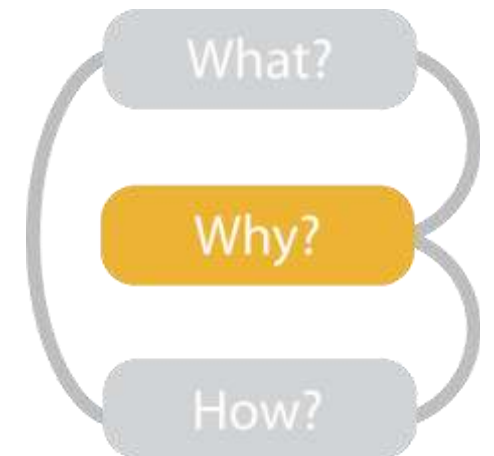


Introduction

Why

Data Visualization

Dr. Claudius Zelenka
Kiel University
cze@informatik.uni-kiel.de



Data Explosion

- Progress in hardware technology allows computers to store an increasing amount of large data
- Computers and the Internet let people consume and produce vast amounts of data
- 2012, Oct 25th: Archive.org „ten petabyte Party“
(<http://blog.archive.org/2012/10/10/the-ten-petabyte-party/>)
 - 10,000,000,000,000,000 bytes

Stunning Data Growth Statistics

<https://techjury.net> : It's possible to see how much data is created every day, as well as how much data we consume regularly. You might be surprised to find out that:

- In 2021, people created [2.5 quintillion bytes](#) of data every day.
- In 2022, [91%](#) of Instagram users **engage with brand videos**.
- In 2022, users send around **650 million Tweets** per day.
- In 2022, [333.2 billion emails](#) are sent every day.

<https://techjury.net/blog/how-much-data-is-created-every-day/#gref>

By 2025, **200+ zettabytes** of data will be in **cloud storage** around the globe.

- **How much data is created every day in 2022? By the end of the year, the world would produce 94 zettabytes of data.** *(Source: Finances Online)*

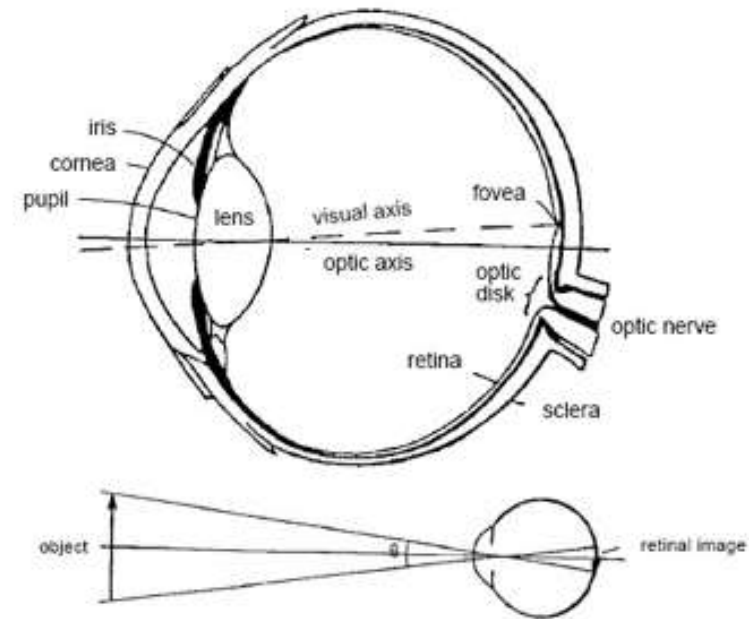
Data Overload

- Data, not information!

Visualization: allow information to be derived from data

- How to transfer information to the user?
- Use human vision:

- Provides high bandwidth sense: human retina can transmit data at roughly 10 million bits per second (Koch et al., 2006) (Ethernet connection: 10 to 1000 million bits per second)
- Can not be used at this theoretical limit for raw data!!!
- Pattern recognition
- Pre-attentive perception
 - *is the subconscious accumulation of information from the environment“*
 - realizing something before you think
- Extends memory and cognitive capacity
- People think visually



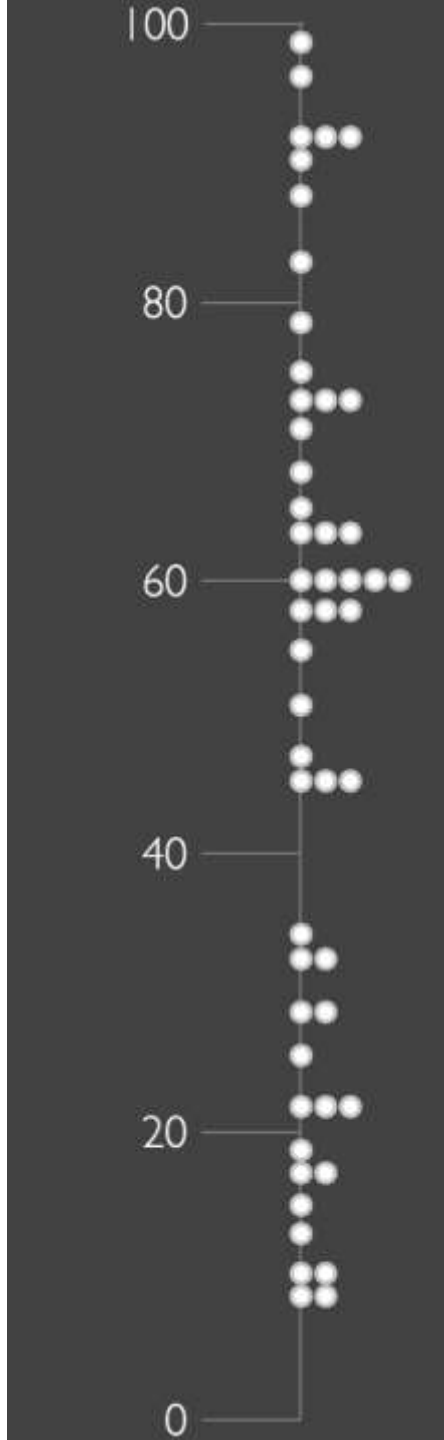
MTHIVLWYADCEQGHKILKMTWYN
ARDCAIREQGHVLMFPSTWYARN
GFPSVCEILQGKMFPSENDRCEQDIFP
SGHLMFHKMVPSTWYACEQTWRN

Count the Vs

MTHI**V**LWYADCEQGHKILKMTWYN
ARDCAIREQGH**L**KMFPSTWYARN
GFPS**V**CEILQGKMFPNSNDRCEQDIFP
SGHLMFHKM**V**PSTWYACEQTWRN

15	19	60
33	11	75
57	34	79
18	51	92
73	22	13
71	60	22
17	10	68
73	18	55
65	46	29
60	73	22
46	92	97
10	58	46
57	17	83
26	99	33
88	92	60
91	29	57
96	12	47

Which number appears
most often?



Which number appears most often?

Anscombe's Quartet

Anscombe's quartet

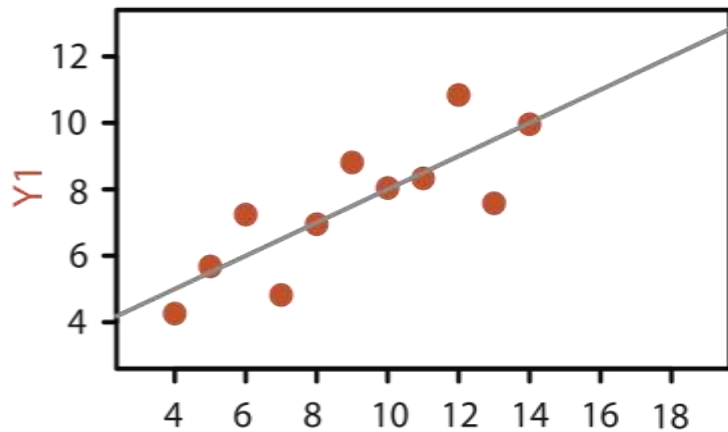
I		II		III		IV	
x	y	x	y	x	y	x	y
10.0	8.04	10.0	9.14	10.0	7.46	8.0	6.58
8.0	6.95	8.0	8.14	8.0	6.77	8.0	5.76
13.0	7.58	13.0	8.74	13.0	12.74	8.0	7.71
9.0	8.81	9.0	8.77	9.0	7.11	8.0	8.84
11.0	8.33	11.0	9.26	11.0	7.81	8.0	8.47
14.0	9.96	14.0	8.10	14.0	8.84	8.0	7.04
6.0	7.24	6.0	6.13	6.0	6.08	8.0	5.25
4.0	4.26	4.0	3.10	4.0	5.39	19.0	12.50
12.0	10.84	12.0	9.13	12.0	8.15	8.0	5.56
7.0	4.82	7.0	7.26	7.0	6.42	8.0	7.91
5.0	5.68	5.0	4.74	5.0	5.73	8.0	6.89

Identical statistics

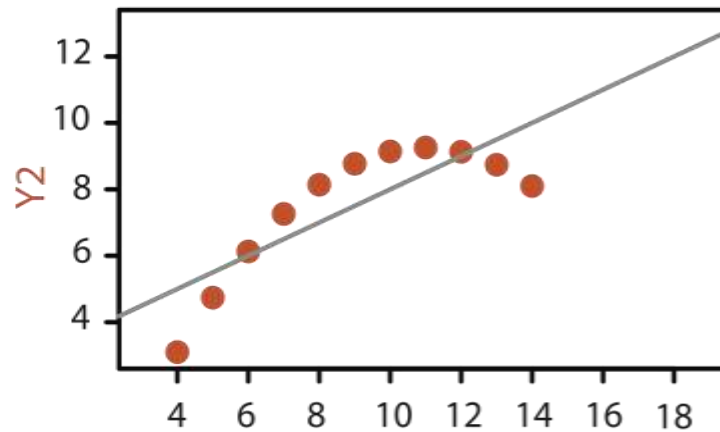
x mean	9
x variance	10
y mean	7.5
y variance	3.75
x/y correlation	0.816

Frank Anscombe, "Graphs in Statistical Analysis"

