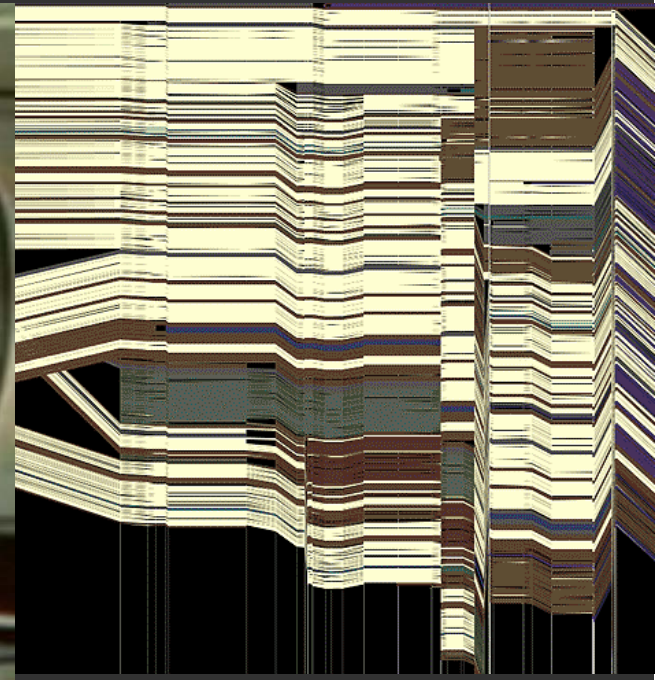
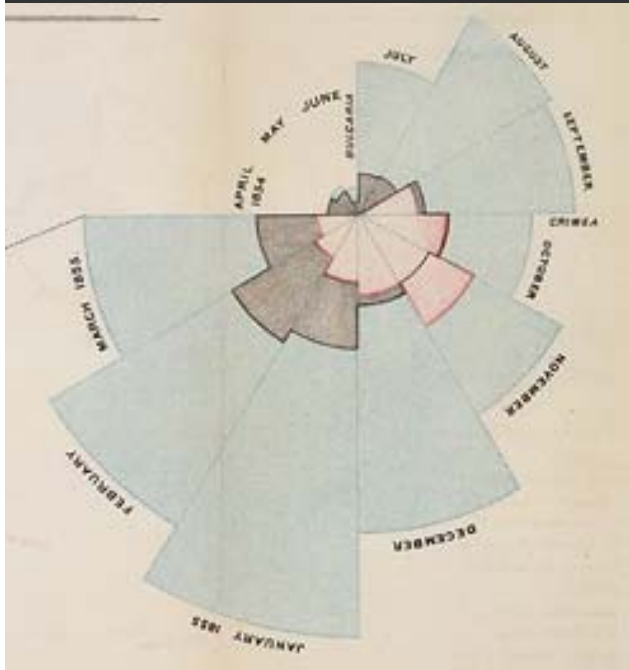


