687

Sources of Error

Sampling errors

Random Errors due to the sample forming only part of the population

Systematic Bias in sampling

Bias During Data Collection

Demand Characteristics

Illusory Superiority

Data entry / processing errors

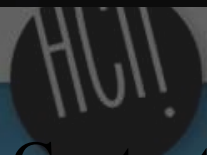
Data is generated accurately but errors introduced during recording or processing



Example: Demand Characteristics

An experimental artifact where participants form an interpretation of the experiment's purpose and unconsciously change their behavior to fit that interpretation





What can we do to minimize demand characteristics in HCI?

Be aware that response bias affects all studies

Dissociate from a particular design or solution

Minimize the differences in social status between investigators and participants

Use triangulation to validate data collected

Tricks for asking sensitive questions



Measurement Errors

Badly Designed Questions

Badly Chosen Sensors

Bad Administration of Measurement
Instrument

Inaccurate Measurements



Sources of Error

Sampling errors

Random Errors due to the sample forming only part of the population

Systematic Bias in sampling

Bias During Data Collection

Demand Characteristics

Illusory Superiority

Data entry / processing errors

Data is generated accurately but errors introduced during recording or processing



Four C's of Data Quality

Is your data *Complete*?

Is your data *Coherent*?

Is your data *Correct*?

Is your data *Accountable*?



Questions about Completeness

Appropriate Data: Does the data you have match the questions you want to answer?

Missing Data: Data does not exist because it was never obtained or was lost

Reporting error: The sensor (or respondent) is incorrect



Is your Data Coherent?

Does the data “add up”?

Does it make sense relative to itself?

Are there extreme values?

Examples

- Non number in a numeric field

- Month field has something other than a month

- Email has no @

- Hourly data adds up to 24 hrs per day

- Etc.*



Is your data Correct?

Itemize aspects of your data that are easy to verify

Compare (collect twice or find alt. source)

Analyze the data collection strategy and look for sources of bias

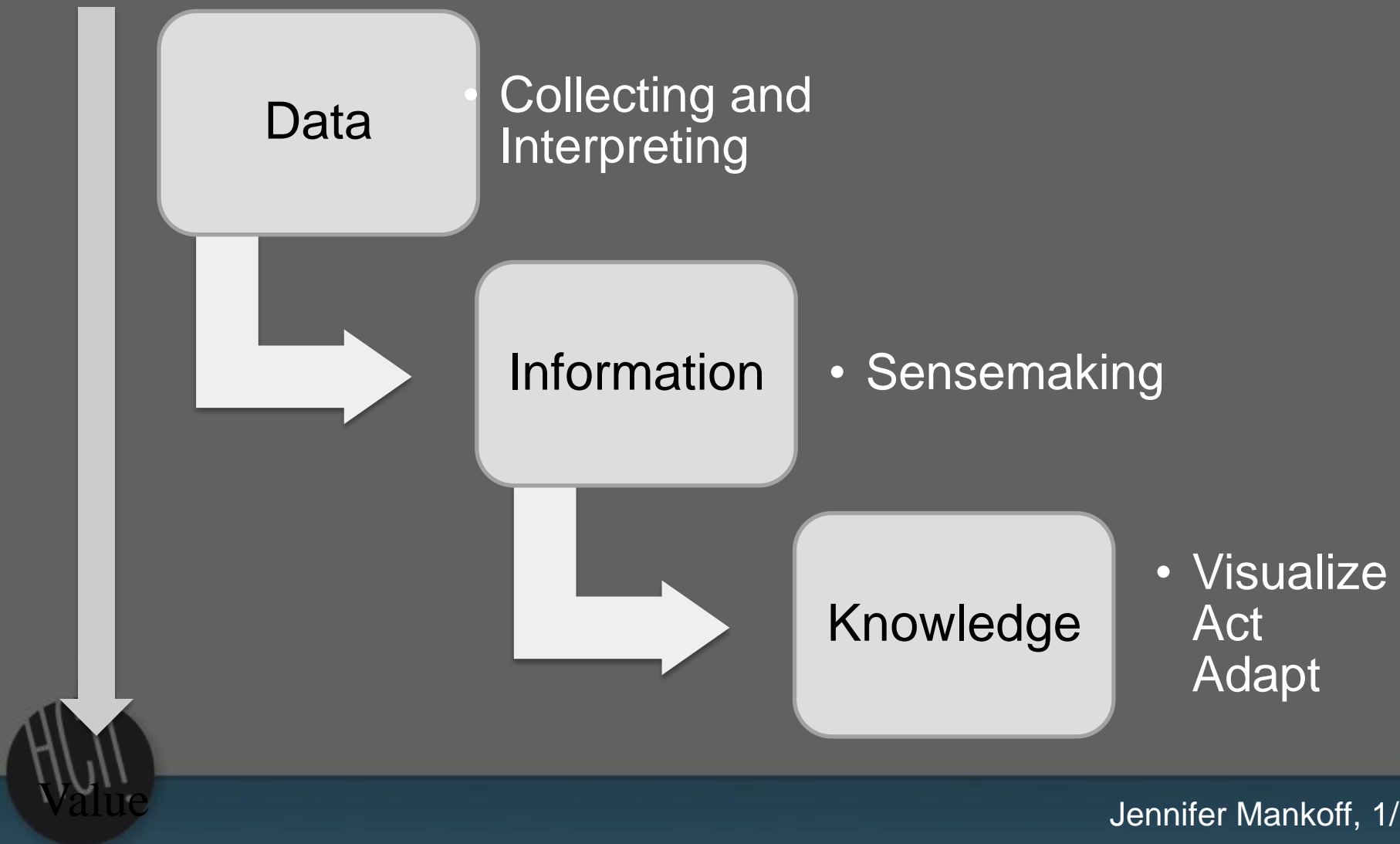
- Within the population

- across variables (surveys with only round values; people who report everything in round numbers)

Determine how much is bad



Making Data Actionable



Understanding Humans

Activities

Routines

Intent

Causality



Understanding Humans

Activities

Routines

Intent

Causality



Understanding Humans

Activities

Routines

Intent

Causality



Why are routines important?

Develop routines to reduce cognitive effort

Deviations and anomalies cause stress and extra effort

→ *at least as* actionable as inferred activities

Longer term, more built-in

Opportunities for change

Barriers for change

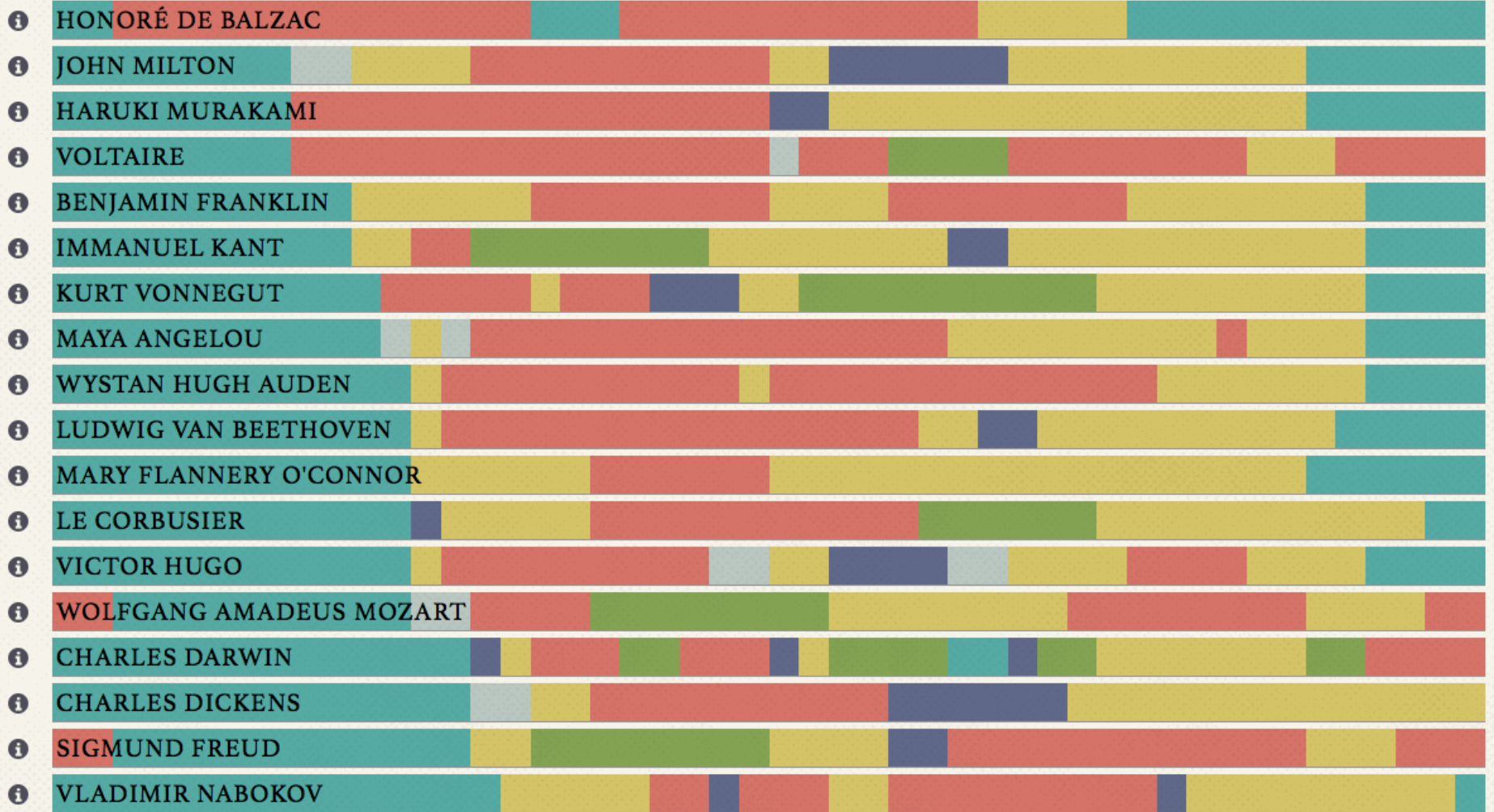
Leverage point for understanding human intent

Untapped as a resource: sensing and using



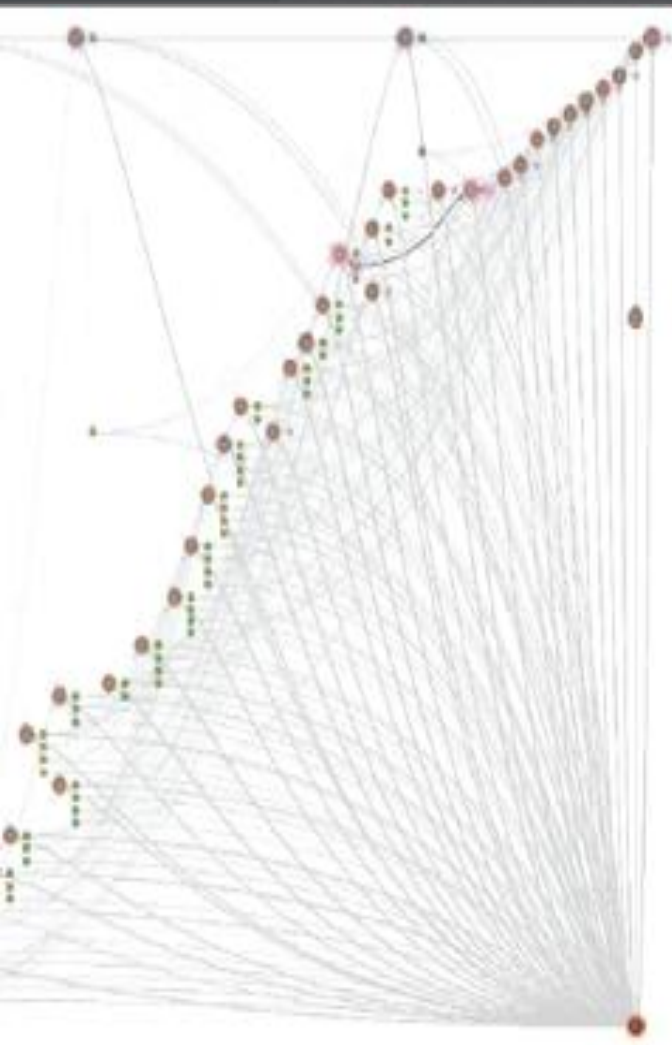
■ SLEEP
 ■ CREATIVE WORK
 ■ DAY JOB/ADMIN
 ■ FOOD/LEISURE
 ■ EXERCISE
 ■ OTHER

12 AM —————> 12 PM —————> 1 2 3 4 5 6 7 8 9 10 11 12 1 2 3 4 5 6 7 8 9 10 11 12

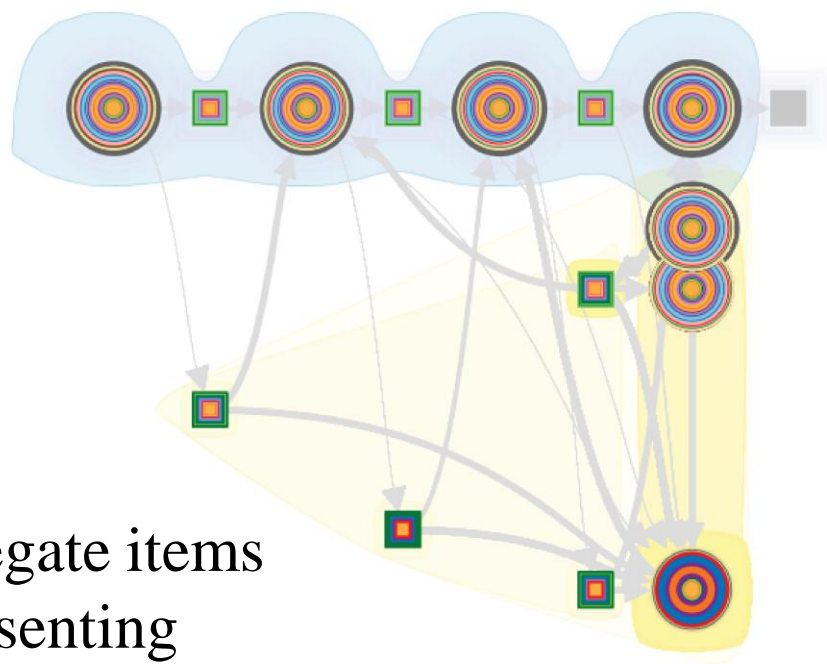
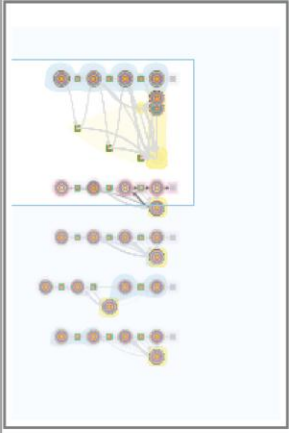


automatically extract *patterns* of human behavior that form routines and *deviations* from demonstrated behavior

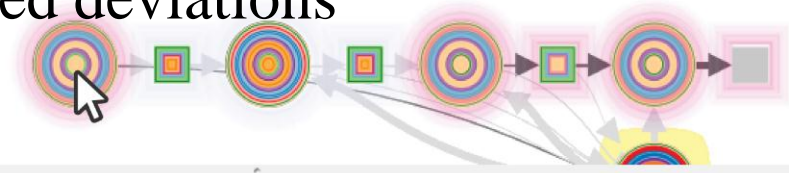




Study



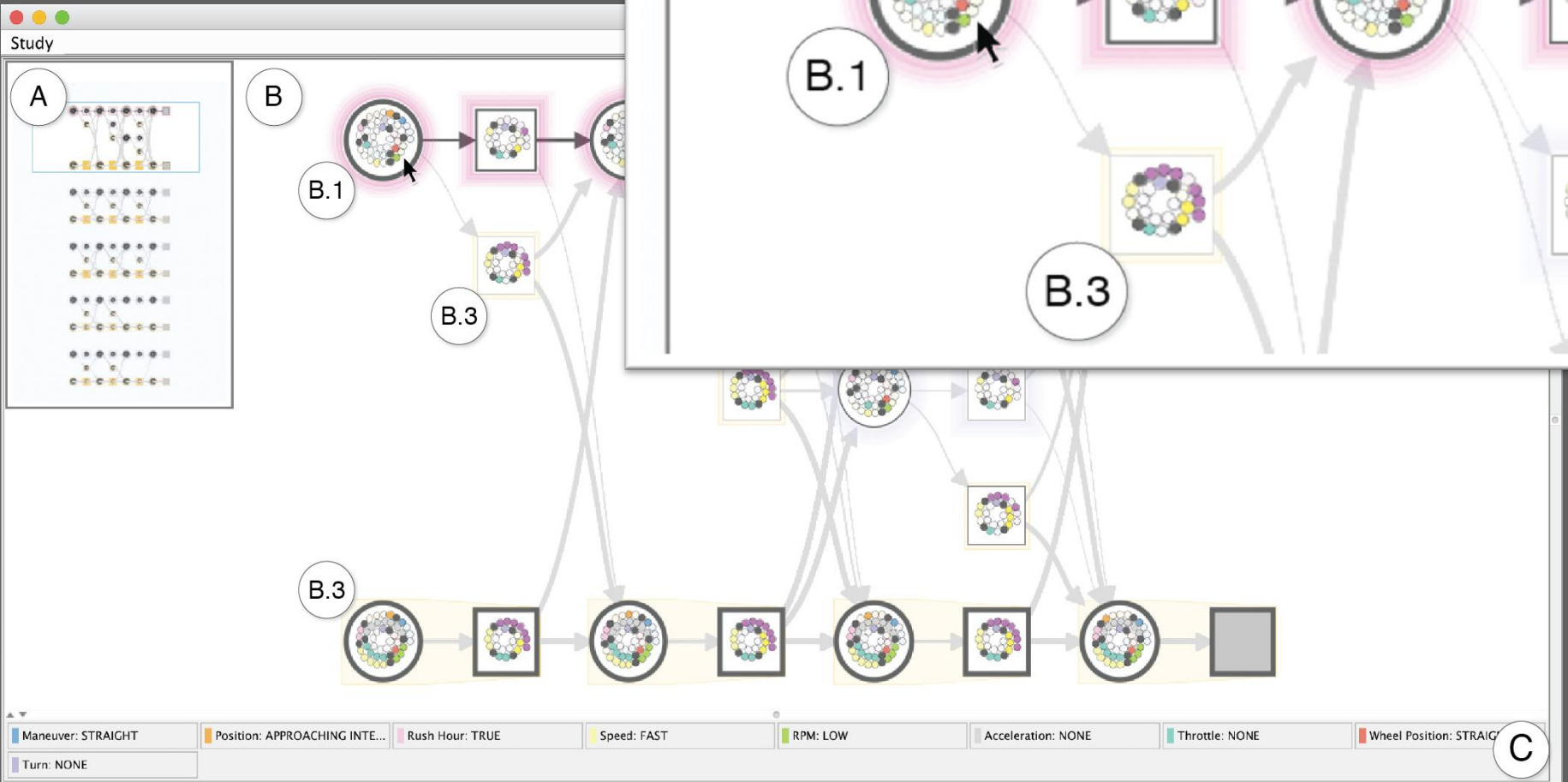
aggregate items
representing
extracted deviations



Maneuver: RIGHT TURN	Position: EXITED INTERSECTION	Rush Hour: TRUE	Speed: SLOW	RPM: LOW
Acceleration: NONE	Throttle: NONE	Wheel Position: RETURNING	Turn: SMOOTH	



Visualization



Now automated extraction (CHI 2017?)

Aggressive (Novice) Drivers

(CHI 2014, CHI 2016, CHI 2017 submission)

