

# Introduction to Big Data Analytics

## Big Data Visualization

Yanwei Fu

School of Data Science, Fudan University



## Part I : Introduction

What is visualization ?

Why do we create visualisation ?

Existing Visualisation Techniques

## Part II : Big Data Visualization

Challenges

Techniques

## Part III : How can we visualize big data

Key techniques

Open source tools

Examples

## Part IV : Visual Analysis of Big Data



## Part I : Introduction



## What is visualization?



# How can we acquire information?

Listen



Taste &  
Smell



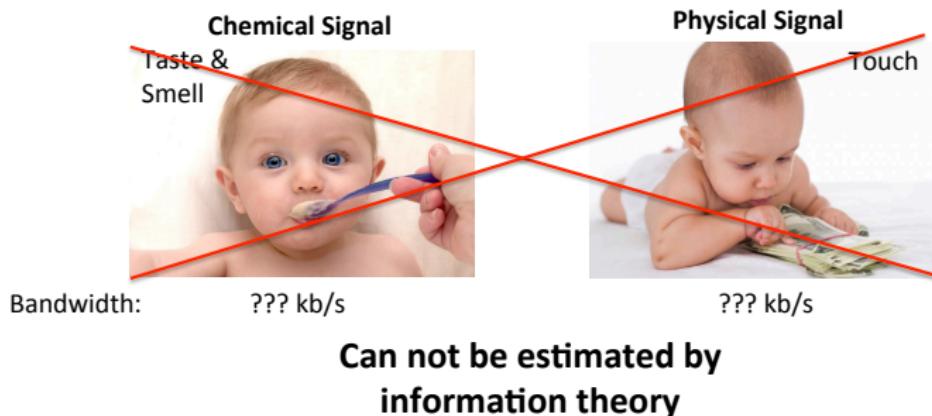
Touch



Look



# Do they effective ?



# Do they effective ?

Sound Signal

Listen



Electronic / Light Signal

Look

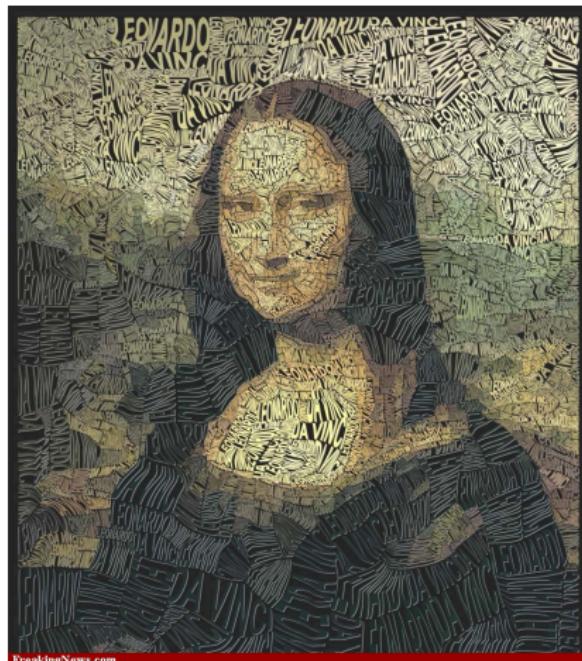


Bandwidth: **about 0.1 KB / s**

**> 100 MB/s**

“Information Visualization, Perception for Design” 3<sup>rd</sup> Edition, by Colin Ware

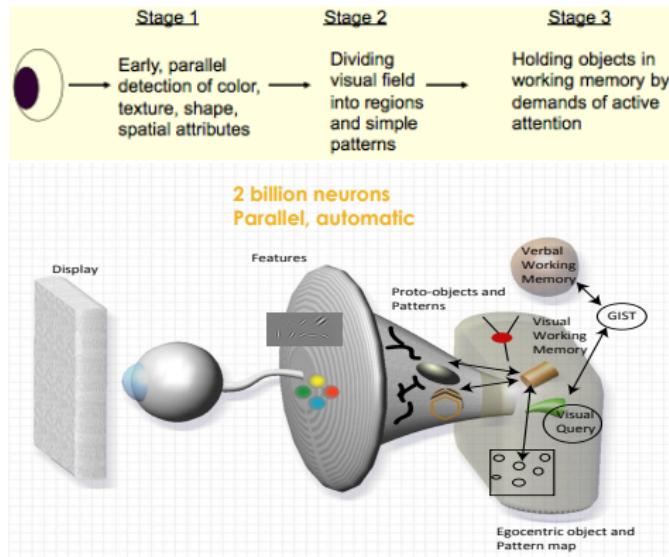
# A picture worth 1000 Words



## The Visual Thinking Pipeline

Parallel Processing to Extract Low-Level Visual Properties such as color, shape, etc

Sequential Goal-Oriented Processing







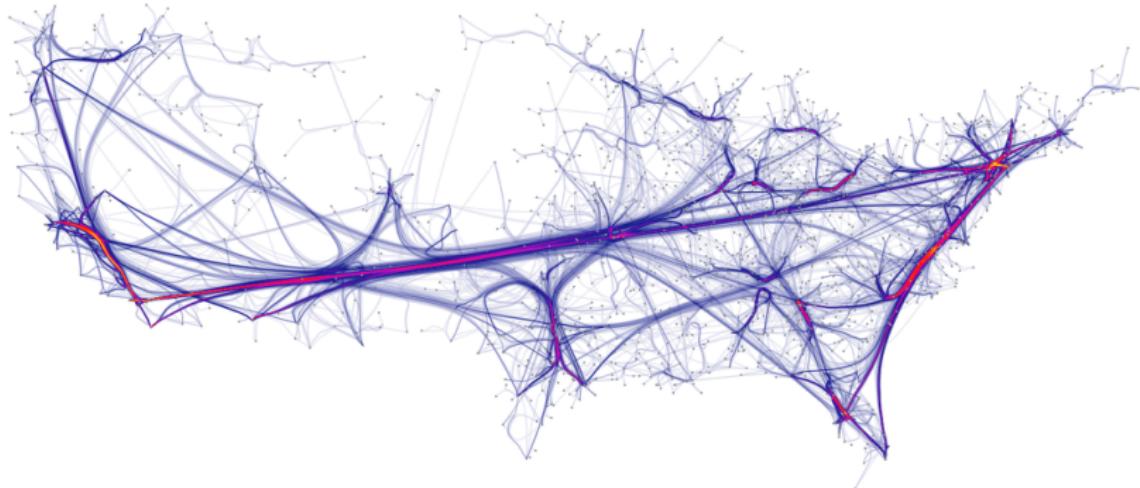
Ching-Yung Lin, Nan Cao, Shixia Liu, Spiros Papadimitriou, Jimeng Sun, and Xifeng Yan. SmallBlue: Social Network Analysis for Expertise Search and Collective Intelligence. ICDE 2009

**Record of human activities to help find data patterns**



**Visualization is used for help reasoning and decision making**

## Summarization of Airlines in United States



Holten, Danny, and Jarke J. Van Wijk. "Force-Directed Edge Bundling for Graph Visualization." Computer Graphics Forum. Vol. 28. No. 3. Blackwell Publishing Ltd, 2009.

# What is Information Visualization ?

“The action or fact of visualizing; the power or process of forming a mental picture or vision of something not actually present to the sight; a picture thus formed.”

-- Oxford English Dictionary

“... finding the artificial memory that best supports our natural means of perception.”

-- Bertin, 1983

**The use of computer-supported, interactive, visual representations of abstract data to amplify cognition**

-- Cart,Mackinlay, Shneiderman, 1999



Why do we create visualization?

Why do we create visualization?



# Why do we create visualization ?

Counting the number of 3s in the  
following Text:

1235693234870452973467  
0378937043679709102539

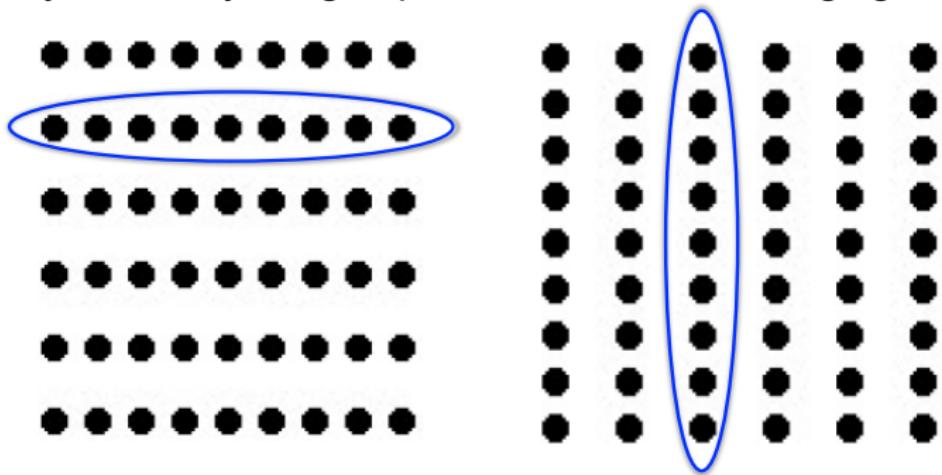


Counting the number of 3s in the following Text:

1235693234870452973467  
0378937043679709102539



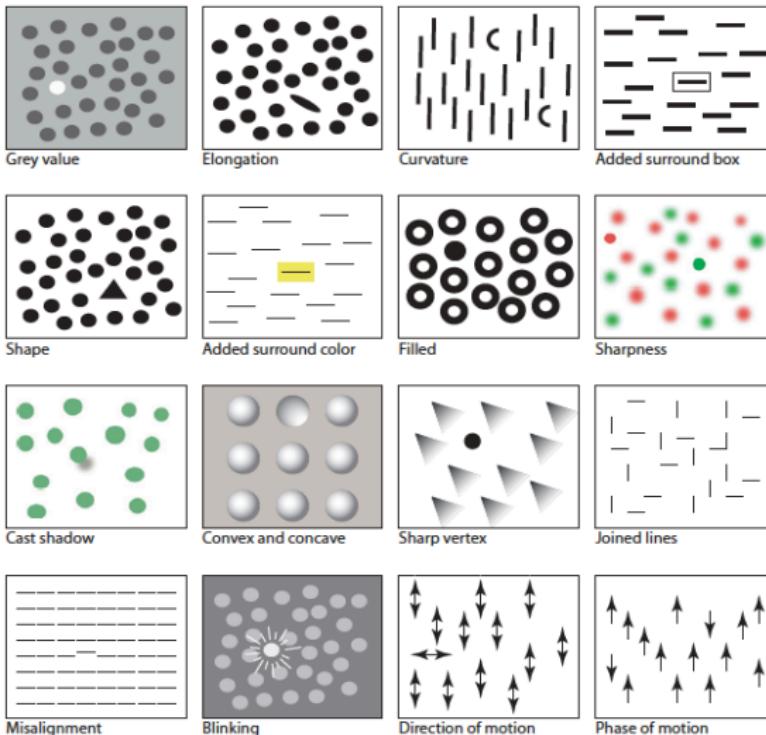
Can you identify the groups of dots in the following figures ?



*Law of Proximity*

*we tend to group elements that are closest to each other*

# Find Patterns: Pre-Attentive Visual Channels

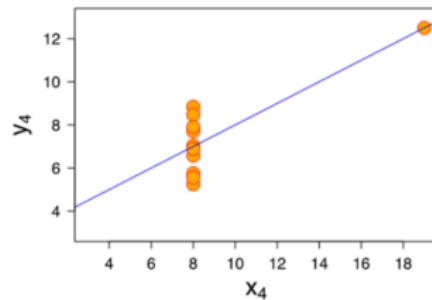
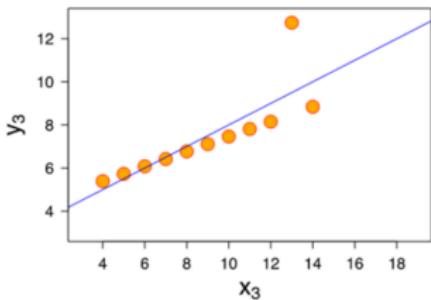
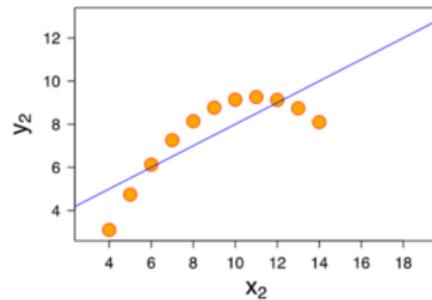
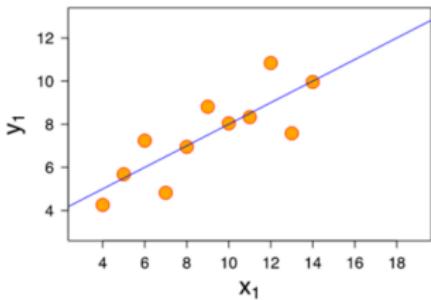


# Why do we create Visualization?

	Set A		Set B		Set C		Set D	
	X	Y	X	Y	X	Y	X	Y
0	10	8.04	10	9.14	10	7.46	8	6.58
1	8	6.95	8	8.14	8	6.77	8	5.76
2	13	7.58	13	8.74	13	12.74	8	7.71
3	9	8.81	9	8.77	9	7.11	8	8.84
4	11	8.33	11	9.26	11	7.81	8	8.47
5	14	9.96	14	8.10	14	8.84	8	7.04
6	6	7.24	6	6.13	6	6.08	8	5.25
7	4	4.26	4	3.10	4	5.39	19	12.50
8	12	10.84	12	9.13	12	8.15	8	5.56
9	7	4.82	7	7.26	7	6.42	8	7.91
10	5	5.68	5	4.74	5	5.73	8	6.89
mean	9.00	7.50	9.00	7.50	9.00	7.50	9.00	7.50
std	3.32	2.03	3.32	2.03	3.32	2.03	3.32	2.03
corr	0.82		0.82		0.82		0.82	
lin. reg.	$y = 3.00 + 0.500x$		$y = 3.00 + 0.500x$		$y = 3.00 + 0.500x$		$y = 3.00 + 0.500x$	



# Seeing data in context



# Why do we create visualization ?

A picture is worth a thousand words

## News illustrated

The sudden explosion of a South Korean entertainer called Psy, has given the world Gangnam Style. It is setting the music and dance world on fire and has a set sequence. We simplify them for your perusal

**★ When to use the steps during the chorus ★**

**Step 1** ahhhhh... Gangnam style...  
Oppa is Gangnam style... ahhhhh... Gangnam style... Oh, oh oh oh, Oppa is Gangnam style... ahhhhh... Gangnam style... Oh, oh oh oh, Oppa is Gangnam style... Eeeeeh- Sexy Lady...  
**Step 2** ahhhhh... Gangnam style...  
**Step 3 or Step 1** (in the last chorus) Oh, oh oh oh, Oppa is Gangnam style... Eeeeeh- Sexy Lady...  
**Step 4** ahhhhh... Gangnam style...  
**Step 5** (in the end) Oh, oh oh oh, Oppa is Gangnam style... Eeeeeh- Sexy Lady oh oh oh, Oppa is Gangnam style.

**Step 1** Riding the horse

Dress classy and dance cheezy!  
Cross your hands like taking the horse reins and move up and down  
Do small jumps with your legs, as if you are riding a horse  
Footsteps: R L R R R L R L L

**Step 2** Lassooing the sexy lady

Lassoing motion with your right arm  
Continue with the horse-riding movement  
Footsteps: R L R R R L R L L

**Step 3** Now everybody is looking at me

A Hands in pockets or waist and small movements combined with the foot steps  
B Finish this move dragging the right leg to the left leg.  
Footsteps A: L R L R L R L R  
Footsteps B: L R L R L R L R

**Step 4** Combine a few 'sexy' moves

C Now move your hips to the outside with quick movements twice  
D Spread and flex your legs. Move your body up and down three times  
Footsteps: L R L R L R L R

**Step 5** Finish with a cool pose

A Cross your hands over your left leg  
B "L" shape with your fingers  
C Spread your arms and raise your right leg (position A). Now get down quickly on your right leg and flex the left one. Now rotate your right arm and with your hand touch your chin doing a "L" shape with your thumb and index finger (position B)

Source: You Tube

HUGO D. SANCHEZ/Digital News

A better communication method



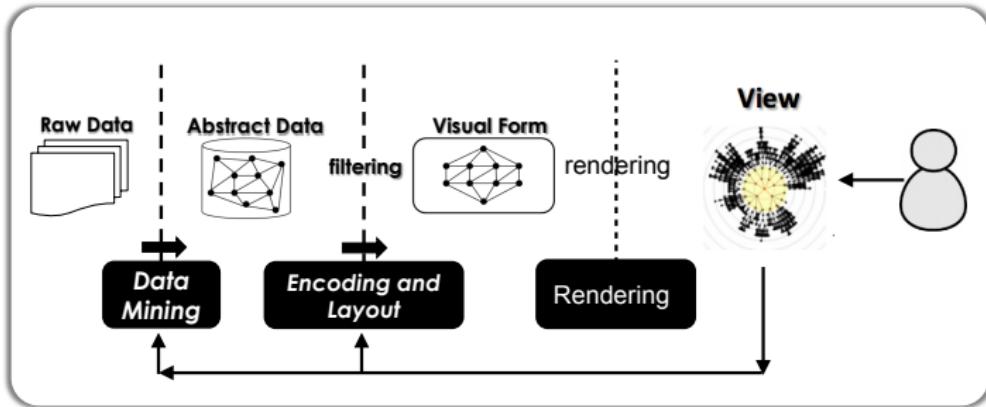
- See data in context
- Find patterns
- Telling a story
- Attract attentions
- Communicate information with others
- Summarization and interpretation
- Graphical calculation
- Expend memory
- Inspire people