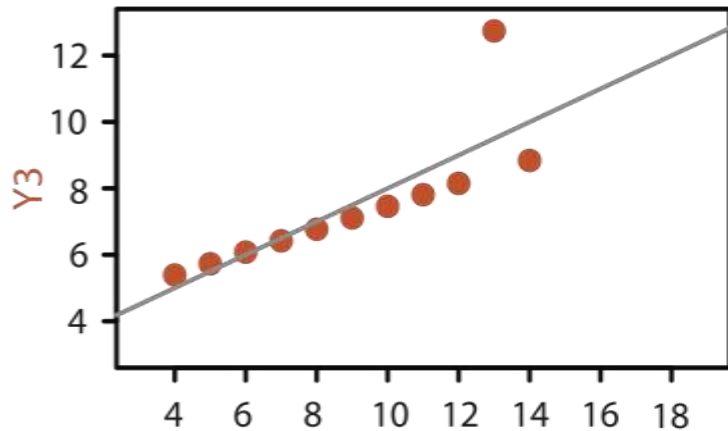
Anscombe's Quartet



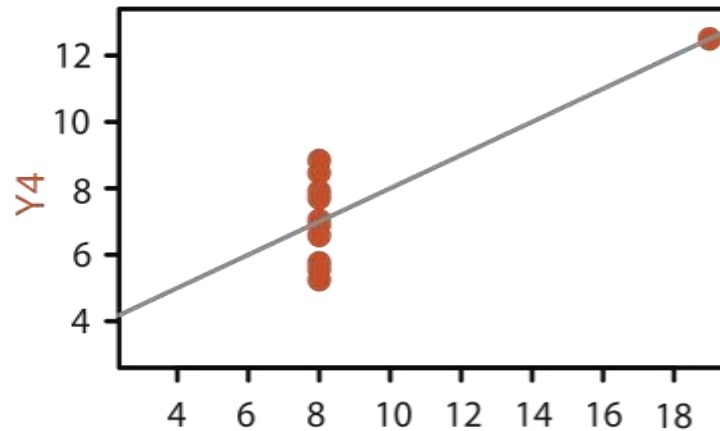
X1



X2



X3



X4

Anscombe's Quartet

Identical statistics

x mean	9
x variance	10
y mean	7.5
y variance	3.75
x/y correlation	0.816

Why vision?

- human visual system is high-bandwidth channel to brain
 - overview possible due to background processing
 - subjective experience of seeing everything simultaneously
 - significant processing occurs in parallel and pre-attentively
- sound: lower bandwidth and different semantics
 - overview not supported
 - subjective experience of sequential stream
- touch/haptics: impoverished record/replay capacity
 - only very low-bandwidth communication thus far
- taste, smell: no viable record/replay devices

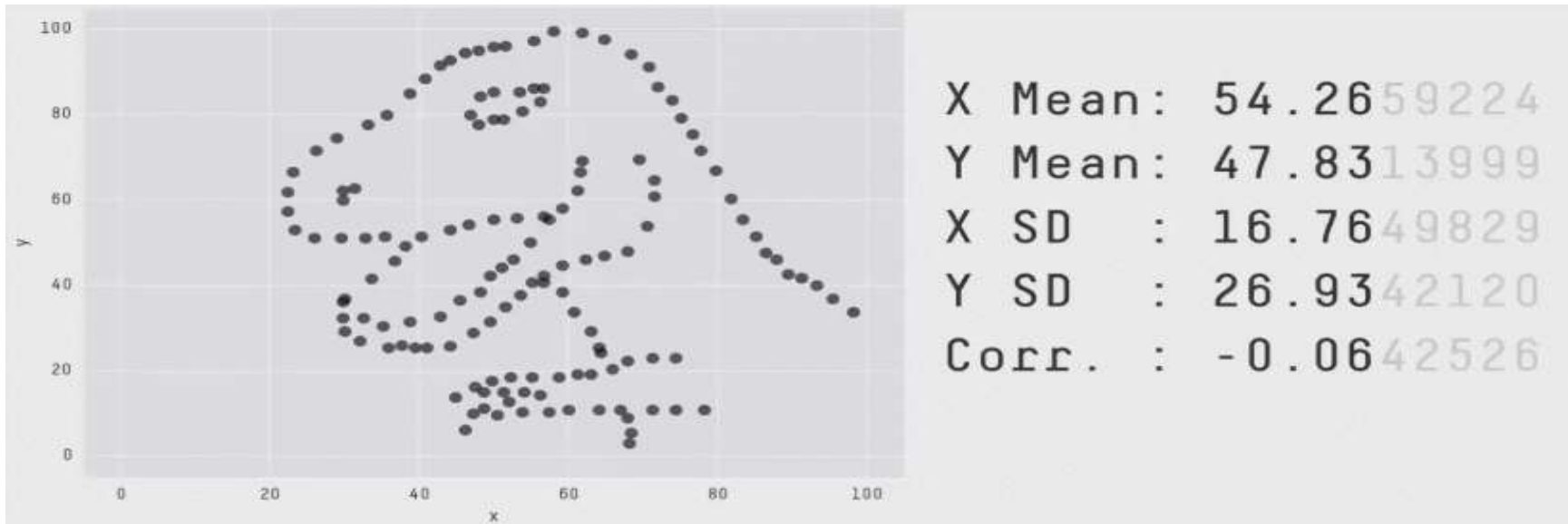
What resource limitations are we faced with?

- computational limits
 - computation time, system memory
- display limits
 - pixels are precious & most constrained resource
 - **information density**: ratio of space used to encode info vs unused whitespace
 - tradeoff between clutter and wasting space
 - find sweet spot between dense and sparse
- human limits
 - human time, human memory, human attention

Anscombe's Quartet

- summaries lose information, details matter
 - confirm expected and find unexpected patterns
 - assess validity of statistical model

Exercise: Dinosaur Dozen



Defining Visualization

- Computer-based visualization systems provide visual representations of **datasets** designed to help **people** carry out tasks more effectively
 - *Tamara Munzner, University of British Columbia*

Alternative Definitions

- "The use of computer-supported, interactive, visual representations of abstract data to amplify cognition" - *Stuart Card*
- an accessible way to see and understand trends, outliers, and patterns in data." - *Tableau*
- "Data visualization is the creation and study of the visual representation of data" - wiki

What is data visualization?

- “Data visualization is the creation and study of the visual representation of data” - wiki
- Input: **data** Output: **visual form** Goal: **insight**



<http://paulbutler.org/archives/visualizing-facebook-friends/>

Data Visualization – part art / part science

- Data visualization is **part art** and **part science**. The challenge is to get the art right without getting the science wrong and vice versa. A **data visualization first and foremost has to accurately convey the data**. It must not mislead or distort. If one number is twice as large as another, but in the visualization they look to be about the same, then the visualization is wrong. At the same time, **a data visualization should be aesthetically pleasing**. Good visual presentations tend to enhance the message of the visualization. If a figure contains jarring colors, imbalanced visual elements, or other features that distract, then the viewer will find it harder to inspect the figure and interpret it correctly.

Claus O. Wilke - Fundamentals of Data Visualization

<https://clauswilke.com/dataviz/introduction.html>

Defining Concepts

- computer supported
- Interactive visual representations
- For Abstract data
- Helping people
- to see and understand
- trends, outliers, and patterns in data,
- and carry out tasks
- more effectively
- Through amplifying cognition

Goals of visualization

Presentation (4,000 years)

- starting point: facts to be presented
- goal: visualization which makes the facts apparent
- “You do not really understand something unless you can explain it to your grandmother.” (Albert Einstein ?)

Confirmative analysis (200 years)

- starting point: hypothesis about the data
- goal: confirmation or rejection of the hypothesis

Explorative analysis (\approx 20 years)

- starting point: no hypothesis about the data
- goal: hypothesis about the data

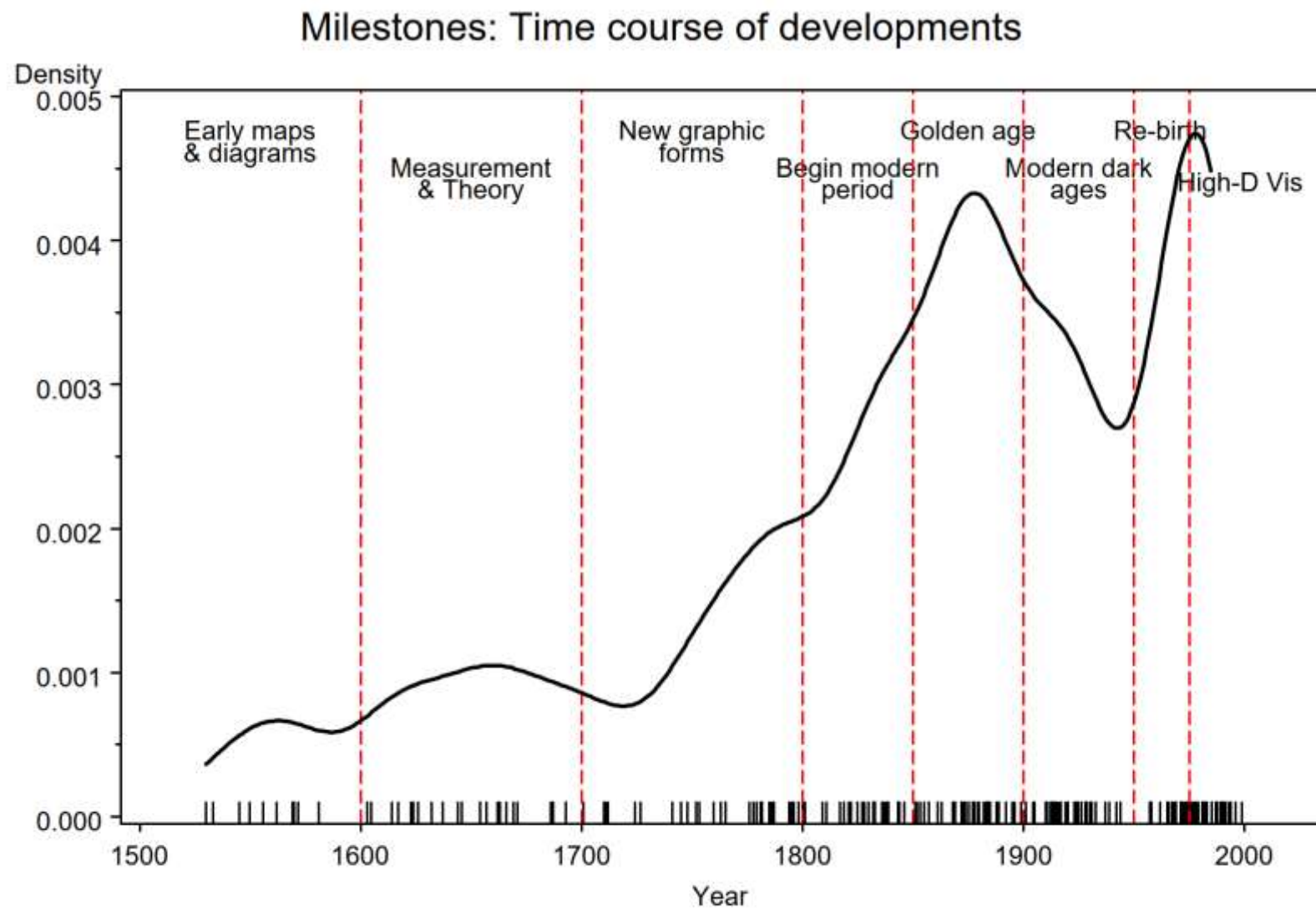


Figure 1: The time distribution of events considered milestones in the history of data visualization, shown by a rug plot and density estimate.