US: 1500 deaths/year

Cost of \$40 billion from crashes



Aggressive (Novice) Drivers: Interventions

Rush Hour

SPEED LIMIT 25

STOP

SPEED LIMIT 25

SELECT DRIVER

REPLAY

PREVIOUS

PLAY NON-AGGRESSIVE

NEXT

SPEEDOMETER: 0, 20, 40, 60, 80, 100, 120

GAS: HI, LOW

BRAKE: LOW, HI



Understanding Humans

Activities

Routines

Intent

Causality



Causality



Hear noise

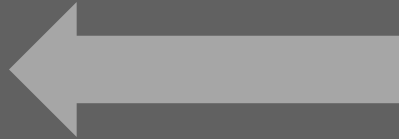
Past experience
leads to expectation



Go to window
to verify

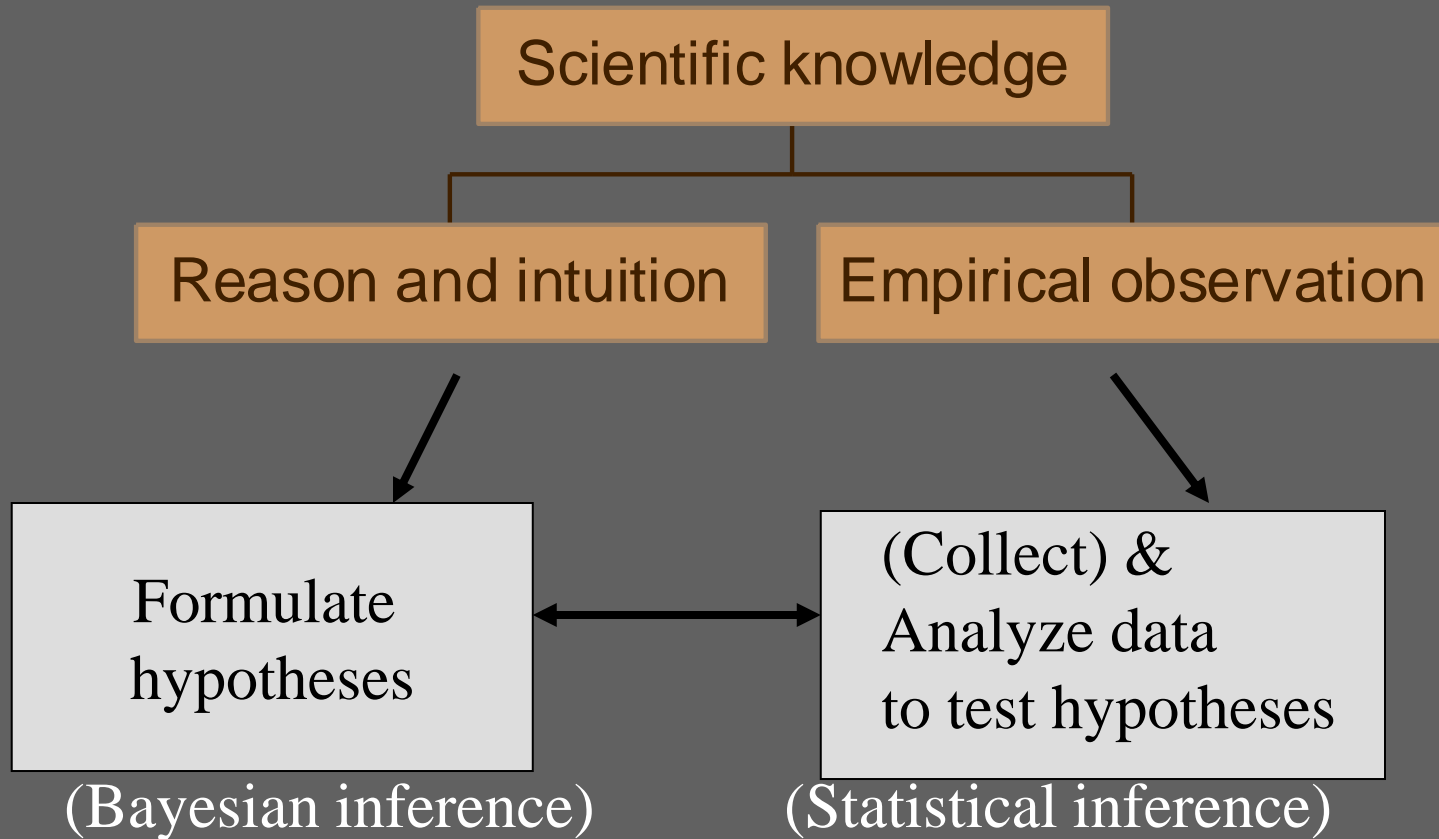


Prediction was false

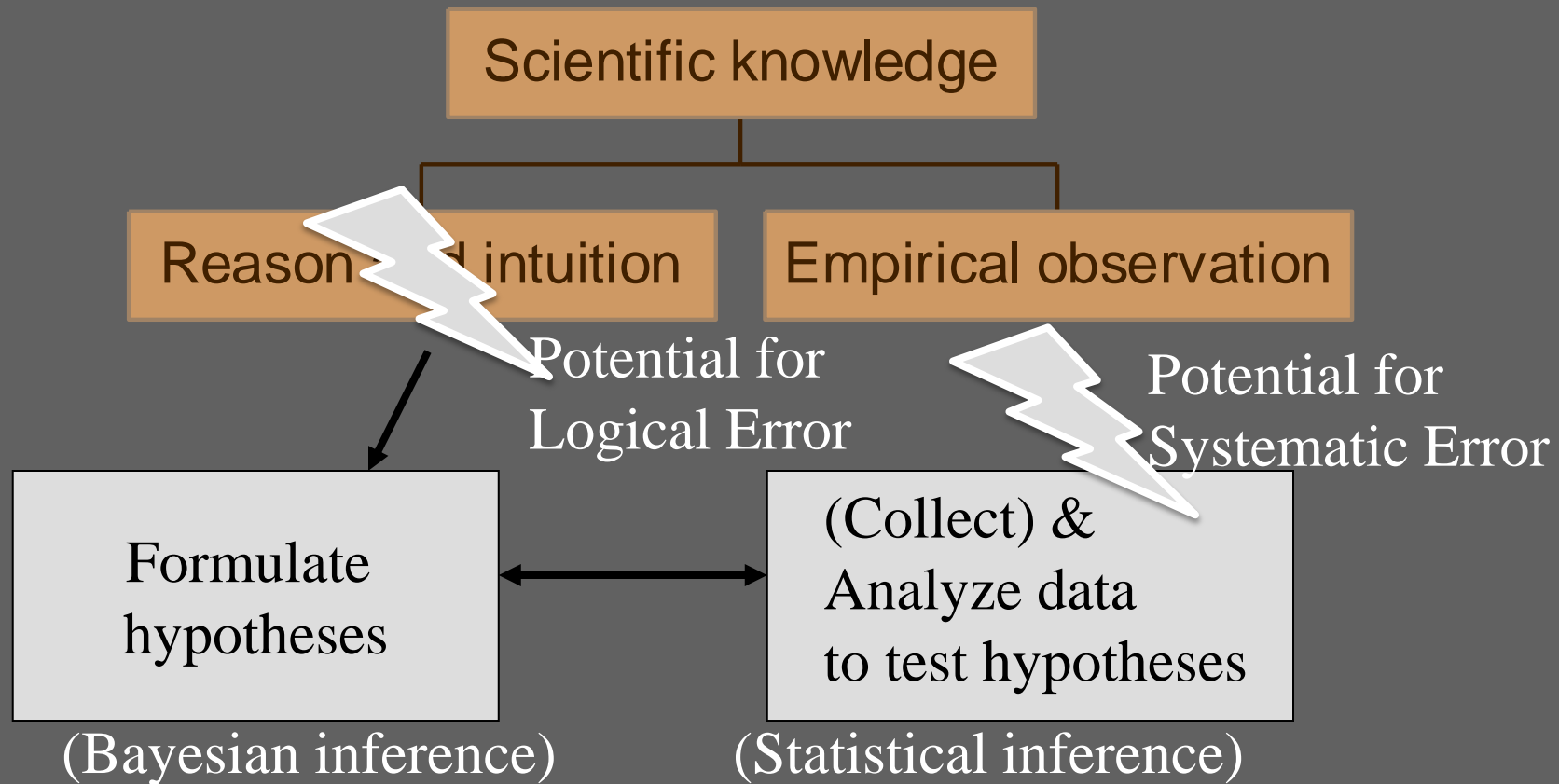


Only with empirical evidence, can we make and test predictions

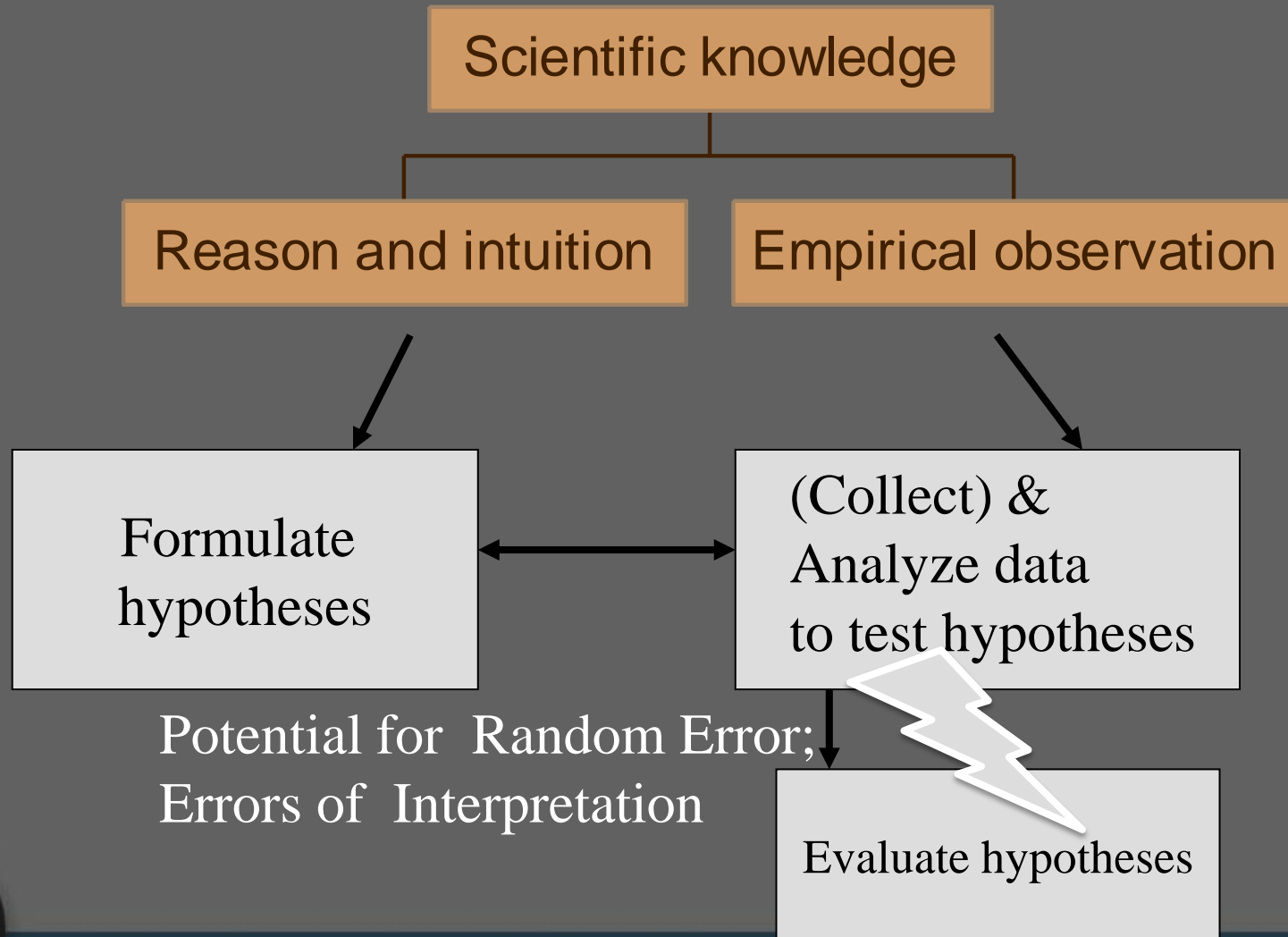
Process vs Frequency data



Process vs Frequency data



Process vs Frequency Data



Correlation is not Causation

There is a 0.91 correlation between ice cream consumption and drowning deaths.

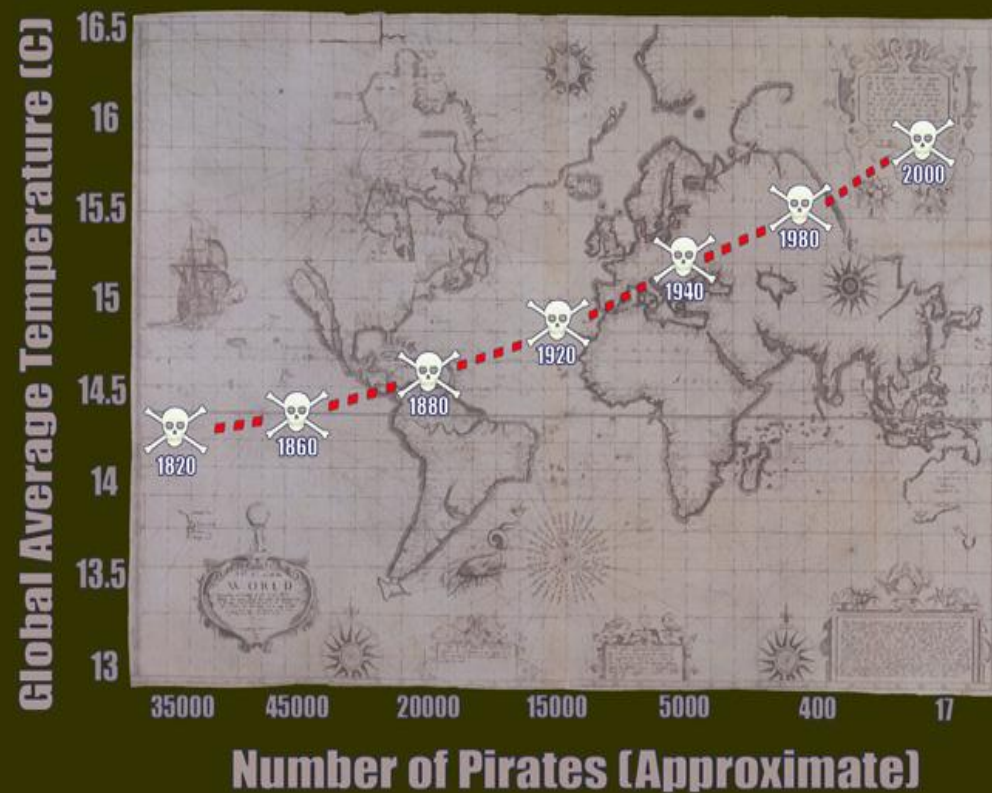
Does eating ice cream cause drowning?

Does grief cause us to eat more ice cream?

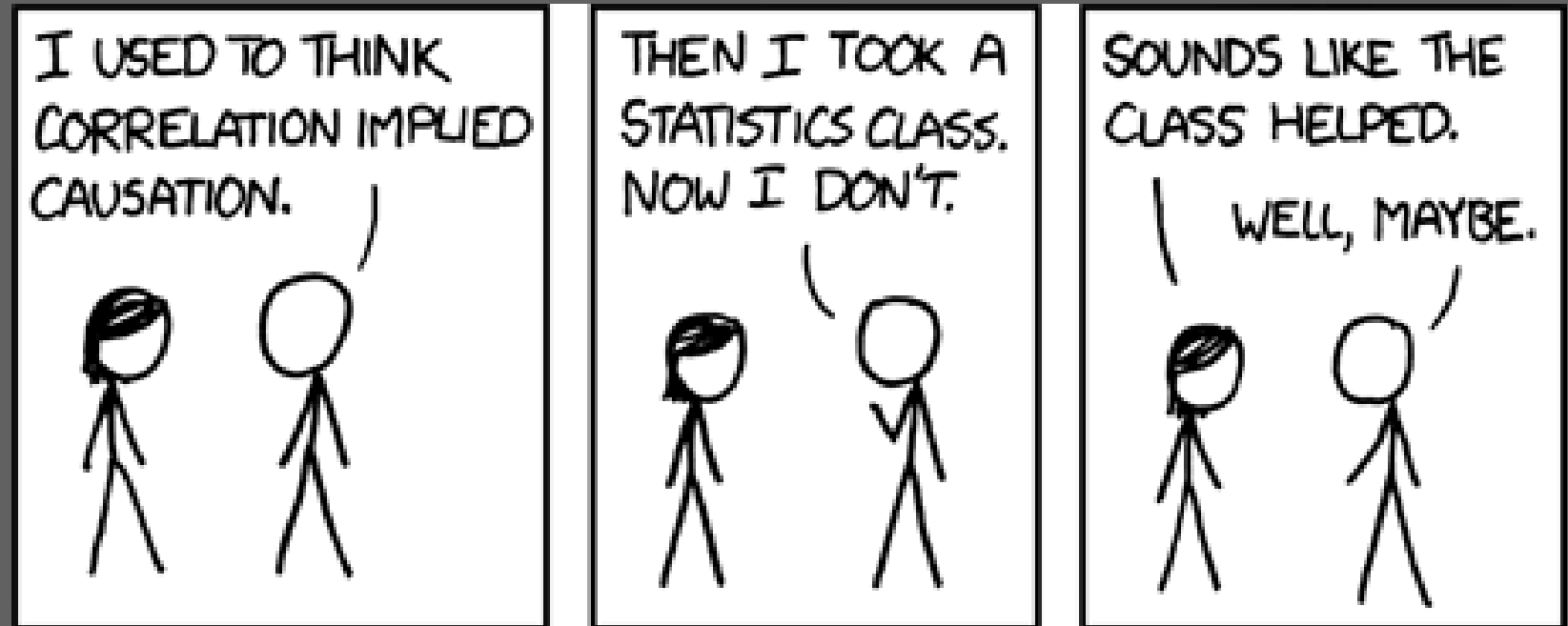


Correlation without causation (1)

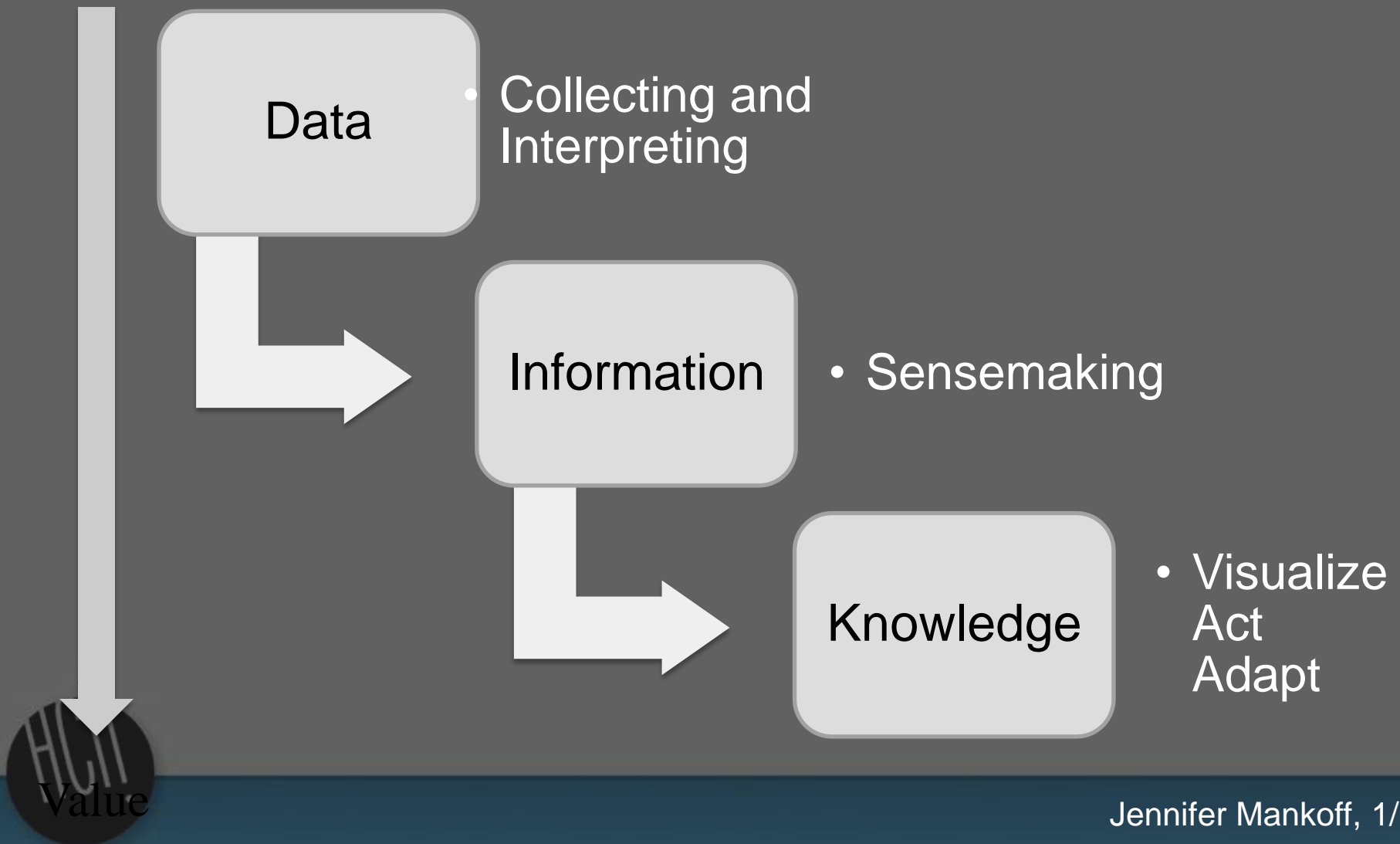
Global Temperature Vs. Number of Pirates

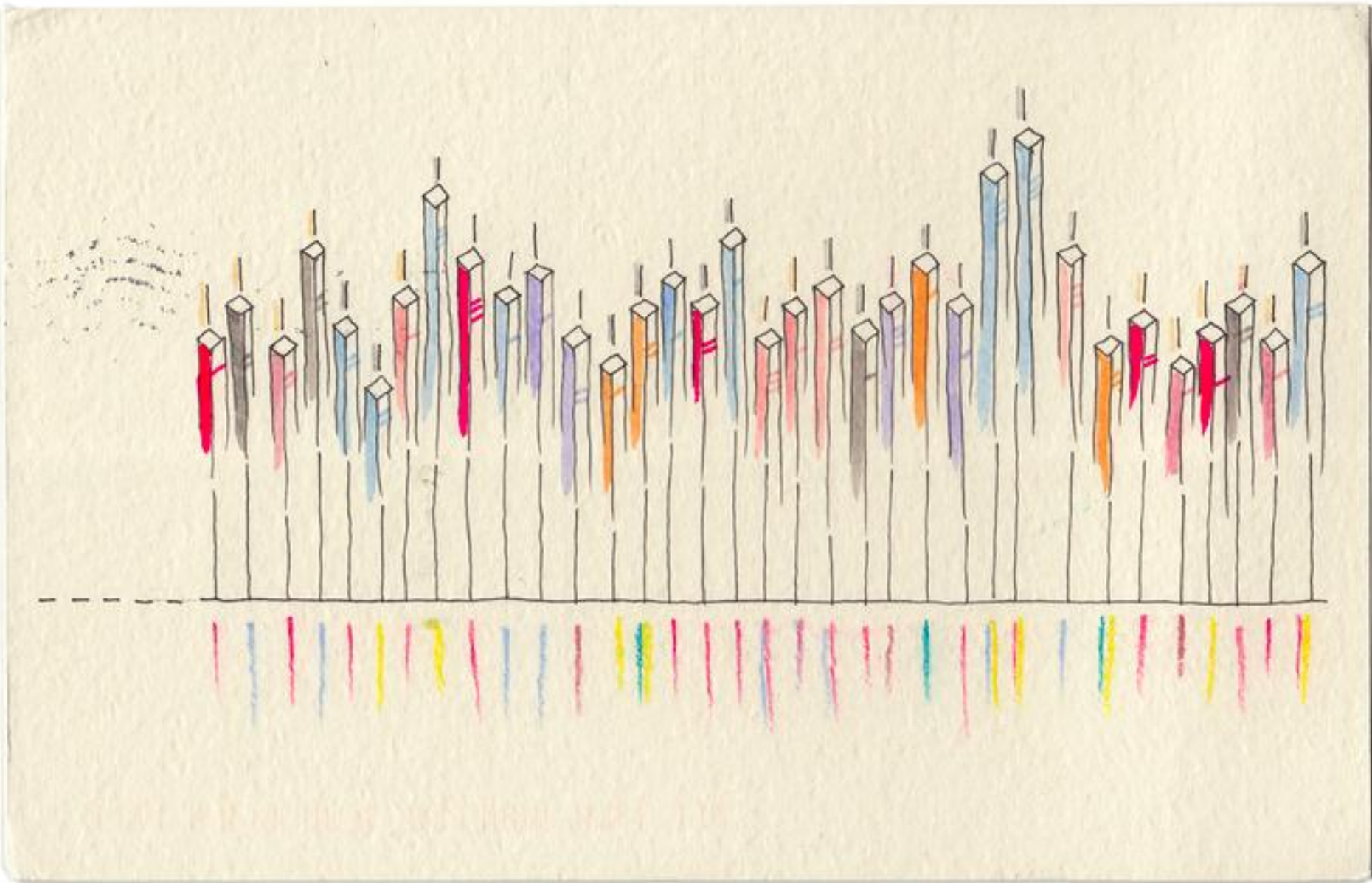


Correlation \neq Causality



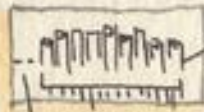
Making Data Actionable





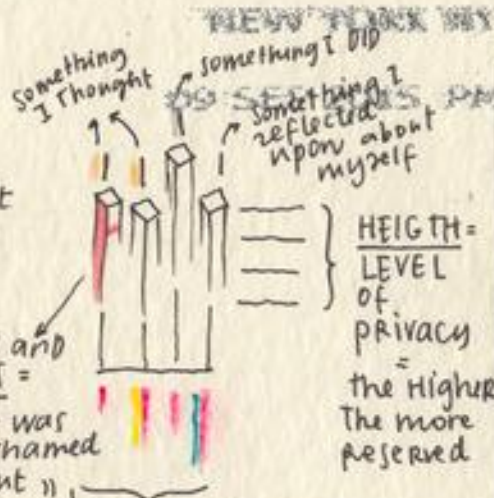
66 Dear Data WEEK 51: Privacy (please!)

HOW TO READ IT:



This week I tracked everything I DID, thought or reflected upon that I would/wanna tell, in chronological order.

--- = dashed
 lines = documenting the FIRST MASSIVE DATA VOID in Dear Data = I forgot to track the whole Monday morning ☺



COLOR and LINES = what was I "ashamed about"

Bottom lines = why did I want it to be SECRET?
 - generally ashamed
 - fear people's judgement
 - somebody would be hurt
 - I am a terrible person
 - I am scared what would have happened

- work
- my attitude
- boyfriend/partners
- thoughts about the future
- Dear Data
- the project's end
- the project
- my behavior with my boyfriend
- my body/physical things
- sensations
- my aspect
- friends/family
- my selfishness
- my attitude
- my habits
- my obsessions
- drinks
- money
- me as a person
- being grumpy
- being irritated
- being mean ☹️

from:
G. LUPI



BROOKLYN
-NY- USA

SEND TO:

STEFANIE POSAVEC

LONDON

-UK-

ENGLAND

DEAR DATA - WEEK 51

A WEEK OF PRIVACY

ABOUT THE DATA: ORIGINALLY I WAS TRACKING EVERY MOMENT I WOULDN'T WANT TO SHARE WITH YOU (OR ANYONE ELSE) BUT MY PHONE DIED + I LOST MY DATA. SO, THIS IS A LIST OF MOMENTS THAT I WOULD PREFER TO KEEP PRIVATE FROM LAST WEEK, BUT MADE FROM MEMORY.

HOW TO READ IT: EACH ^{'CONSIDERED'} SYMBOL IS ONE MOMENT FROM THE WEEK THAT I WOULD PREFER TO KEEP PRIVATE.



SYMBOLS ARE ORDERED BY HOW EMBARRASSED / UNCOMFORTABLE I WOULD FEEL IF YOU KNEW WHAT THESE MOMENTS WERE!