Interactive Tools for Data Transformation & Visualization

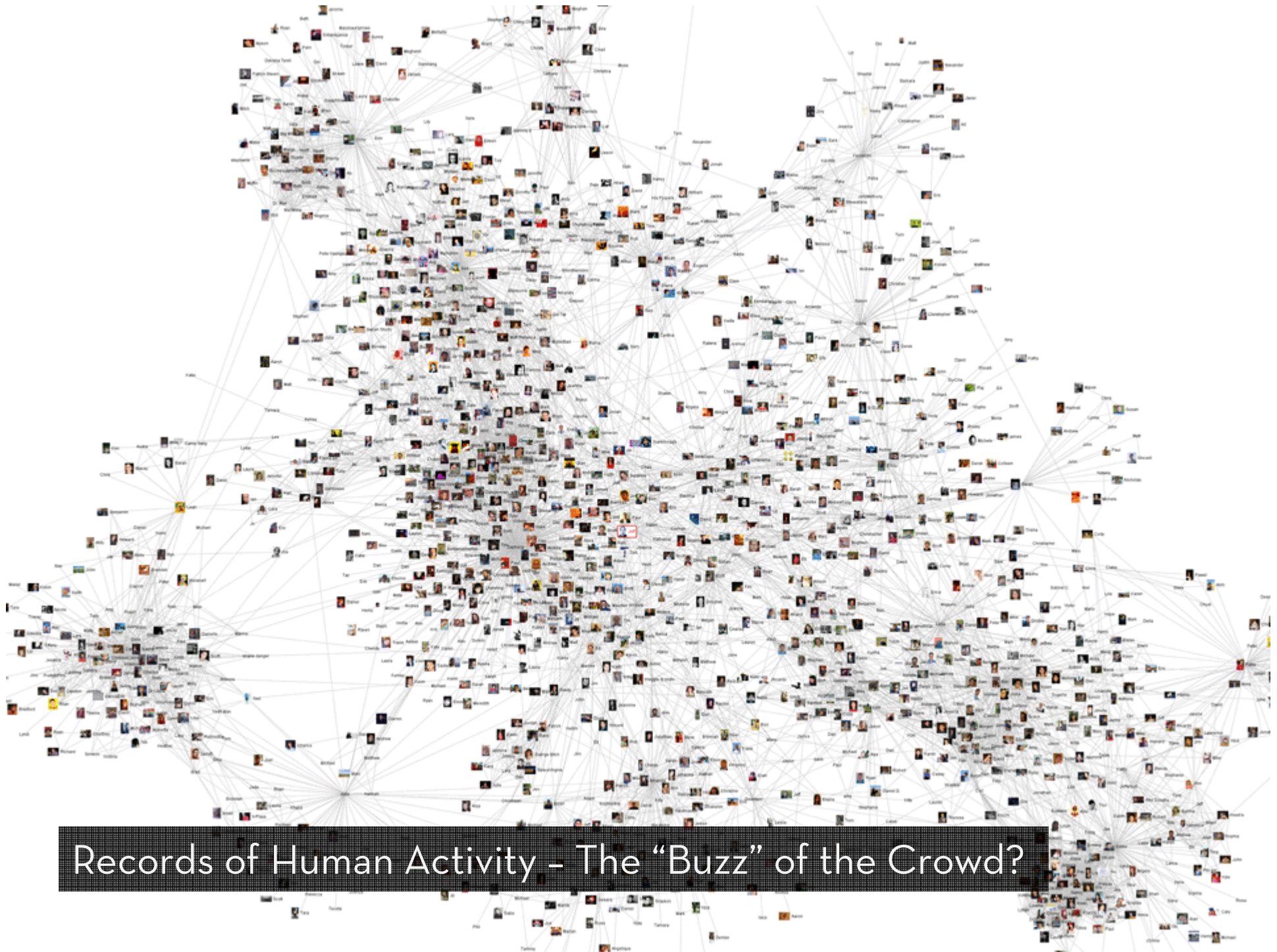


Jeffrey Heer Stanford University

**How much data (bytes)
will we produce in 2010?**

2010: 1,200 exabytes
10x increase over 5 years

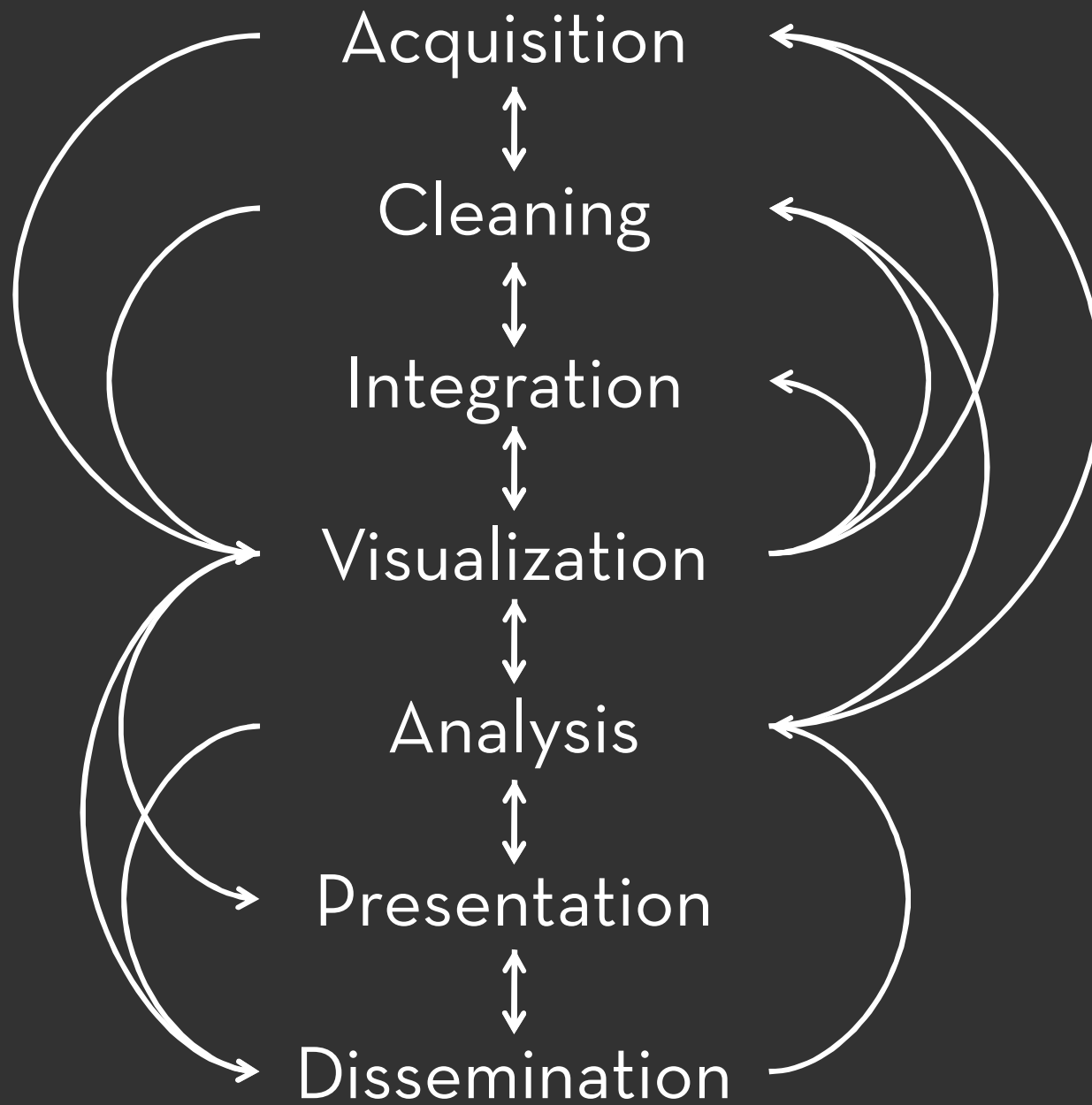
Gantz et al, 2008, 2010

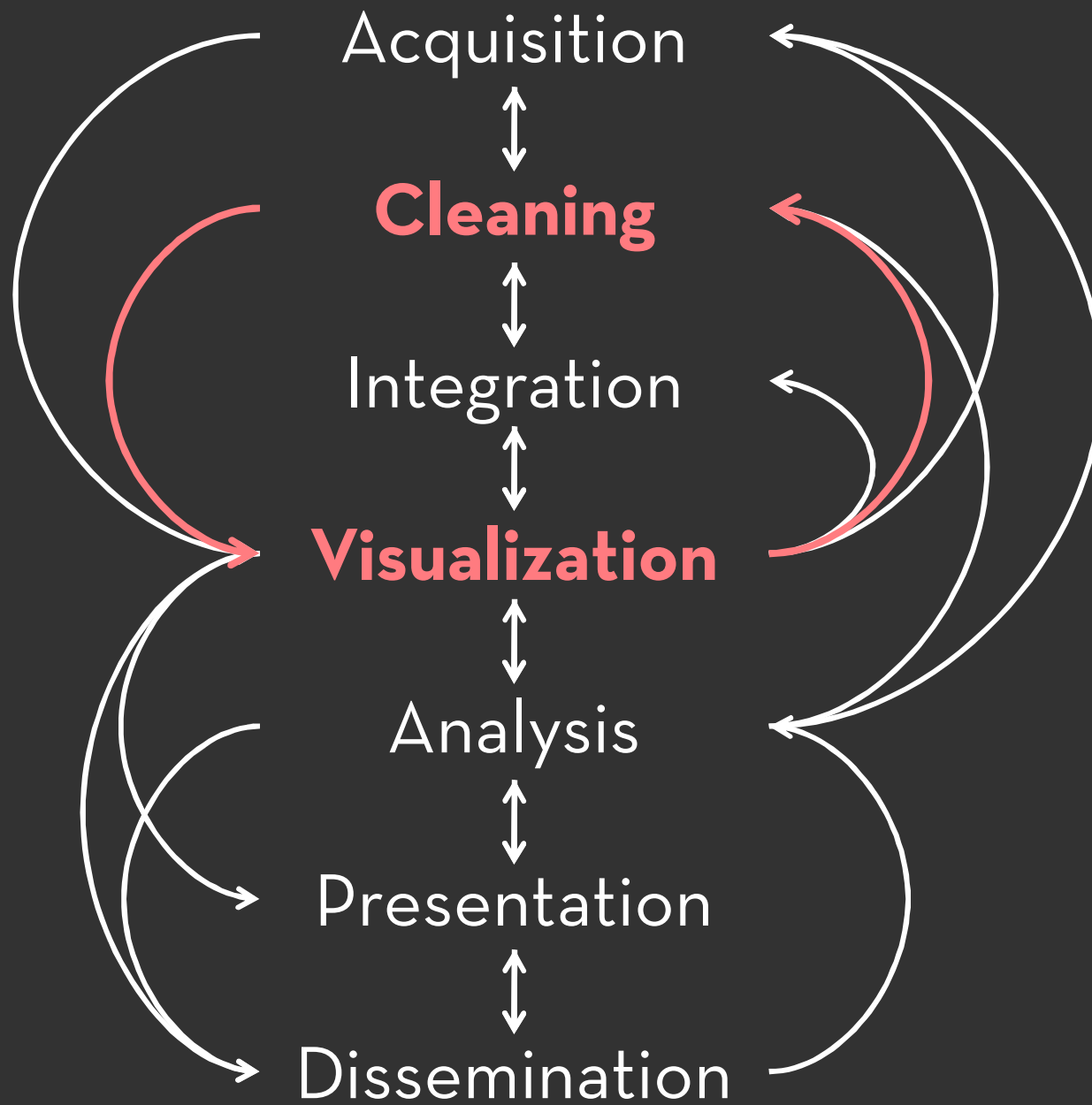


Records of Human Activity – The “Buzz” of the Crowd?

The ability to take data—to be able to **understand** it, to **process** it, to **extract value** from it, to **visualize** it, to **communicate** it—that's going to be a hugely important skill in the next decades, ... because now we really do have **essentially free and ubiquitous data**. So the complimentary scarce factor is the ability to understand that data and extract value from it.

Hal Varian, Google's Chief Economist
The McKinsey Quarterly, Jan 2009





Data Wrangler

Transform Script

Import Export

► Split **data repeatedly** on **newline** into rows

► Split **split repeatedly** on **,** into **columns**

► Promote row **0** to header

Text

Columns

Rows

Table

Clear

Delete rows **7,9**

Delete empty rows

Fill rows **7,9** in **all columns** by **copying** values from **above**

	Year	#	Property_crime_rate
0	Reported crime in Alabama		
1			
2	2004		4029.3
3	2005		3900
4	2006		3937
5	2007		3974.9
6	2008		4081.9
7			
8	Reported crime in Alaska		
9			
10	2004		3370.9
11	2005		3615
12	2006		3582
13	2007		3373.9

with **Sean Kandel**, Andreas Paepcke & Joe Hellerstein

From UI to running code...

```
split('data').on(NEWLINE).max_splits(NO_MAX)
split('split').on(COMMA).max_splits(NO_MAX)
columnName().row(0)
delete(isEmpty())
extract('Year').on(/.*\/).after(/in /)
fill('extract').method(COPY).direction(DOWN)
delete('Year starts with "Reported crime in"')
columnName('extract').to('State')
```

Data Wrangler

Declarative data transformation language

- **Tuple mapping** – split, merge, extract, delete
- **Lookups and joins** – e.g., FIPS code to US state
- **Reshaping** – e.g., cross-tabulation
- **Sorting, aggregation, etc.**
- Informed by prior work in databases, namely Potter's Wheel & SchemaSQL

Data Wrangler

Declarative data transformation language

+

Mixed-initiative interface for data transforms

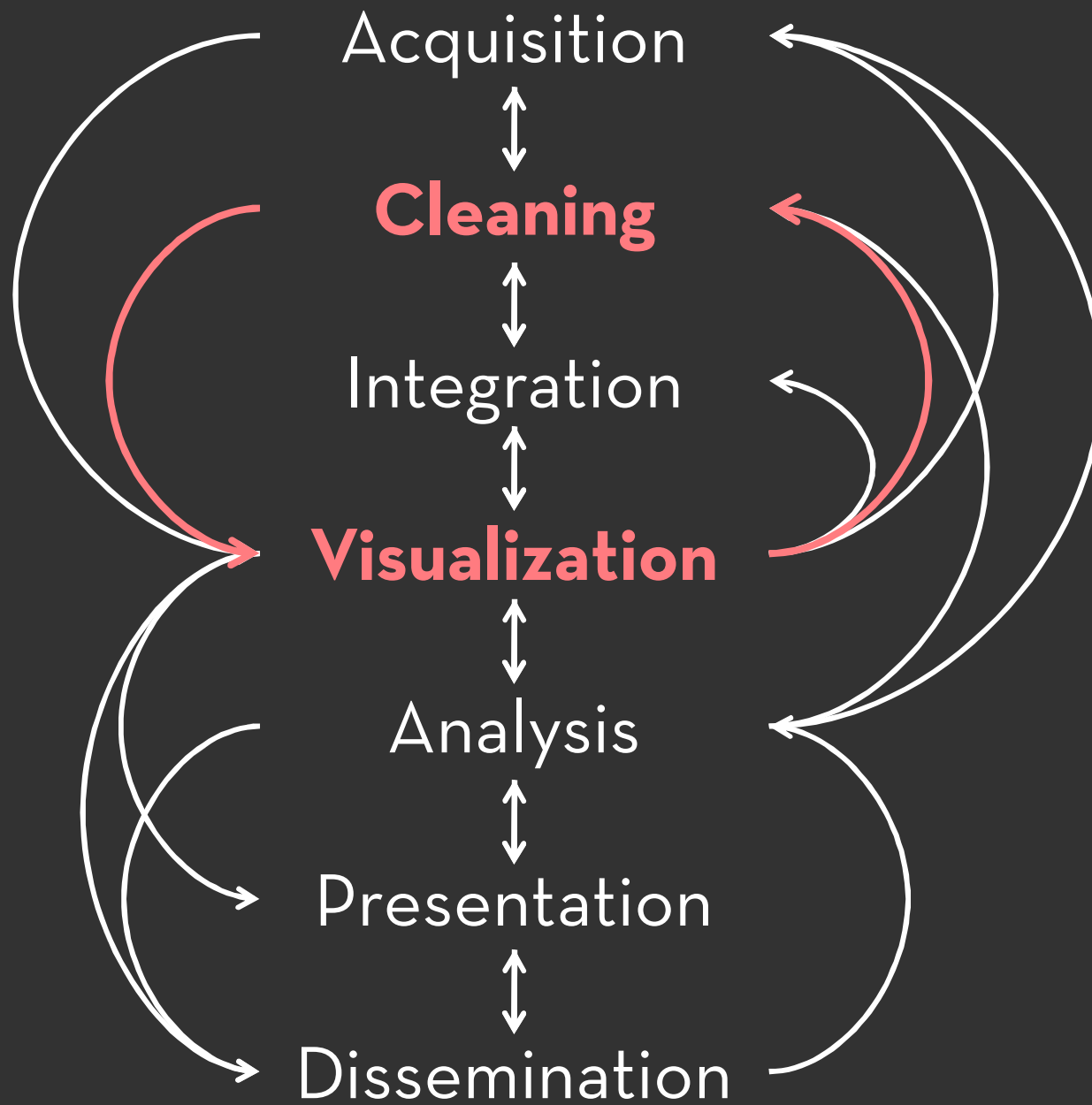
- **Select** data elements of interest
- **Suggest** applicable transforms
- Enable rapid **preview and refinement**

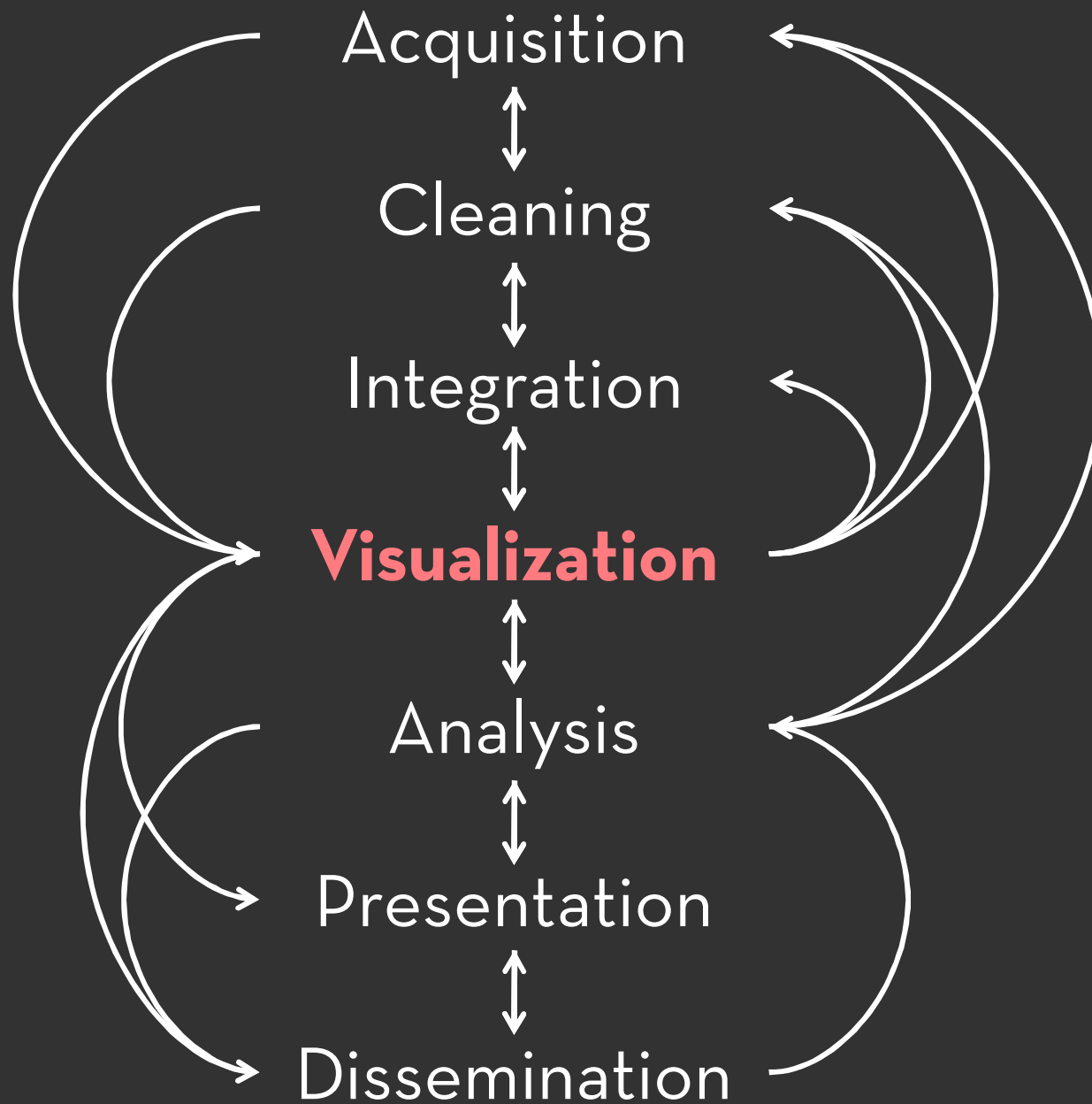
Comparative Evaluation

Compared Wrangler performance to Excel with 3 data cleaning tasks on small data sets.