A Brief History of Data Visualization

Michael Friendly, in: Handbook of Computational Statistics: Data Visualization

Visual analytics vs. Data mining

- let computers do what computers are good at

Data mining

- let humans do what they're good at

Visual analytics

- Data mining focuses more on automatic algorithms
- Visualization keeps human in the loop and focuses more on interactive analysis
- **Visualization is suitable when there is a need to augment human capabilities rather than replace people with computational decision-making methods. (Tamara Munzner)**

Analysts: Data visualization tools key to 'big data' analytics success

Mark Brunelli, Senior News Editor

Published: 30 Nov 2011



Demand for [data visualization tools](#) is rising sharply, partly as a result of more companies seeking to gain valuable business insights through “big data” analytics initiatives. But achieving success with data visualization often requires fresh thinking about how to present information to business users, especially in big-data environments, according to data management analysts.

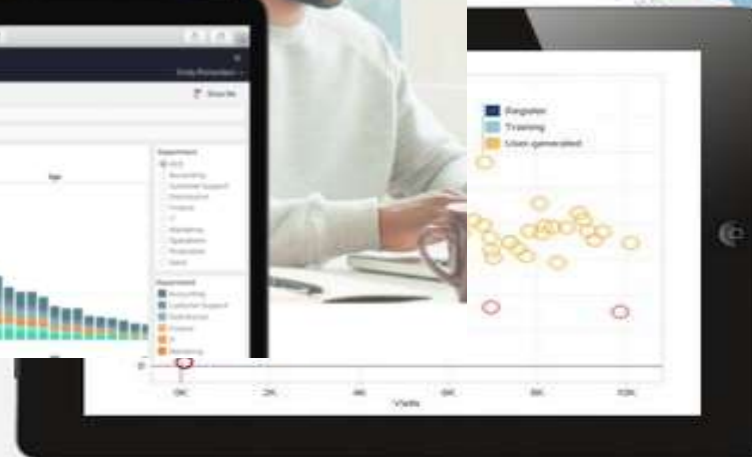
Visualization industry and business applications - Tableau software



Changing the way you think about data

THE TABLEAU PLATFORM

SEE IT IN ACTION



Make analytics easy.
Analysts, executives, IT, everyone.



THE WALL STREET JOURNAL.
 ASIA EDITION Monday, February 24, 2014 As of 10:08 AM rd

Home World Asia China India Japan SE Asia Business Markets Tech

TOP STORIES IN MARKETS

1 of 12 Big Oil Returns to Returns

2 of 12 Groupon Clips Its Own Wings

AXA Grows

WSJ BLOGS

Deal Journal

An up-to-the-minute take on deals and deal makers:

April 3, 2013, 10:57 AM

Tableau Software Plots Latest Big Data IPO

- Big Data Analytics Specialist Tableau Software Raises \$254M In IPO, Shares Close 64% Up; Marketo's First Day Up 78% To \$23.10
- Tableau was acquired by Salesforce for \$15.7B on Jun 10, 2019

Tableau Software, Inc.



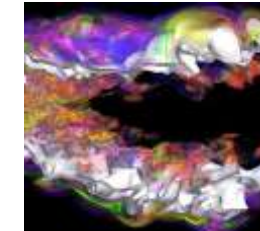
Type	Subsidiary
Industry	Software
Founded	Mountain View, California, U.S. (2003)
Founders	Christian Chabot Chris Stolte Andrew Beers Pat Hanrahan
Successor	Salesforce
Headquarters	Seattle, Washington, U.S.
Key people	Mark Nelson (CEO) Christian Chabot (Chairman)
Products	Business intelligence Data visualization Analytics
Revenue	877,000,000 United States dollar (2017)
Net income	5,873,000 United States dollar (2014)
Number of employees	4,181 (2019) 40

Scientific Subfields

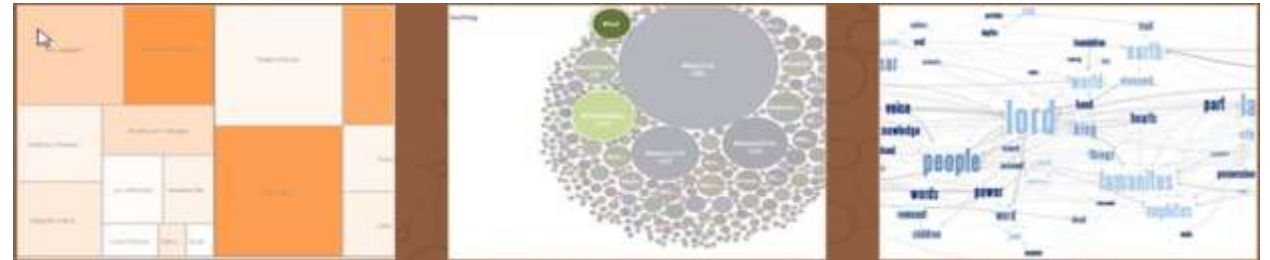
- Scientific Visualization (SciVis) – Spatial data



Analytical reasoning



- Information Visualization (InfoVis)– Abstract data



- Visual Analytics (VAST) –



Explainable AI and Visual Analytics

AI: non-linear function with many parameters from input to output.

Text - > Category

Text classification

Image - > Category

Image classification

Text -> Text

Translation/ Summary / Text completion

Image - > Text

Image Captioning

Text -> Image

Image Generation

How to make a function with 200 billion parameters understandable?