TYPE OF MOMENT I WANTED TO KEEP PRIVATE +

- A THOUGHT I HAD
- ONE OF MY ACTIONS
- MY INTERACTIONS WITH OTHERS
- MY BEHAVIOUR

① = I WOULDN'T BE TOO EMBARRASSED / UNCOMFORTABLE
⑤ = I WOULD BE VERY EMBARRASSED / UNCOMFORTABLE

IF I HAD TO SHARE THIS MOMENT, WHO I WOULD BE WILLING TO SHARE IT WITH:

- MORE PRIVATE: NO ONE BUT ME.
- MY HUSBAND
- MY FRIENDS (+ ALL ABOVE)
- MY PARENTS (+ ALL ABOVE)
- LESS PRIVATE

FROM: S. POSAVEC



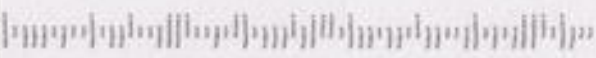
Royal Mail
Mount Pleasant
Mail Centre
02-09-2015
34102357

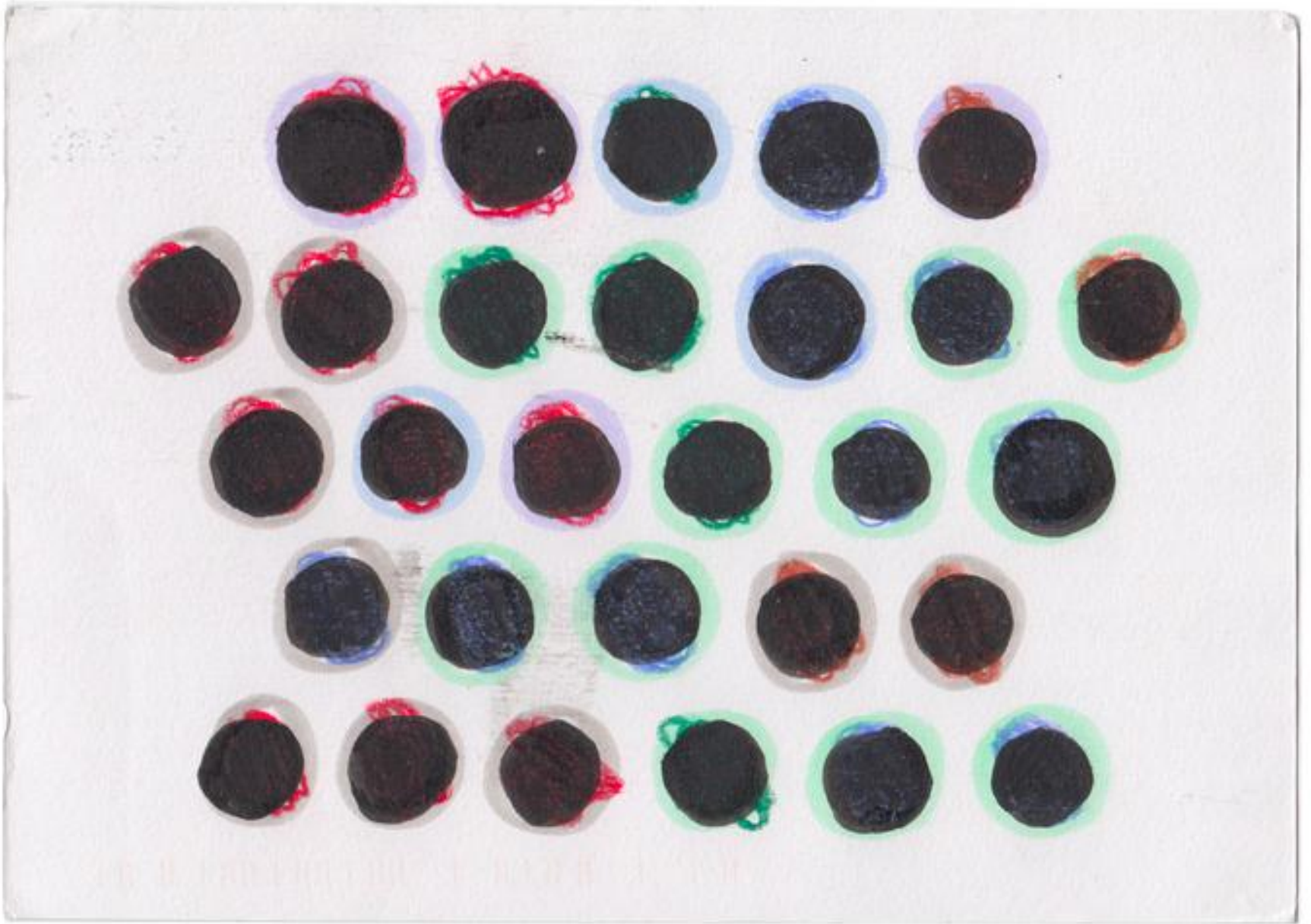


TO: GIORGIA LUPI

BROOKLYN, NY 11249
USA

BY AIR MAIL
par avion
Royal Mail®





What is Information Visualization?

Visualize: to form a mental image or vision of

...

Visualize: to imagine or remember as if actually seeing.

American Heritage dictionary, Concise Oxford dictionary



The Power of Visualization

1. Start out going Southwest on ELLSWORTH AVE Towards BROADWAY by turning right.
- 2: Turn RIGHT onto BROADWAY.
3. Turn RIGHT onto QUINCY ST.
4. Turn LEFT onto CAMBRIDGE ST.
5. Turn SLIGHT RIGHT onto MASSACHUSETTS AVE
6. Turn RIGHT onto RUSSELL ST.

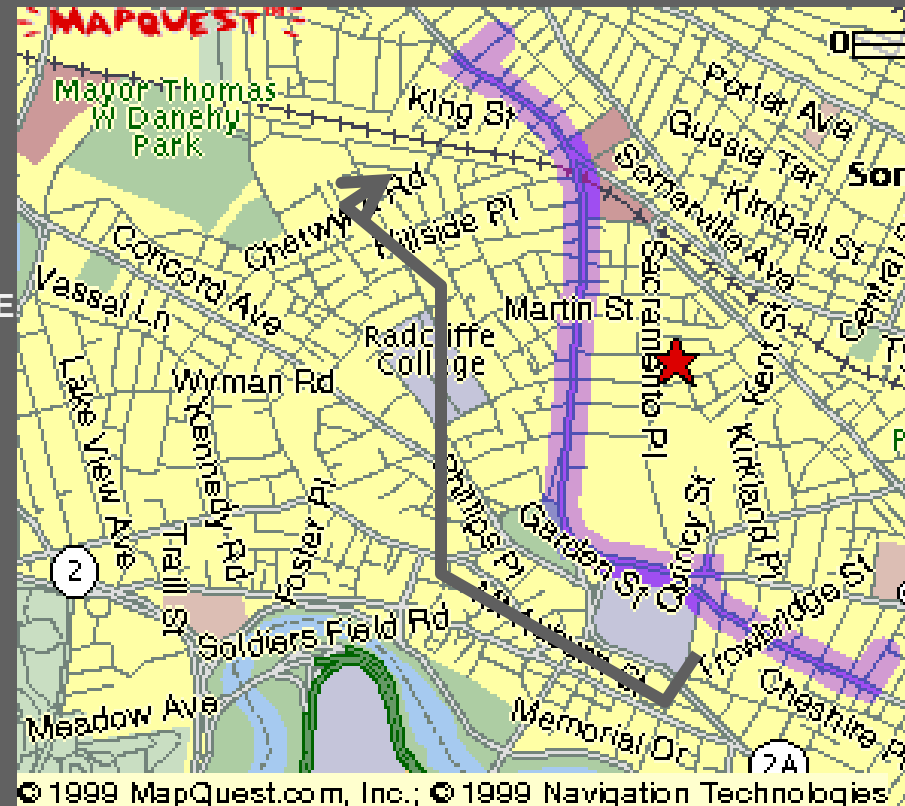
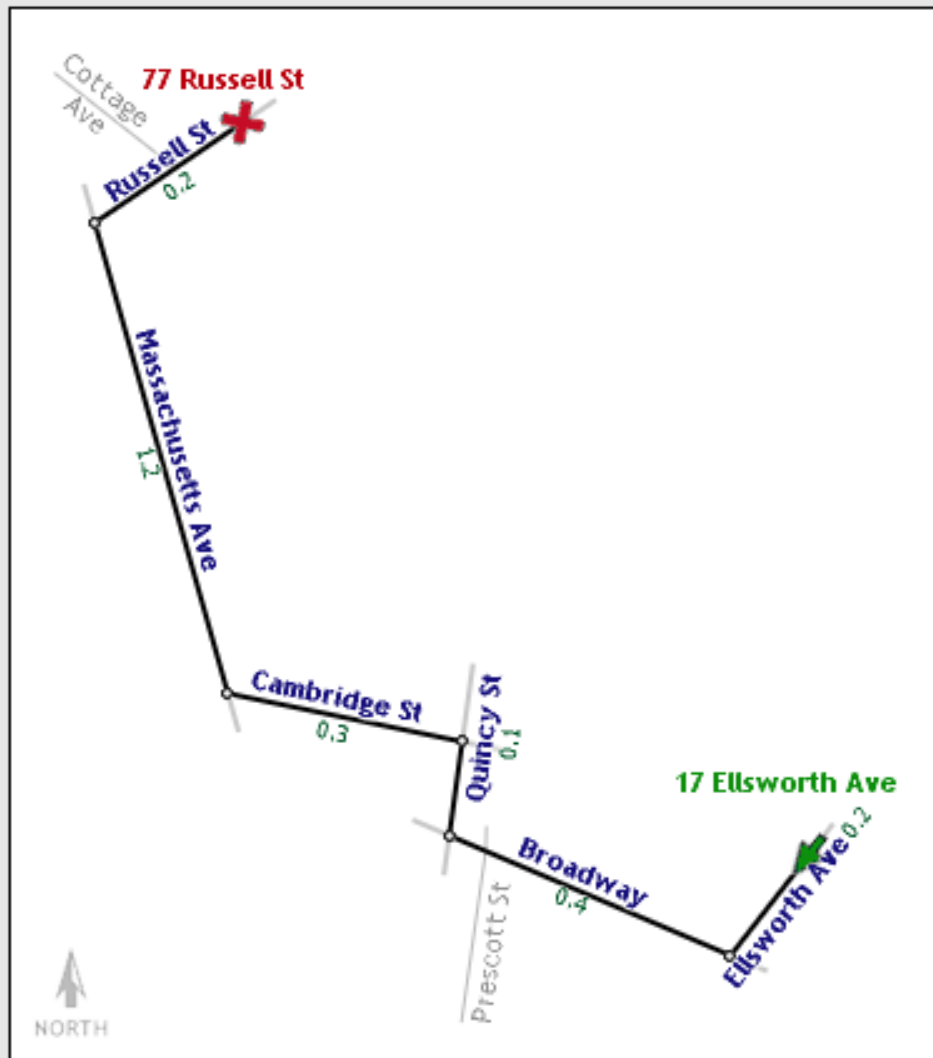


Image
from
mapquest.
com

The estimated travel time is 5 minutes for 2.16 miles of travel, total of 6 steps.

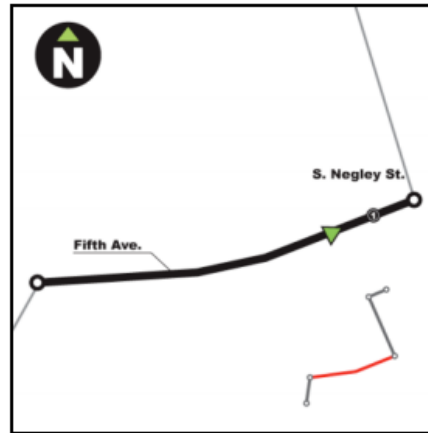
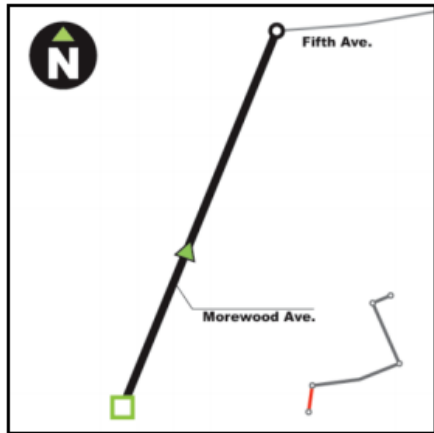


Directions		Elapsed Distance
1	Begin at 17 Ellsworth Ave on Ellsworth Ave and go Southwest for 500 feet	0.1
2	Turn right on Broadway and go Northwest for 0.4 miles	0.5
3	Turn right on Quincy St and go North for 200 feet	0.5
4	Turn left on Cambridge St and go West for 0.3 miles	0.8
5	Bear right on Massachusetts Ave, Mass Ave, RT-2A and go North for 1.2 miles	2.0
6	Turn right on Russell St and go Northeast for 1000 feet to 77 Russell St	2.2

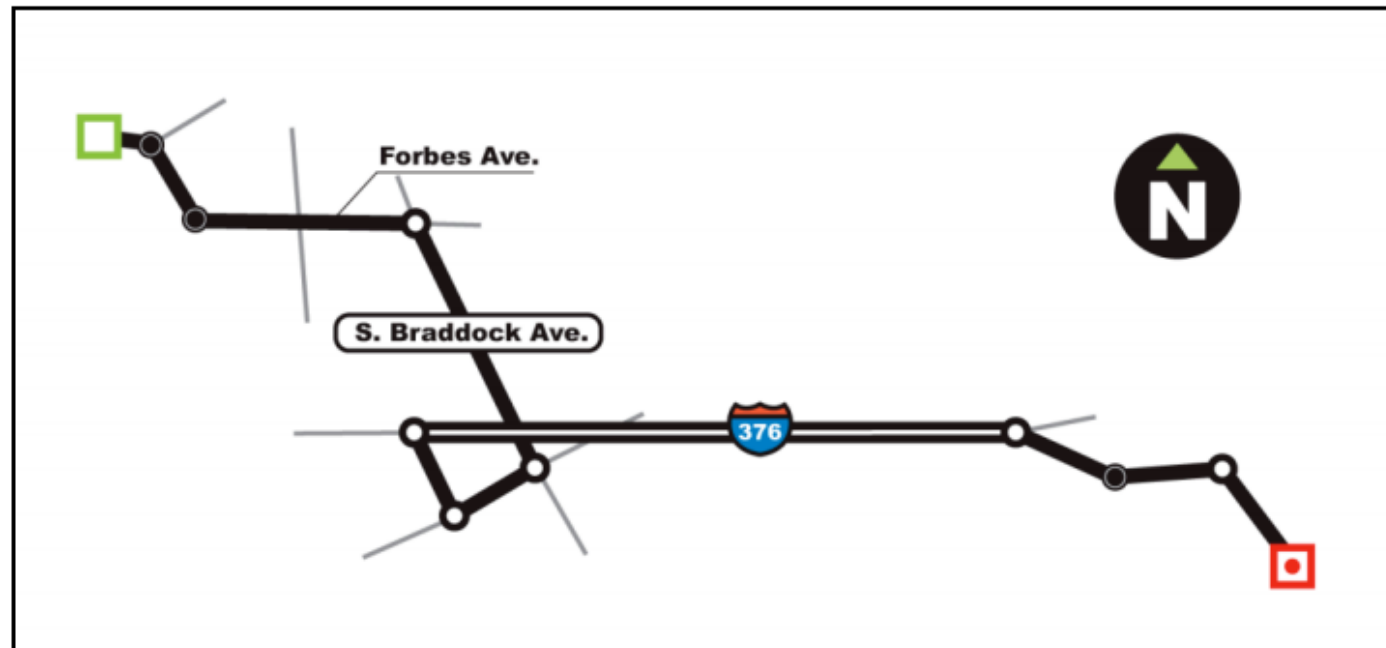


Line drive tool by Maneesh Agrawala <http://graphics.stanford.edu>

The Power of Design *and* interaction in Visualization



Lee, Forlizzi & H



Planning a Visualization

1. What is its goal?
2. What visual queries does it support?
3. What are some compelling, useful examples? [COPY COPY COPY!]
4. Could it have been done more simply?



Making Queries

Define the query[ies] you wish to support

“The special skill of designers ... [is] the talent to analyze a design in terms of its ability to support the visual queries of others...”

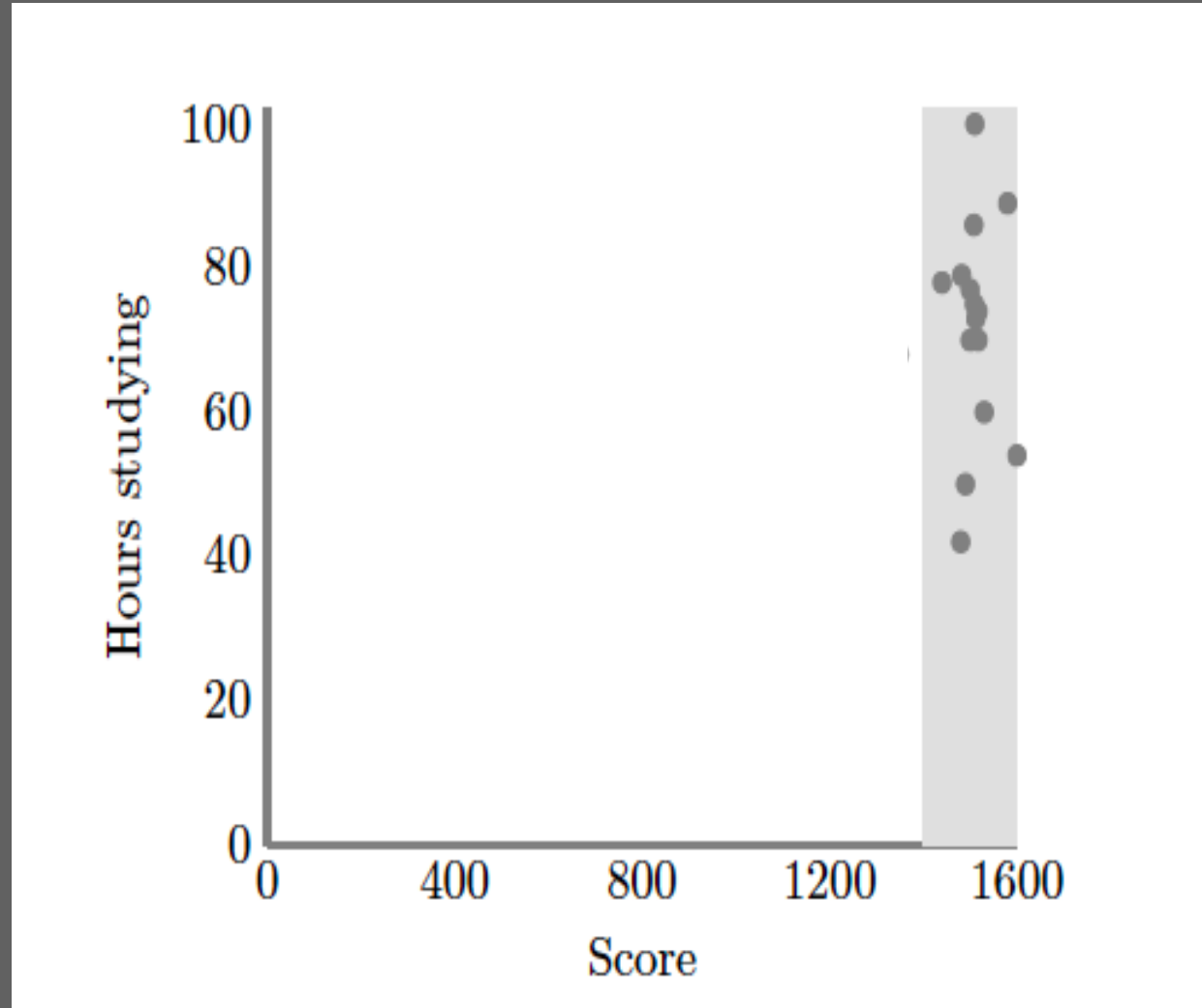
Patterns ⇔ visual system

Cognitive process prediction

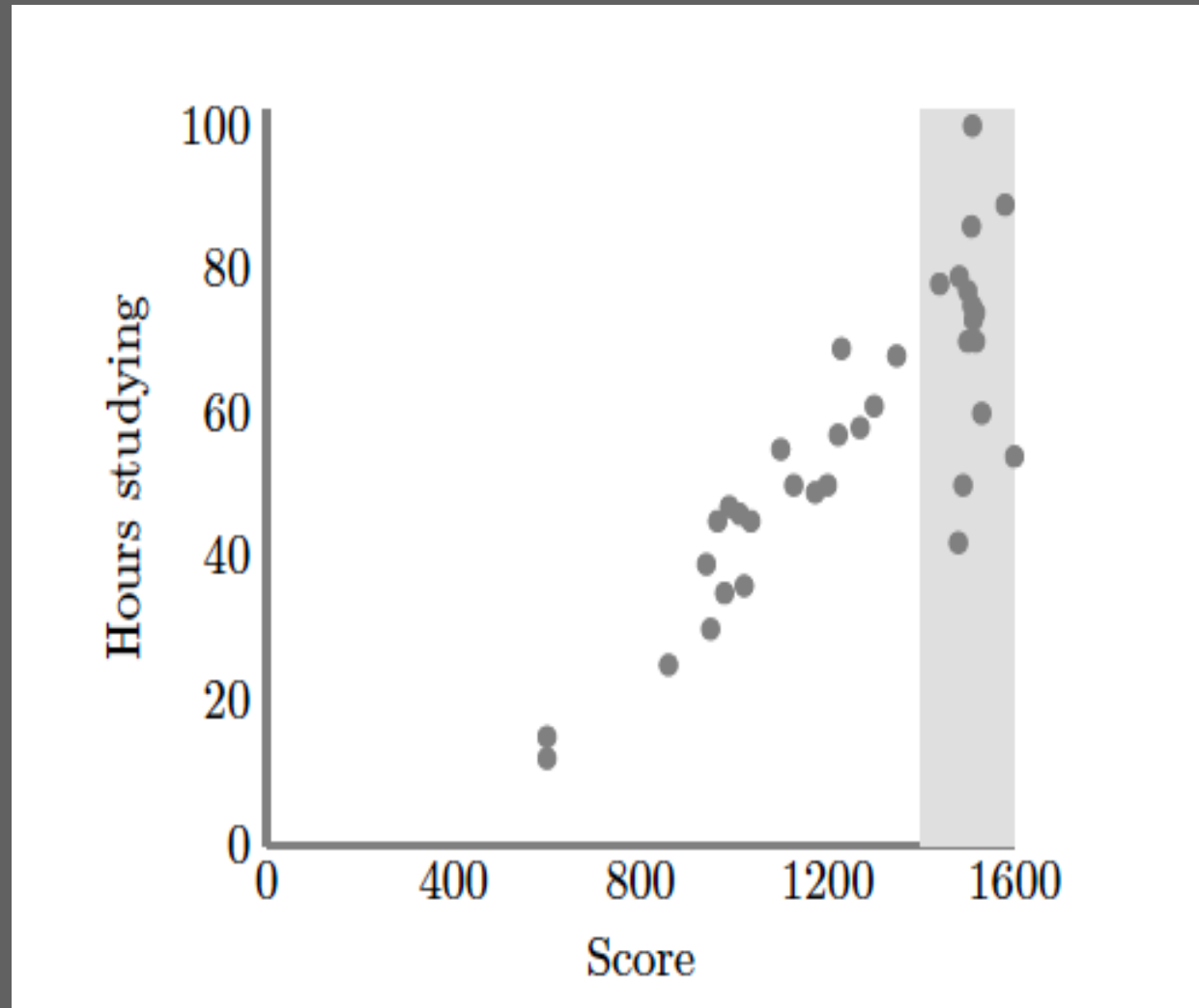
A continual fresh eye



Visualizations frequently casue us to draw conclusions



Which is why visualization choices are so important...



Planning a Visualization

1. What is its goal?
2. What visual queries does it support?
3. What are some compelling, useful examples? [COPY COPY COPY!]
4. Could it have been done more simply?



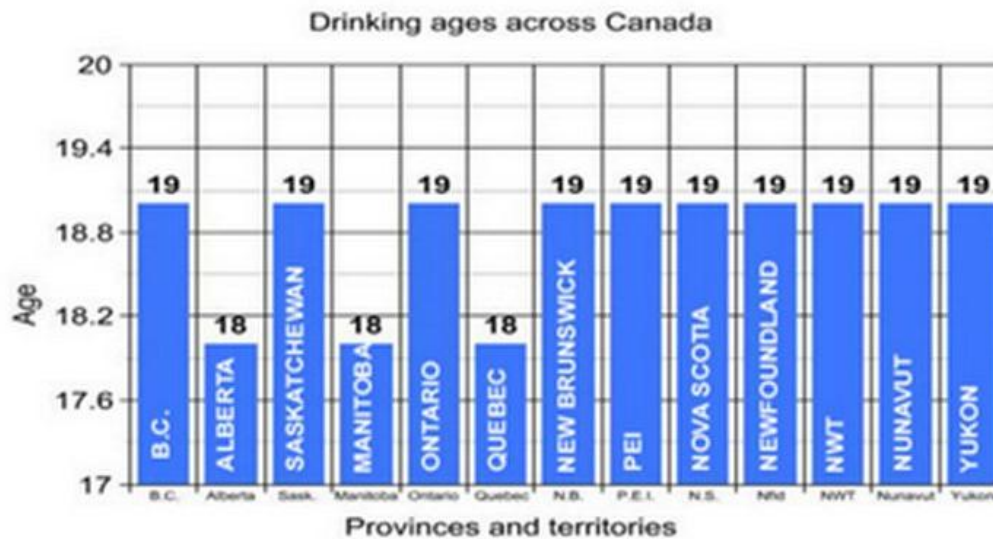
Area vs Size



Confusing Axes

Drinking age will remain 19 in Saskatchewan

CBC News Posted: Mar 4, 2013 11:59 AM CST | Last Updated: Mar 4, 2013 11:55 AM CST 25



Canadian Centre on Substance Abuse

You have to be 19 in Saskatchewan to have a drink, while in Alberta and Manitoba, the drinking age is 18. (CBC)

The Saskatchewan Party government has ruled out lowering the drinking age, four months after party members put the issue in the public eye.