Median completion time for Wrangler at least twice as fast in all tasks.

Skilled Excel users benefit disproportionately!





How do people create visualizations?

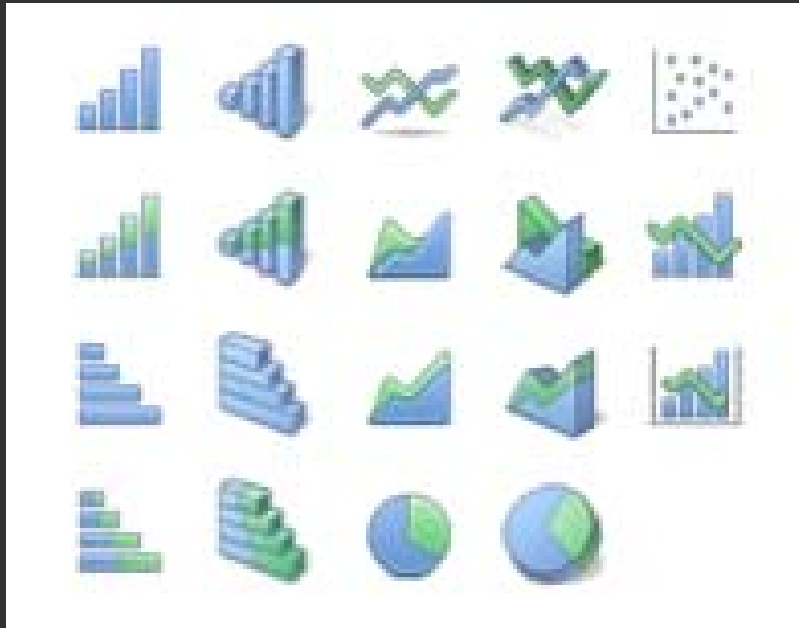
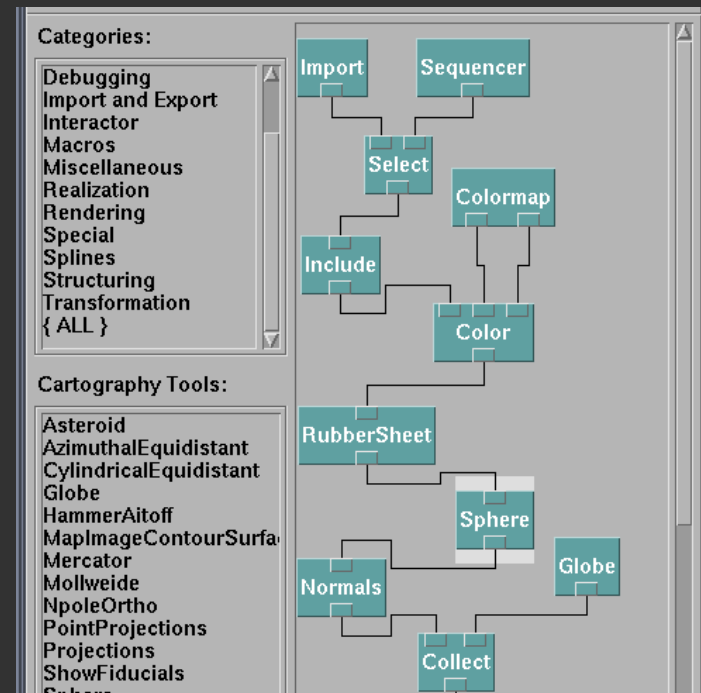


Chart Typology

Pick from a stock of templates
Easy-to-use but limited expressiveness
Prohibits novel designs, new data types



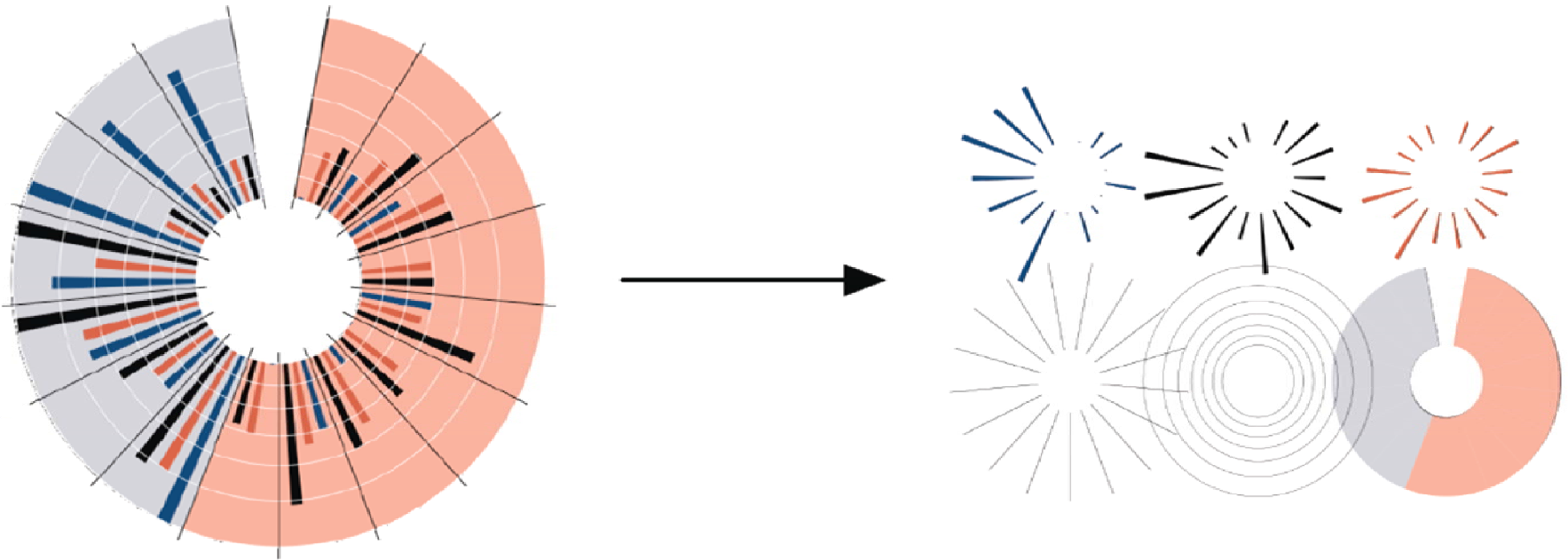
Component Model Architectures

Compose common high-level operations
Permits more combinatorial possibilities
Novel views require new operators, in turn requiring software engineering.

Today's first task is not to invent wholly new [graphical] techniques, though these are needed. Rather we need most vitally to recognize and reorganize the essential of old techniques, to make easy their assembly in new ways, and to modify their external appearances to fit the new opportunities.

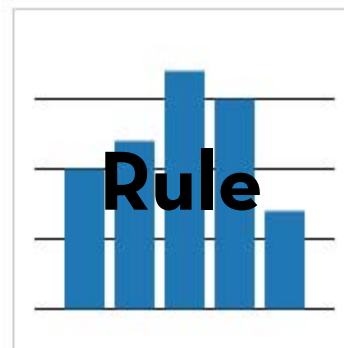
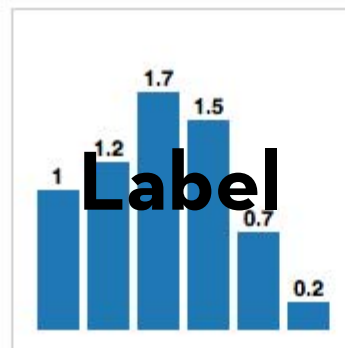
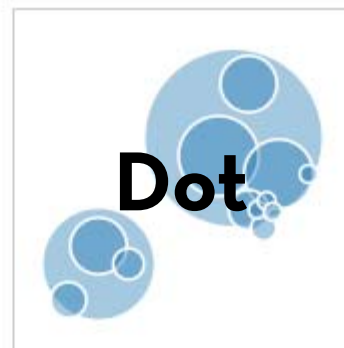
J. W. Tukey, *The Future of Data Analysis*, 1962.

Protovis: A Declarative Language for Visualization



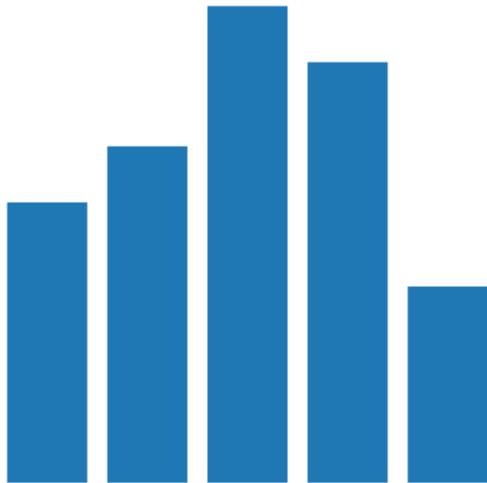
A graphic is a composition of data-representative marks.

with **Mike Bostock** & **Vadim Ogievetsky**



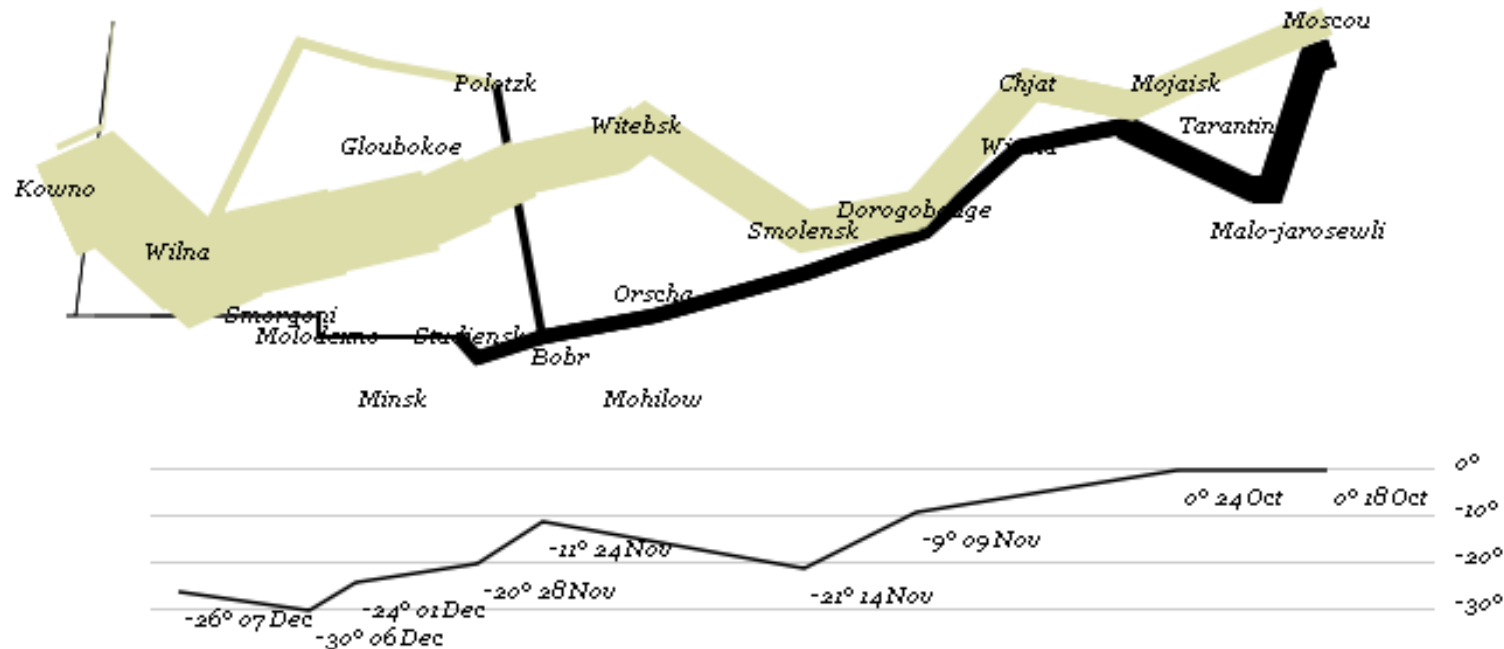
Protovis

Create customized visualizations using a declarative specification language.



```
var vis = new pv.Panel();  
vis.add(pv.Bar)  
  .data([1, 1.2, 1.7, 1.5, .7])  
  .bottom(10)  
  .width(20)  
  .height(function(d) d * 70)  
  .left(function() this.index * 25 + 20);  
vis.render();
```