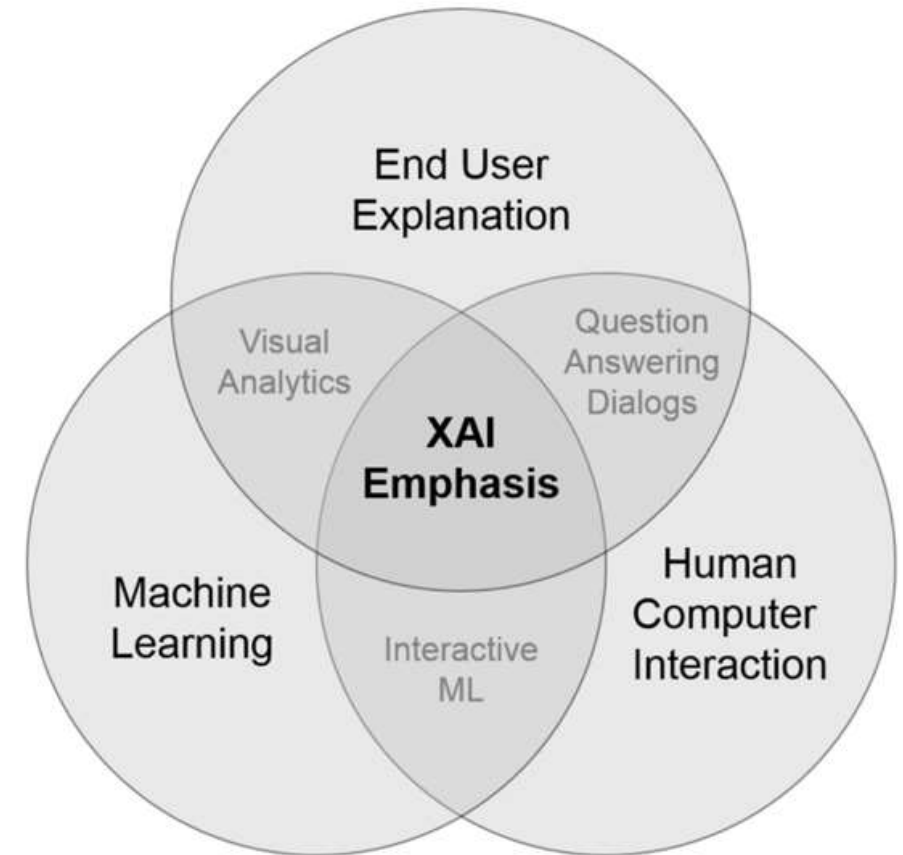
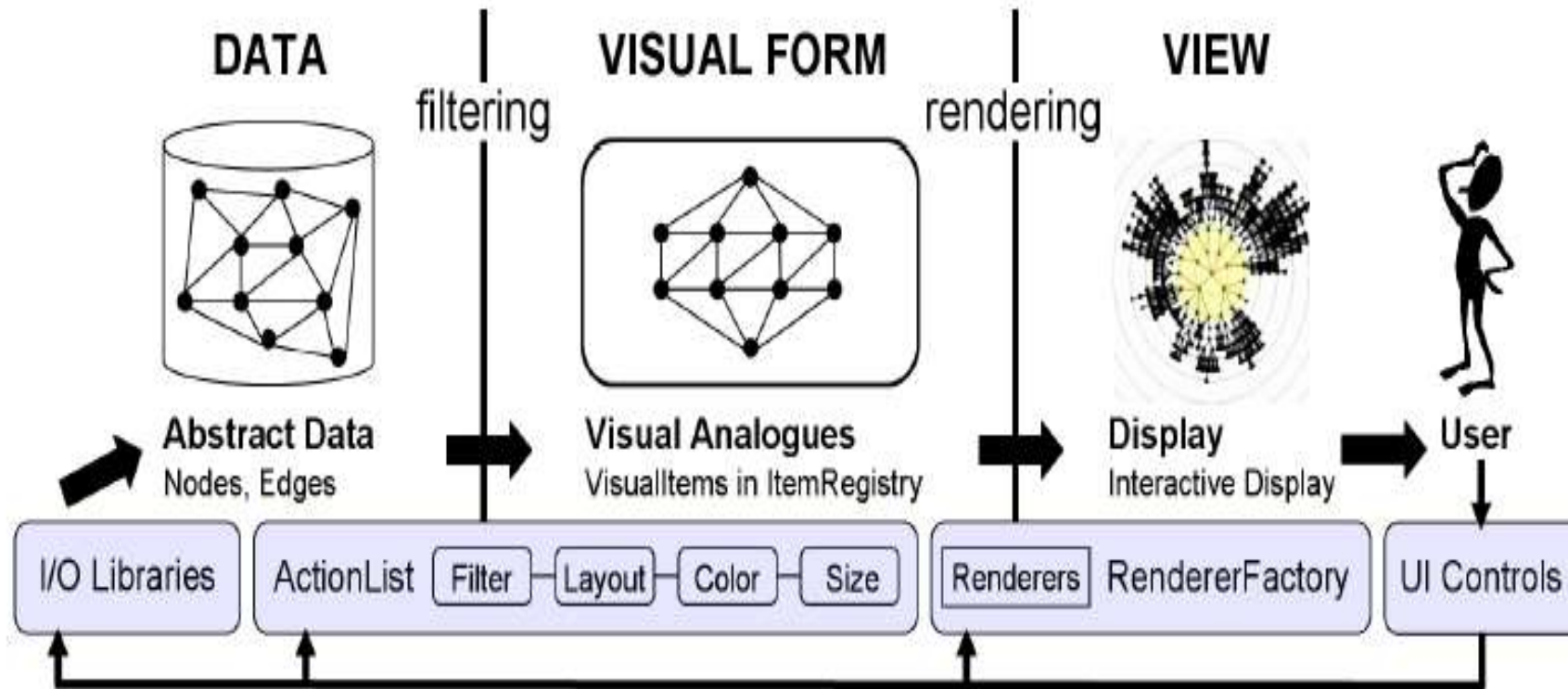
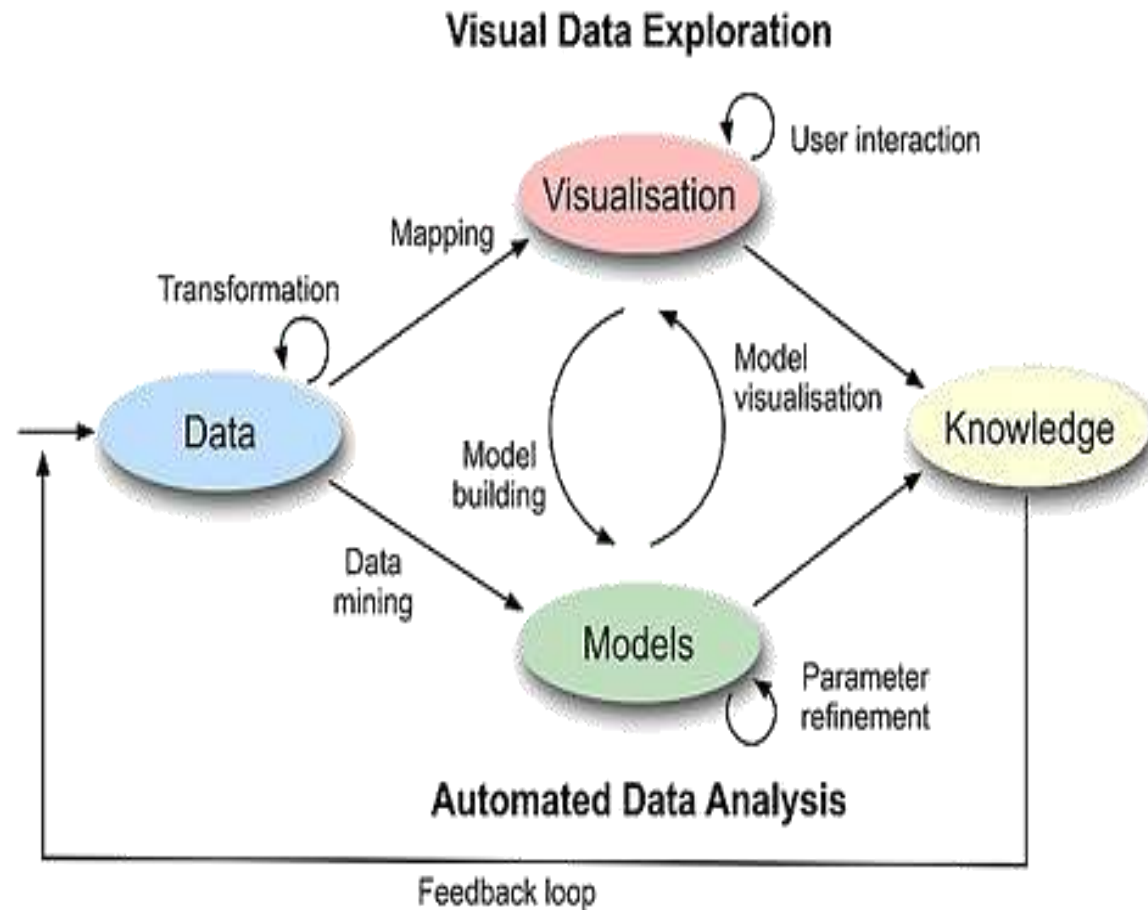


Figure 2: XAI Emphasis

Visualization Reference Model: 1990s

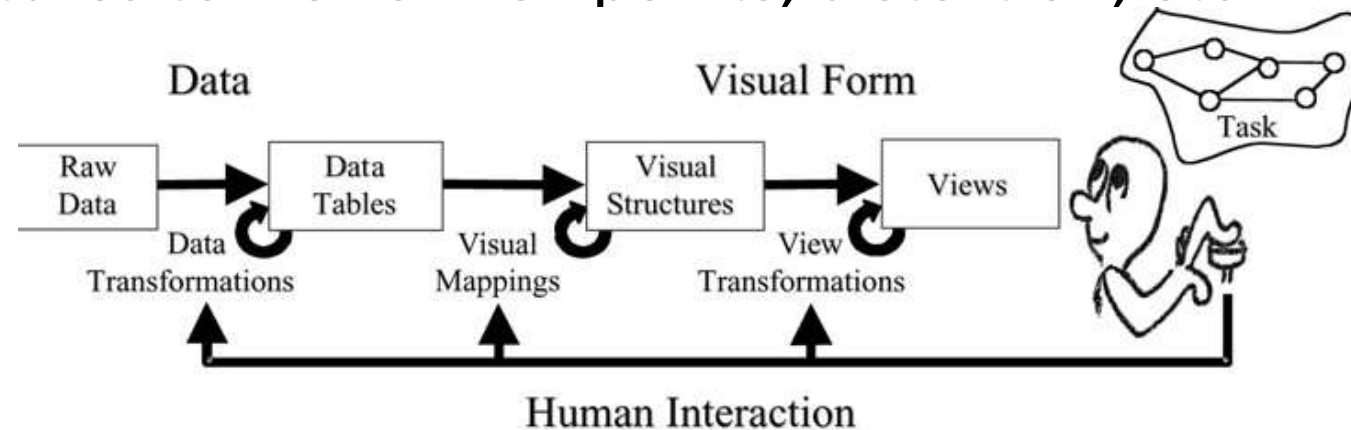


Visualization Reference Model: 2000s



Reference Models

- Raw table to data table: filtering, data cleaning
- Data table to visual structures: pick mappings
- Visual structures to views: viewpoints, distortion, etc.



Raw Data: idiosyncratic formats
Data Tables: relations (cases by variables) + meta-data
Visual Structures: spatial substrates + marks + graphical properties
Views: graphical parameters (position, scaling, clipping, ...)