Facebook

Twitter

Share

Email

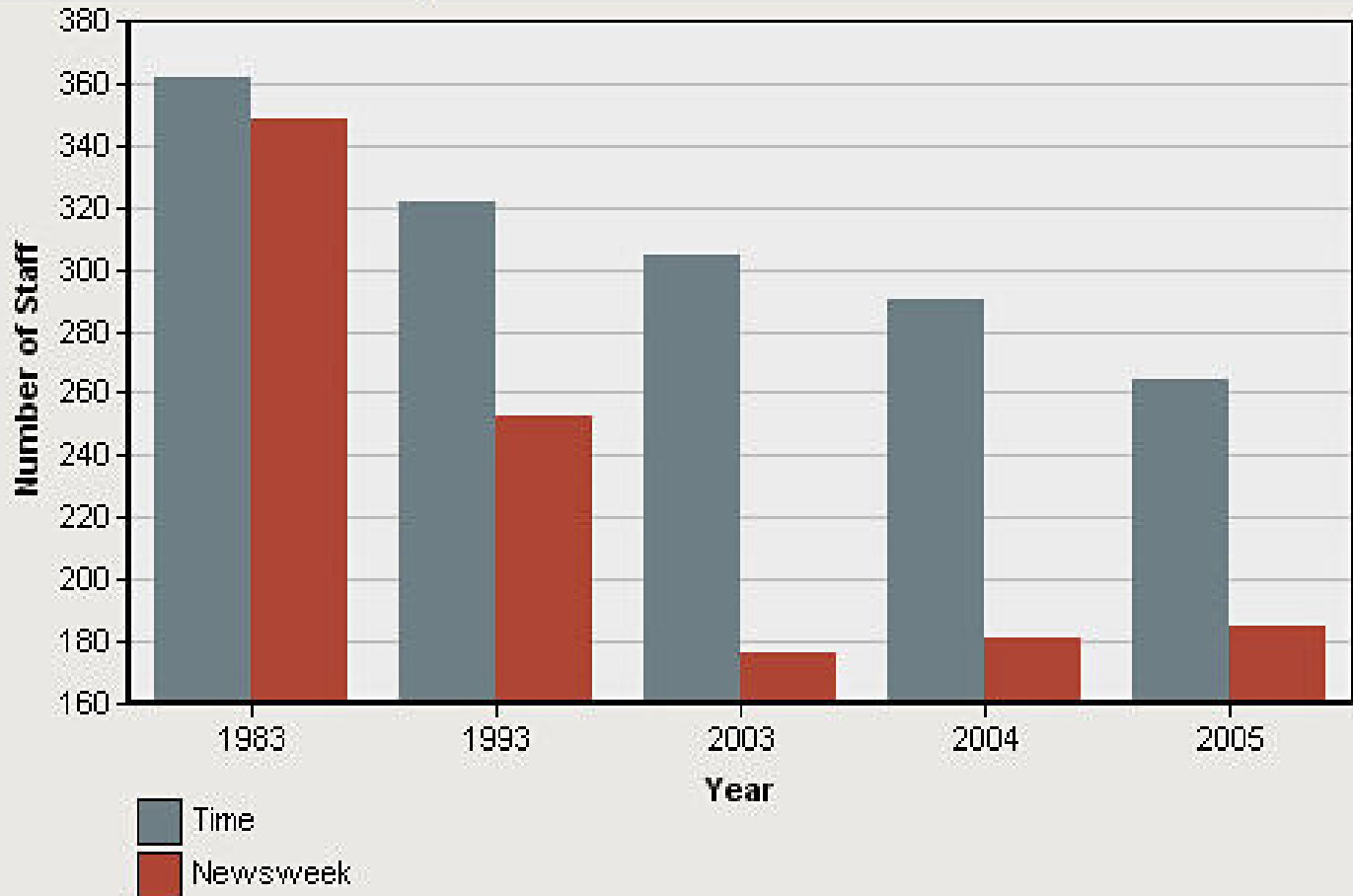
Stay Co

Mobile Fac

Misleading Axis

NEWS MAGAZINE STAFF SIZE OVER TIME

Time and Newsweek select years 1983 - 2005



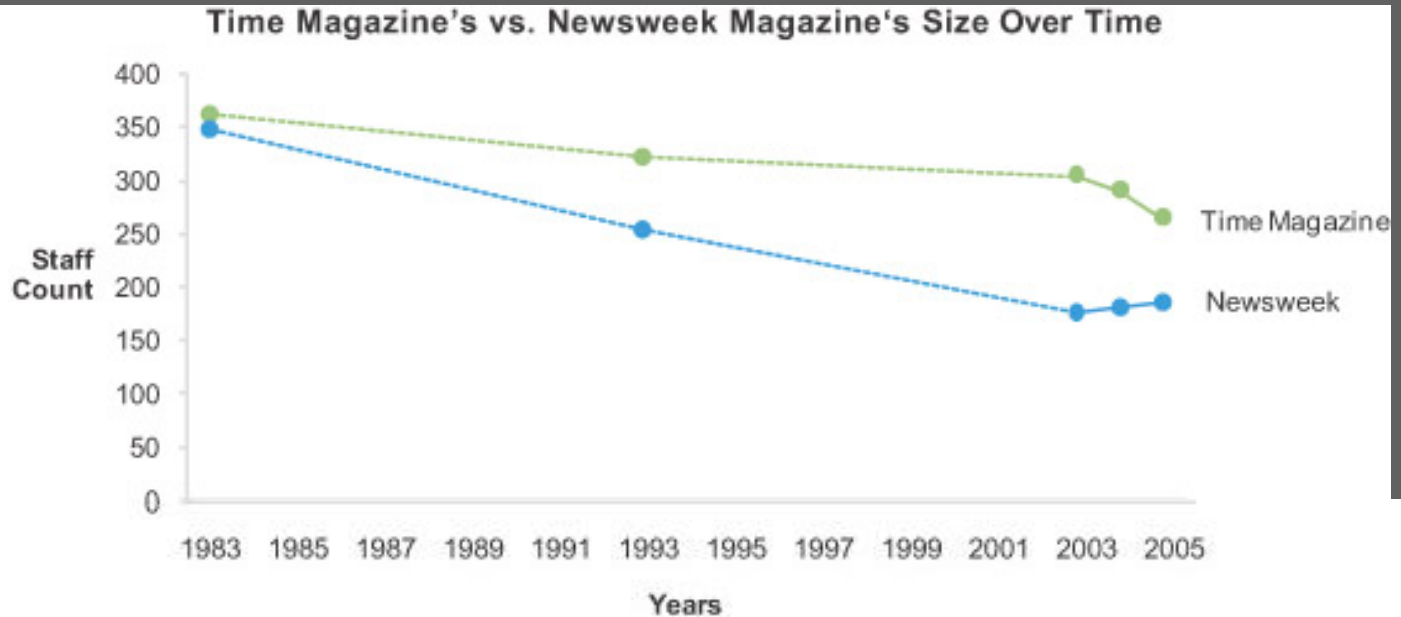
60

11/

Jennifer Mankoff, 6/12

16/

Improved

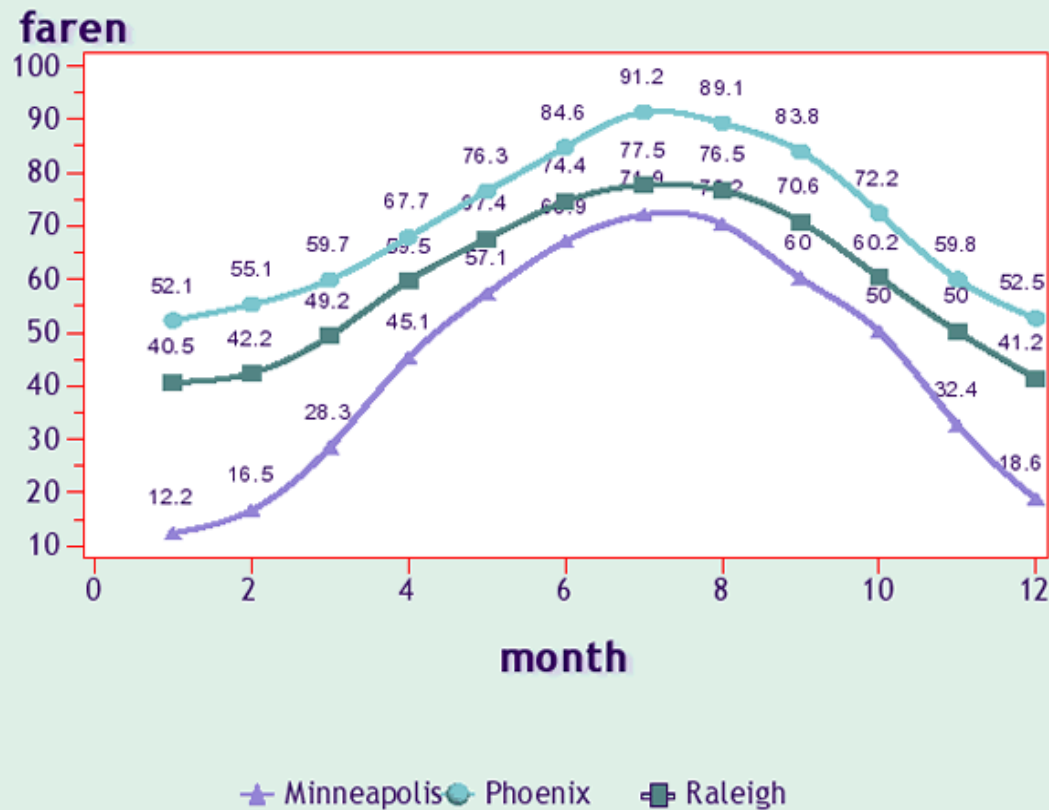


Note: A dashed line connecting two points indicates that there are years between the points for which values were not available. If the values were available, the shape of the lines might vary significantly.

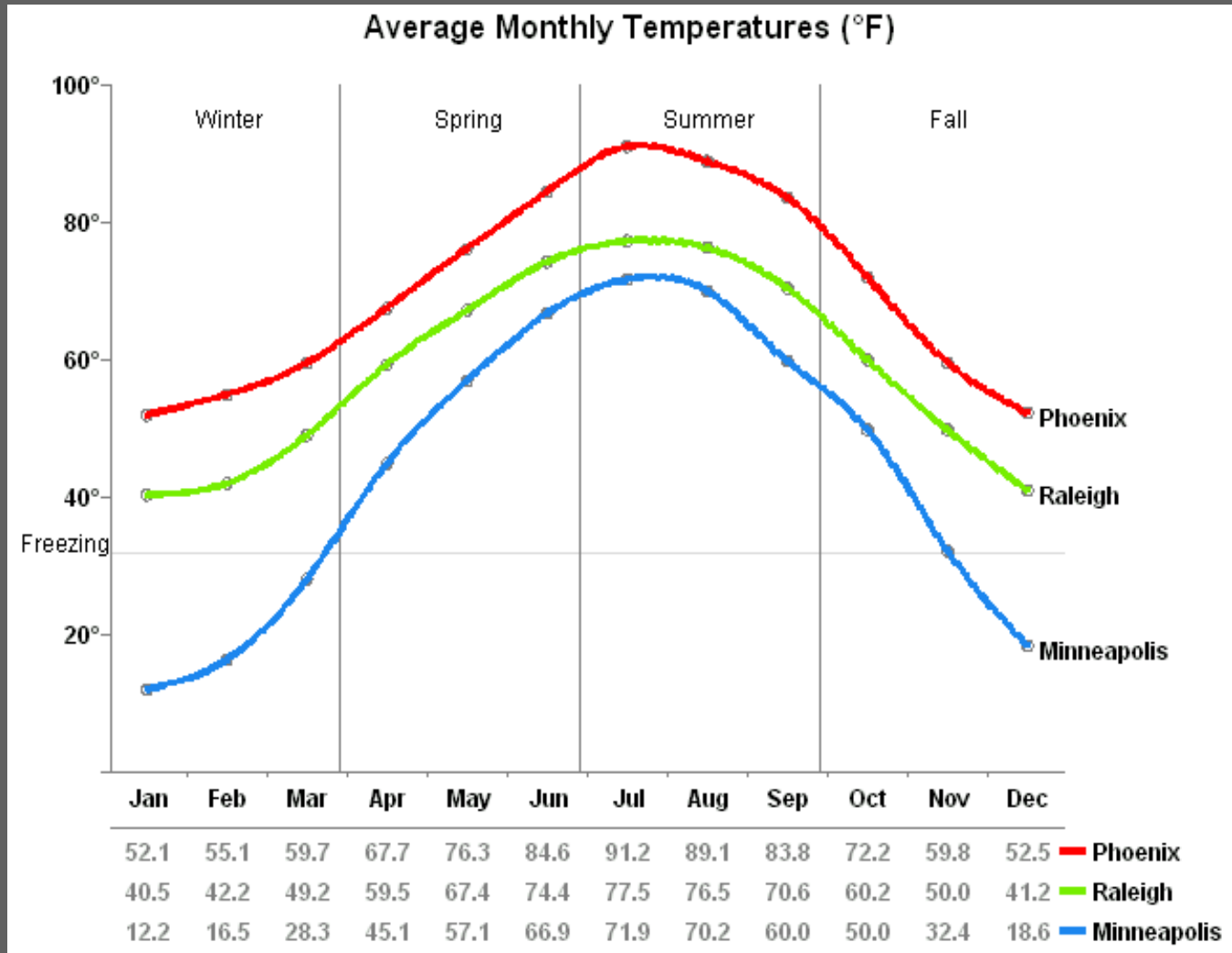


Overly Complex

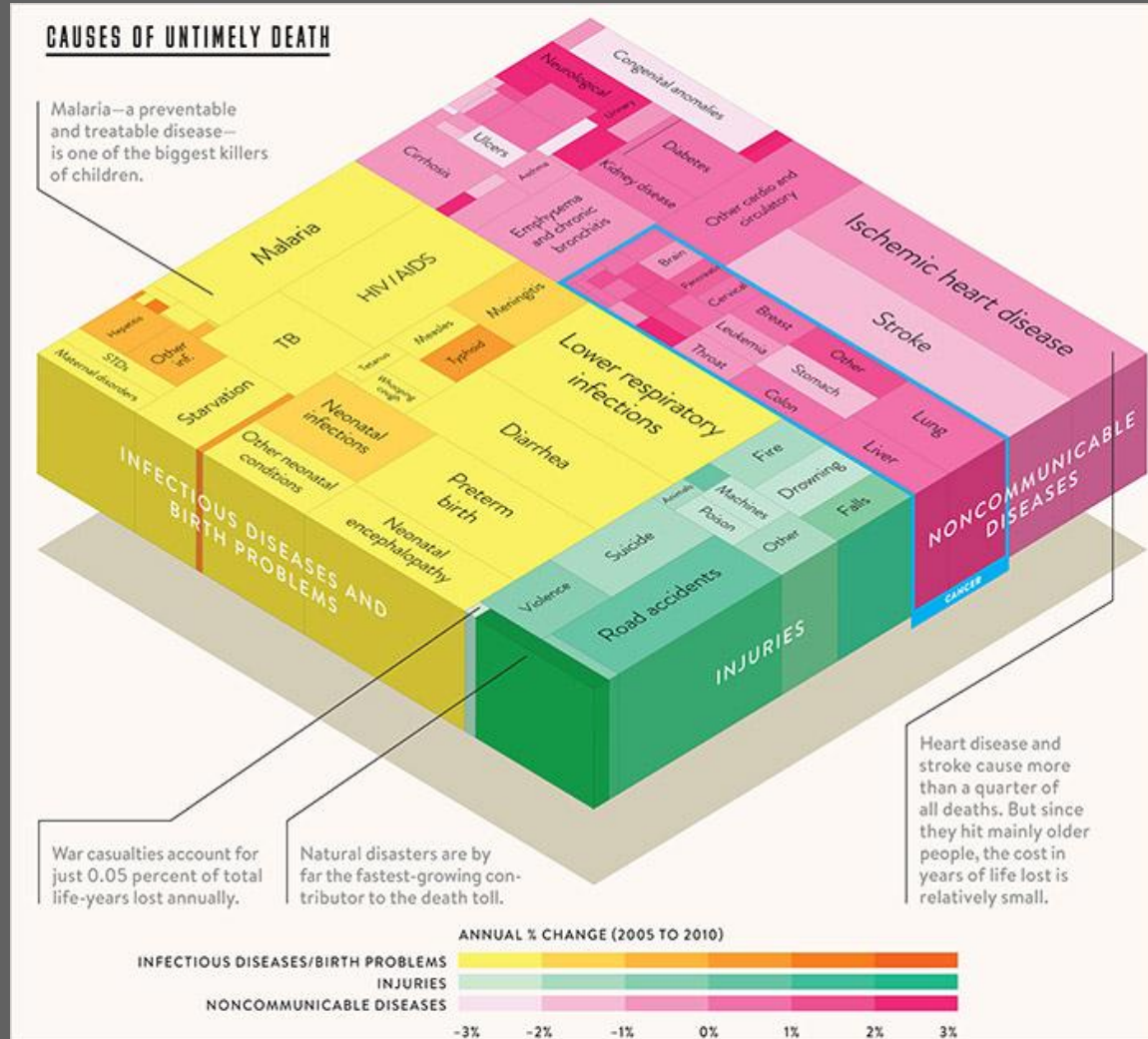
Average Monthly Temperature



Improved



Hard to read

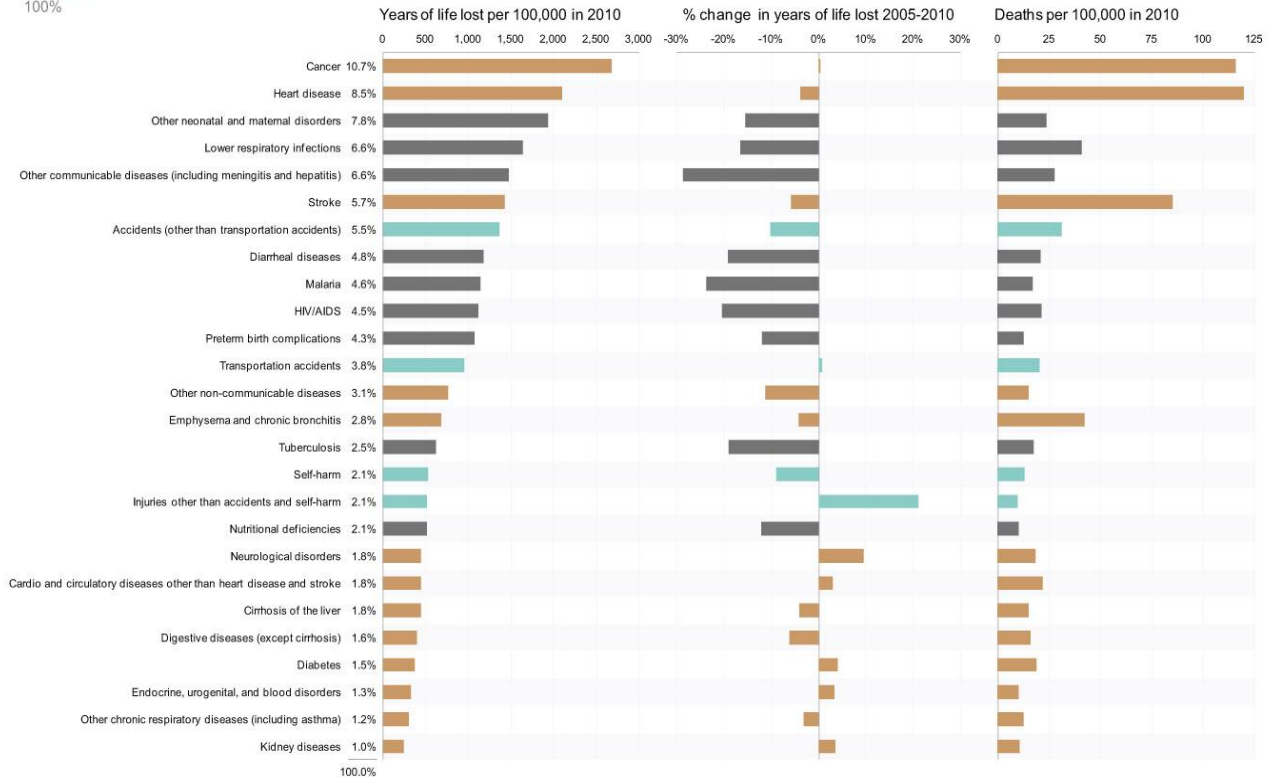


Improved

Global Causes of Lost Life

100% ■ Communicable, maternal, neonatal, and nutritional disorders
 43% ■ Non-communicable diseases
 13% ■ Injuries

Comparing the number of deaths alone, as shown in the rightmost graph below, doesn't tell the entire story. Some causes of death have a greater effect on the young, which can be seen when comparing years of life lost in the leftmost graph.



Some causes of death contribute disproportionately to years of life lost because of their effect on the young. For example, malaria, while not huge in the number of deaths, is much more significant in the number of years that are lost.

Two interesting changes reside in "Injuries other than accidents and self-harm." War, which accounted for only 0.05% of years of life lost, decreased since 2005 by 31.5% in years of life lost per 100,000 people. Natural disasters, which accounted for 0.65% of years of life lost, increased by 217% in years of life lost per 100,000.

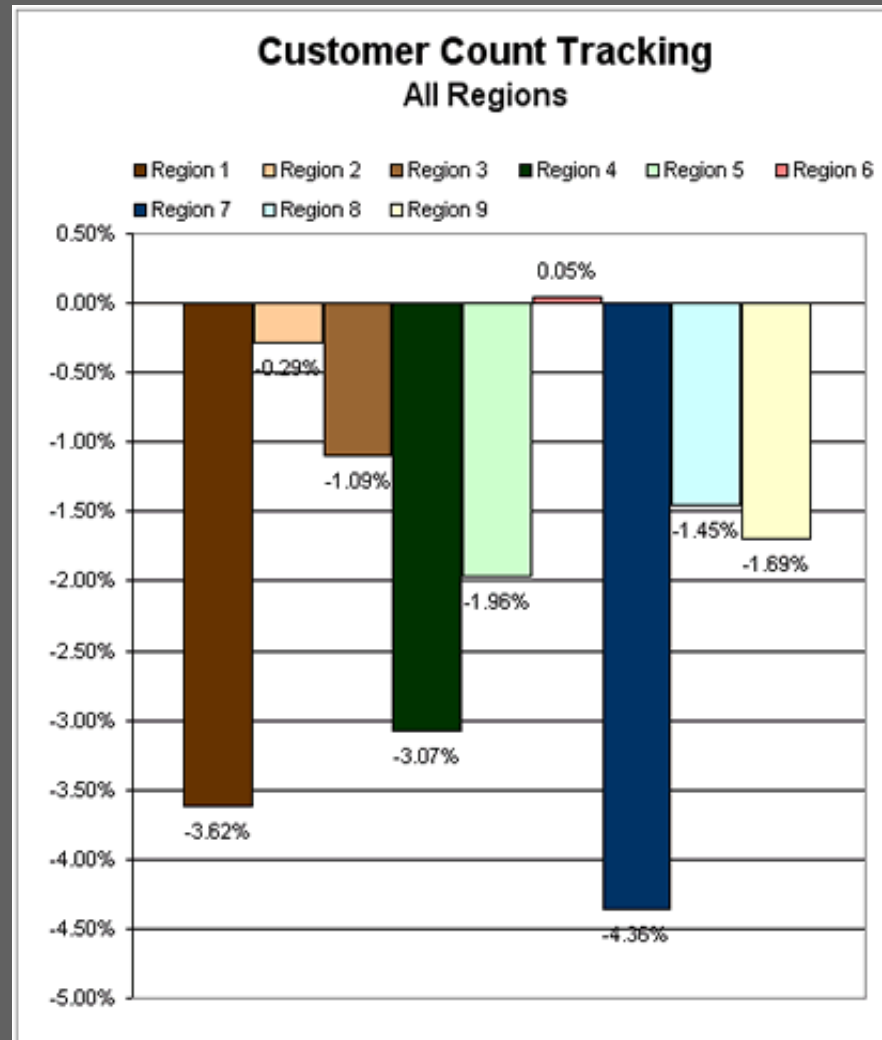
Communicable, maternal, neonatal, and nutritional disorders (the gray bars) are often easier to prevent through healthcare than other causes of death. This reveals itself in the graph above by the fact that all of these disorders have decreased during this five year period.

The five forms of cancer that cause the most deaths are trachea/bronchus/lung (2.9%), stomach (1.4%), liver (1.4%), colon/rectum (1.4%), and breast (0.8%).