Protovis (<http://protovis.org>) – Declarative Visualization Specification



```

var army = pv.nest(napoleon.army, "dir", "group");
var vis = new pv.Panel();

var lines = vis.add(pv.Panel).data(army);
lines.add(pv.Line)
  .data(function() army[this.idx])
  .left(lon).top(lat).size(function(d) d.size/8000)
  .strokeStyle(function() color[army[panelIndex][0].dir]);

vis.add(pv.Label).data(napoleon.cities)
  .left(lon).top(lat)
  .text(function(d) d.city).font("italic 10px Georgia")
  .textAlign("center").textBaseline("middle");

```

```

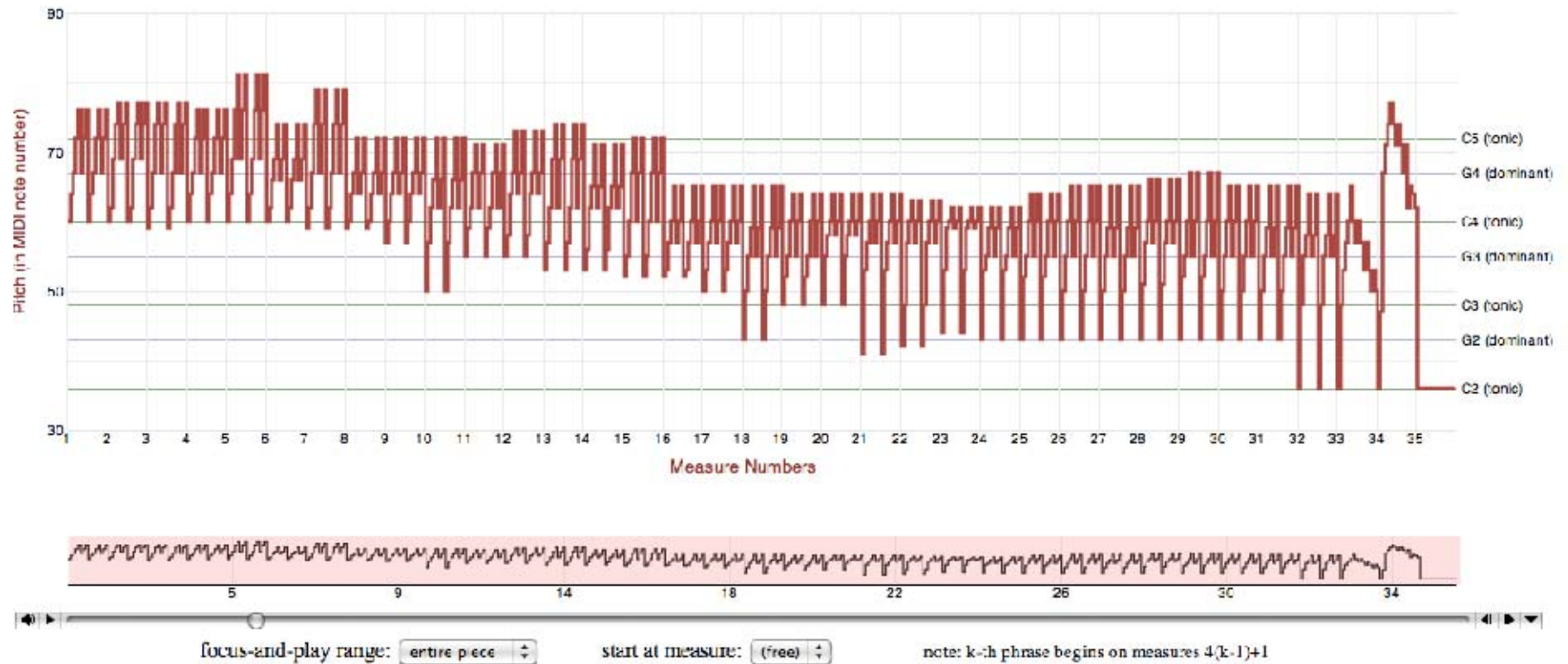
vis.add(pv.Rule).data([0,-10,-20,-30])
  .top(function(d) 300 - 2*d - 0.5).left(200).right(150)
  .lineWidth(1).strokeStyle("#ccc")
  .anchor("right").add(pv.Label)
  .font("italic 10px Georgia")
  .text(function(d) d+"°").textBaseline("center");

vis.add(pv.Line).data(napoleon.temp)
  .left(lon).top(tmp) .strokeStyle("#0")
  .add(pv.Label)
  .top(function(d) 5 + tmp(d))
  .text(function(d) d.temp+"° "+d.date.substr(0,6))
  .textBaseline("top").font("italic 10px Georgia");

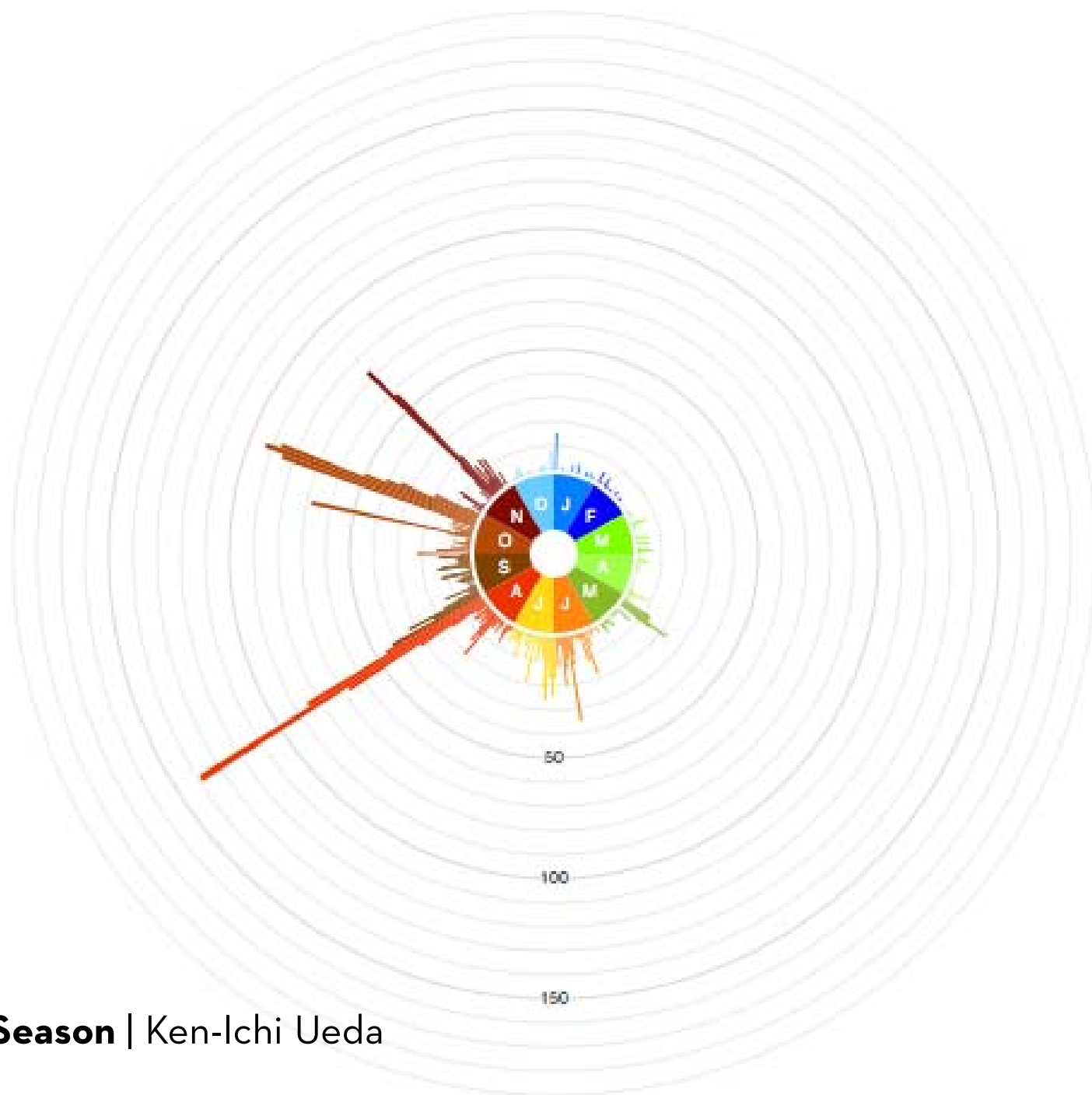
```

PRELUDE NO.1 IN C MAJOR, BWV 846
(FROM WELL-TEMPERED CLAVIER, BOOK 1)