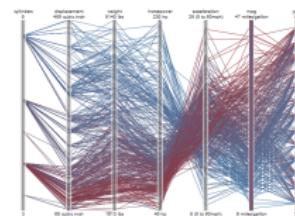
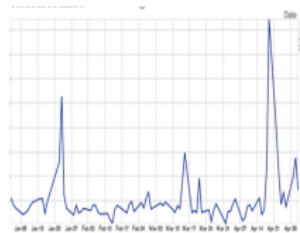
## Existing Visualisation Techniques



# Visualization & Visual Analysis Reference Model



# Taxonomy by data types

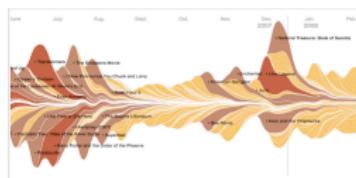


1D

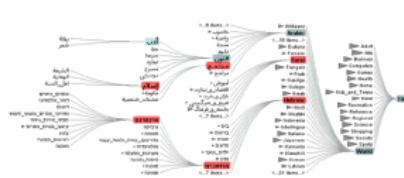
2D

3D

Multi-D



Temporal



Tree

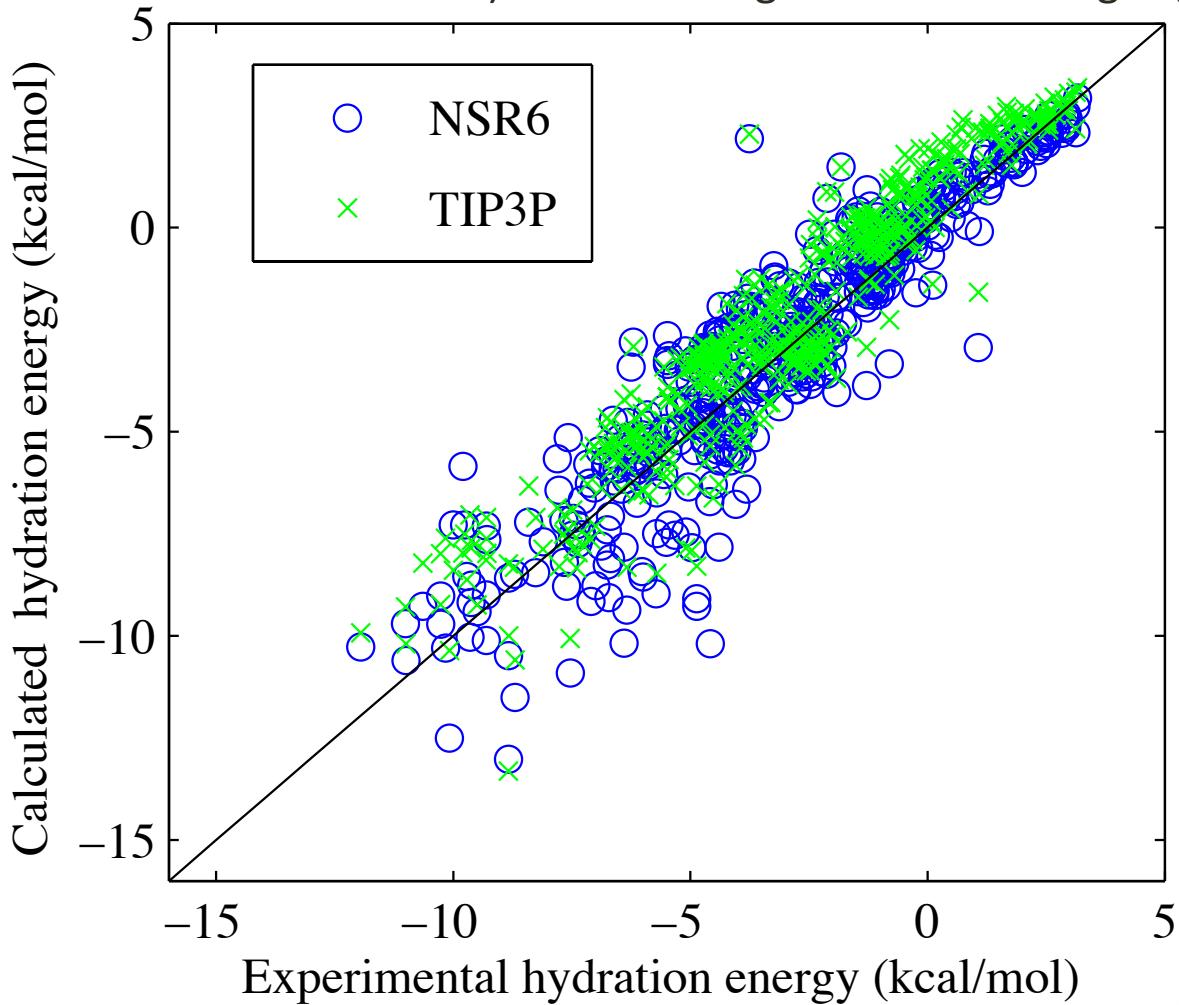


Graph

# Example: scatter plot

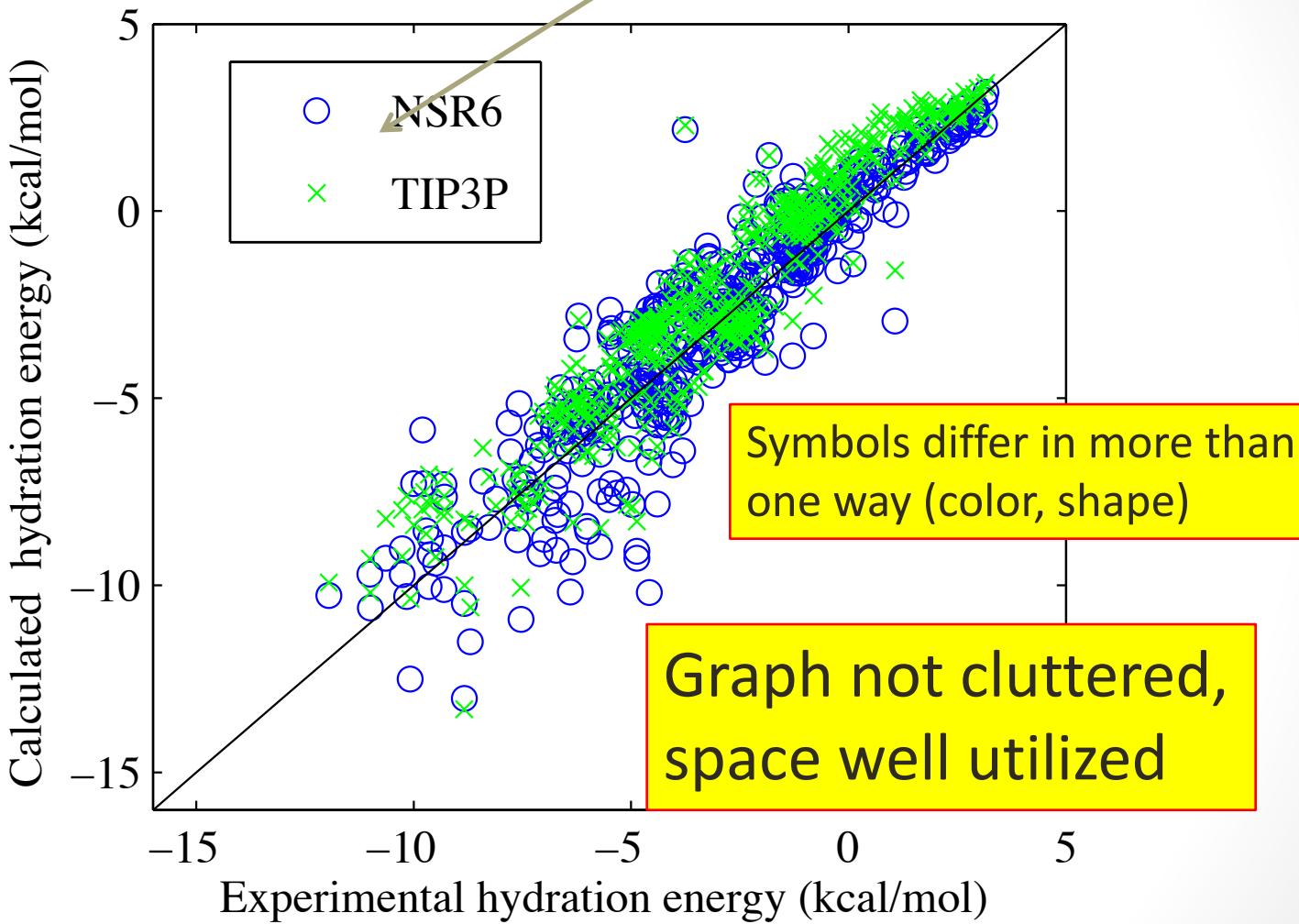
Immediately tells:

- a) How good the over-all agreement is
- b) Are there serious outliers
- c) Is the agreement uniform?
- d) Where disagreement is strongest/weakest



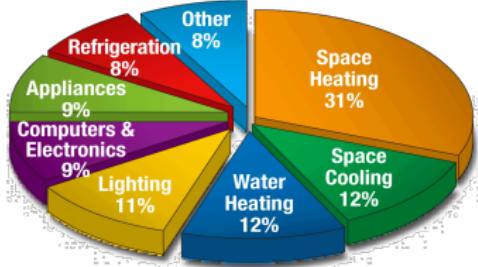
# Elements of a good graph:

Large enough font, thick lines.