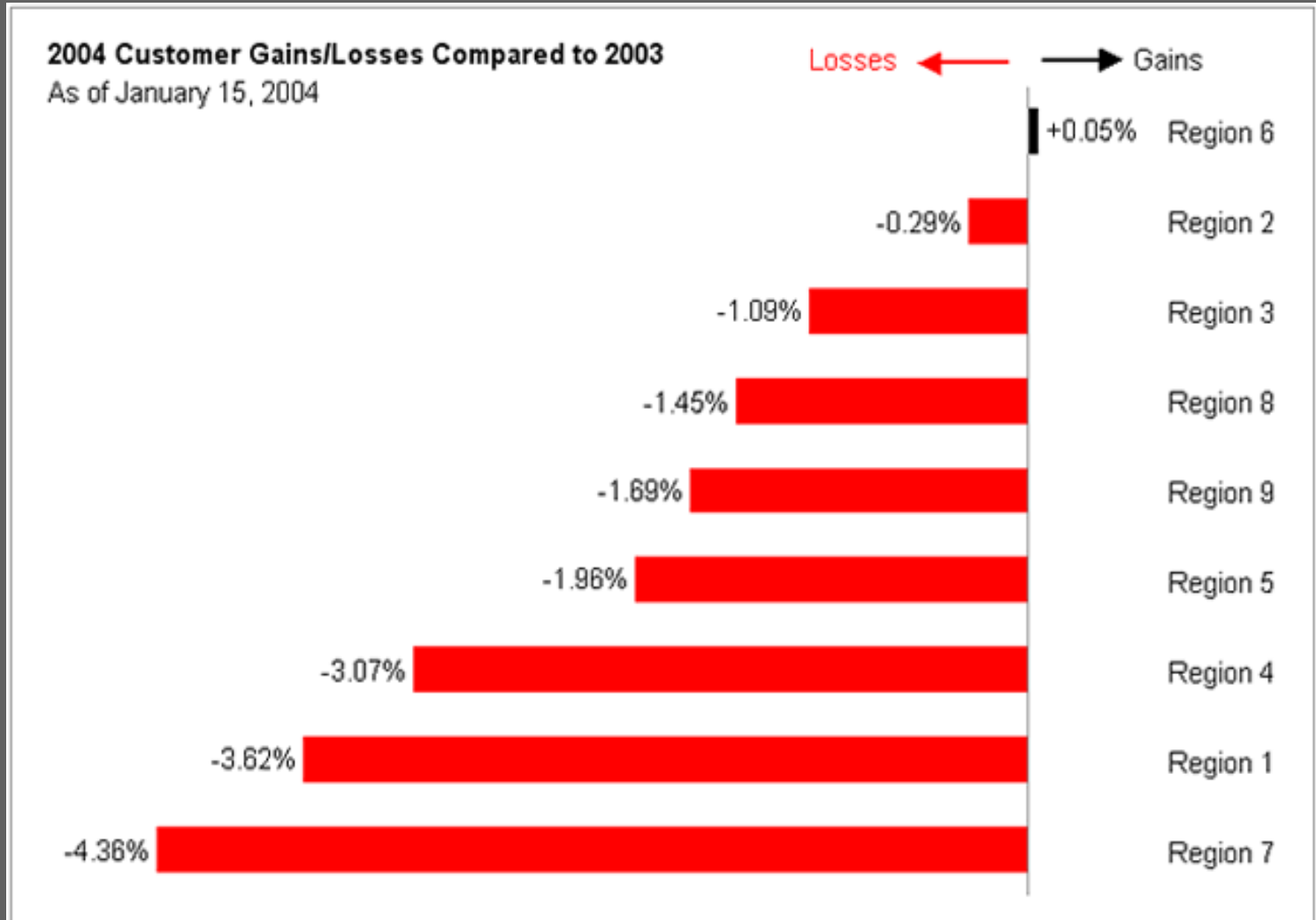
All cardiovascular and circulatory diseases combined account for 30% of deaths.



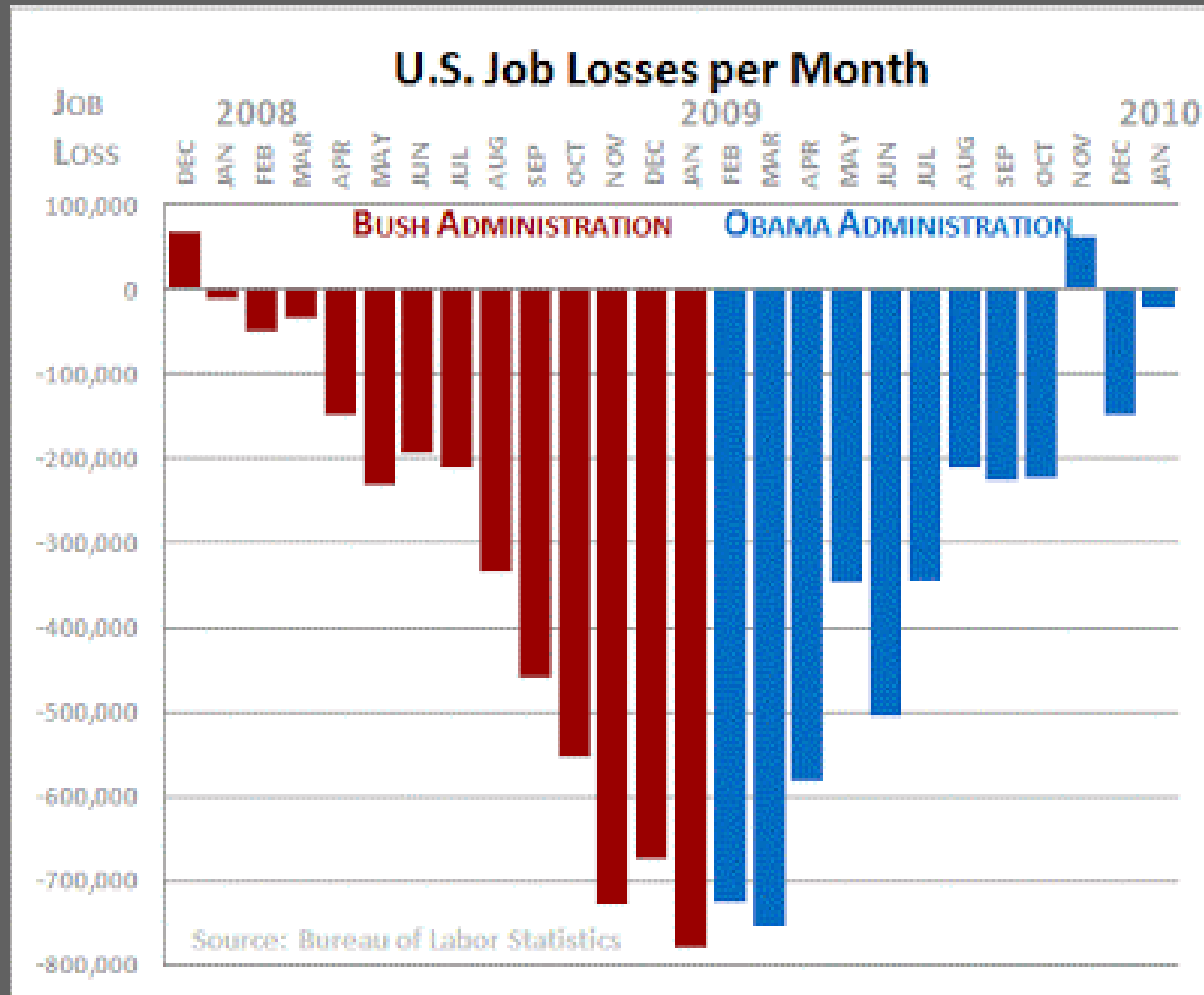
Comparisons Difficult



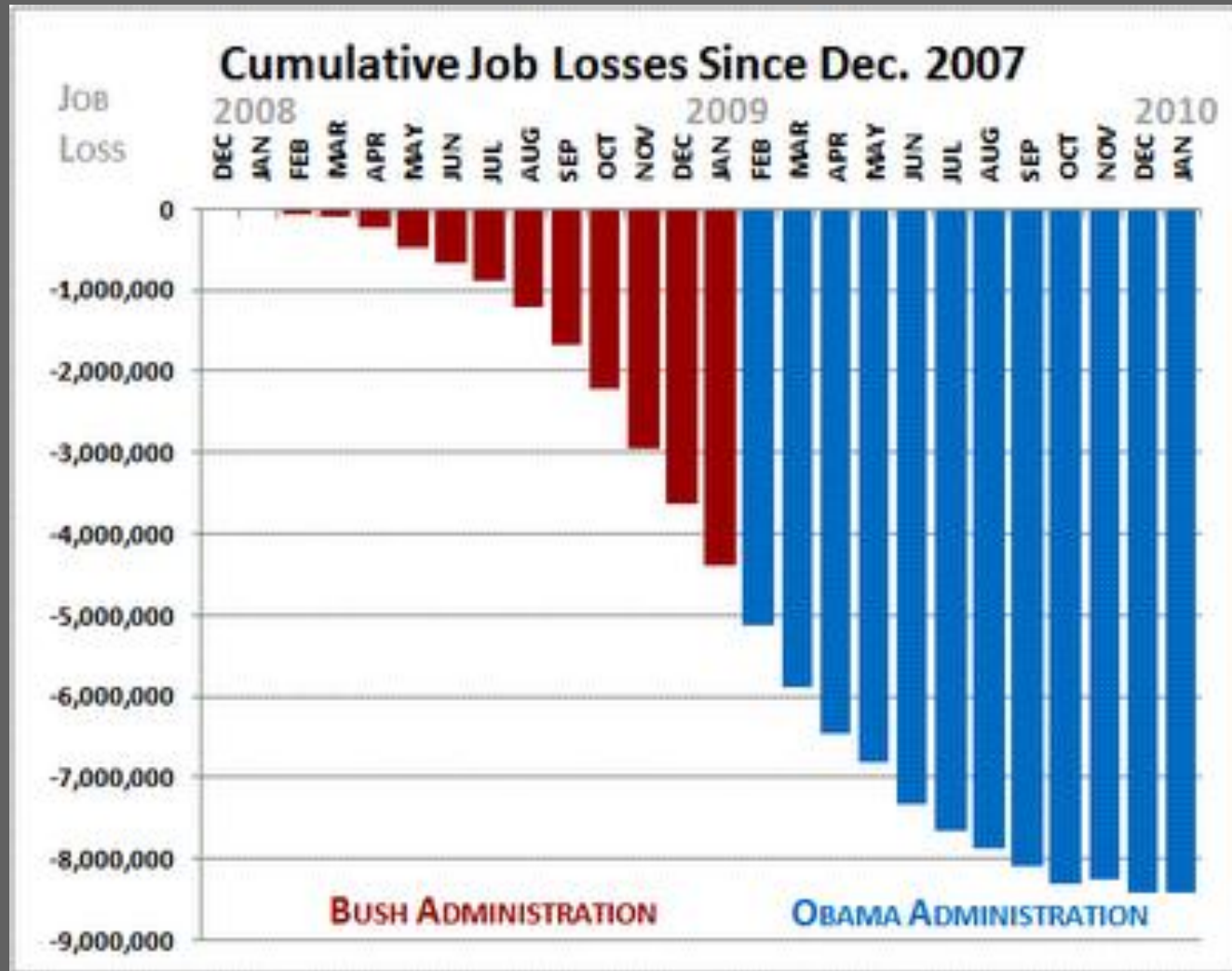
Improved



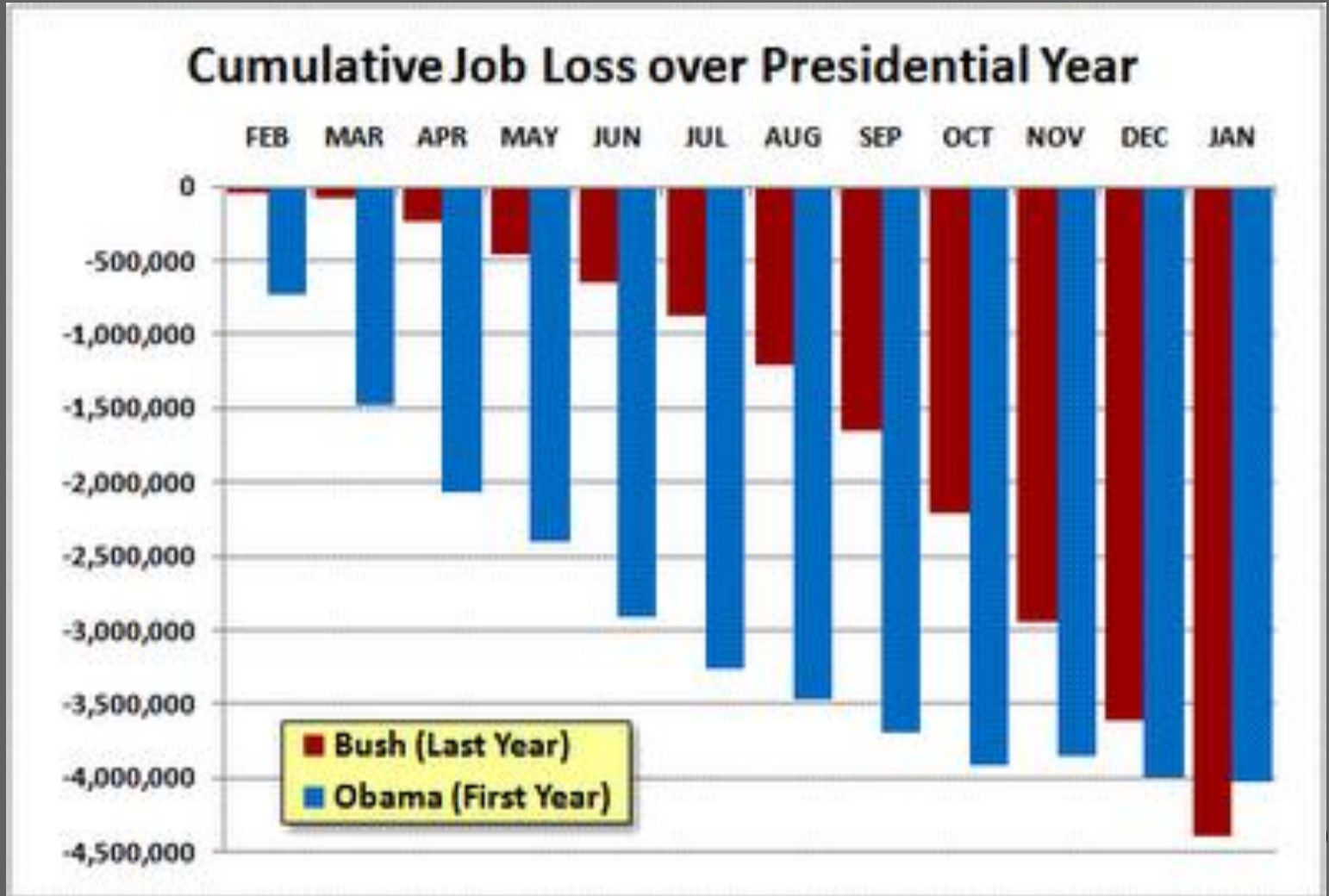
Highlighting wrong aspect of data?



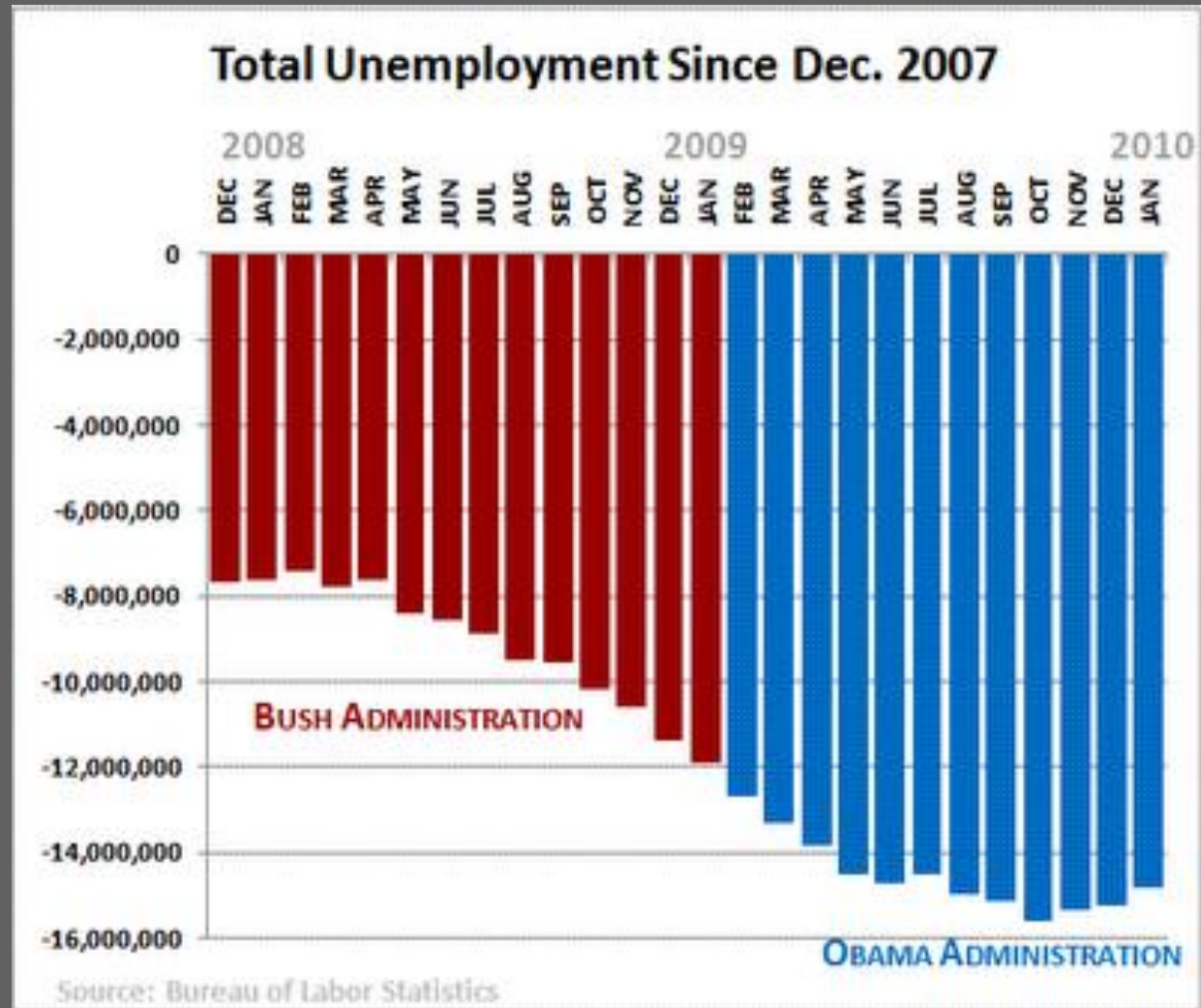
An alternative view



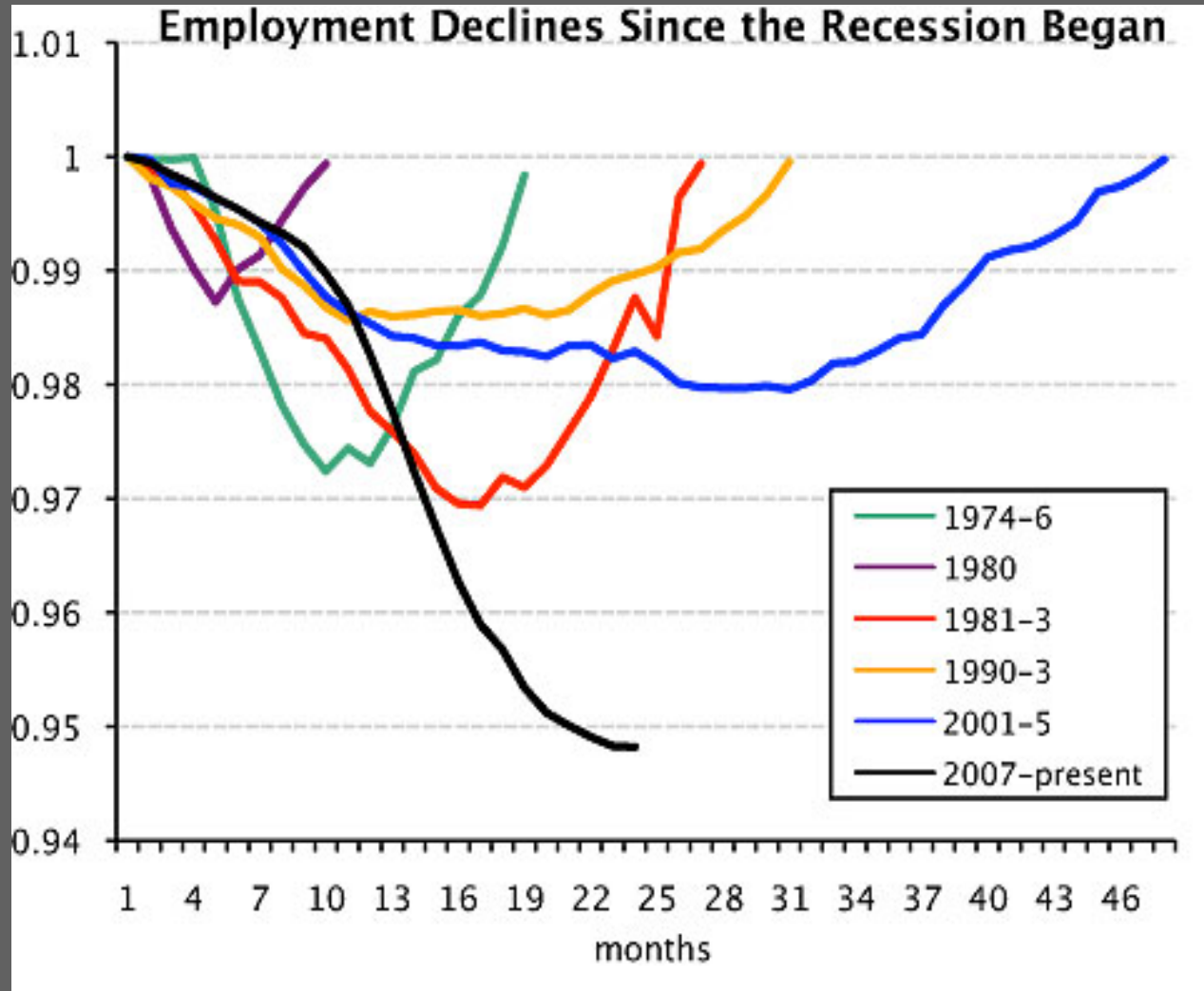
An alternative view



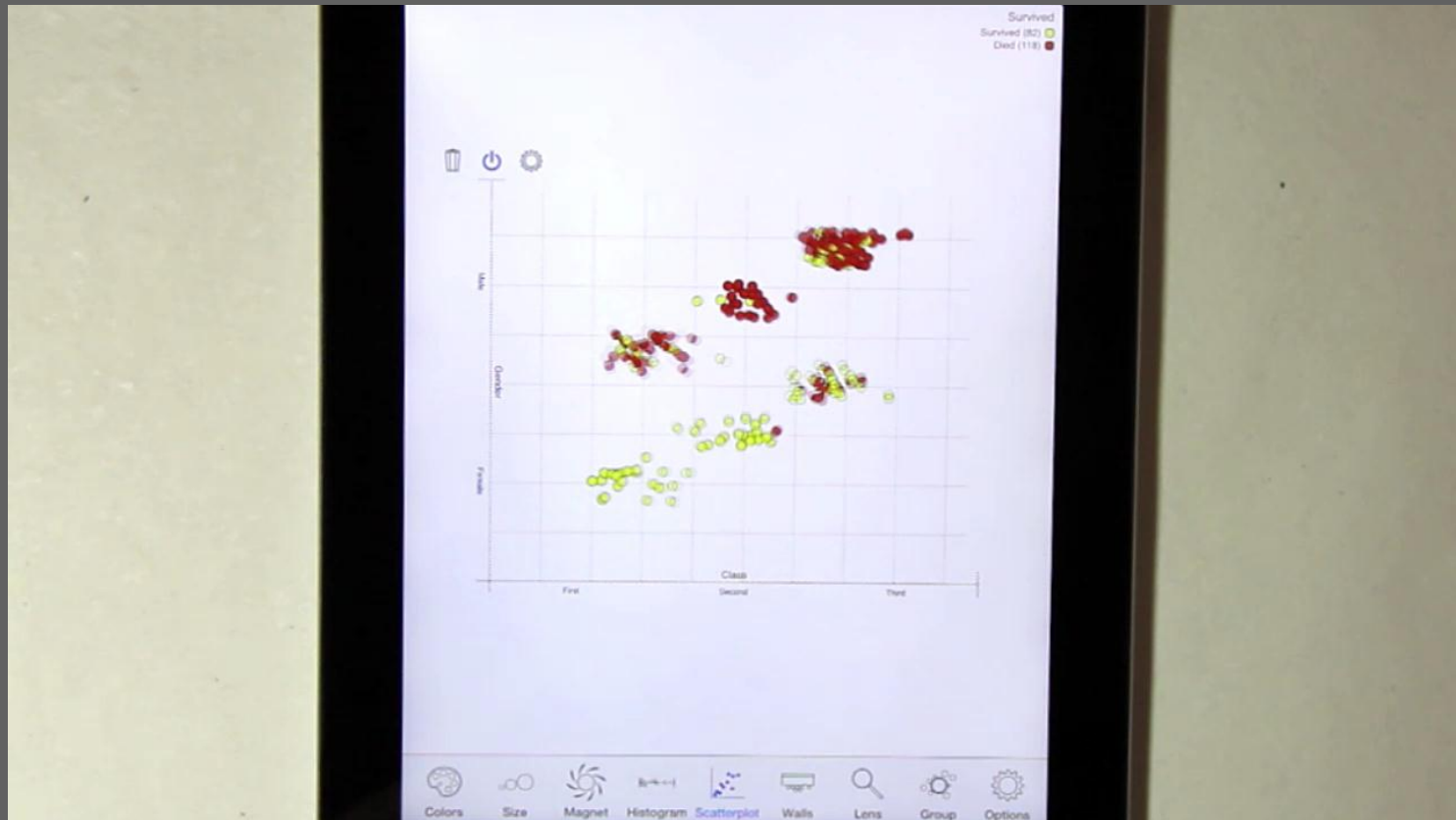
An alternative view



An alternative view



Interactive Visualization



Jeffrey M. Rzeszotarski and Aniket Kittur. 2014. Kinetica: naturalistic multi-touch data visualization. In Proceedings of the 32nd annual ACM conference on Human factors in