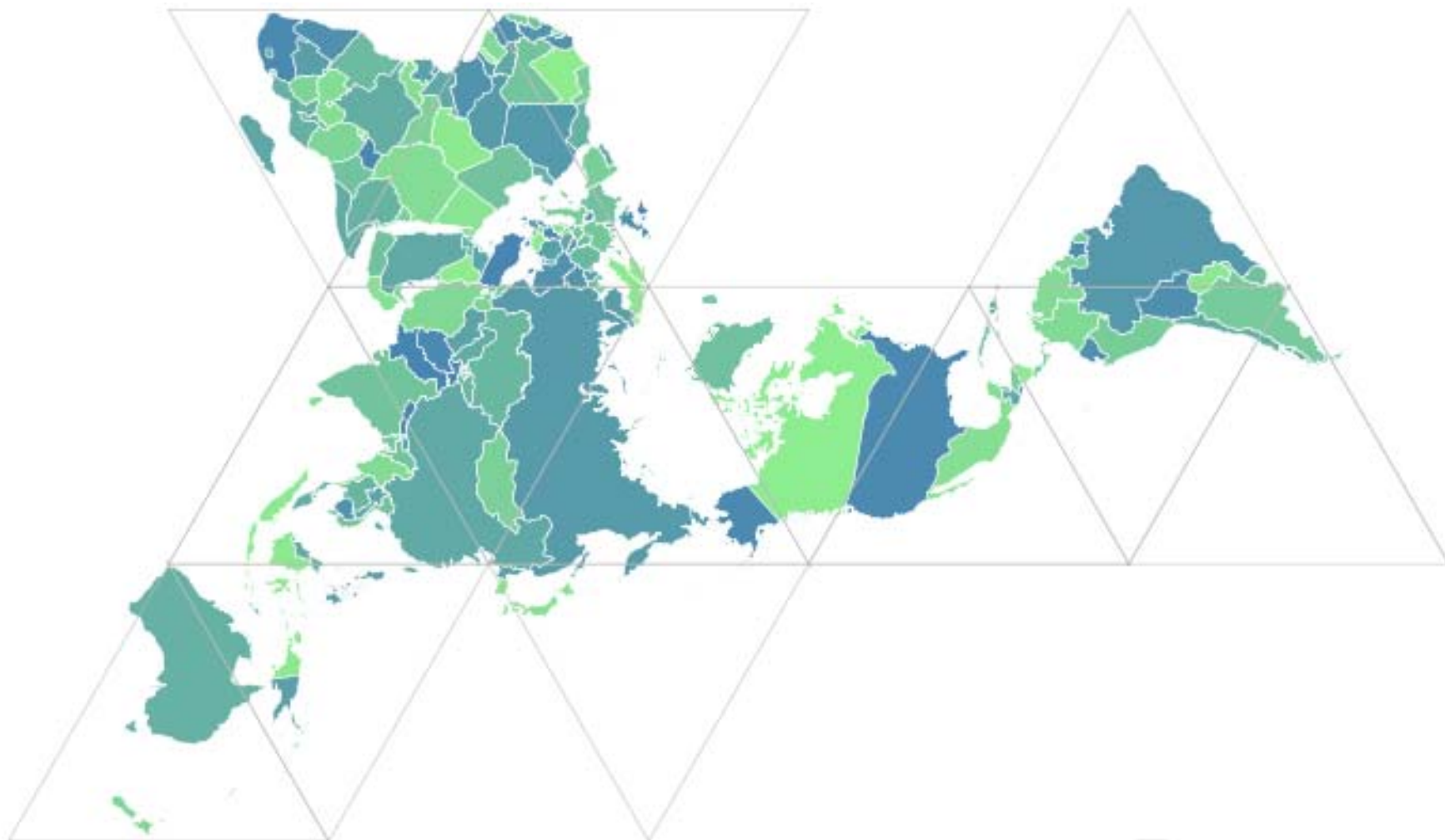
BY J.S. BACH



Bach's Prelude #1 in C Major | Jieun Oh



FlickrSeason | Ken-Ichi Ueda



Dymaxion Maps | Vadim Ogievetsky

Exploiting Declarative Specification

Protovis has led to faster designs, less code

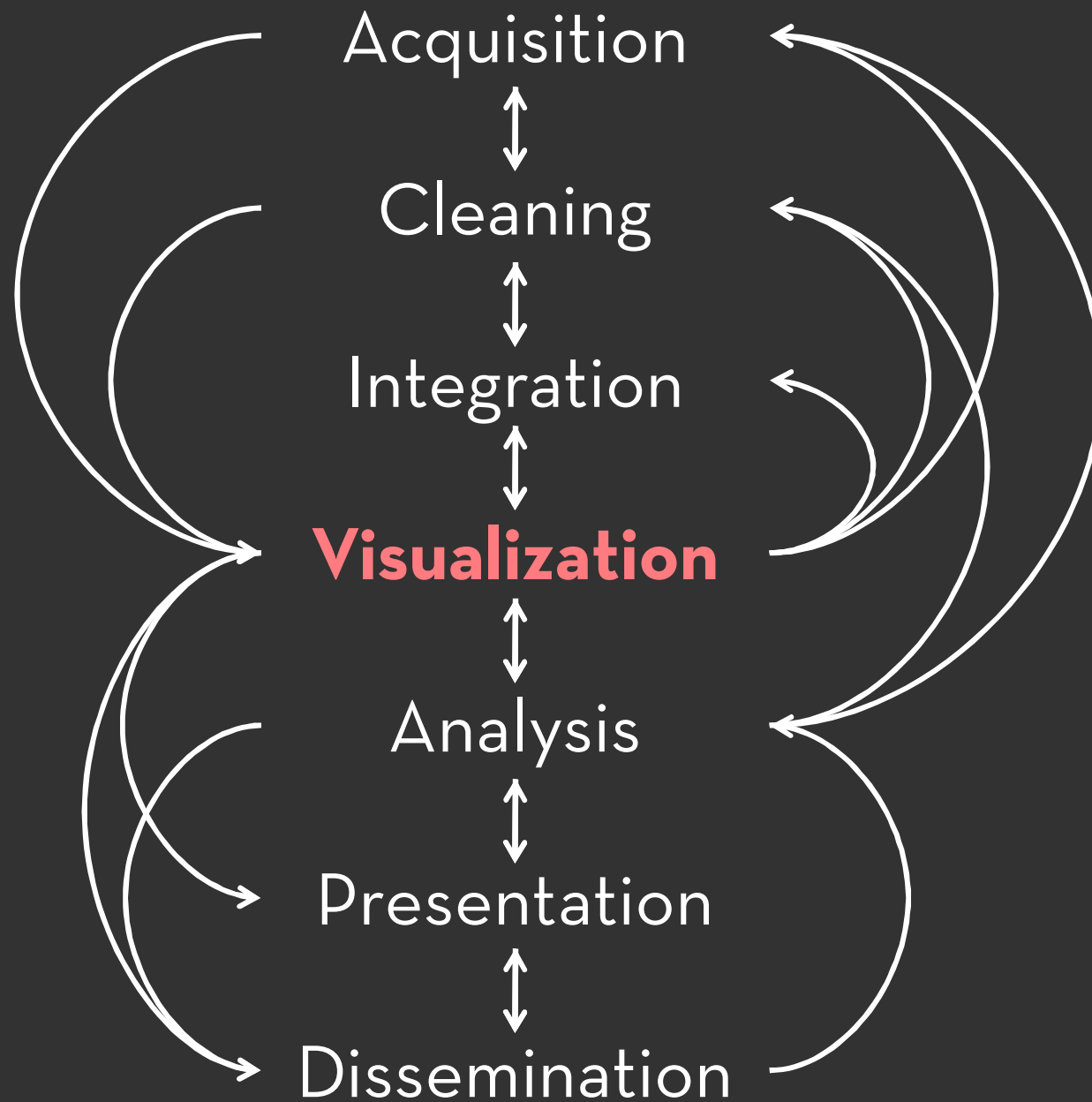
Job Voyager: 5x less code, 10x less dev time

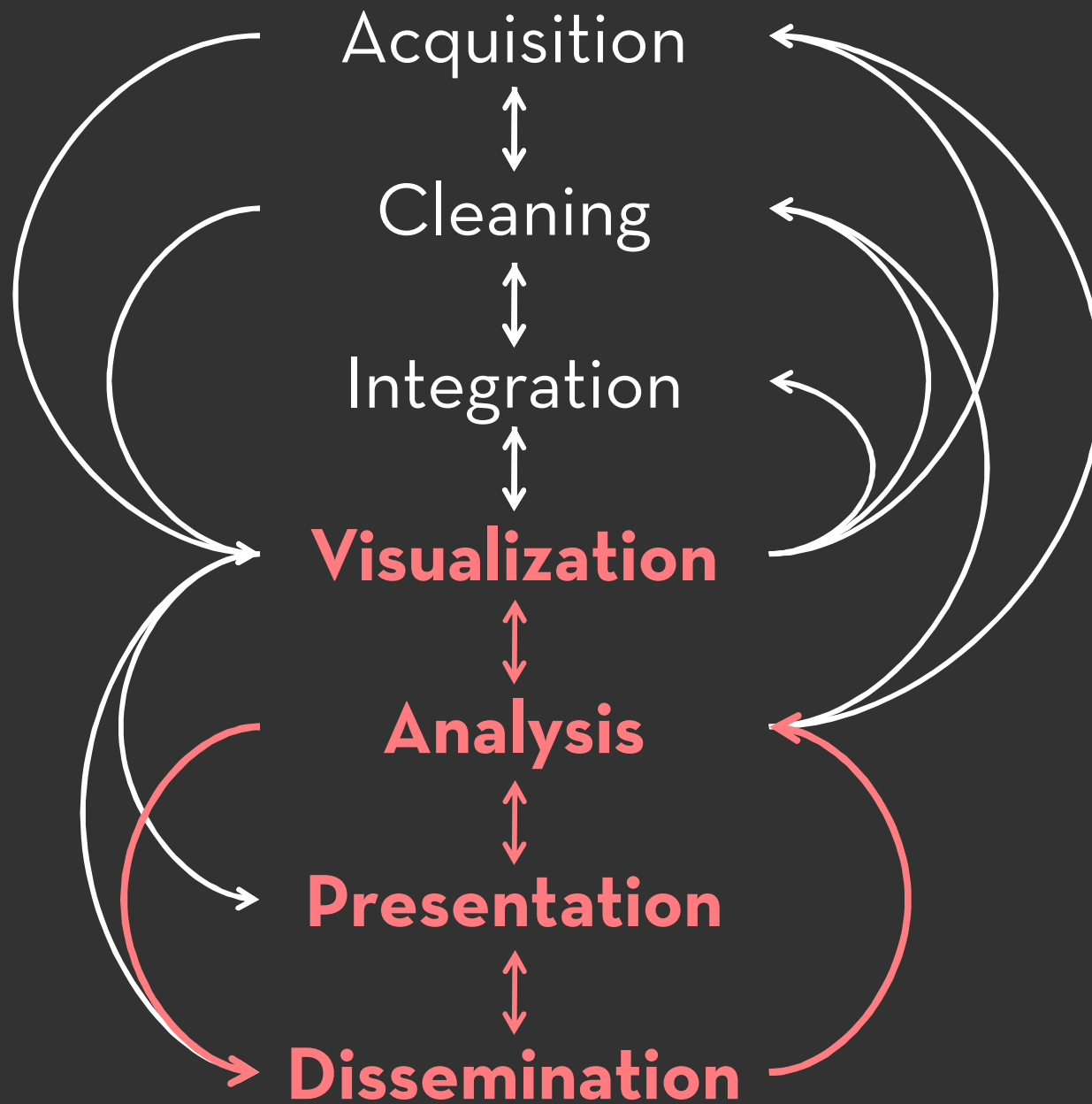
Over 20,000 downloads and widely in use

Multiple implementations: JavaScript & Java

Behind-the-scenes optimization & parallelization

20x scalability over prior systems (in Java)





sense.us

A Web Application for Collaborative
Visualization of Demographic Data

with **Fernanda Viégas** and **Martin Wattenberg**

sense.us > social data visualization - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://sense.us/

sense.us

Log in to add comments! ibm e-mail pass login help

Tags: >>

Tap & hold to see the view

No. of Bookmarks +

Dogear Tags

Dogear Tag Usage, May 2005 to August 2006
source: IBM Dogear
10 comments

TAGS: >>

Tap & hold to see the view

Total Bookmarks +

Dogear People

Dogear Bookmarking by Person, May 2005 to August 2006
source: IBM Dogear
1 comment

>>

Male | Green | Women

% of Work Force +

Job Voyager

Reported Occupations of U.S. Labor Force, 1850-2000
source: <http://ipums.org>
139 comments

>>

Europe +

% of U.S. Population +

Birthplace Voyager

Reported Birthplace of U.S. Residents, 1850-2000
source: <http://ipums.org>
10 comments