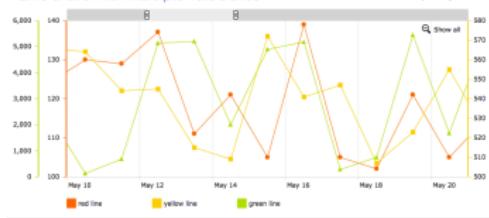


All axes labeled, units shown

# Examples: Visualizing 1D Numerical Data



► Line chart with multiple value axes

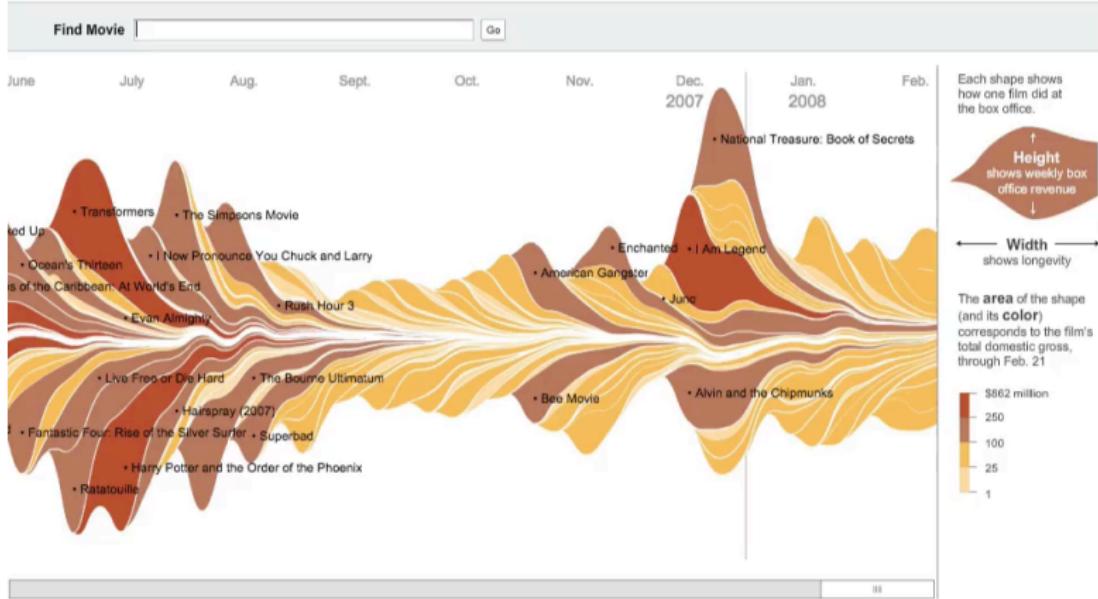


Open in jsFiddle

NFC



# Example : Visualizing 1D Ordinal Data



Sources: Baseline StudioSystems; Box Office Mojo

Mathew Bloch, Lee Byron, Shan Carter and Amanda Cox

[http://www.nytimes.com/interactive/2008/02/23/movies/20080223\\_REVENUE\\_GRAPHIC.html](http://www.nytimes.com/interactive/2008/02/23/movies/20080223_REVENUE_GRAPHIC.html)



# Examples: 2D Data

S&P 500 AUG 29 2008 04:00 PM

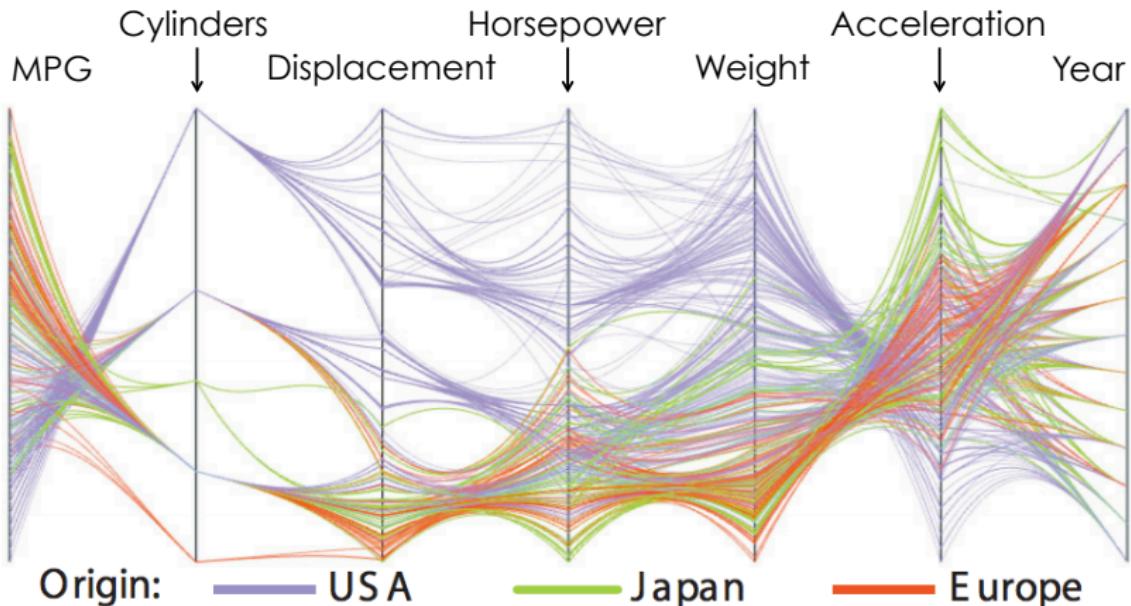
[finviz.com](http://finviz.com)



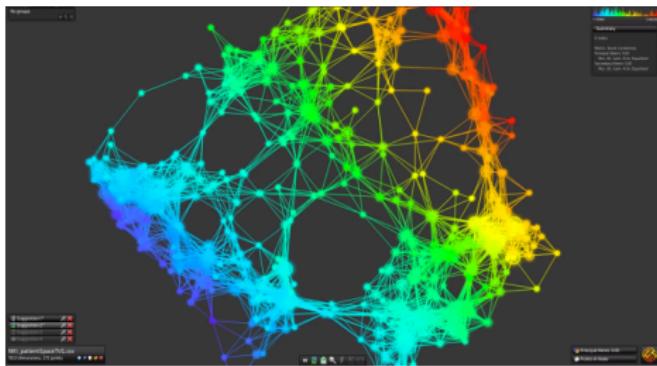
Size of each Cell: Stock Market Value  
Color: Stock Change



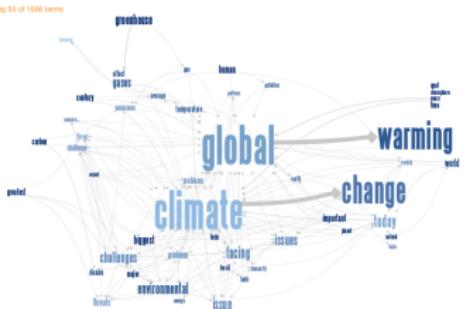
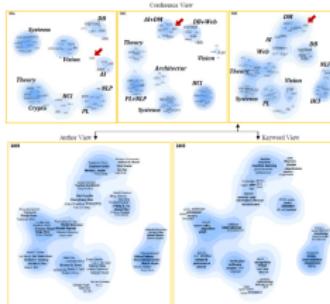
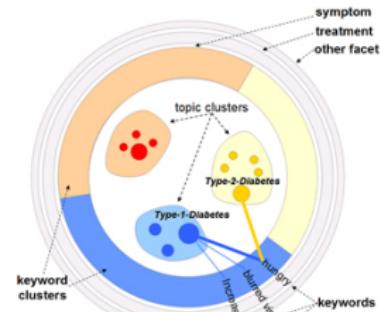
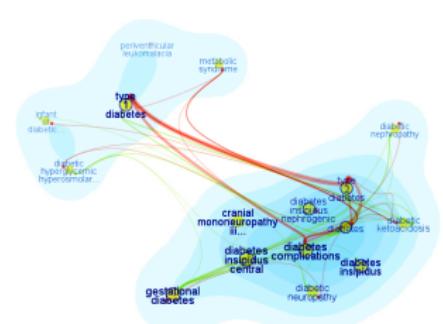
# Example: Multi-Dimensional Data