73

11/

Jennifer Mankoff, 1/12

16/

Interactive Visualization

<http://queue.acm.org/detail.cfm?id=2146416>



Visualizing Big Data

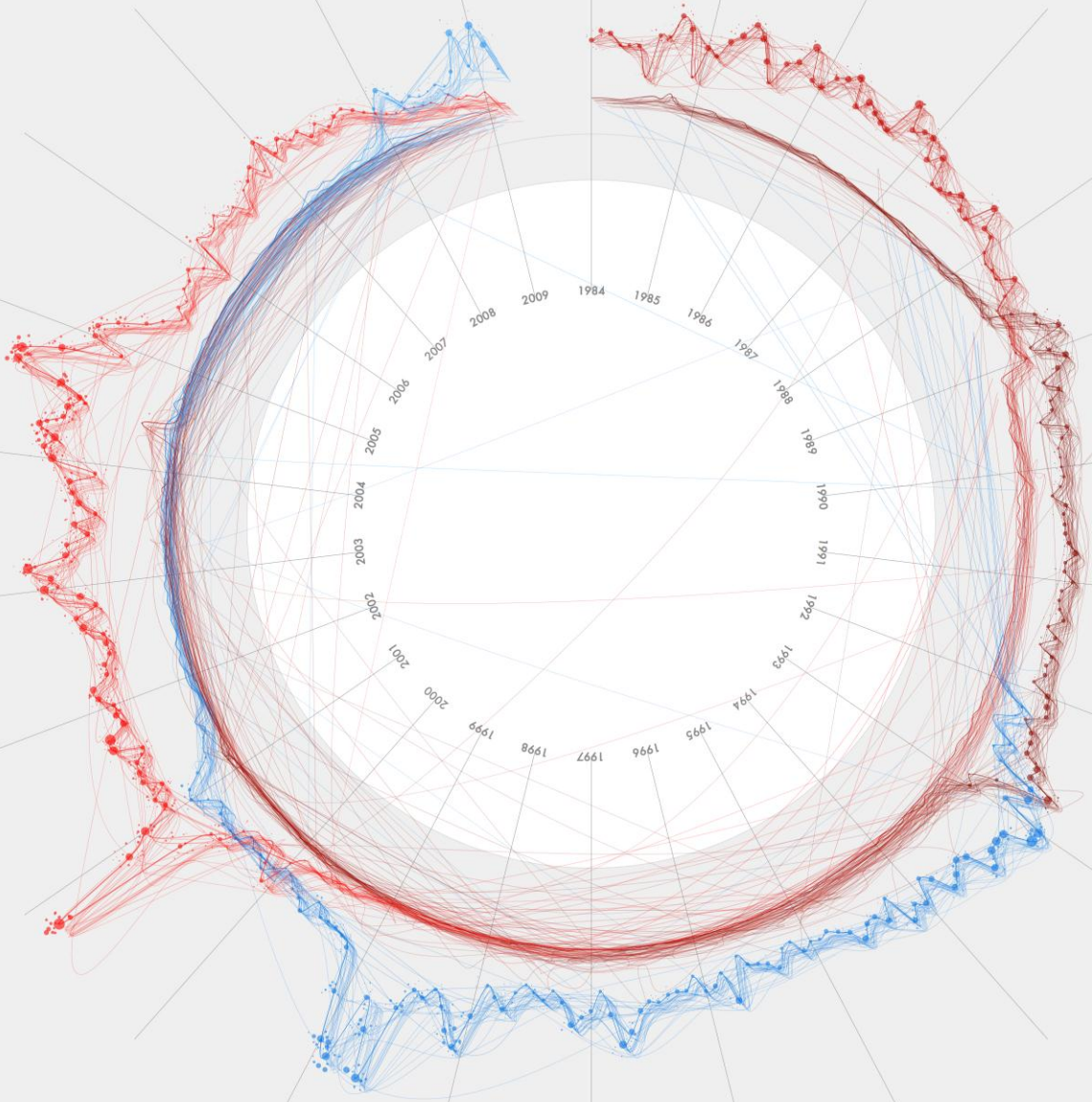
First a tour

Then some techniques

<http://www.smashingmagazine.com/2007/08/02/data-visualization-modern-approaches>



REAGAN, RONALD WILSON
BUSH, GEORGE
CLINTON, BILL
BUSH, GEORGE W
OBAMA, BARACK
1984-2009



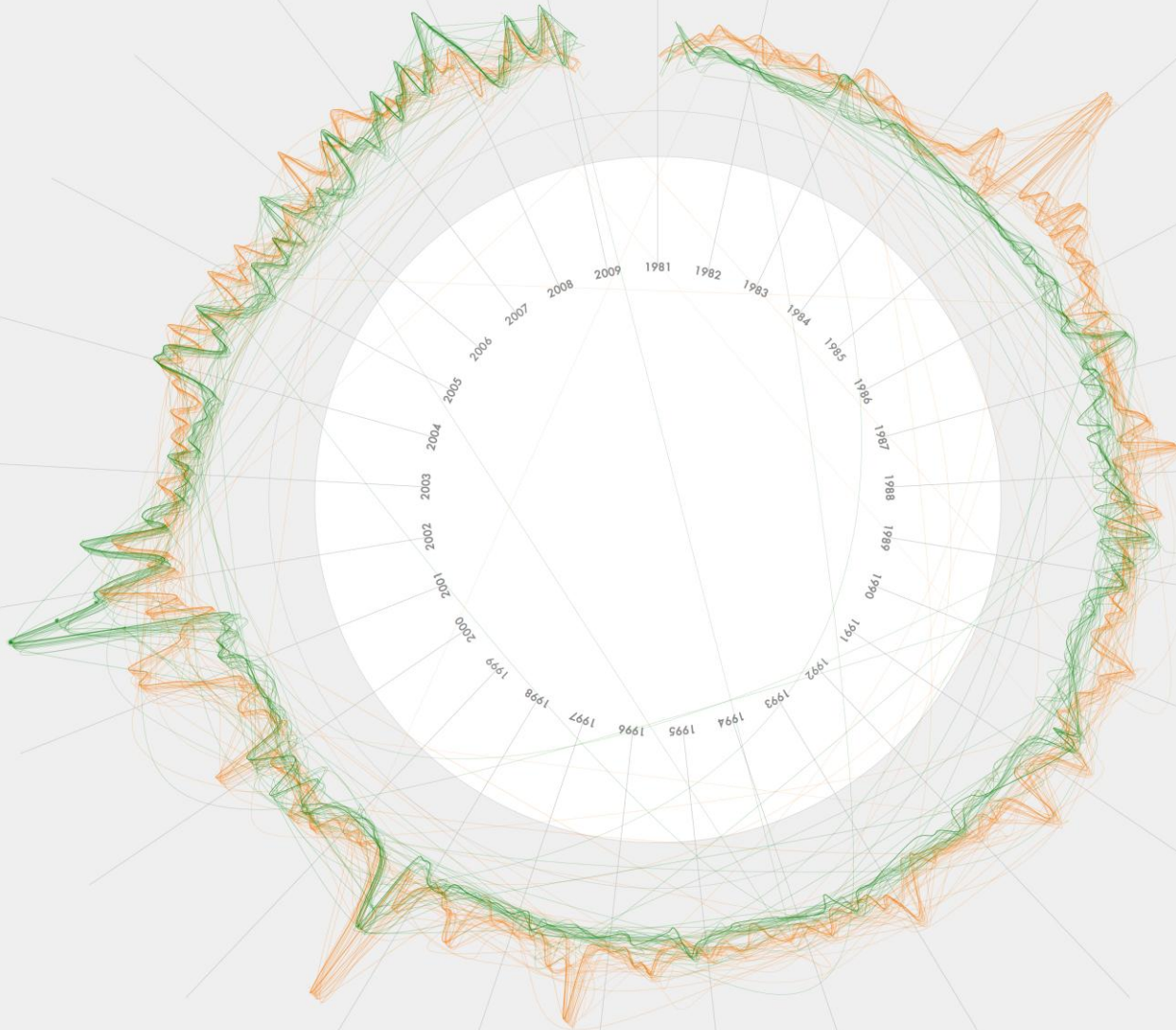
Jer Thorp Data Artist

This graph charts the frequency of mention of the five US Presidents between 1984 and 2009. It also indicates weighting of stories - the darkest line shows front page stories while the lighter lines indicate stories buried deeper in the

Jer Thorp Data Artist

NYTimes Threads - India & Pakistan

This graph charts the frequency of articles mentioning India and Pakistan in the NYT between 1981 and 2009. It also indicates weighting of stories - the darkest line shows front page stories while the lighter lines indicate stories buried deeper in the paper

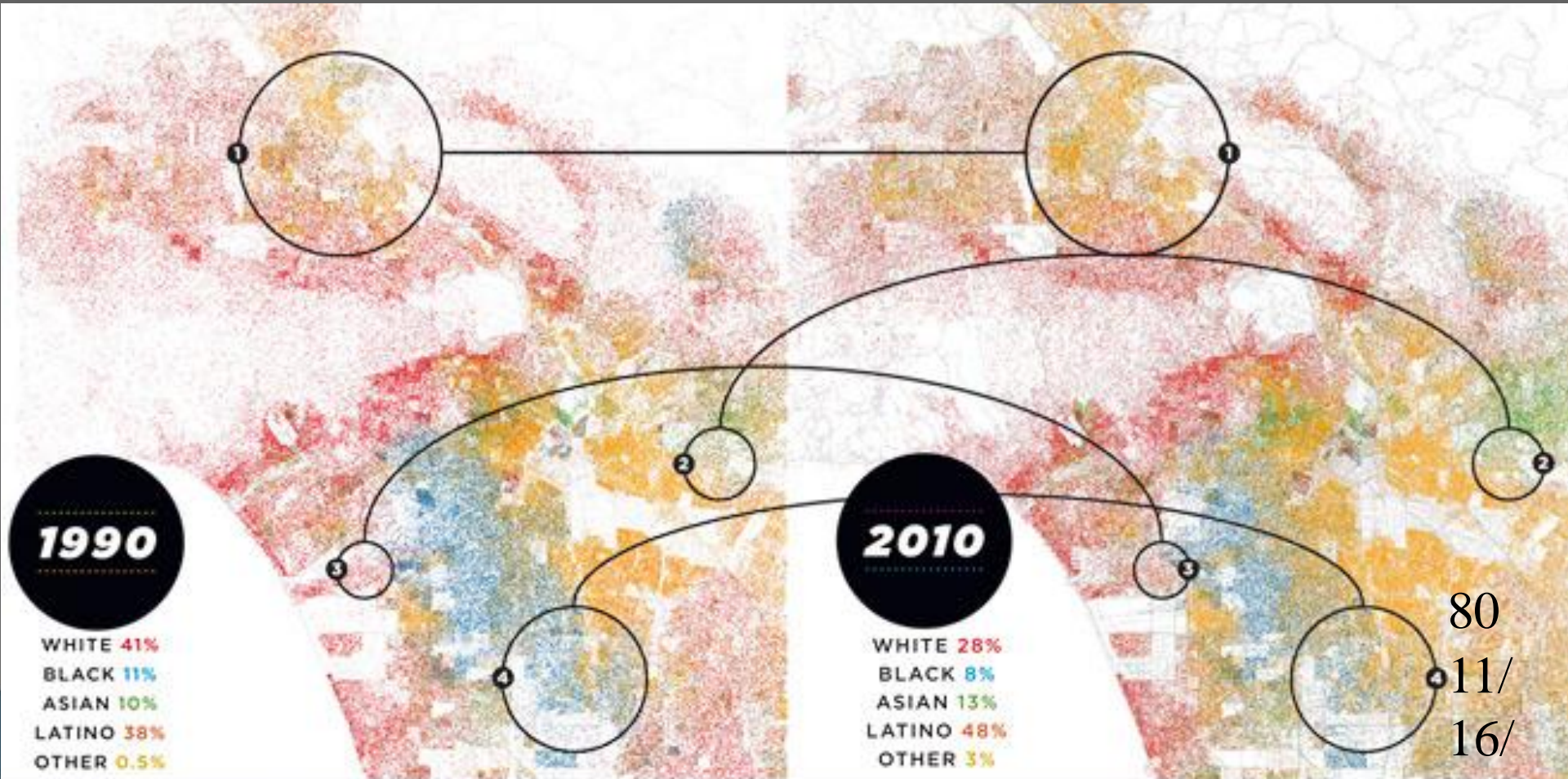


Race in LA (Eric Finkelstein)



Nice narrative @lamag.com

<http://www.lamag.com/features-hidden/race-in-la-see-how-weve-grown/#0412infographic>



World travel and communications recorded on Twitter

Green is physical movement from place to place; purple is @replies from someone in one location to someone in another; combining to white where there is both.

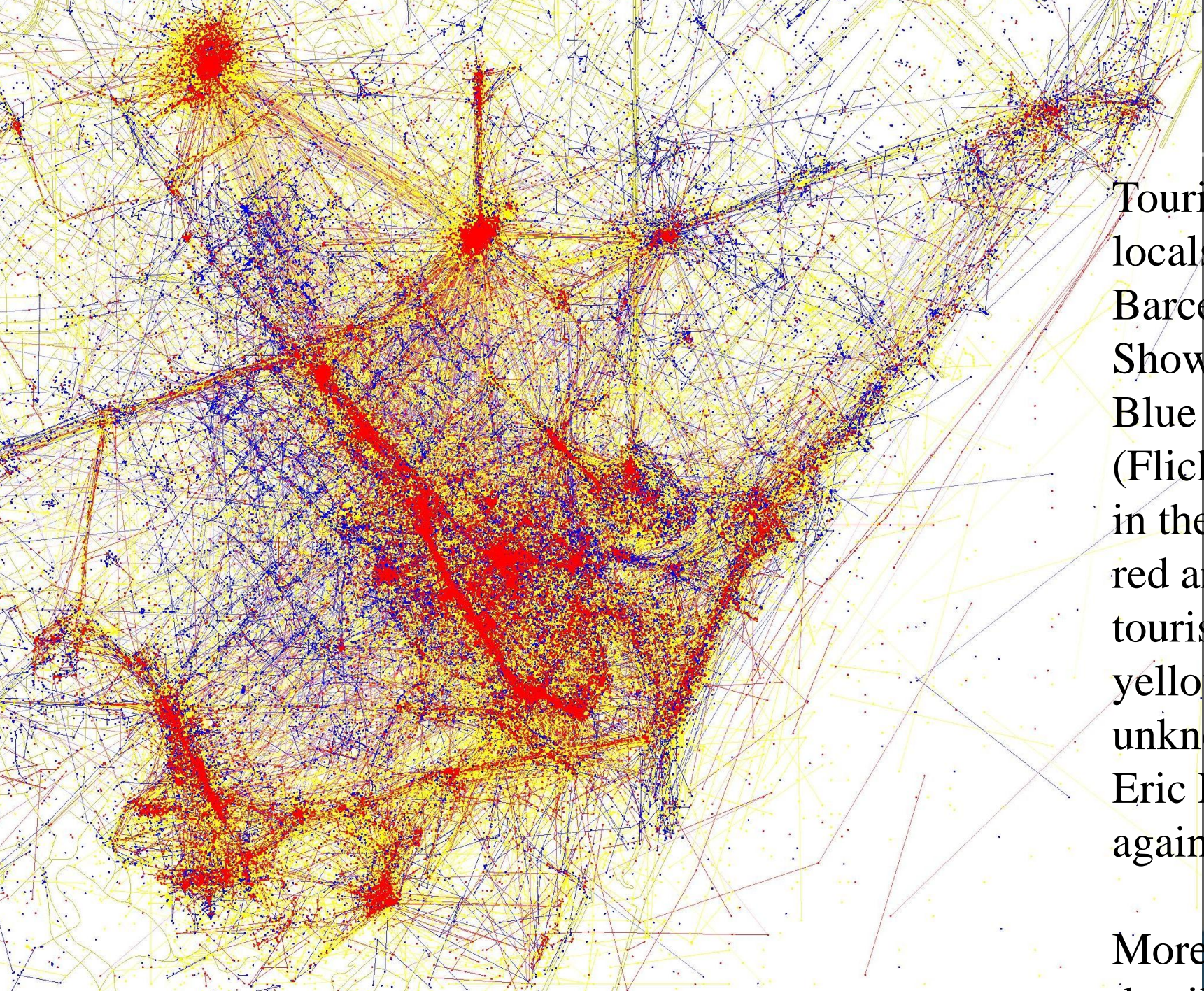
Reported trips to Null Island excluded; all other geotags trusted.
Endpoints of trips are real data; routes in between are fabricated.
Brightness is logarithmic.

Data from the Twitter streaming API through September 1, 2011.

Continent shapes from Natural Earth. Author: Eric Fischer

<https://www.flickr.com/photos/walkingsf/6635655755/in/photostream/>





Tourist vs
locals in
Barcelona.
Shows.
Blue photos
(Flickr) live
in the city,
red are
tourists,
yellow are
unknown.
Eric Fischer
again

82

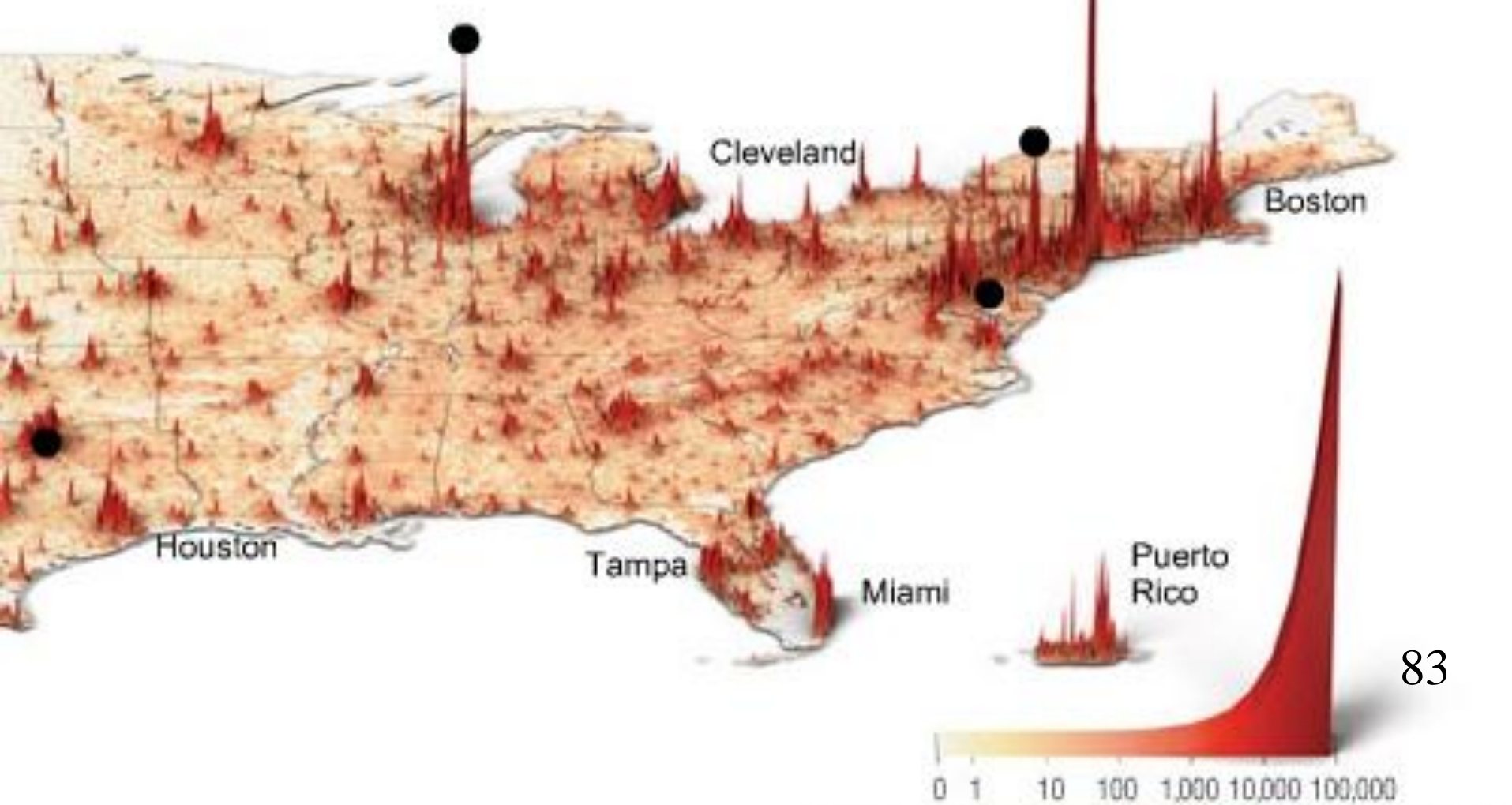
11/

More
details:
koff, 6/12

Time Magazine

bers.