Bachelor's degree or higher pct of persons age 25+ 2000

U.S. Census State Map

State Map of 2000-2005 Census Data
source: U.S. Census Bureau
16 comments

Country Map +

1870 Male +

1870 Female +

1910 Male +

1910 Female +

1950 Male +

1950 Female +

1990 Male +

1990 Female +

1990 Person Count +

Population Pyramid

U.S. Population Demographics, 1850-2000
source: <http://ipums.org>
7 comments

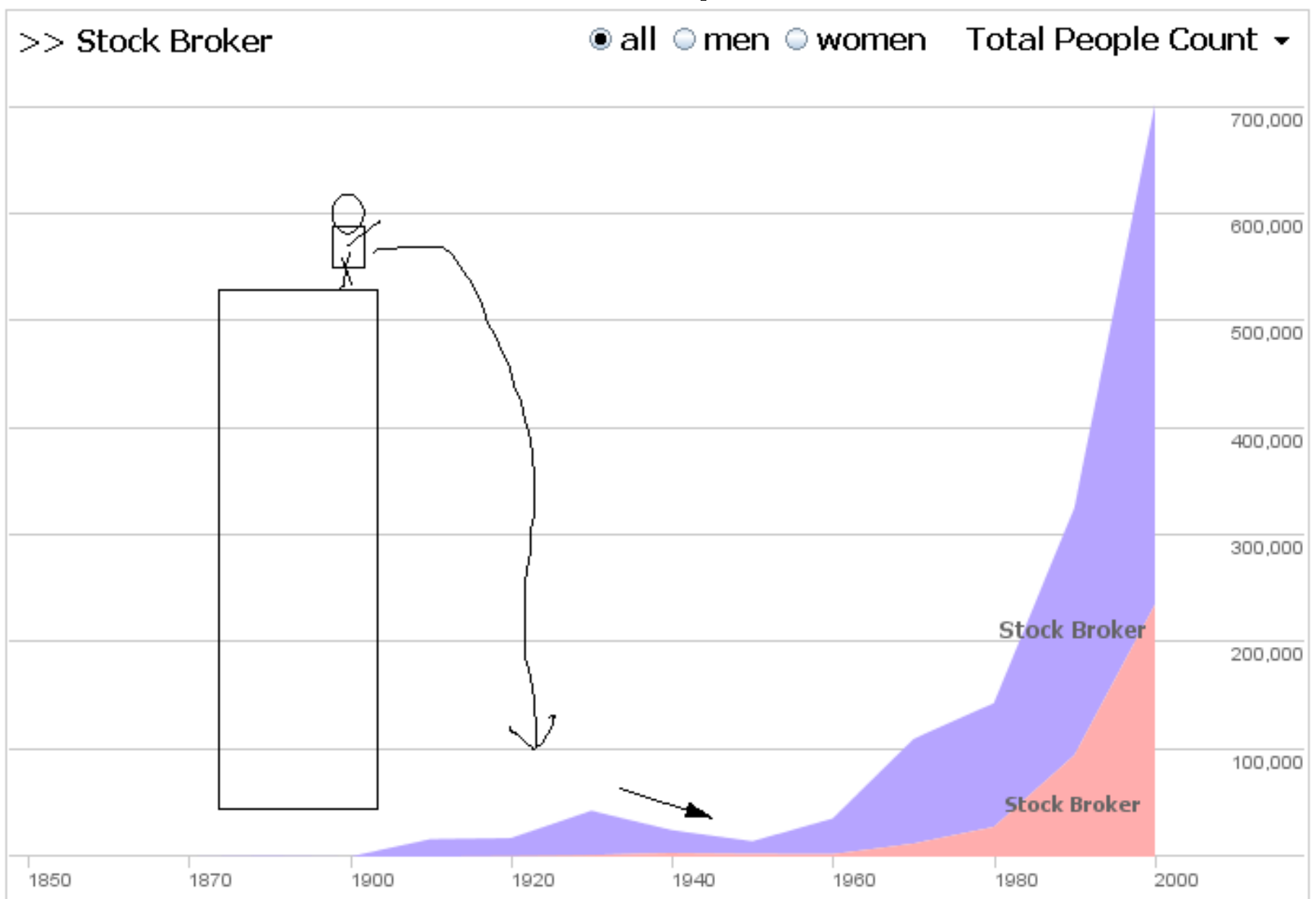
Done

>> mili

● all ○ men ○ women % of Work Force ▼



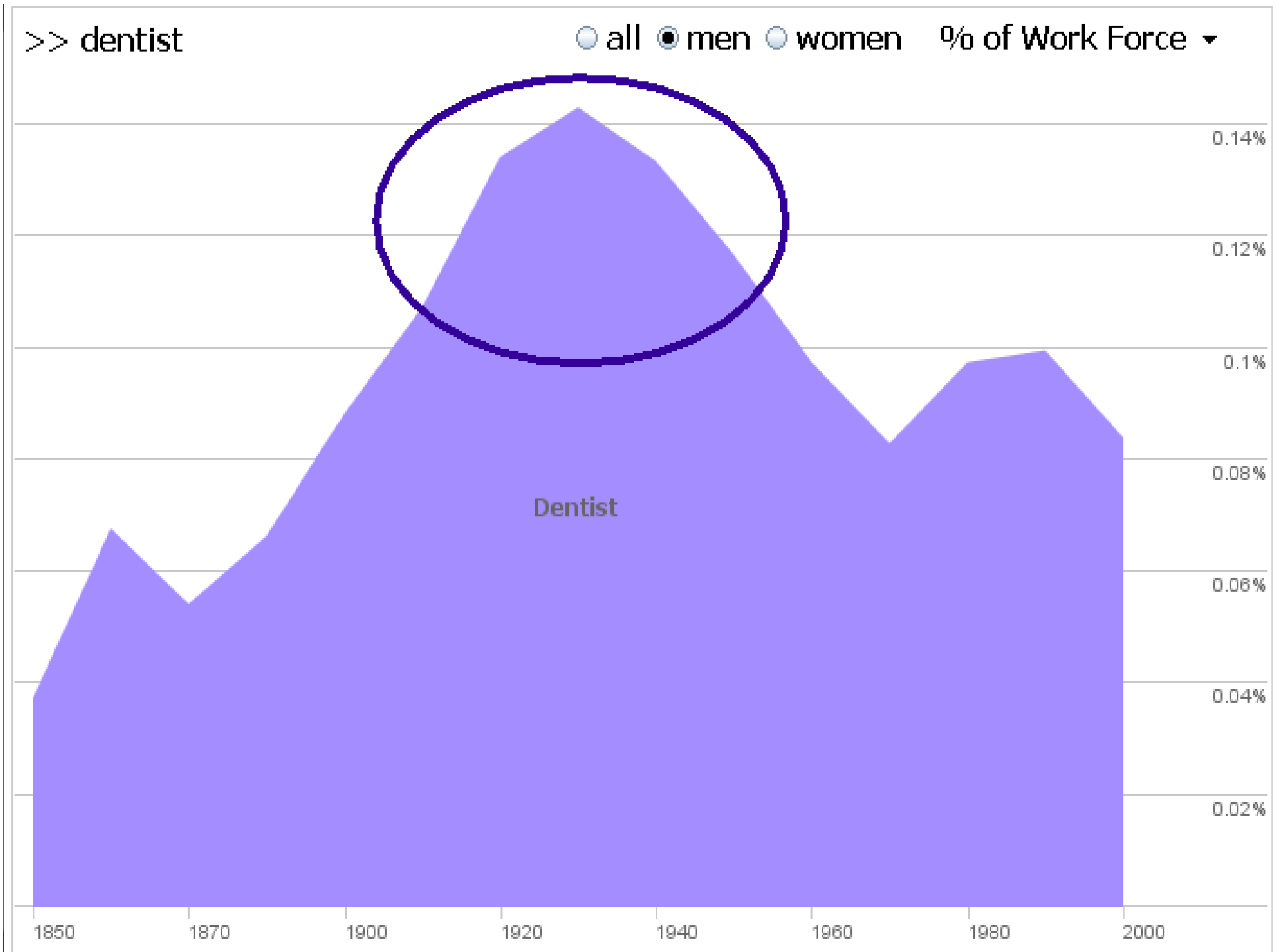
Great depression "killed" a lot of brokers