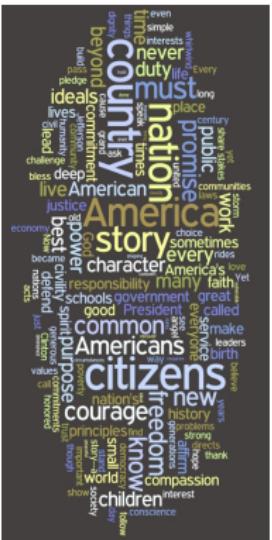
# Examples: Visualizing Structured Data



# Examples: Visualizing Unstructured Data



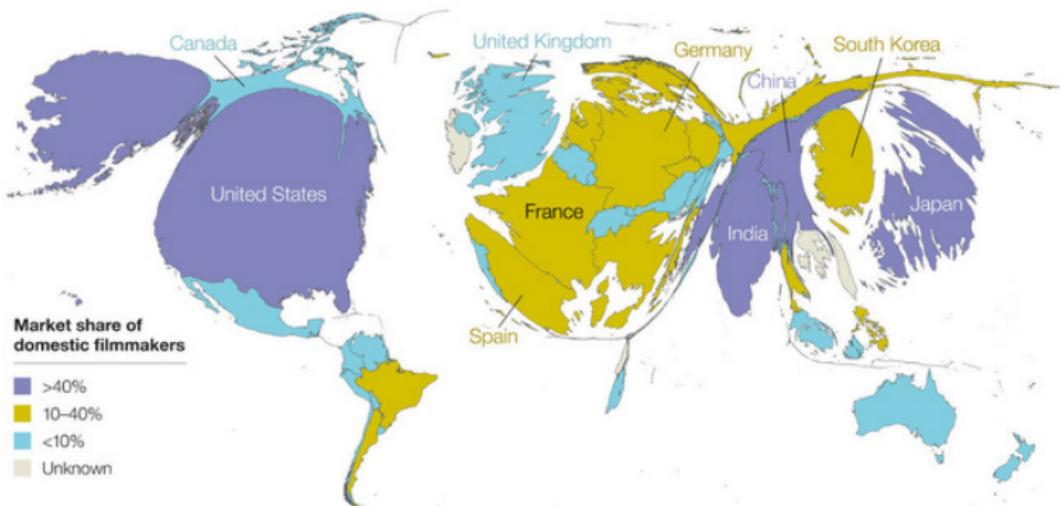
Visualization of Text Documents



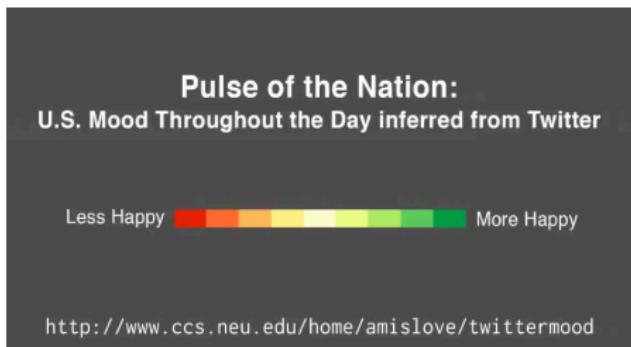
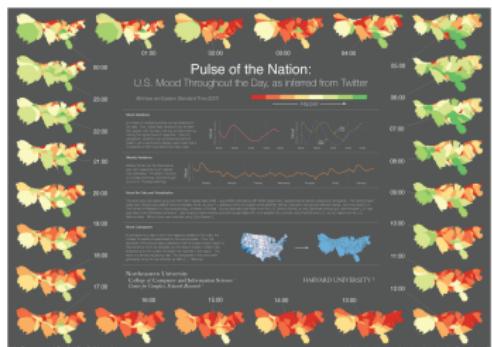
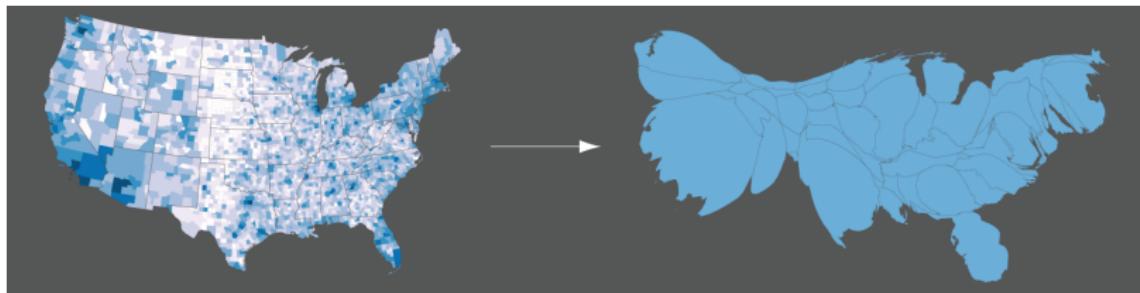
# Examples: Geospatial

Larger cinema markets support stronger domestic film industries.

Countries sized by relative share of worldwide box office revenue, 2009



# Examples: Visualizing Spatial Temporal Data



# Examples: Visualizing Spatial Temporal Data

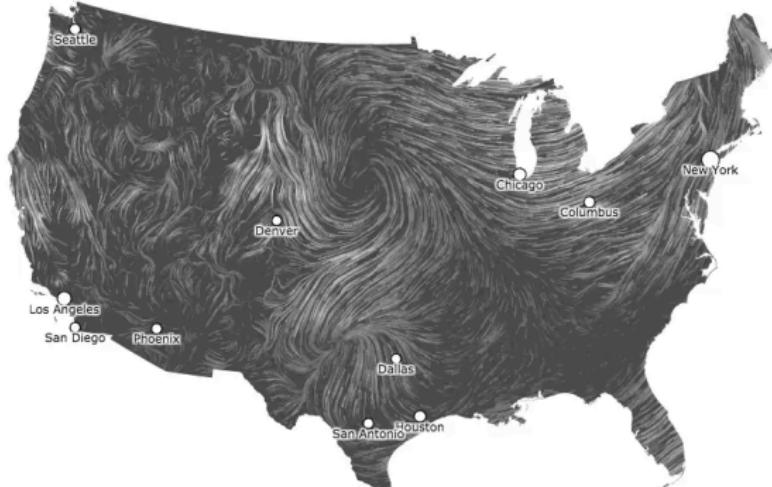
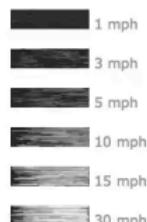
## wind map

Dec. 3, 2014

11:35 am EST

(time of forecast download)