Explainable AI and Visual Analytics

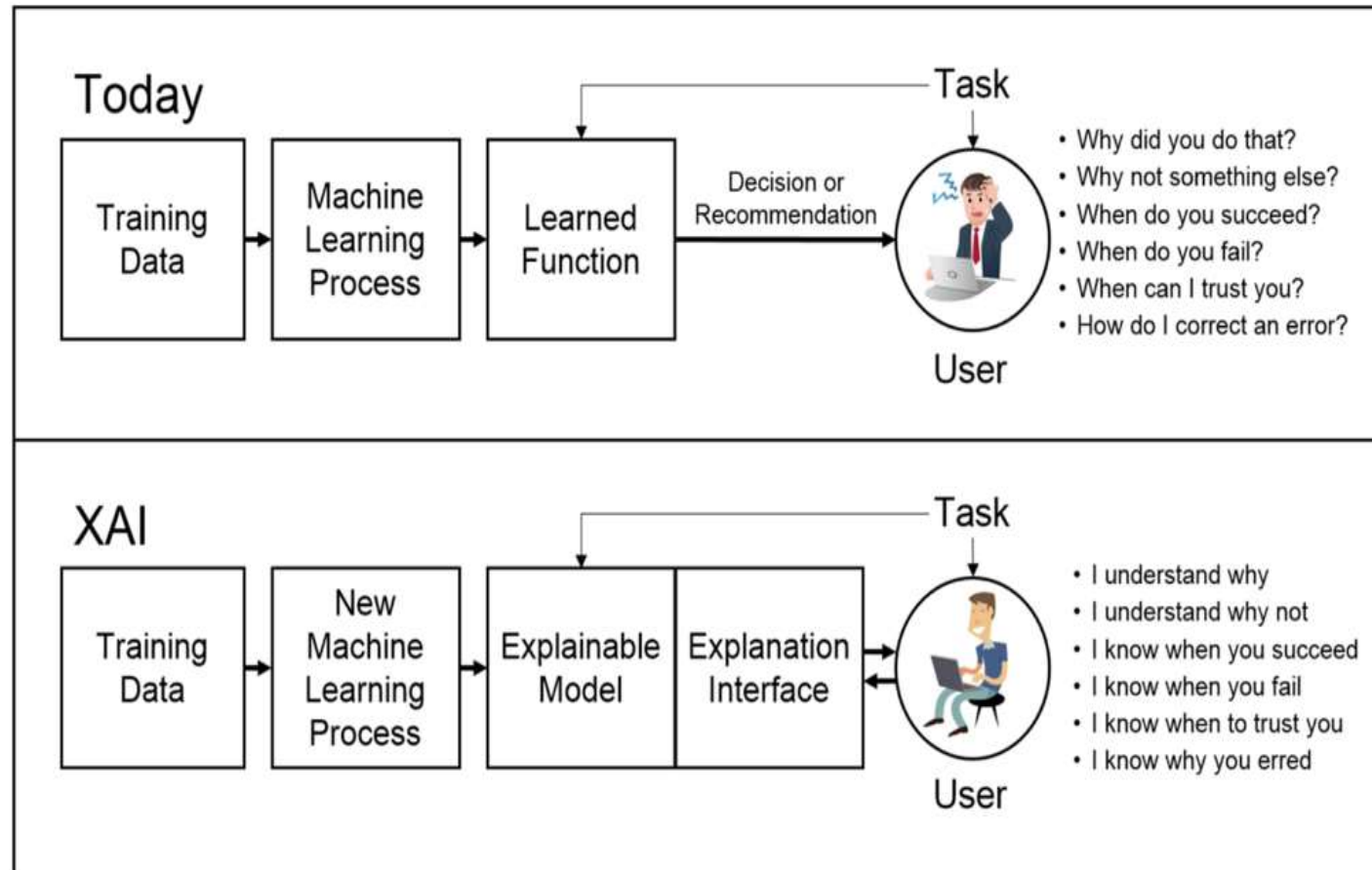


Figure 1: XAI Concept

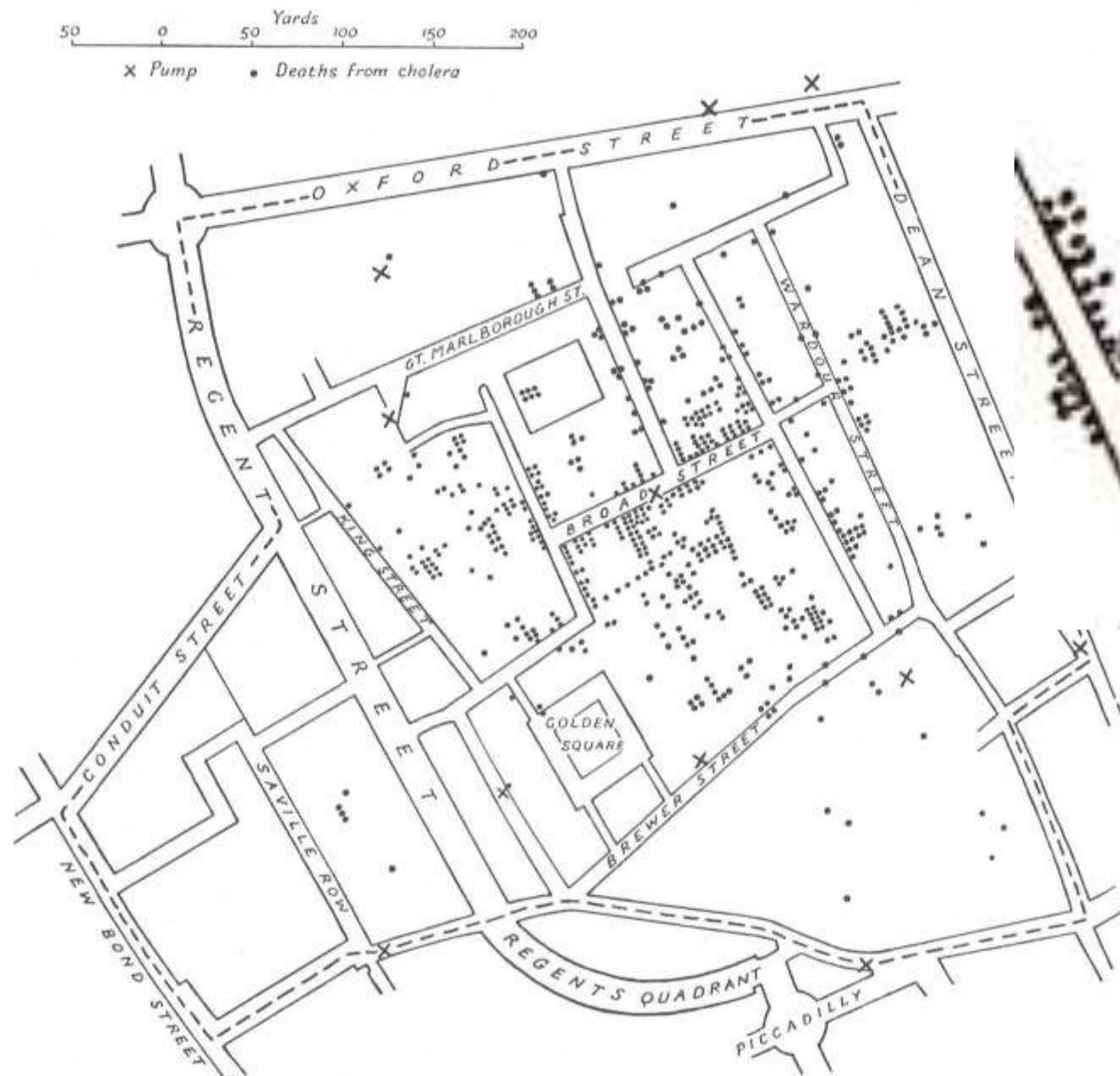
Cholera outbreak in London 1854

An early and worthy use of a map to chart non-geographical patterns

Cholera broke out in Broad Street area in London on 31 August 1854, with over 500 deaths

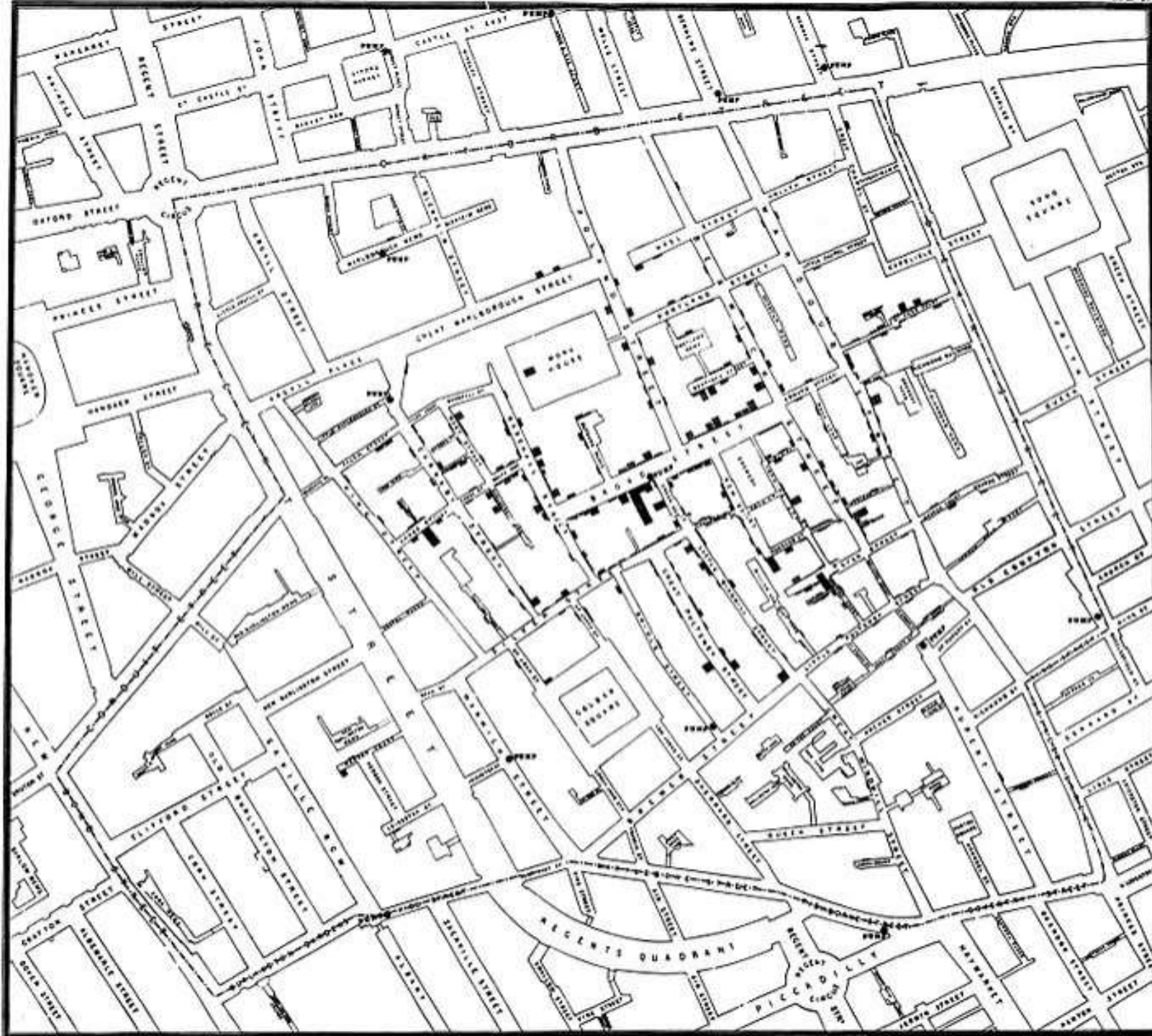
John Snow, M.D., obtained a list of deaths and by persistent case-by-case detective work

(Super famous example, pictures by from Tassu Takala, Aalto University)





The contaminated pump is located at the intersection of Broad Street and Cambridge Street



Cholera outbreak in London 1854

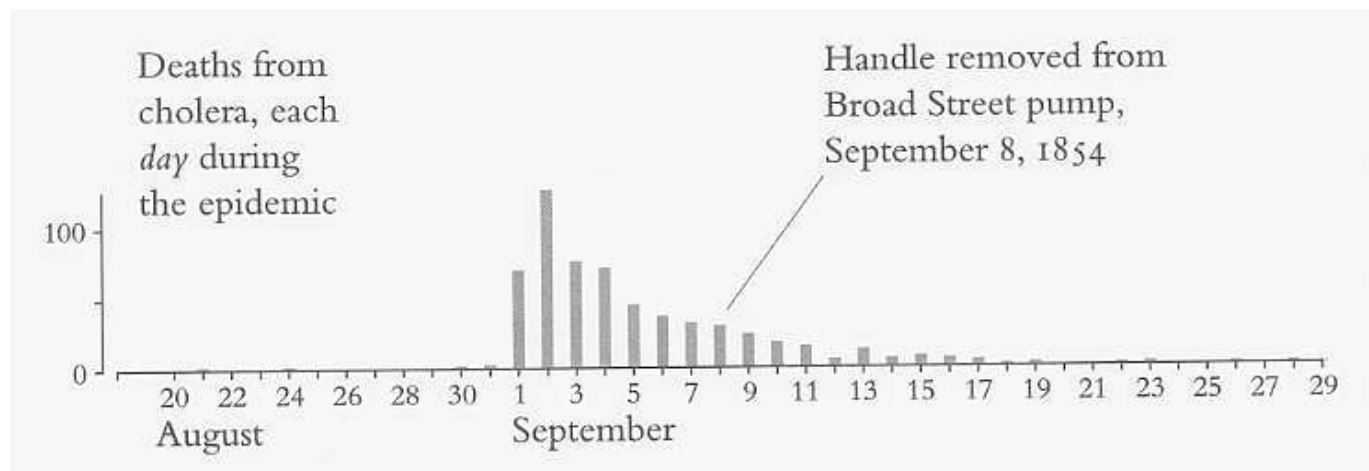
- he discovered the probable cause for the epidemics: A water pump at the Broad Street
 - The pump handle was removed on 7 September and the epidemics ended.
 - Previously it was thought that cholera spread via impure air etc.
-
- Cholera is caused by vibrio bacteria, who only become pathogen in the presence of specific phages
 - vibrio bacteria also occur in shallow water (1m depth) above 20` in the Baltic Sea. Dangerous only for the elderly.