>> dentist

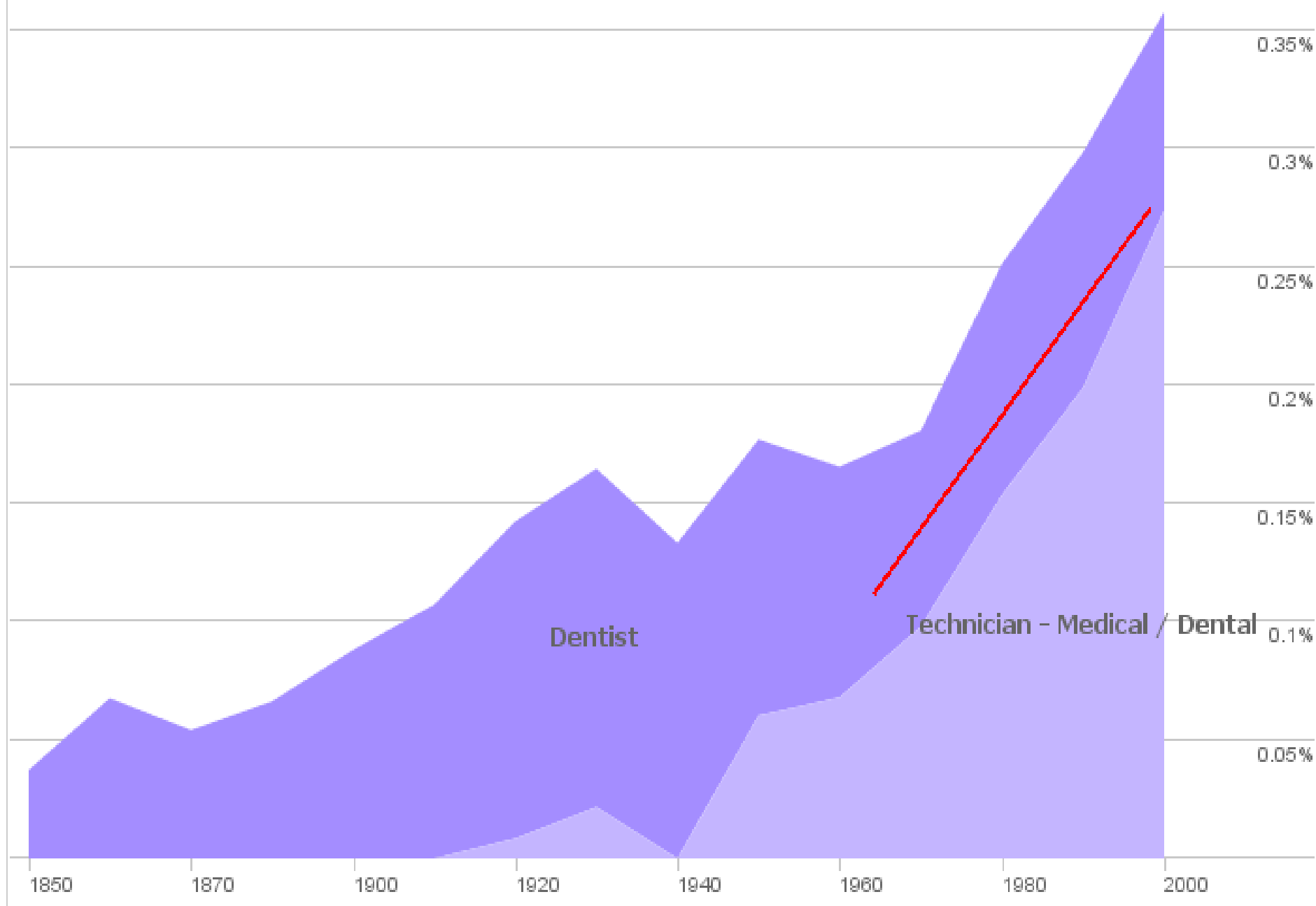
all men women % of Work Force ▼

Dentist



>> dent

all men women % of Work Force ▼



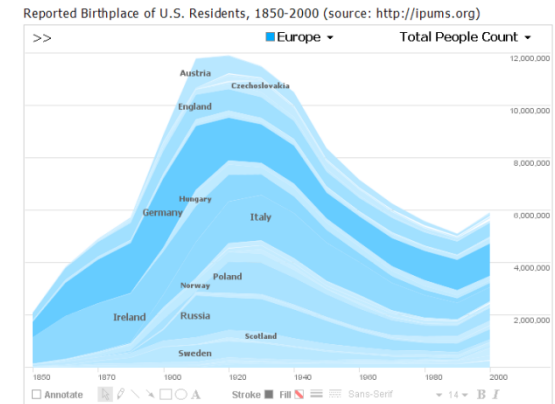
Voyagers and Voyeurs

Complementary faces of analysis

Voyager – focus on visualized data

Active engagement with the data

Serendipitous comment discovery

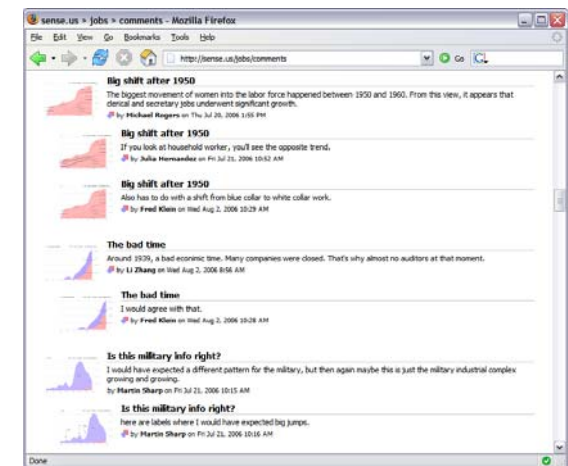


Voyeur – focus on comment listings

Investigate others' explorations

Find people and topics of interest

Catalyze new explorations





many eyes

Sign in

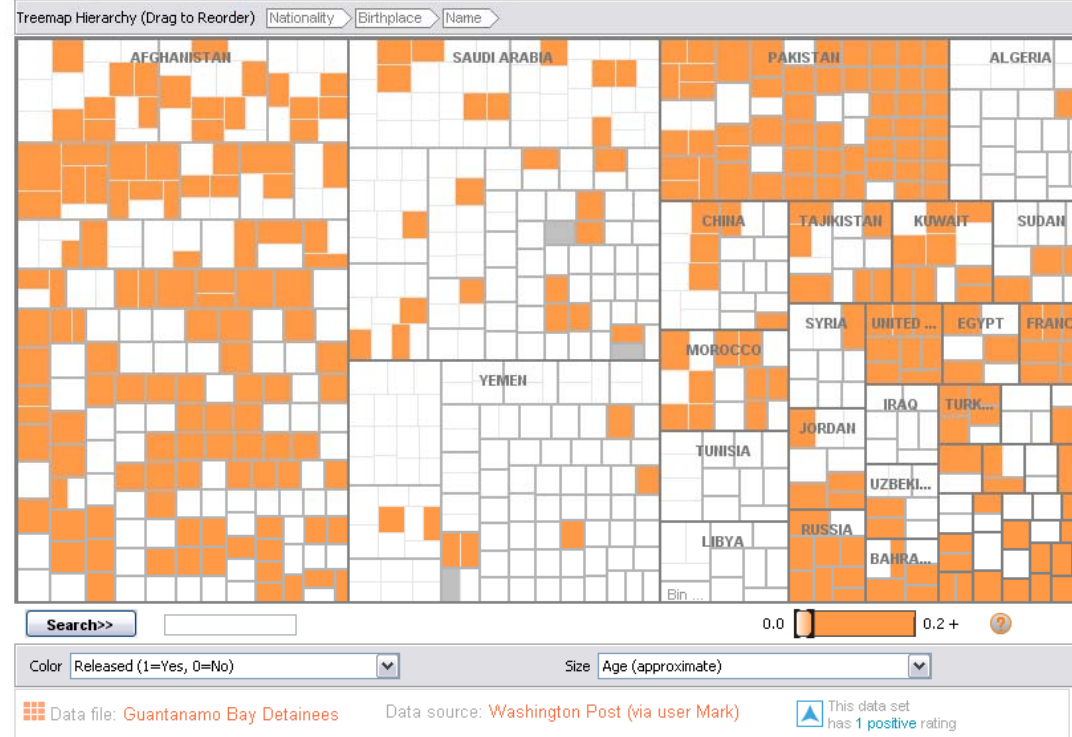
data sets

search

Visualizations : Guantanamo Bay Detainees, release status & age

Can't see the visualization? Download the latest Java plugin [here](#). On Macs: best viewed in Safari.

Created by: [Martin Wattenberg](#) Created on: Saturday February 24, 12:06 PM



share this



watch this



add to topic hub



rate this

Comments (4)



[Martin Wattenberg](#) says:

In this view, orange means released, white means not released. Gray means committed suicide.

Posted Saturday February 24, 12:07 PM

[see view for this comment](#)



[Martin Wattenberg](#) says:

I'd be curious to hear ideas on why various countries have the release proportions that they do.

Posted Saturday February 24, 12:13 PM

[see view for this comment](#)



[Mark](#) says:



This visualization has 1 positive rating



You can add this visualization to a topic hub! [Learn more.](#)



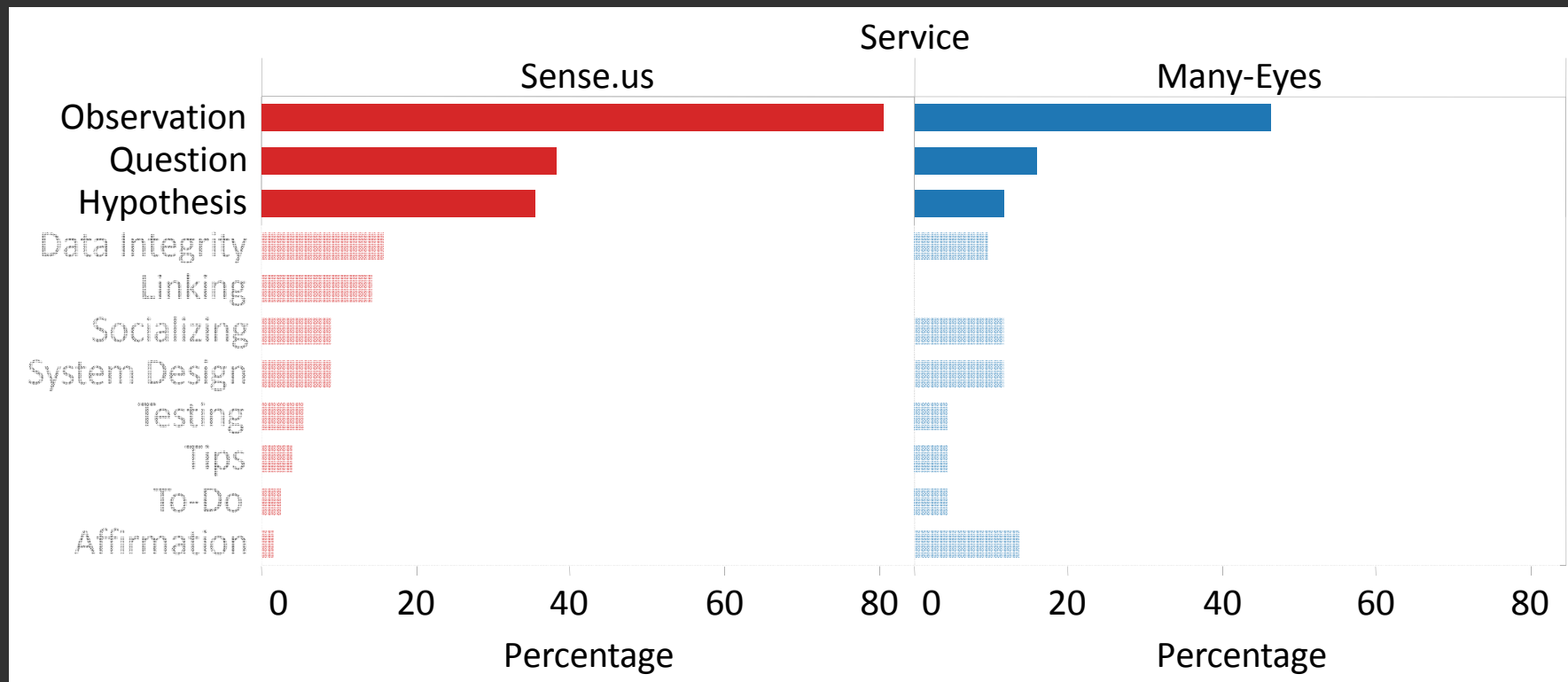
Want to keep track of this visualization? Add it to your watchlist!

Learn more:

[About the Treemap](#)