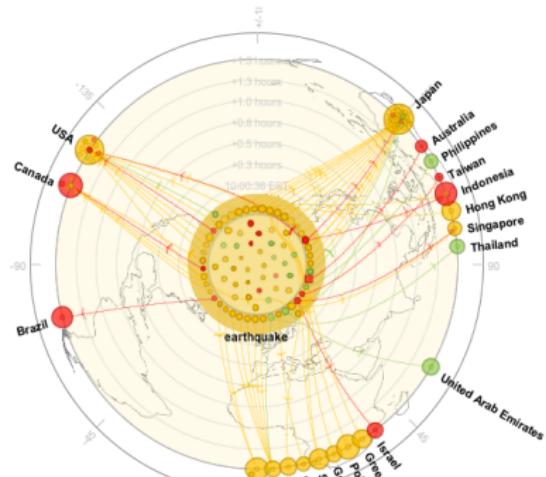
top speed: 31.5 mph  
average: 8.2 mph



<http://hint.fm/wind/>



# Visualization is not just a beautiful picture

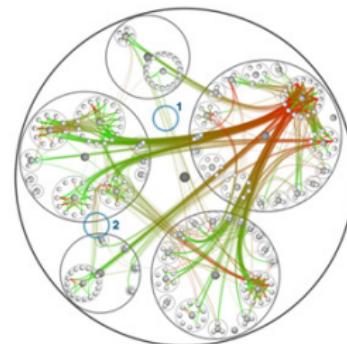


The purpose of visualization is to reveal the insight of the data

# InfoVis vs Computer Graphics

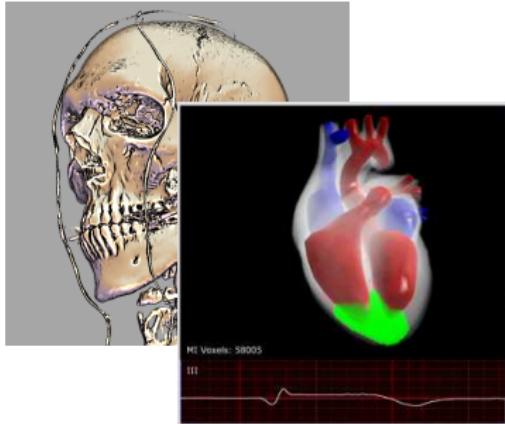


Realism

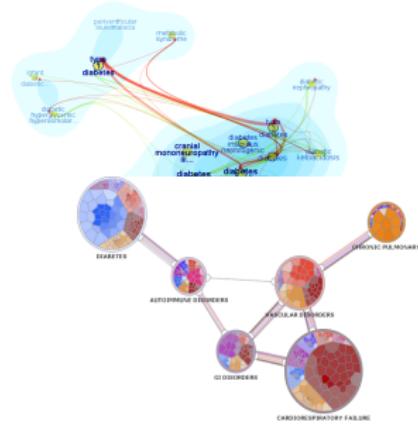


Information

# InfoVis v.s. Scientific Visualization

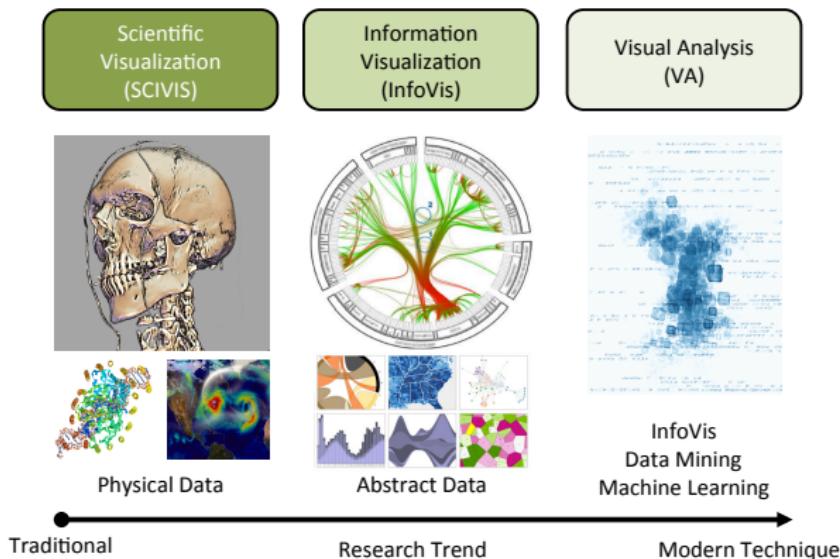


**Physical Data**



**Artificial Data**

# Three Sub-areas



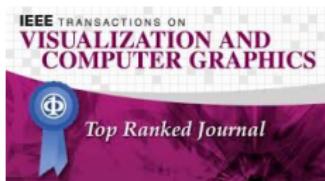
# Major Conferences and Journals

Scientific  
Visualization  
(SCIVIS)

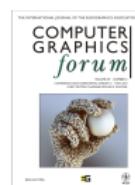
Information  
Visualization  
(InfoVis)

Visual Analysis  
(VA)

VisWeek: IEEE SCIVIS, INFOVIS, VAST



IEEE Transaction on Visualization  
and Computer Graphics



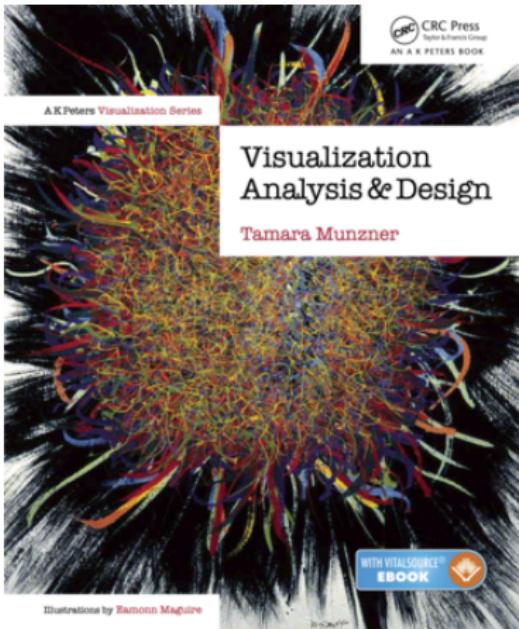
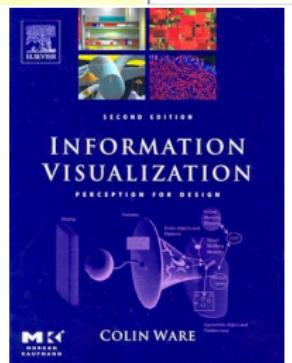
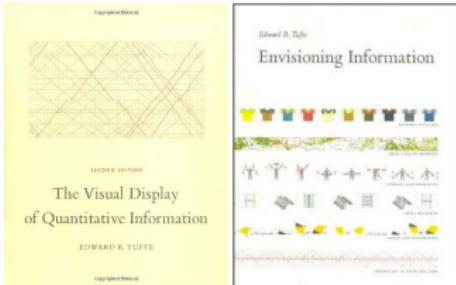
Computer Graphics  
Forum



IEEE Computer Graphics  
and Application

36

# Recommended Books



## Part II: Visualising Big Data

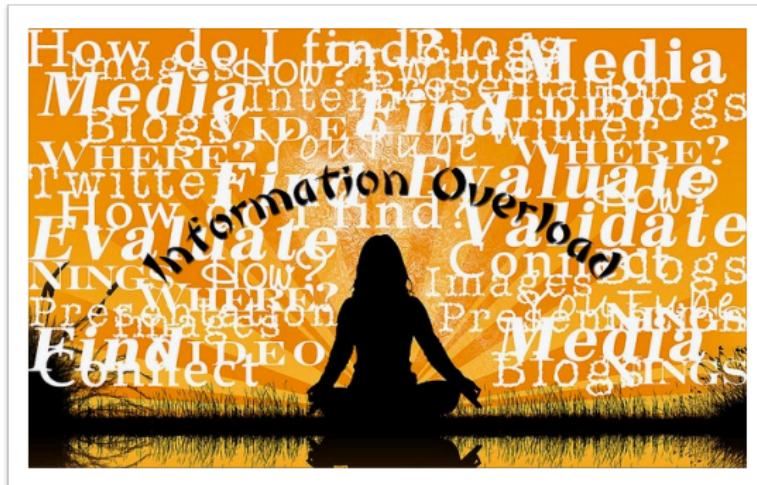


# Information Overload

difficult caused by too much information



# How can we acquire useful information from the overwhelming data



# Big Data Visualization

76425 species



Tree of Life by Dr. Yifan Hu

14.8 million tweets



500 million users



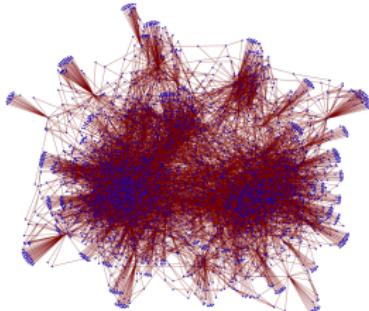
Facebook friendship graph by Paul Butler

## Challenging Task:

Squeezing millions and even billions of records into million pixels  
( $1600 \times 1200 = 2$  million pixels)



# Challenges



Visual clutter

How can we avoid visual  
clutters like overlaps  
and crossings?



Performance issues

How can we render the  
huge datasets in real time  
with rich interactions?



Limited cognition

How can users understand  
the visual representation  
when the information  
is overwhelming?

# Techniques (1) : Pixel Oriented Visualization

data item

attr1 

attr2 

attr3 

attr4 

attr5 

attr6 

- ✿ A multidimensional data item contains 6 attributes

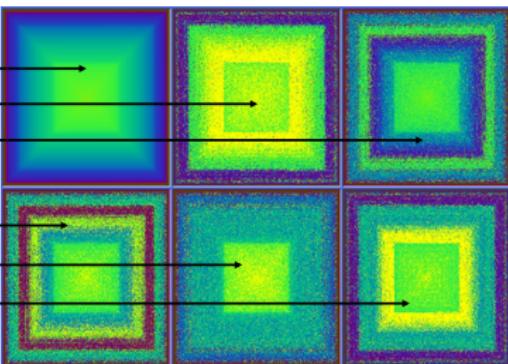


# Technique(1) : Pixel Oriented Visualization

Database visualization (10,000 items, 6 dimensions)