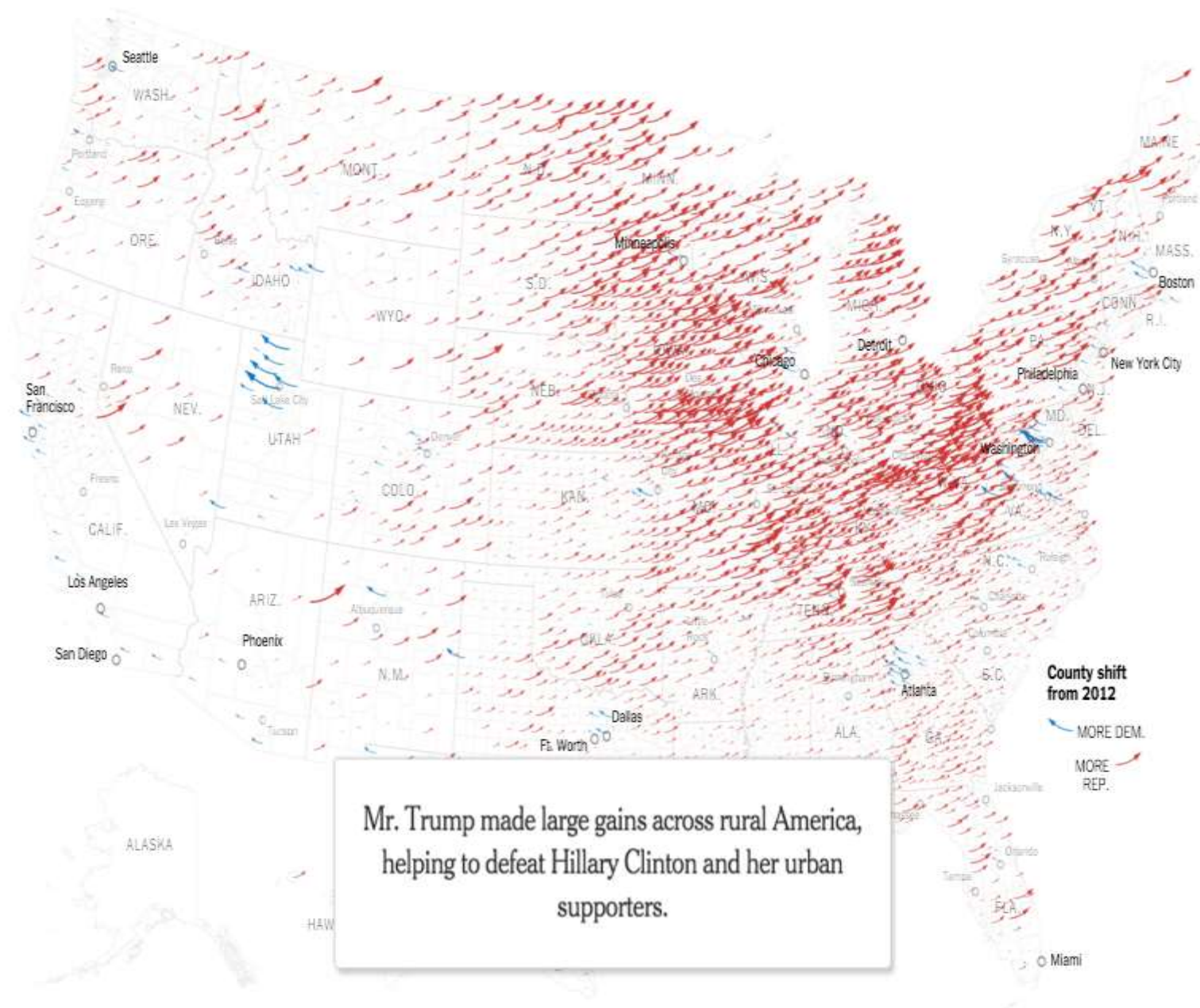
Cholera outbreak in London 1854

- Data was placed in an appropriate context to make the **relation between cause and effect** apparent. Time series, for instance, would not have been useful for finding the cause in this case.
- **Quantitative comparisons** were made. For example, Snow found that the employees of the adjacent brewery were saved because they didn't drink the water from the polluted well. They were saved by the beer(!).
- **Alternative explanations were considered.** Snow also analysed deaths that occurred far away from the Broad Street.

Cholera epidemic

- Famous example of data mapping and data visualization
- Critical view:











World History

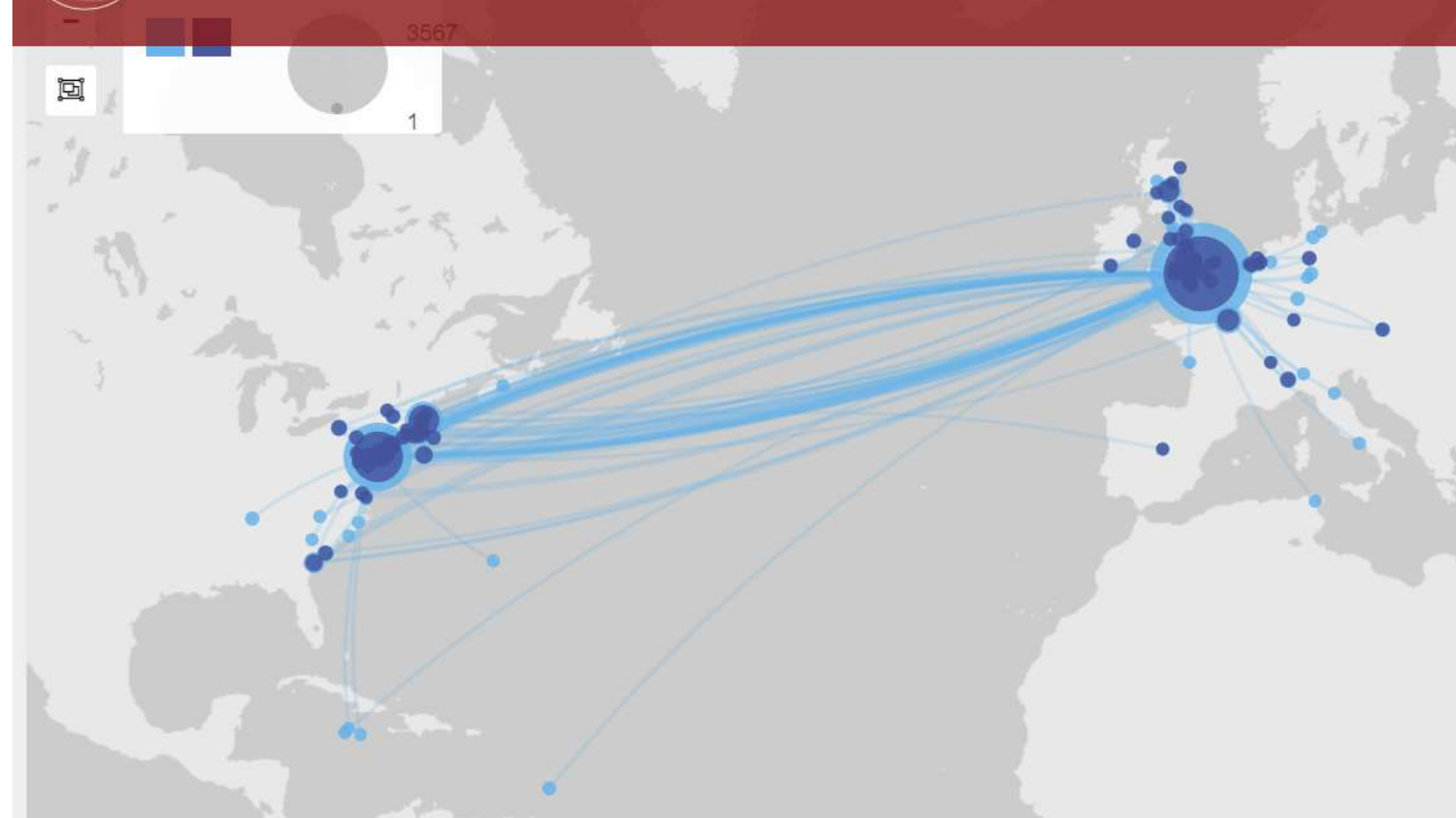


Timeline of historical events



Mapping the Republic of Letters

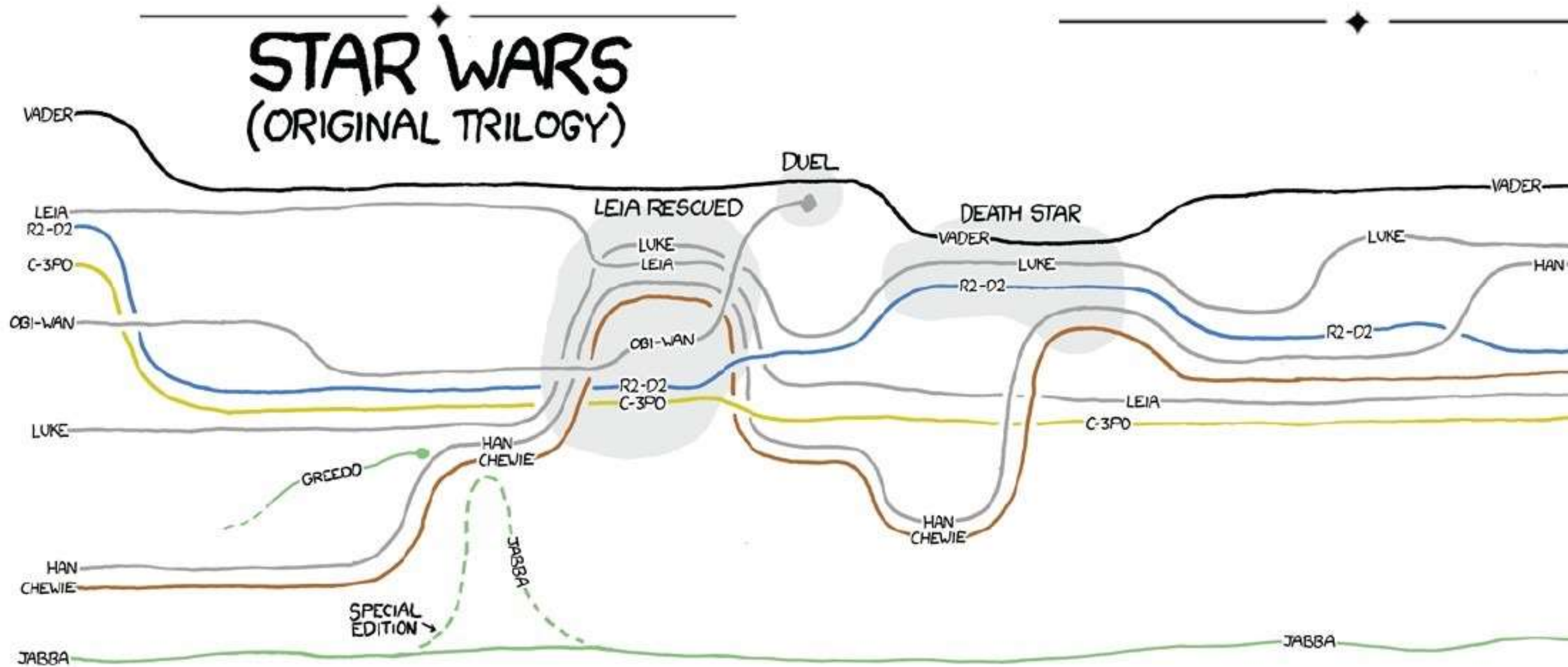
Benjamin Franklin Papers: The London Decades