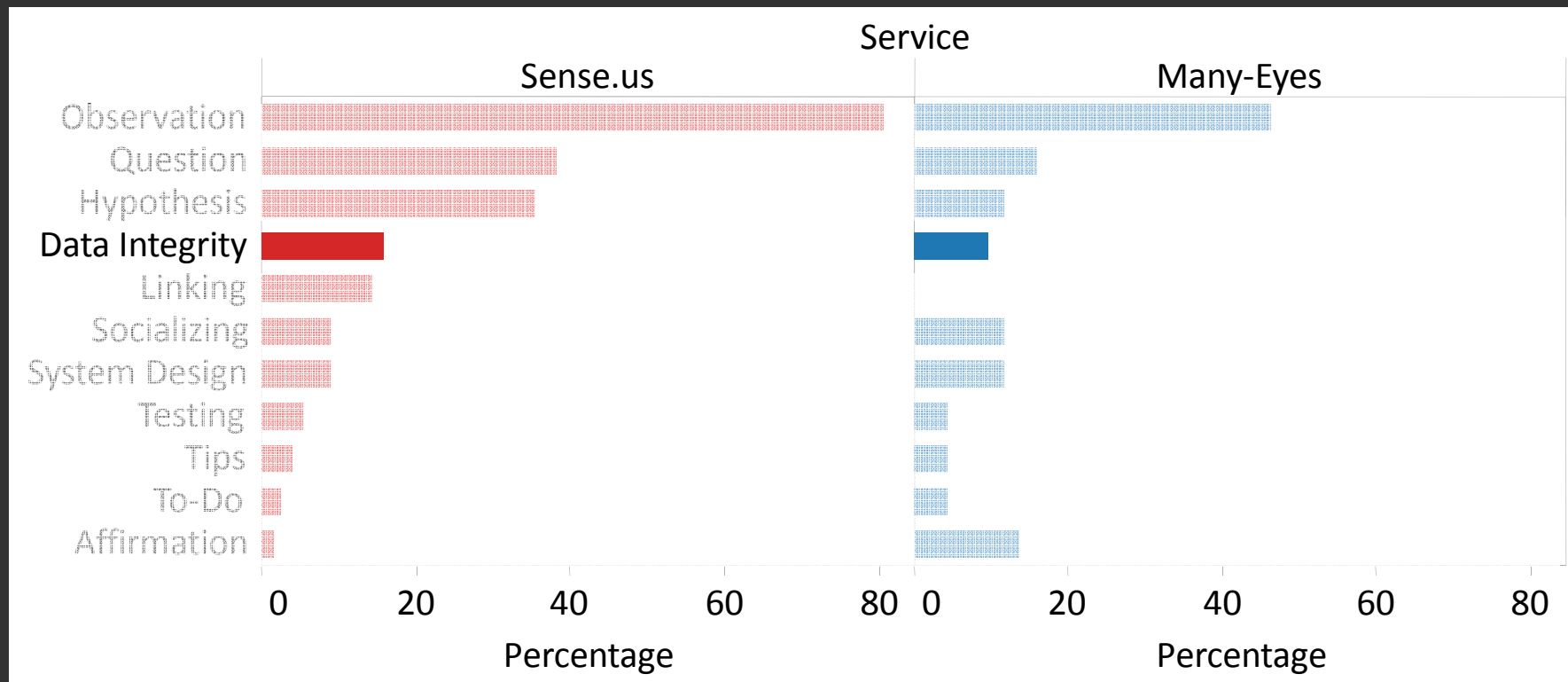
Many-Eyes

Content Analysis of Comments

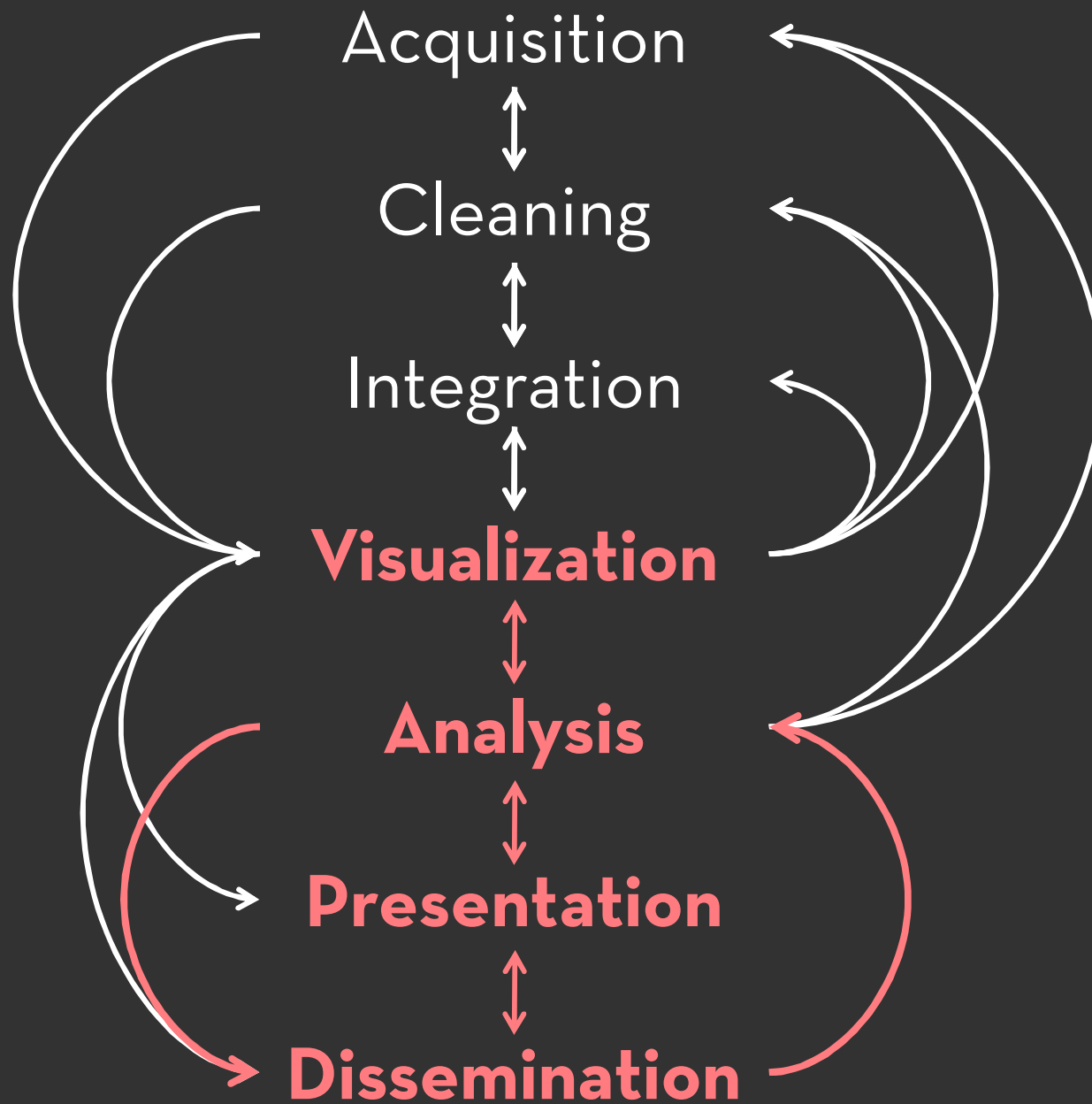


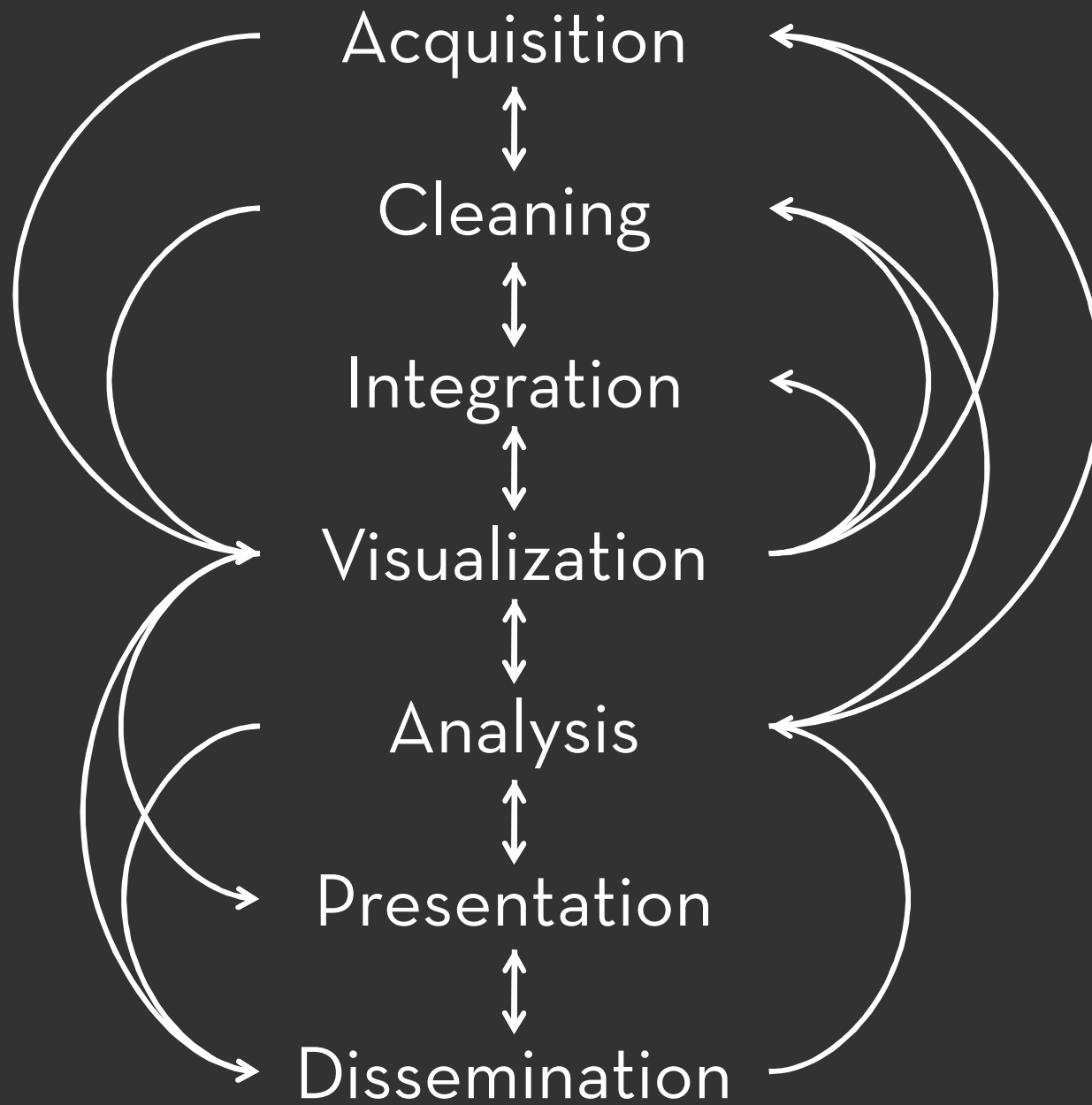
Feature prevalence from content analysis (min Cohen's $\kappa = .74$)
High co-occurrence of Observation, Question, and Hypothesis

Content Analysis of Comments



16% of sense.us comments and **10%** of Many-Eyes comments reference *data integrity* issues.





Students & Collaborators

Mike Bostock

Jason Chuang

Sean Kandel

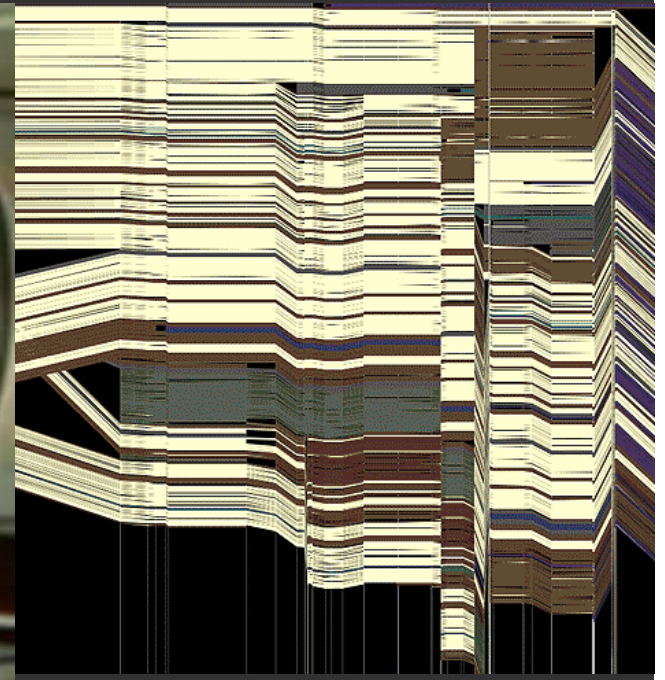
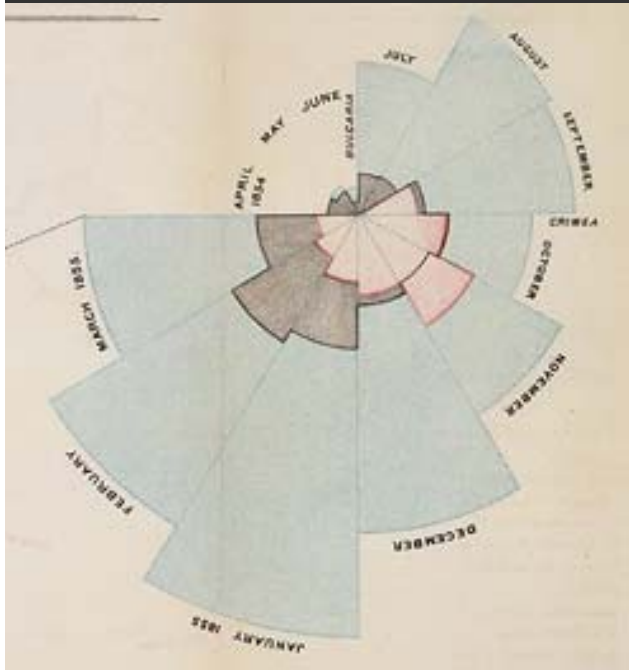
Diana MacLean

Vadim Ogievetsky

Joe Hellerstein, Andreas Paepcke

Fernanda Viégas, Martin Wattenberg

Interactive Tools for Data Transformation & Visualization



Jeffrey Heer <http://vis.stanford.edu>