Jan	Feb	Mar	Apr	May	Jun
-99.99	-99.99	315.7	317.45	317.5	317.26
315.62	316.38	316.71	317.72	318.29	318.16
316.43	316.97	317.59	319.62	320.03	319.59
316.58	317.77	318.34	319.58	320.59	319.77
317.94	318.56	319.68	320.63	321.01	320.55
318.74	319.08	319.86	321.39	322.24	321.47
319.57	-99.99	-99.99	322.24	321.89	
319.44	320.44	320.89	322.13	322.16	321.87
320.62	321.59	322.39	323.87	324.01	323.75
322.06	322.5	323.04	324.42	325	324.09
322.57	323.25	323.89	325.62	325.57	325.36
324.42	325.64	326.89	328.34	329.34	328.76
325.03	325.99	326.87	328.14	328.07	327.66
326.17	326.68	327.18	327.78	328.92	328.57
326.77	327.63	327.75	329.72	330.07	329.09
328.55	329.56	330.3	331.5	332.48	331.07
329.35	330.71	331.48	332.65	333.09	332.25
330.4	331.41	332.04	333.31	333.96	333.6
331.75	332.56	333.5	334.3	334.97	334.25
333.02	333.92	334.7	336.07	336.74	335.27
334.97	335.39	336.64	337.76	338.01	337.89
336.23	336.76	337.98	338.69	339.47	339.29
338.01	338.36	340.08	340.77	341.46	341.17
339.23	340.47	341.38	342.51	342.91	342.25
340.75	341.61	342.7	343.57	344.13	343.35
341.37	342.52	343.1	344.94	345.75	345.32
342.47	343.8	344.95	347.05	347.43	346.79
344.97	346	347.43	348.35	349.39	348.25
346.3	346.96	347.6	349.25	350.22	349.25
348.02	348.47	349.42	350.99	351.84	351.25
350.43	351.73	352.22	353.59	354.22	353.79
352.76	353.67	353.68	355.13	355.67	355.13
353.66	354.7	355.39	356.2	357.16	356.23
354.72	355.75	357.16	358.6	359.33	358.24
355.29	356.22	357.11	358.5	359.66	358.5
356.7	357.16	359.38	369.46	368.28	359.6
358.37	358.91	359.87	361.26	361.68	360.95
359.97	361	361.64	363.45	363.79	363.26
362.05	363.25	364.02	364.72	365.41	364.97
363.18	364	364.56	366.35	366.79	365.62
365.33	366.15	367.31	368.61	369.3	368.87
366.87	368.07	369.59	371.12	371	370.75
369.14	369.46	370.77	371.66	372.82	371
370.28	371.5	372.12	372.78	374.02	373.3
372.43	373.09	373.52	374.86	375.55	375.41
374.68	375.63	376.11	377.65	378.35	378.13
376.79	377.37	378.41	380.52	380.63	379.57
378.37	379.69	380.41	382.1	382.28	382.13
381.38	382.03	382.64	384.62	384.95	384.06
382.45	383.68	384.23	386.26	386.39	385.87
385.07	385.72	385.85	386.71	388.45	387.64

Order by degree of interests max

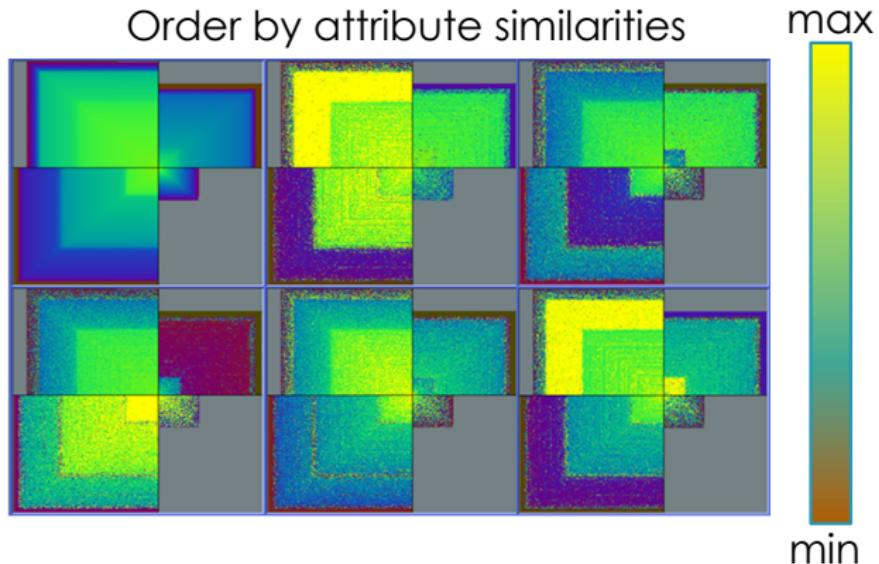


(Keim & Kriegel, 1994; 1996)

min

# Techniques (1): Pixel Oriented Visualization

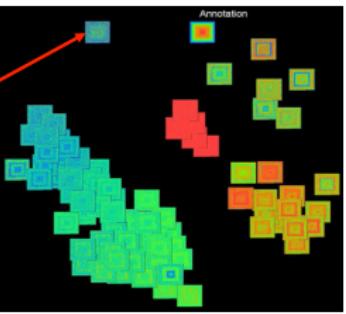
Database Visualization (10,000 items, 6 dimensions)



# Techniques (1) : Pixel Oriented Visualization

## Different Ways for splitting the display region

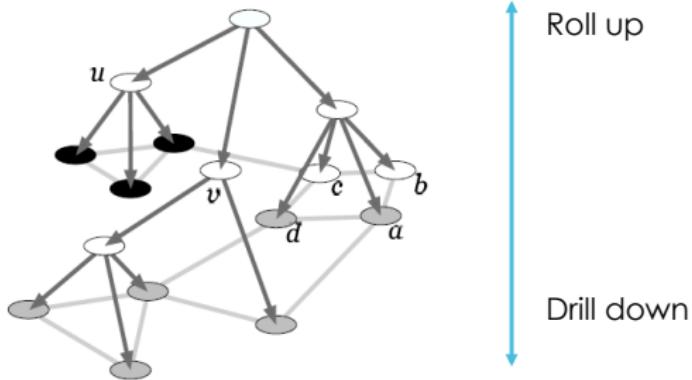
	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec	Annual Average
1998	-99.99	-99.99	-99.99	315.99	317.45	317.59	317.26	315.99	314.93	313.52	312.44	313.33	316.9
1999	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.26	315.26	315.97
2000	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2001	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2002	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2003	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2004	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2005	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2006	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2007	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2008	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2009	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2010	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2011	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2012	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2013	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2014	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2015	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2016	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2017	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2018	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2019	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2020	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2021	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2022	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2023	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2024	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2025	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2026	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2027	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97
2028	316.42	316.42	316.97	317.58	319.02	320.03	317.16	318.98	318.58	318.91	315.19	315.19	315.97



(Yang et al., 2006)

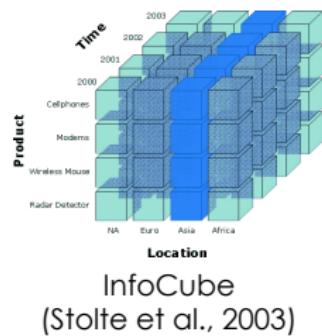
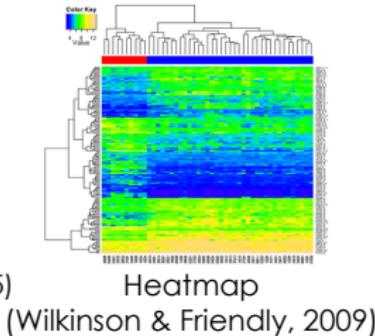
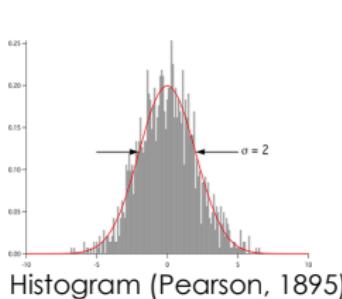
Values above represent monthly concentrations adjusted to represent 2400 hours on the 15th day of each month. Units are parts per million by volume (ppmv) expressed in the 2033A STO meteorite mole fraction scale. The "annual average" is the arithmetic mean of the twelve monthly values where no missing values are present.

## Techniques (2): Aggregation & Level of Details (LOD)



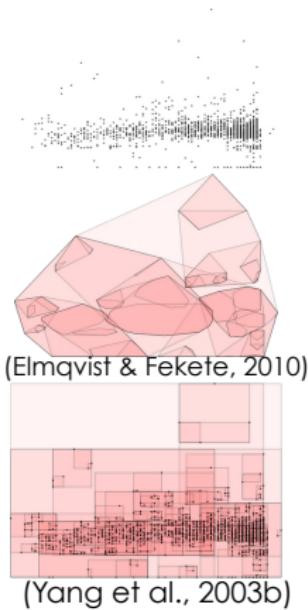
Building a tree for aggregating data items in either a bottom-up or top-down approach