- This map depicts the geographic scope of Franklin's correspondence network during the London Decades, between 1757-1775. It shows the places where correspondence within the Papers of Benjamin Franklin originated and to where it was sent.
- You can hover your cursor over each light blue and dark blue dot to see that name of the place in the form of the city, state, and country.
- The size of the dots corresponds to the number of documents that originate from or are sent to each place providing a visual indication of the places where correspondence was most frequently sent to and from.

<http://republicofletters.stanford.edu/publications/franklin/papers/>

Movie Narrative Charts



Financial data



Bloomberg terminal

Text

Text Visualization Browser

A Visual Survey of Text Visualization Techniques (IEEE PacificVis 2015 short paper)

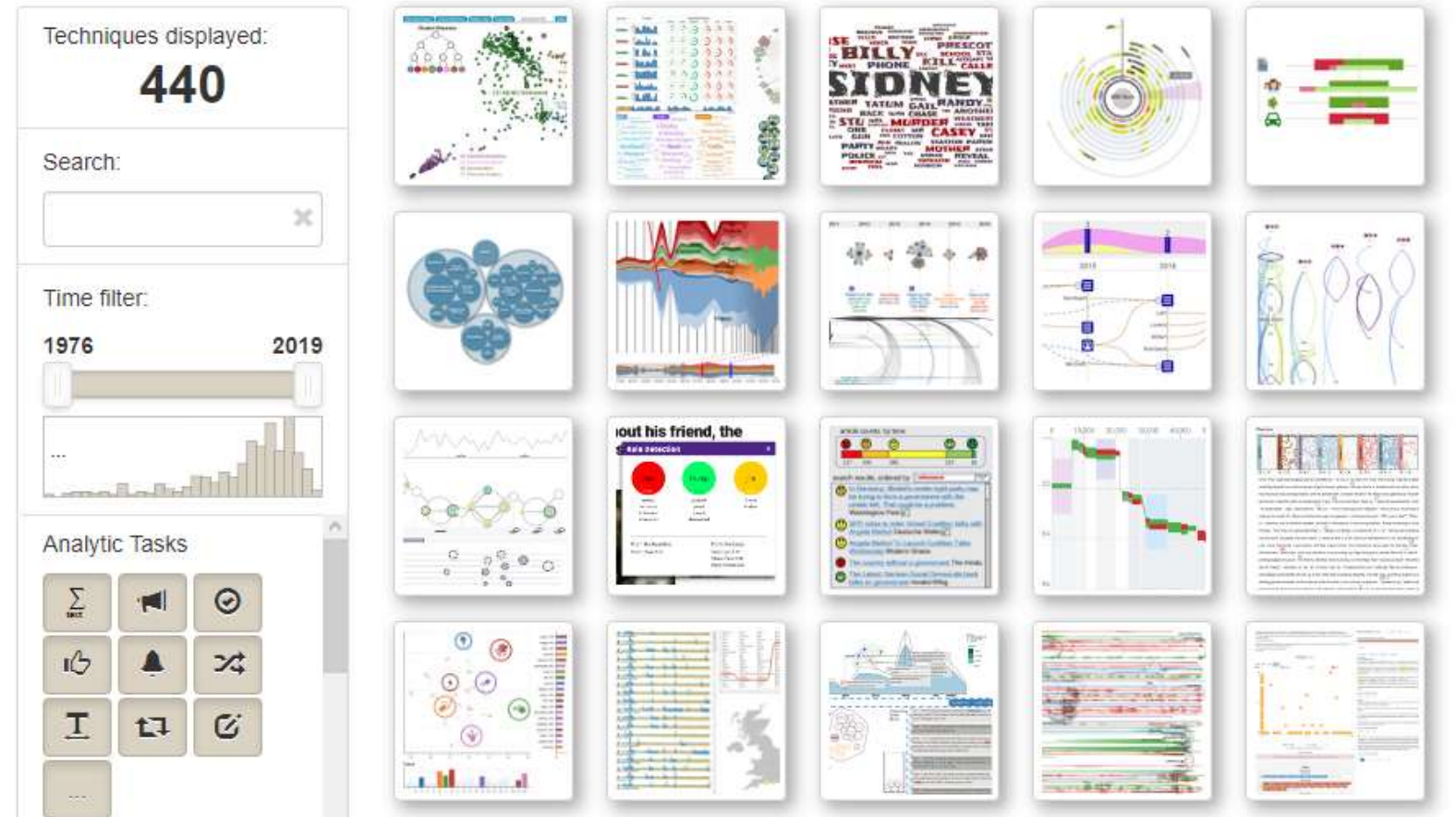
Provided by ISOVIS group

About

Summary

Add entry

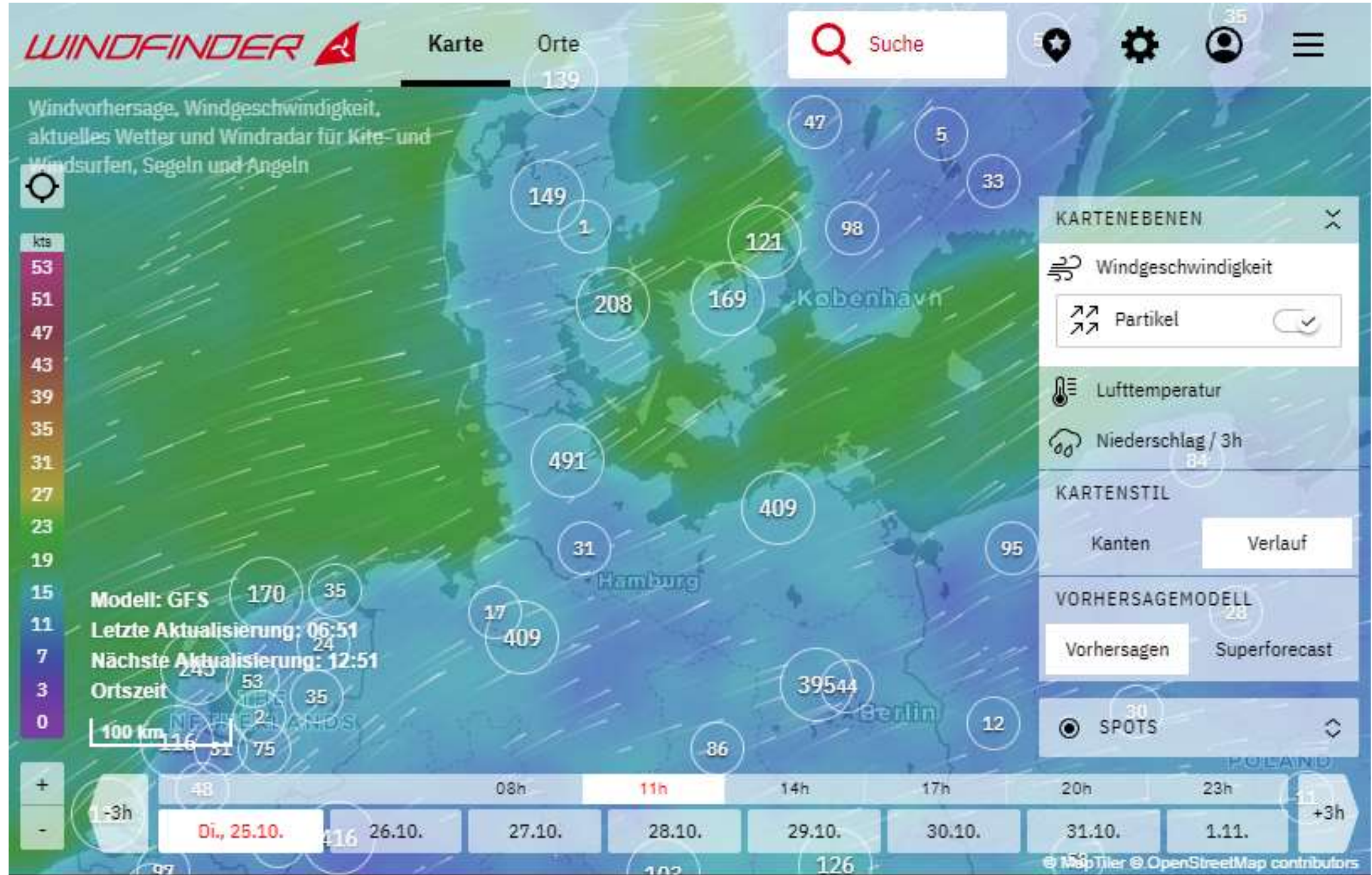
Other surveys

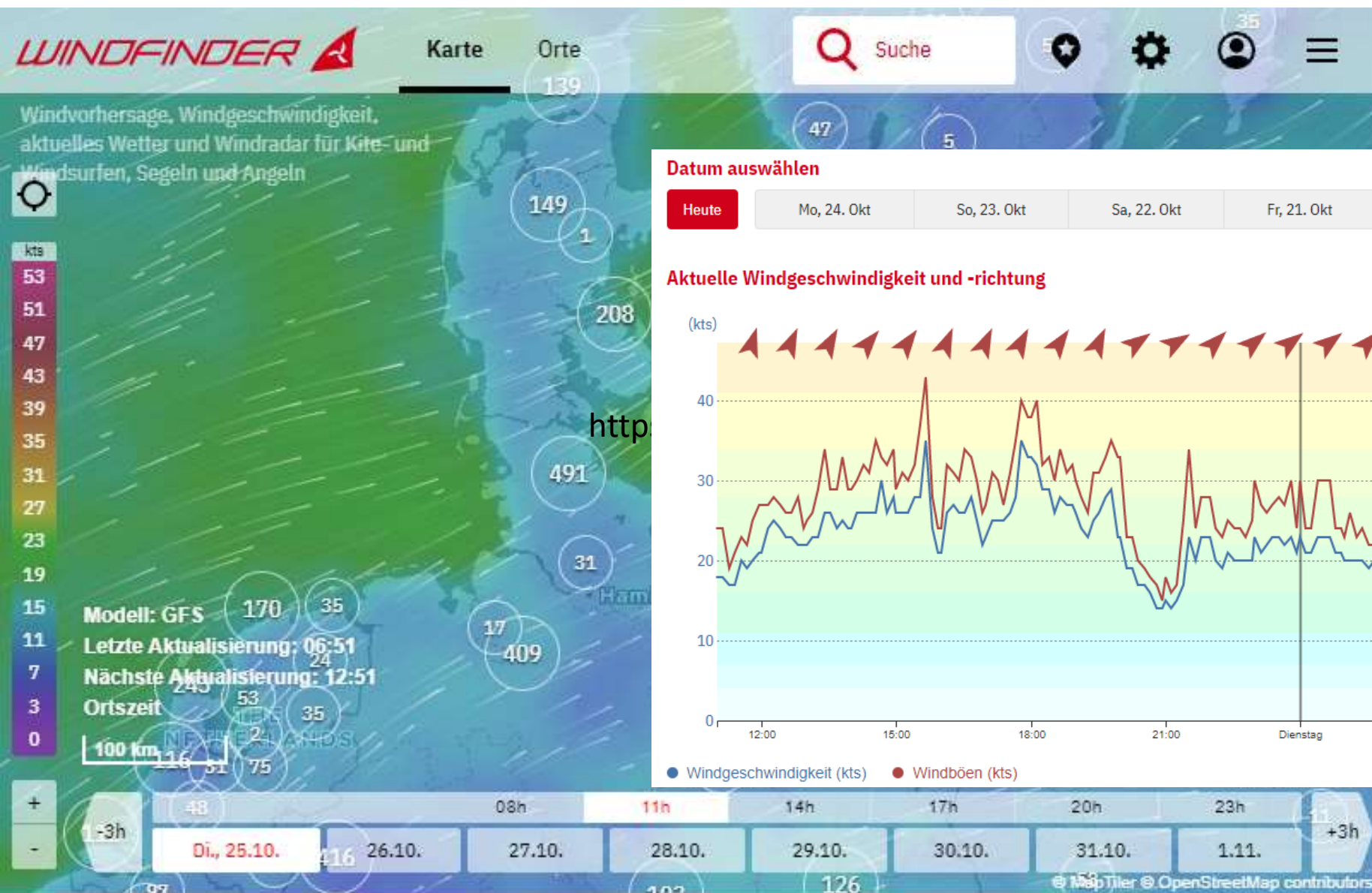


<https://textvis.lnu.se>

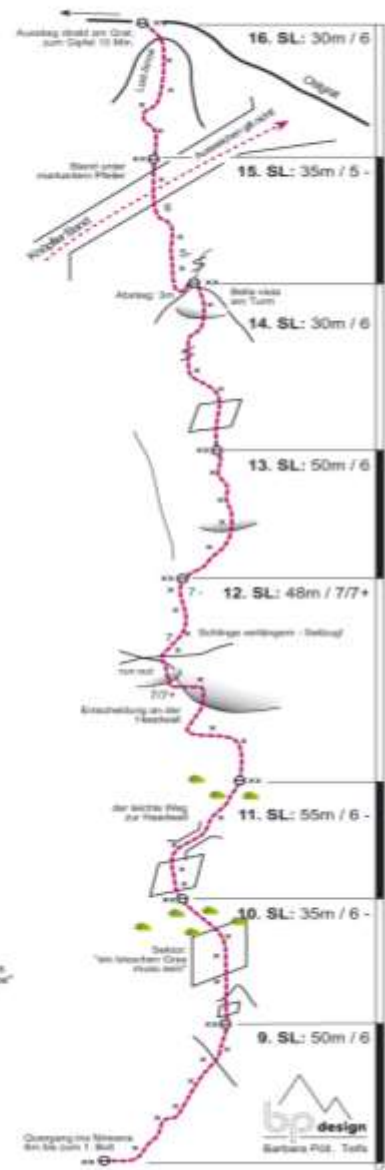
User Interfaces

































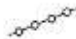




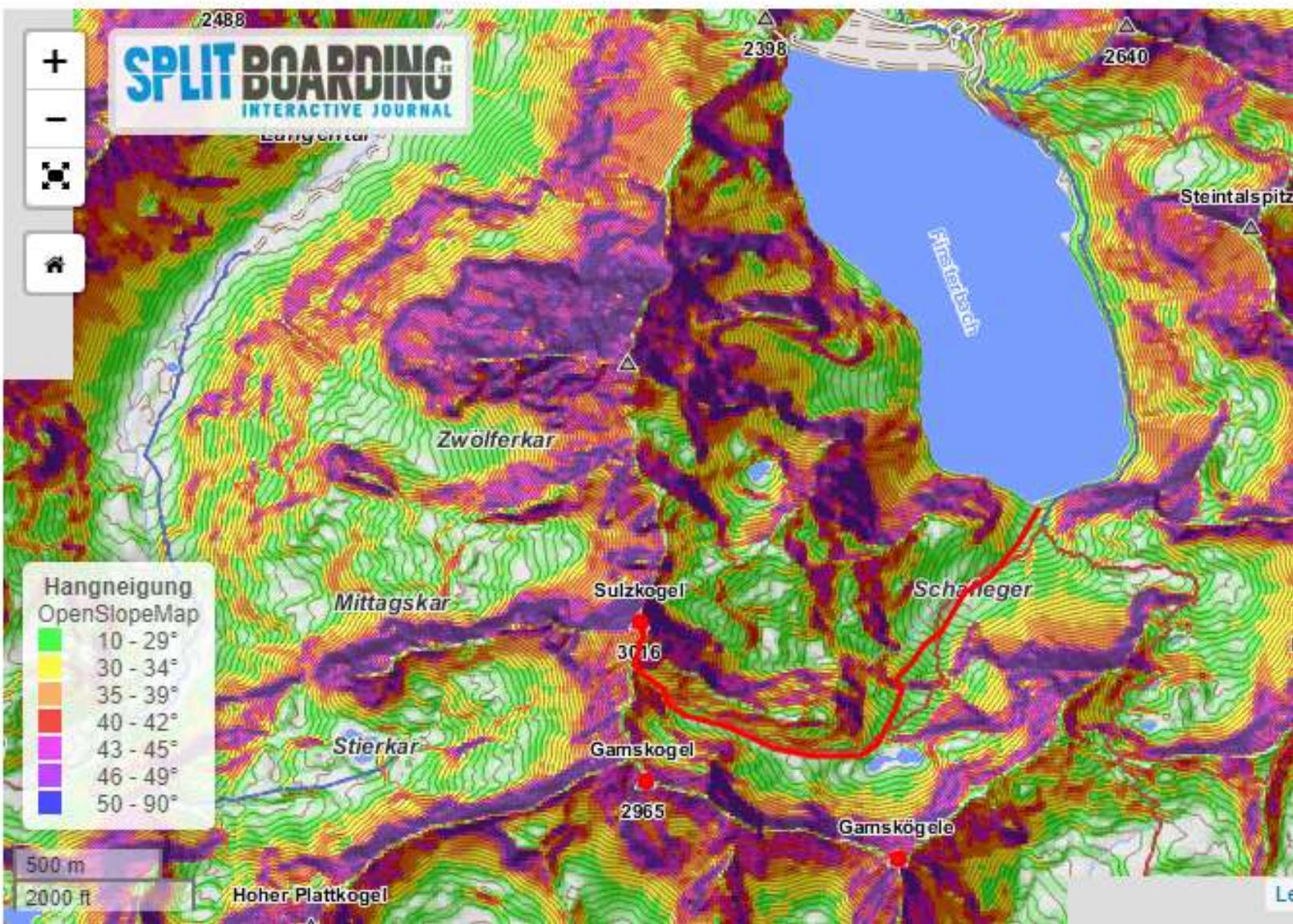


Windfinder.com - powered by IFM Kiel

[illegible]

						