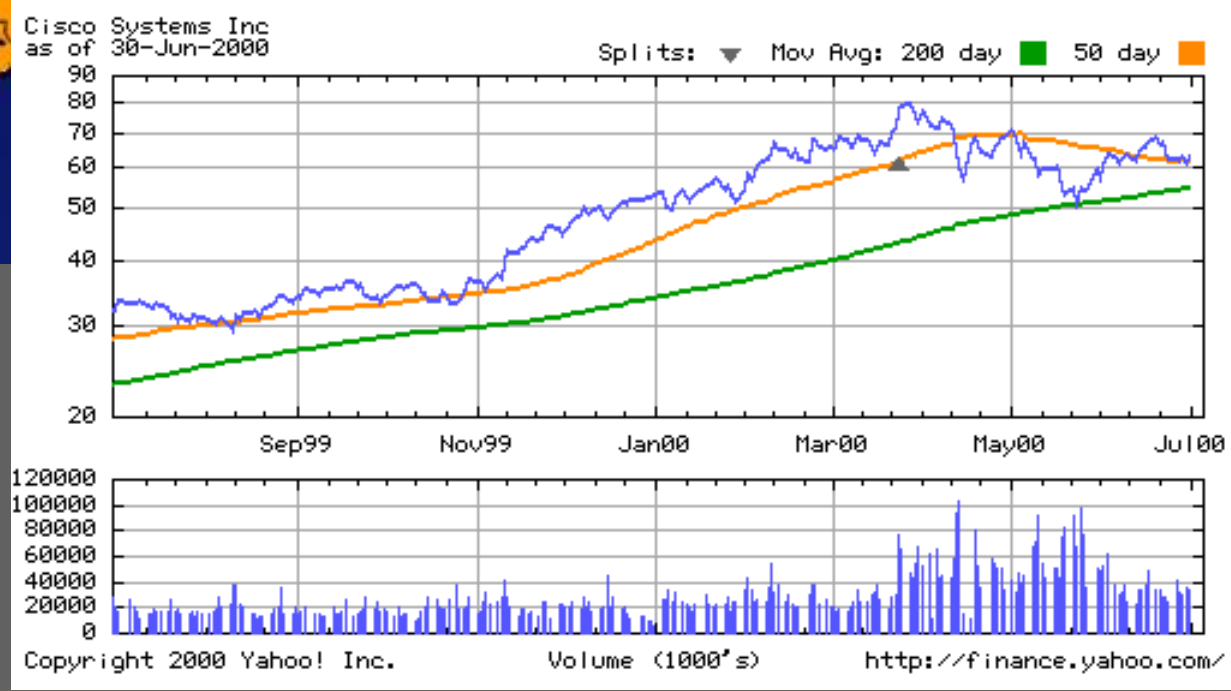
80% of the U.S. population lives in a metropolitan area. Top five population centers are numbered



07.03.2000 -20 -10 0 10 20 30 40 50 60 70 80 90 100 °F

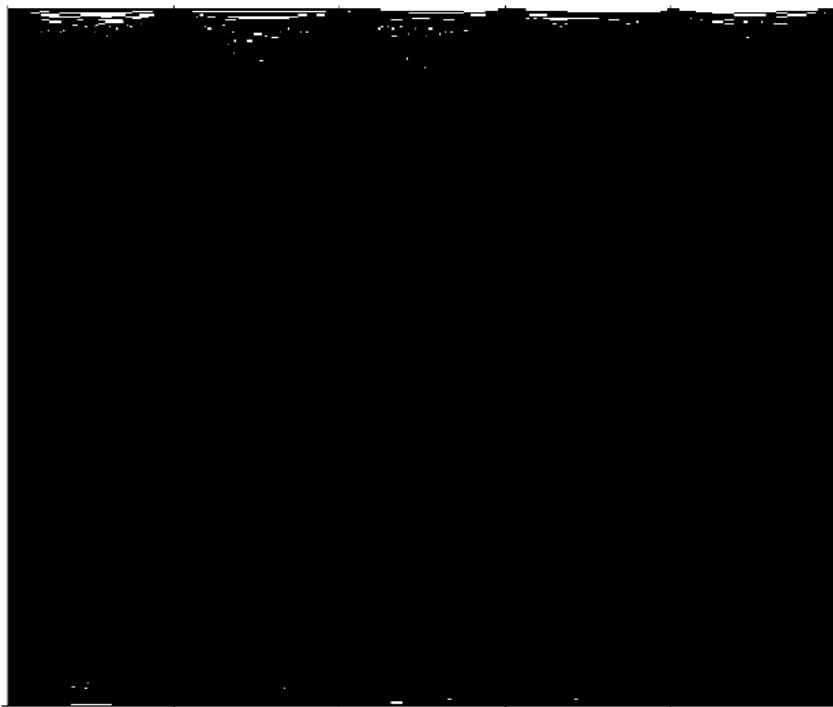


Small Data?

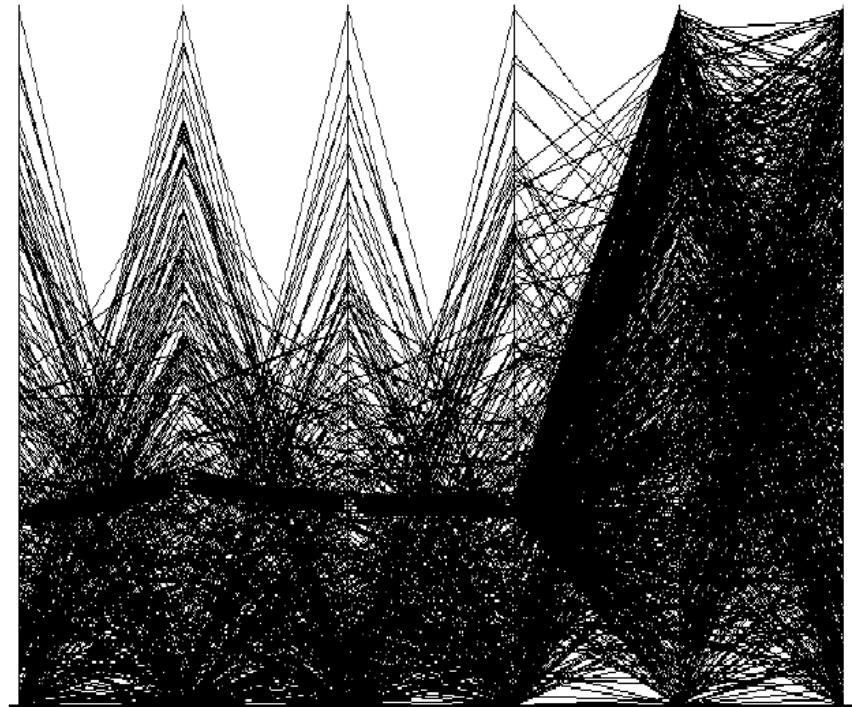


Images from yahoo.com

Visualization Techniques for Big Data



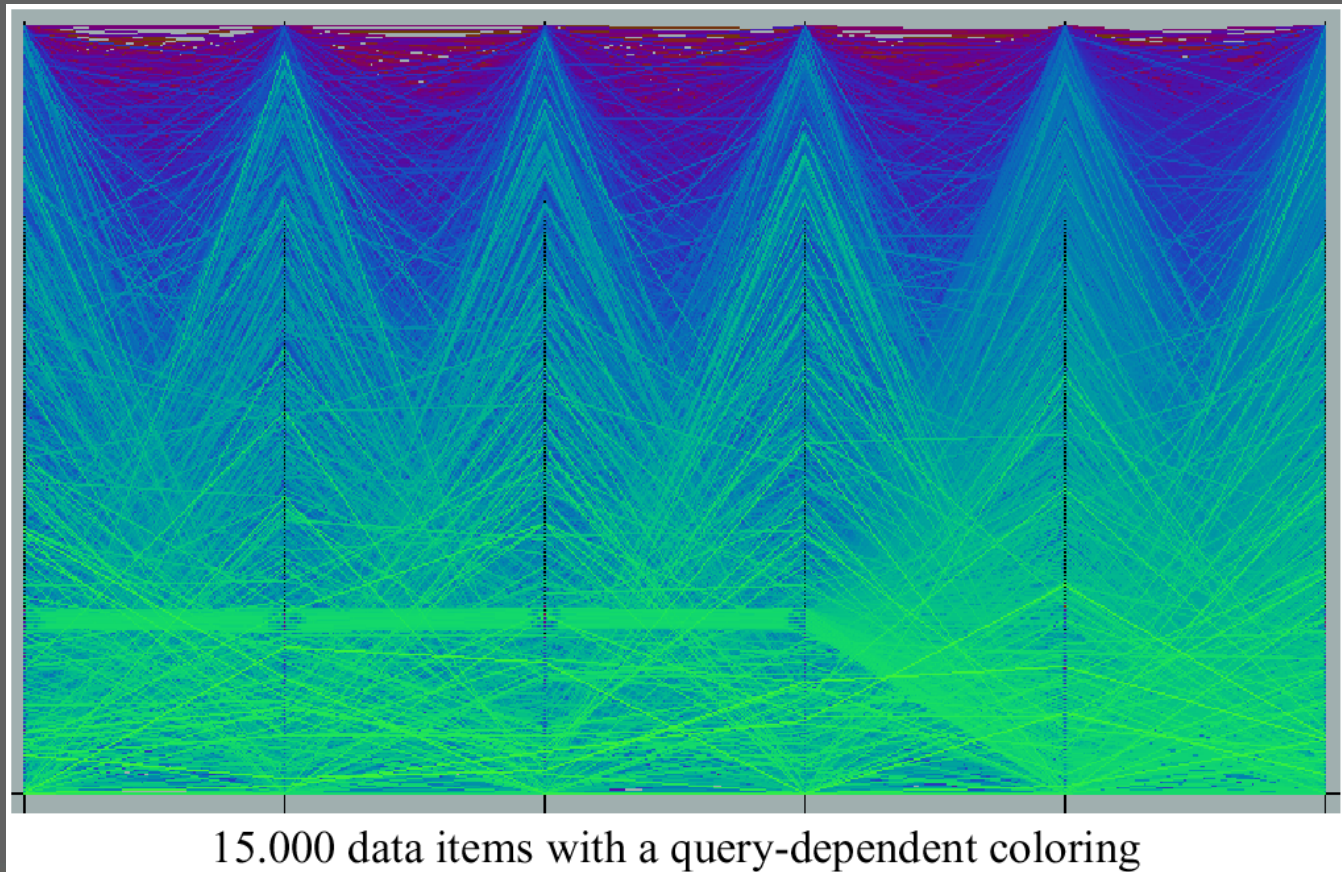
15.000 data items with noise



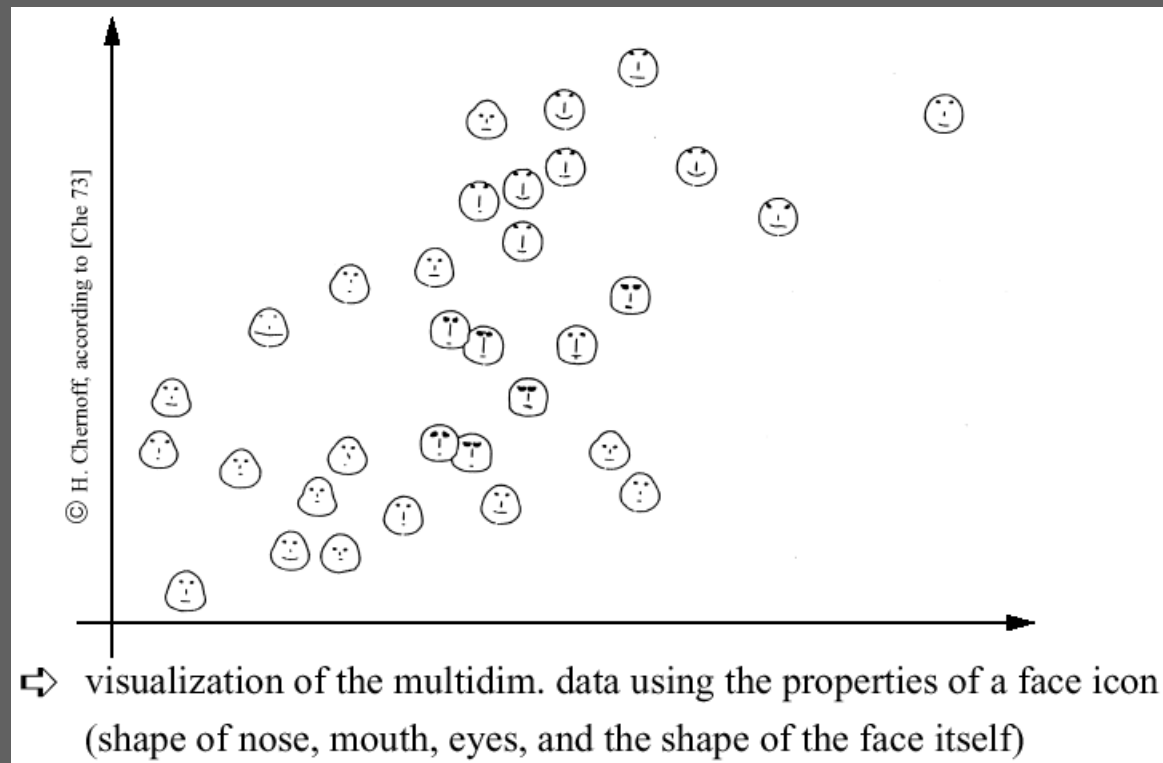
5% of the data (750 data items)



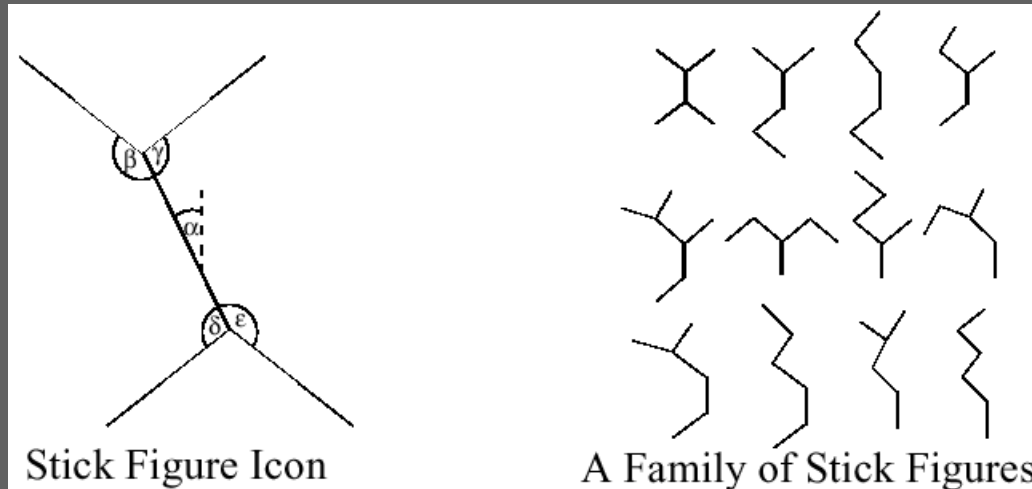
Query Dependent Coloring



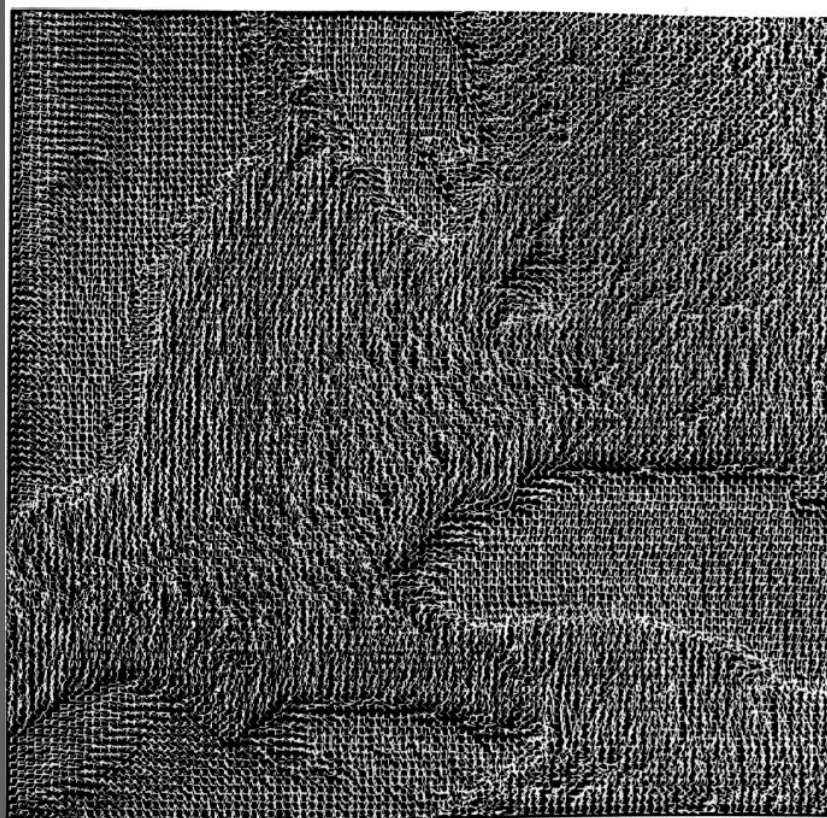
Chernoff-faces



Stick Figures



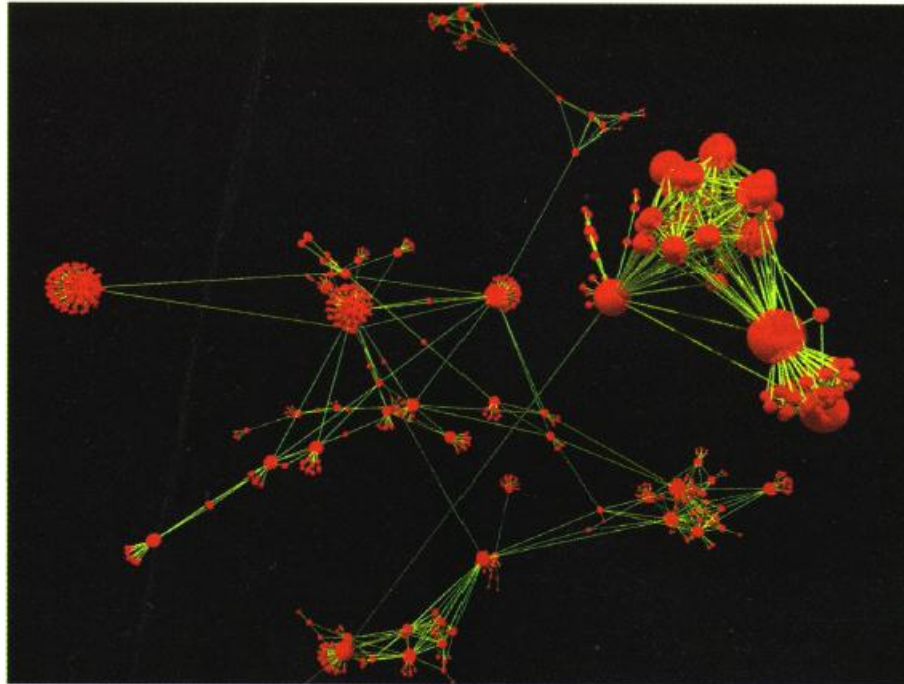
Stick Figures



Graph-based Techniques

Narcissus [HDWB 95]

used by permission of B. Hendley, University of Birmingham



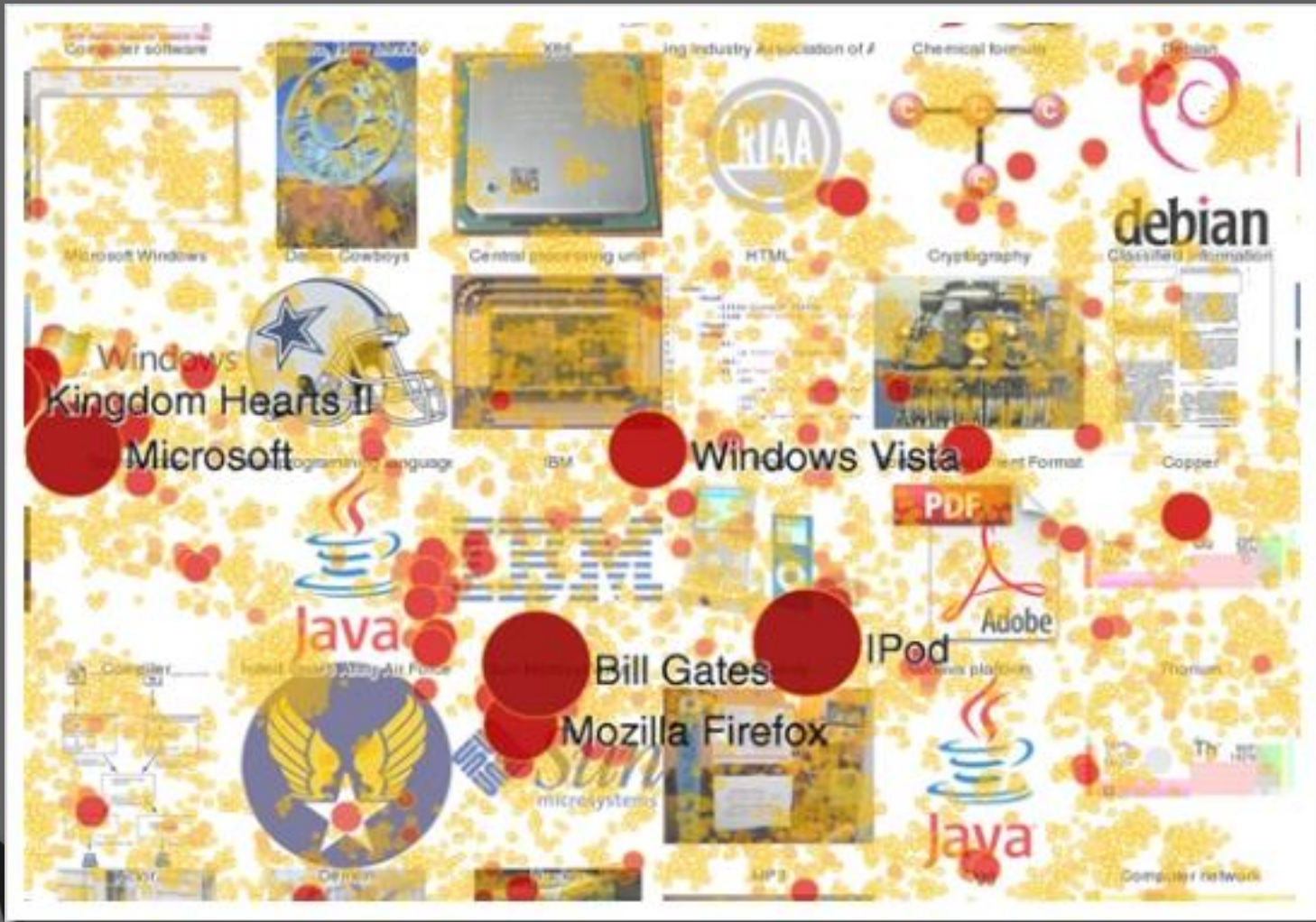
visualization of
a large number
of web pages

⇒ visualization of complex highly interconnected data (e.g., graphs such as the web)



Zoomable UIs

<http://gigapan.com/gigapans/4304>



Distortion Techniques

Basic Idea: Distortion of the image to allow a visualization of larger amounts of data

An alternative to zoomable Uis (or complement)



Distortion Techniques

Basic Idea: Distortion of the image to allow a visualization of larger amounts of data

Simple:

Perspective Wall [MRC91]

Bifocal Displays [SA 82]

TableLens [RC94]

Graph. Fisheye Views [Fur 86, SB94]

Hyperbolic Repr. [LR94, LRP95]

Complex:

Hyperbolic Repr. [LR94, LRP95]

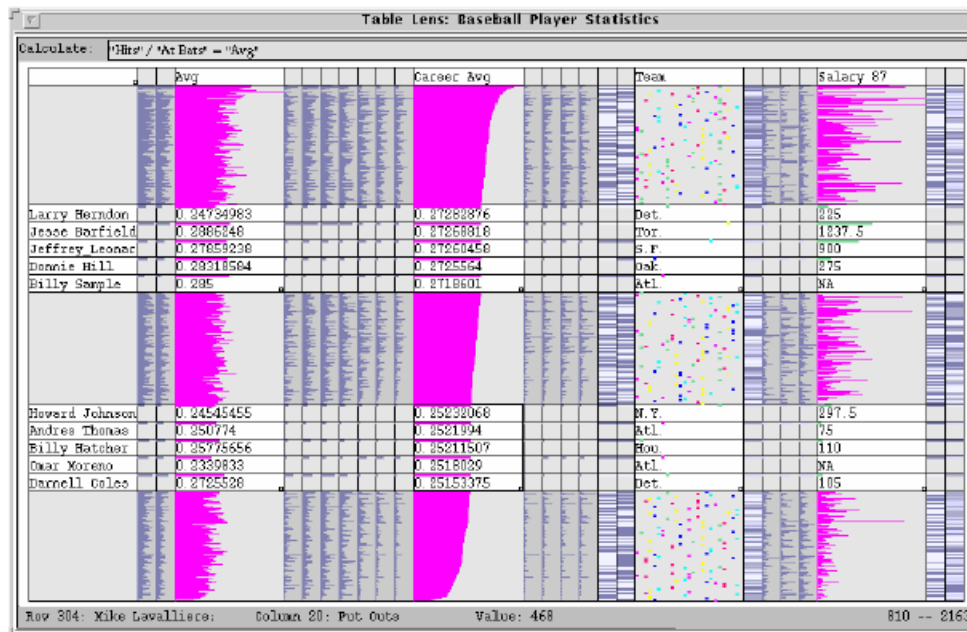
3D-Hyperbolic Repr. [MB95]

Hyperbox [AC91]