Sichtbare Route visible route	Verdeckte Route hidden route	Variante variant	Standplatz belay	Abseilstelle abseiling point/anchor	Biwakplatz bivouac/bivi area	Pendelquergang pendulum
						
Bohrhaken bolt	Normalhaken (normal) piton	Sanduhr thread	Dach roof	Kamin chimney	Klemmblock chockstone	Kante edge
						
Platten slabs	Rampe ramp	Rinne groove/channel/gully	Riss crack	Verschneidung dihedral/V-cleft	Wasserstreifen water streak/runnel	Überhang overhang
						
Band ledge	Absatz terrace/pedestal/platform	Nische niche	Baum tree	Latsche mountain pine	Wechte cornice	Schnee snow
			GT	WB		
Gras grass	Schutt rubble	Schlüsselstelle crux	Gedenktafel plaque	Wandbuch summit register/visitors log	Gipfelkreuz summit cross	Fixes Drahtseil fixed wire rope





Stop or Go

Strategische Lawinenkunde für Tourengeher

EDITION Berg & Steigen
CEAV

Entscheidungsstrategie

CHECK 1

Gefahrenstufe und Hangneigung

1 gering	2 mäßig	3 erheblich	4 groß	5 sehr groß
	Verzicht auf 40° und mehr	Verzicht auf 35° und mehr	Verzicht auf 30° und mehr „Spitzkehrgelände“	Verzicht auf Touren allgemein

CHECK 2

Gefahrenzeichen erkennen

Wahrnehmen	Beurteilen	Handeln
Neuschnee?	GEFÄHRlich für mich?	YES: STOP Ausweichen Abbrechen
Triebsschnee?		
Lawinen?		
Durchfeuchtung?		
Setzungsgeräusche?	NO: GO	

Recap

- What is Data Visualization?
 - Why visualize?
 - Models of Data Visualization
 - Examples
-
- Next Week: More on What and Why