Table Lens

used by permission of R. Rao, Xerox PARC



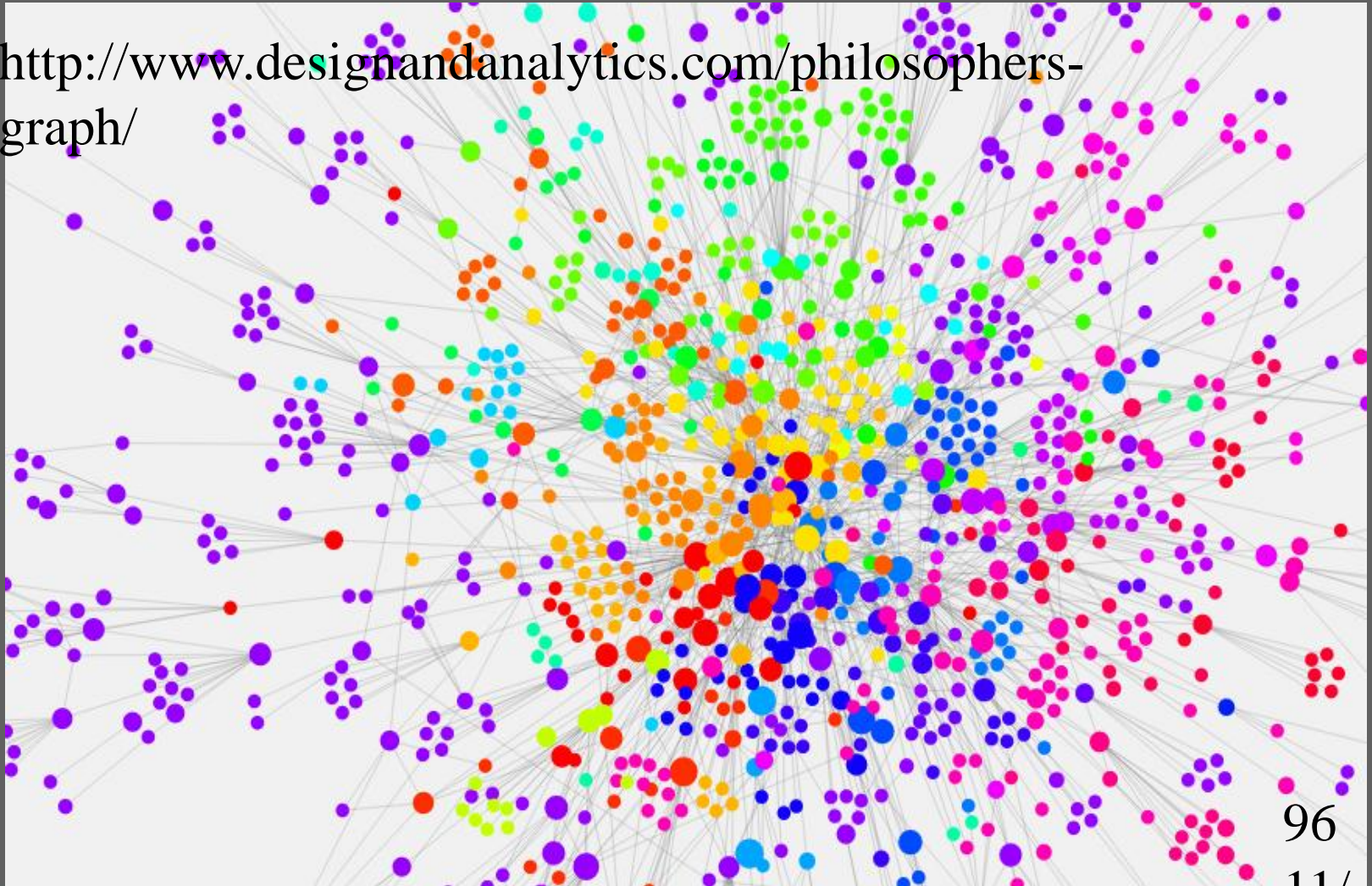
visualization of a baseball database with a few rows being selected in full detail

⇒ compact visualization of a table (spreadsheet / database) with the possibility of viewing portions of the table in more detail



Networks of Information

<http://www.designandanalytics.com/philosophers-graph/>

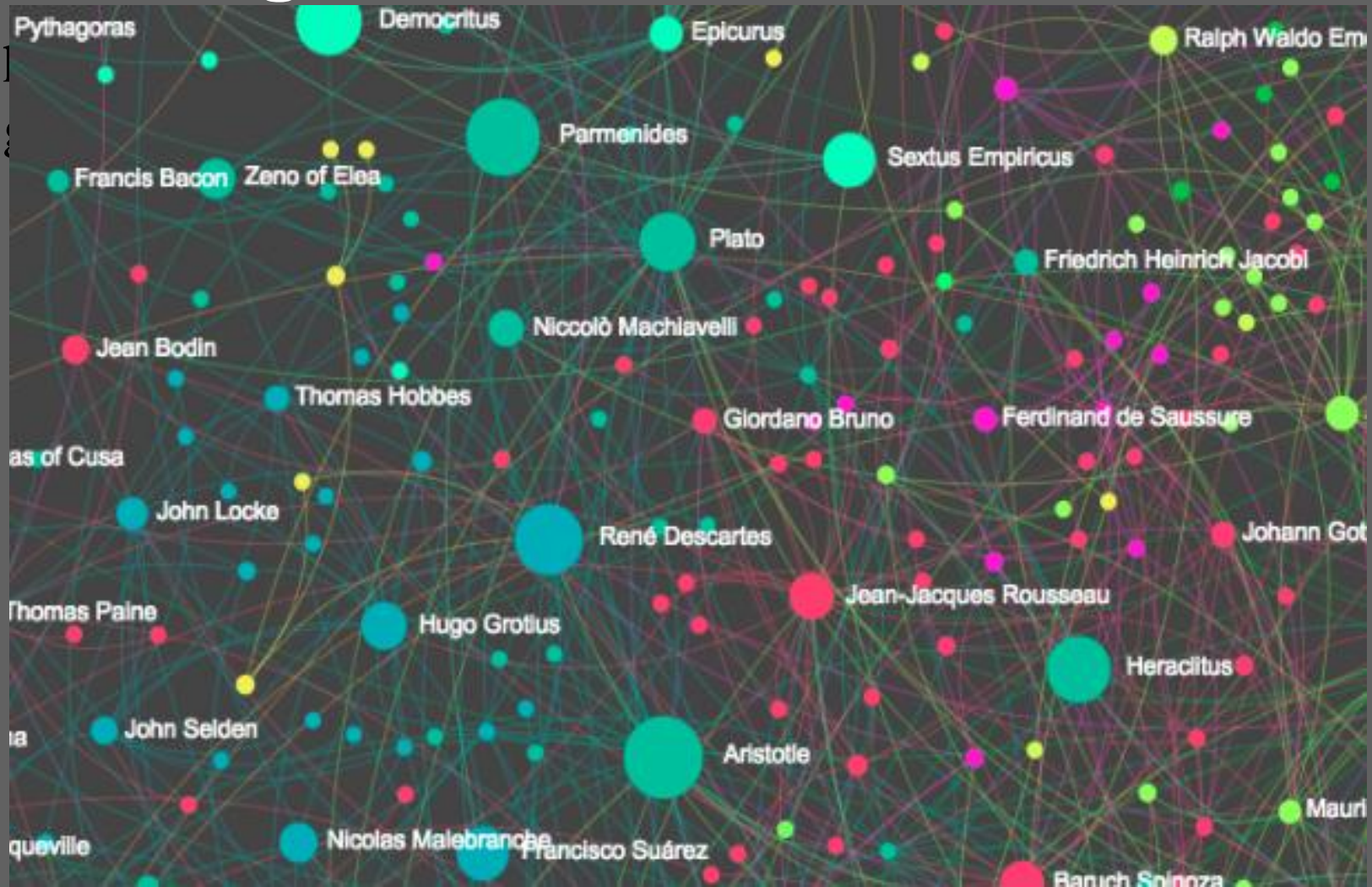


96

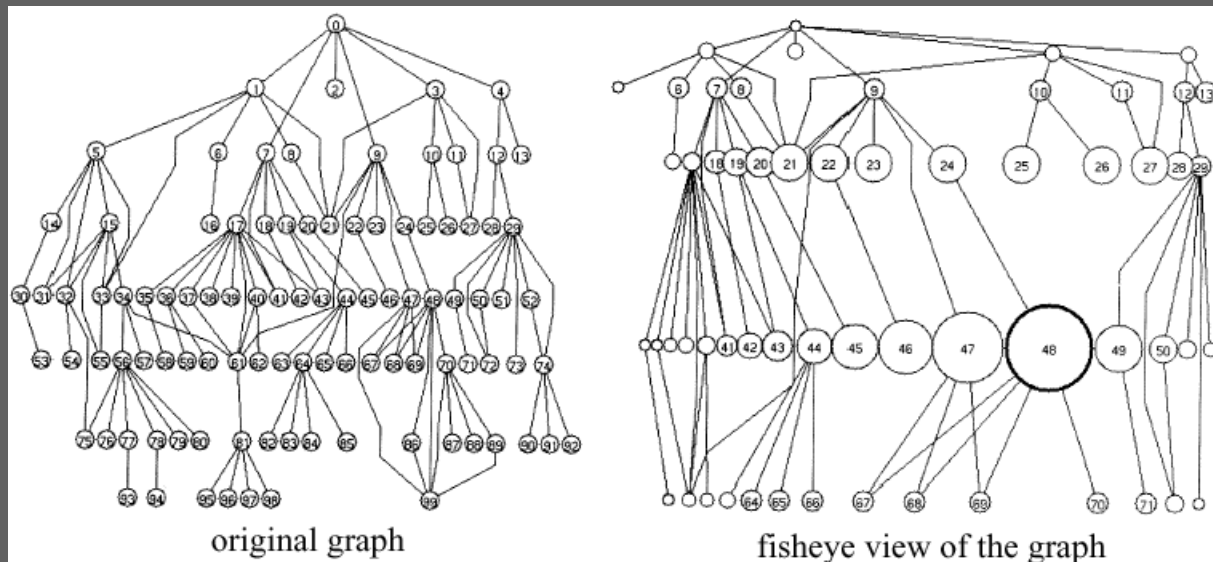
11/



Networks of Information – with zooming

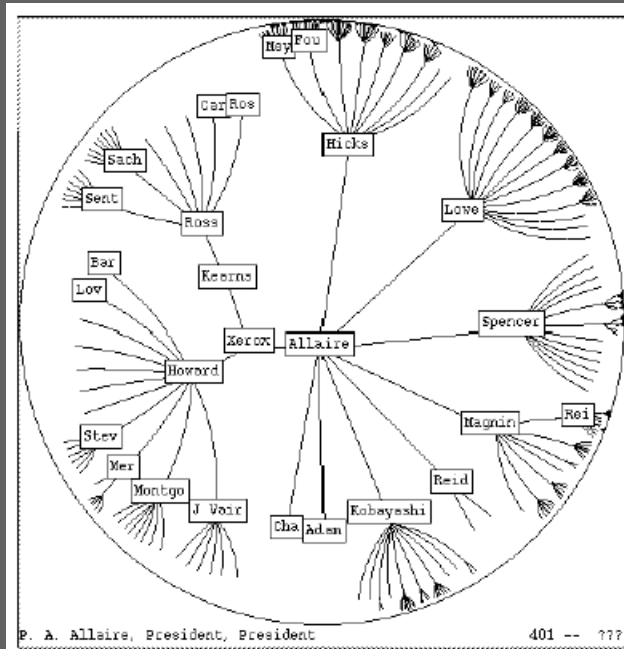


Fisheye View

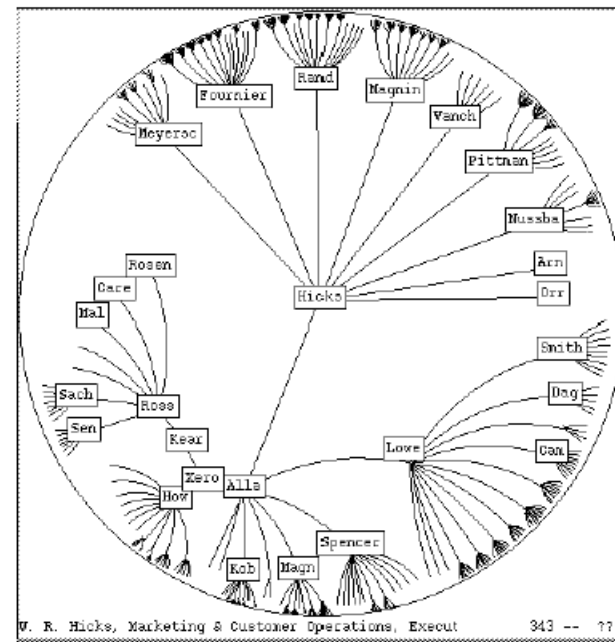


- ⇒ graph visualization using a fisheye perspective
- ⇒ shows an area of interest quite large and with detail and the other areas successively smaller and in less detail

Hyperbolic Tree



used by permission of R. Rao, Xerox PARC



used by permission of R. Rao, Xerox PARC

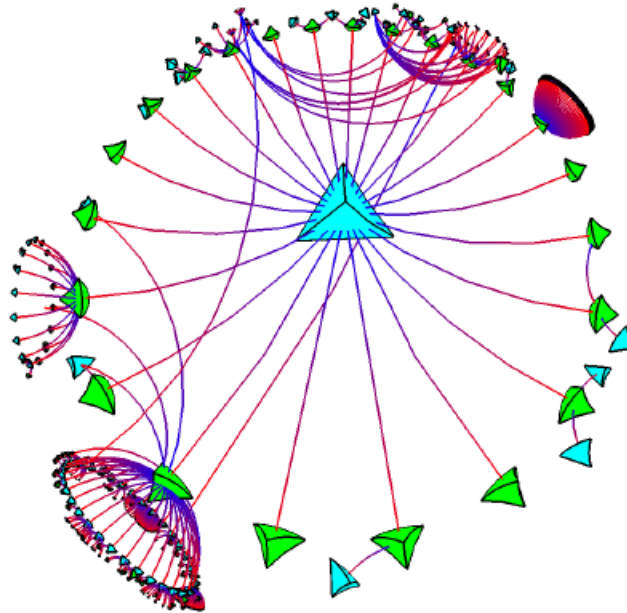
visualization of a large organizational hierarchy

⇒ visualization of a tree structure in hyperbolic space with different foci



3D Hyperbolic Representation

used by permission of T. Munzner, Stanford University



visualization
of a large number
of connected
web-pages

⇨ visualization of a graph in 3D hyperbolic conetree-like representation



Summary

Make use of every pixel

Use color and location to break the data up and allow the viewer to easily filter

Give a sense of things and use zooming for detail

Add a dimension such as time or height for a key variable

Allow exploration through distortion, filtering, highlighting, and linking

Exploit hierarchy and connectivity

101

11/



Applied Data Use: StepGreen

Capture

Self reported data
on behavior

Motion & GPS

Temperature & Energy

Know

Green Actions

Transportation choices

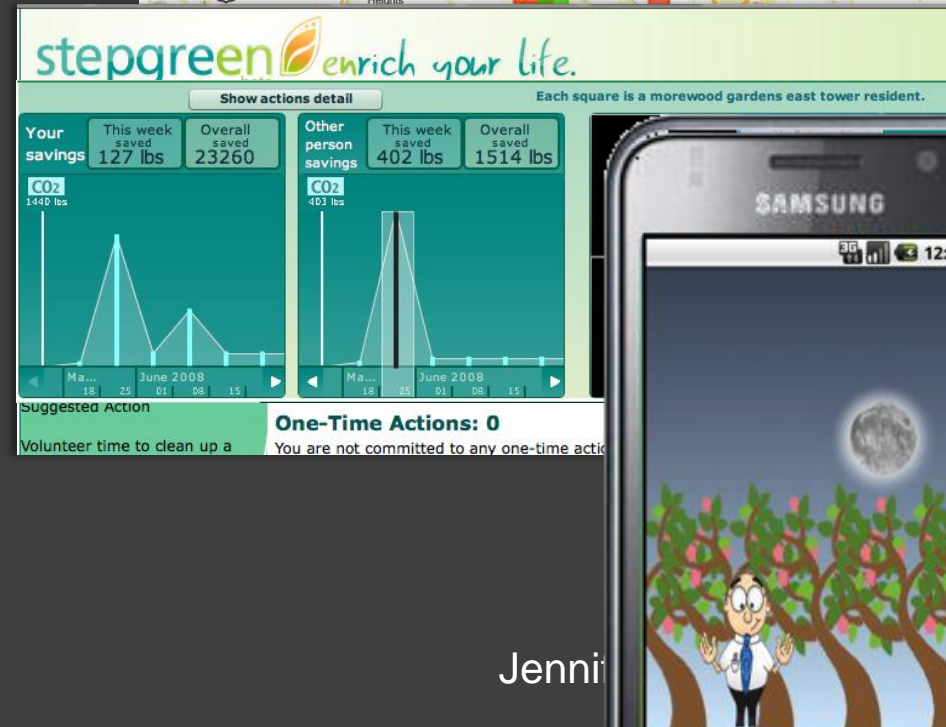
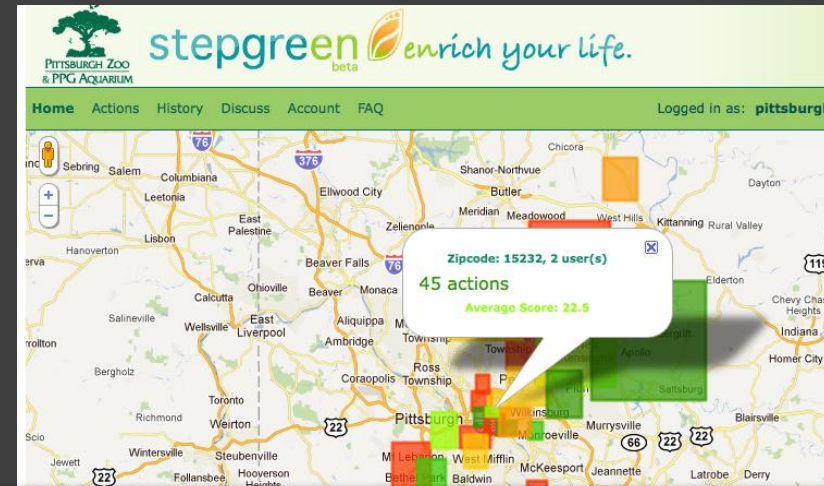
Appliances in Use

Adapt and Act

Visualize

Expose

...



Applied Example: Ubigreen

Capture

Self reported data
on behavior

Motion & GPS

Temperature & Energy

Know

Green Actions

Transportation choices

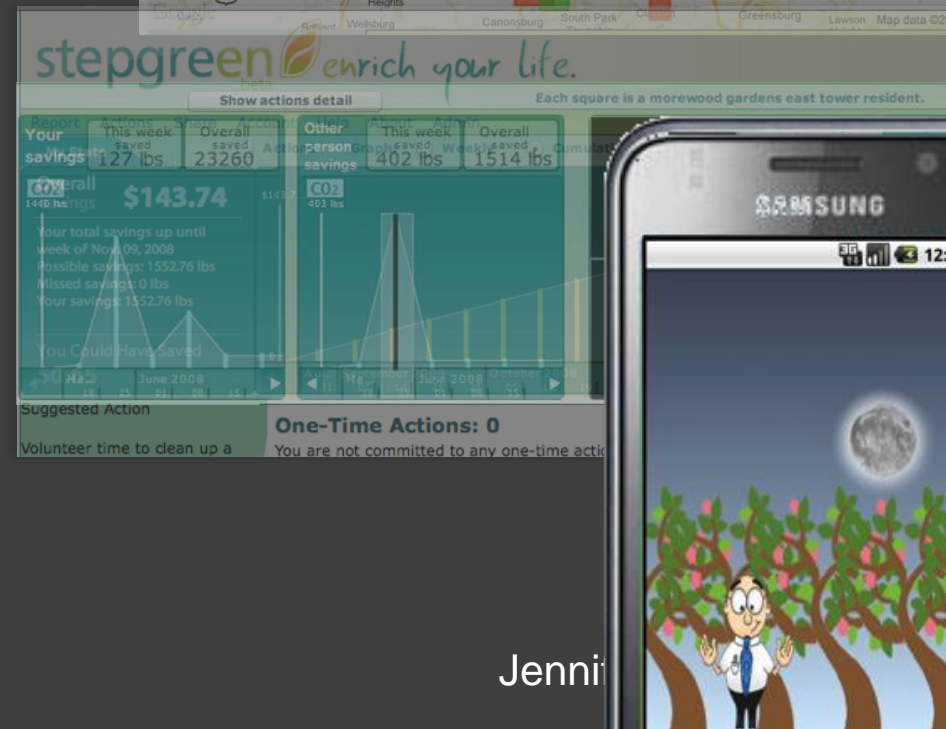
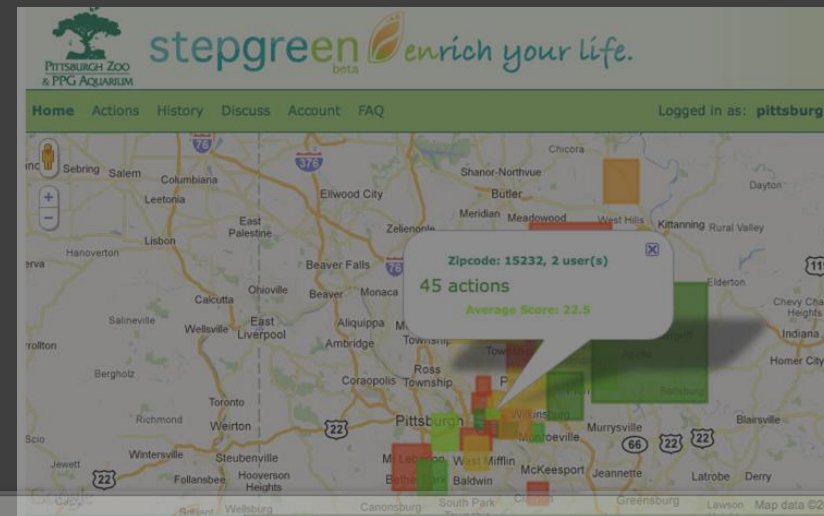
Appliances in Use

Adapt and Act

Visualize

Expose

...



3 week field study [CHI'09]

Current
Activity

Phone
Background
(Wallpaper)

Values
Icon Bar

Evolving
Image



Engagement



“It’s omnipresent”

- Participant 9

“I want to have different stories every week
... to maintain curiosity in the app”

- Participant 8

Real-life game

One participant complained that when a trip hadn’t been automatically recorded, “I felt like I was being cheated out of my ‘points’”

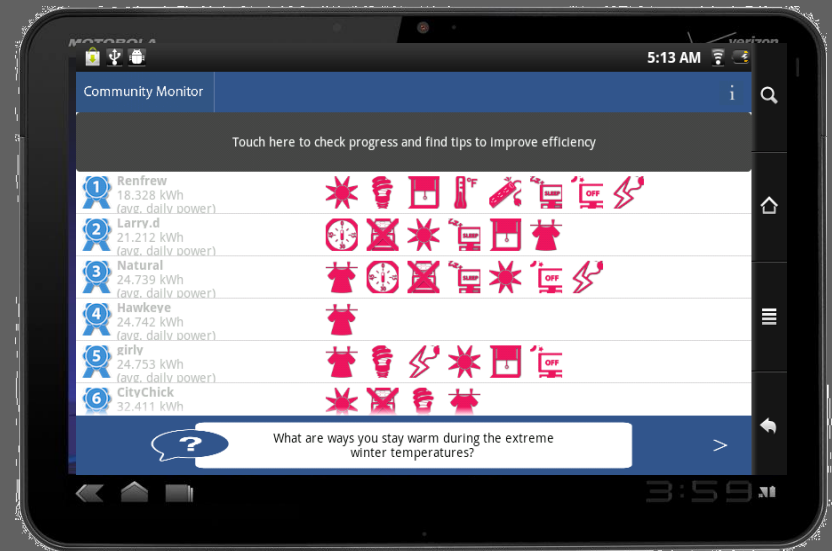
- Participant 15

Social

“Some people at work knew about the polar bear and every day they asked me about it.
‘Did you get a seal today?’”

- Participant 14

Longer Term, Real-World Deployments



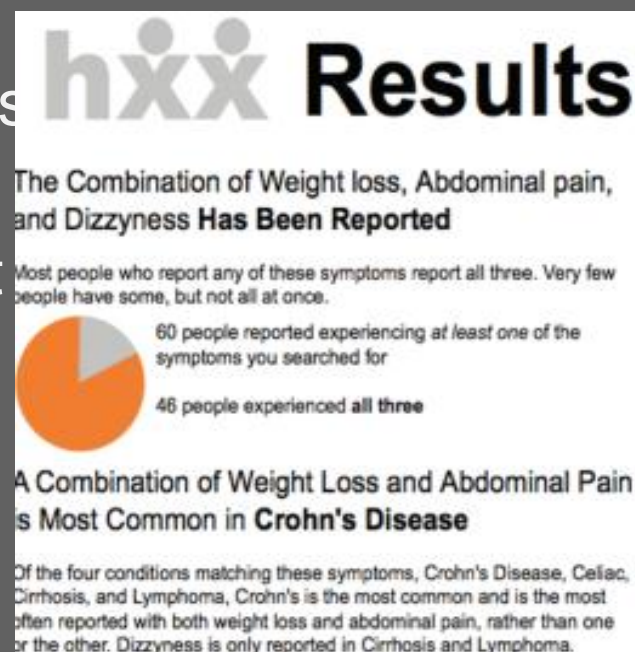
Other Application Areas

Health

Self reported data on symptoms and conditions

Internet data: Extract argument structures & enhance search

Forums: predict expertise, highlight time on site, *etc.*



Summary: Making Data Actionable