# Technique (2) : Aggregation & LOD

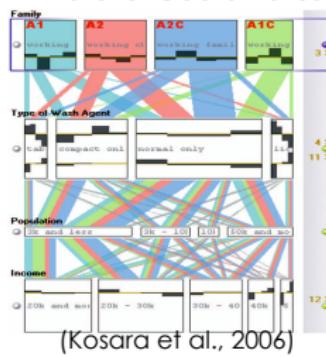


# Techniques (2) : Aggregation

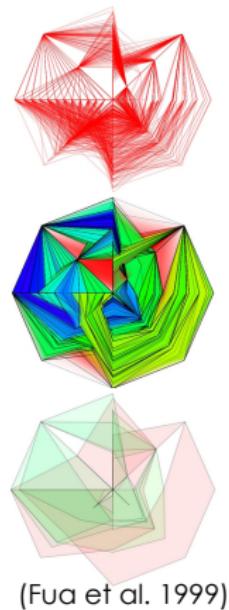
Scatter Plots



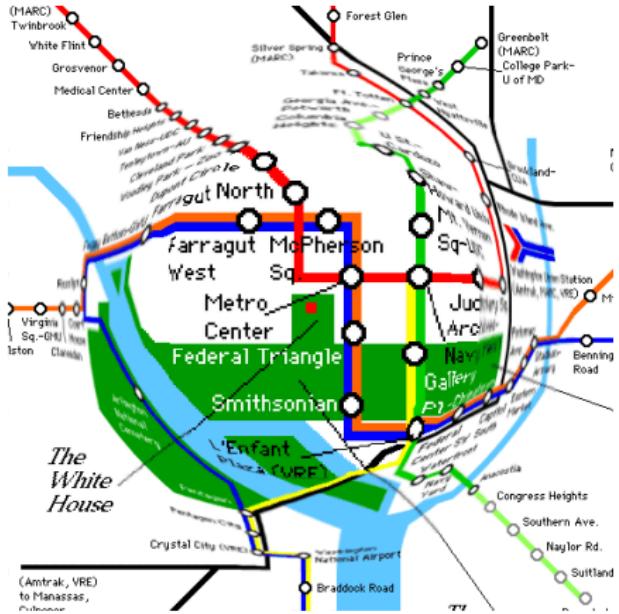
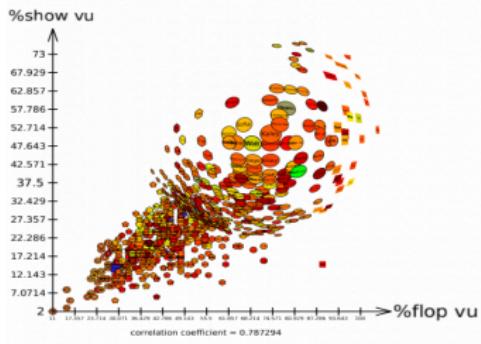
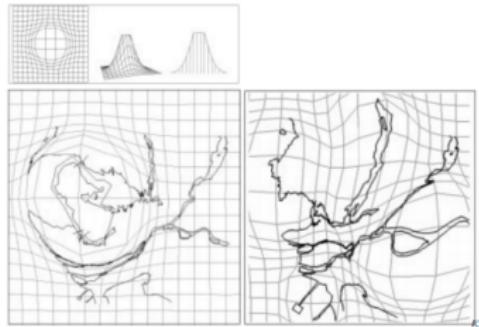
Parallel Coordinates



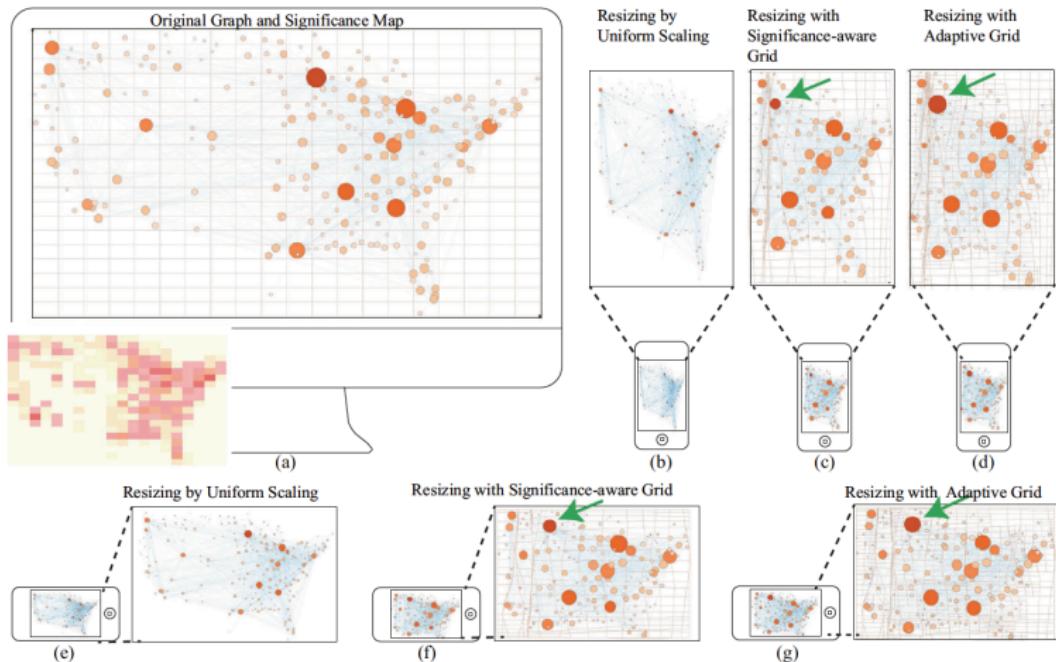
Star Plots



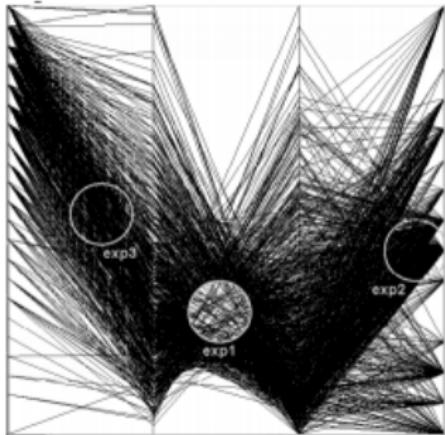
# Technique (3) : Distortion



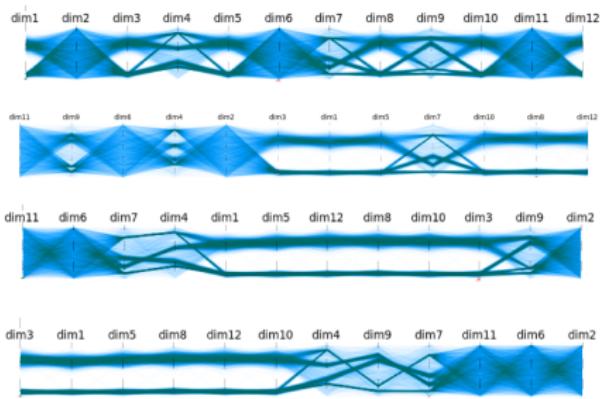
# Techniques (3) : Distortion



# Technique (4) : Clutter Reduction



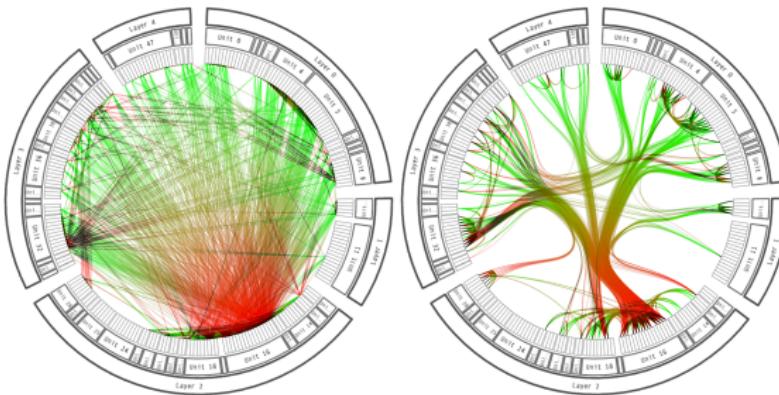
Sampling



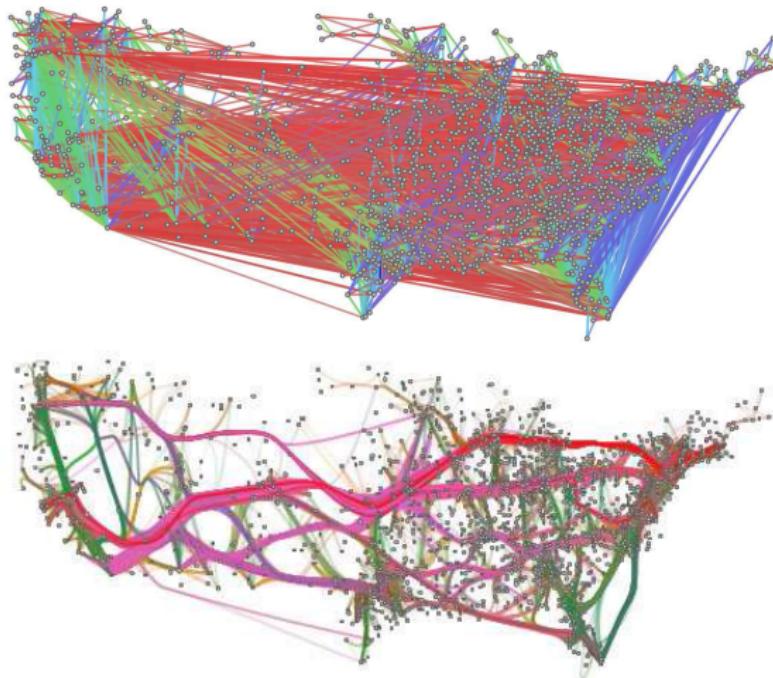
Reordering

# Technique (4): Clutter Reduction

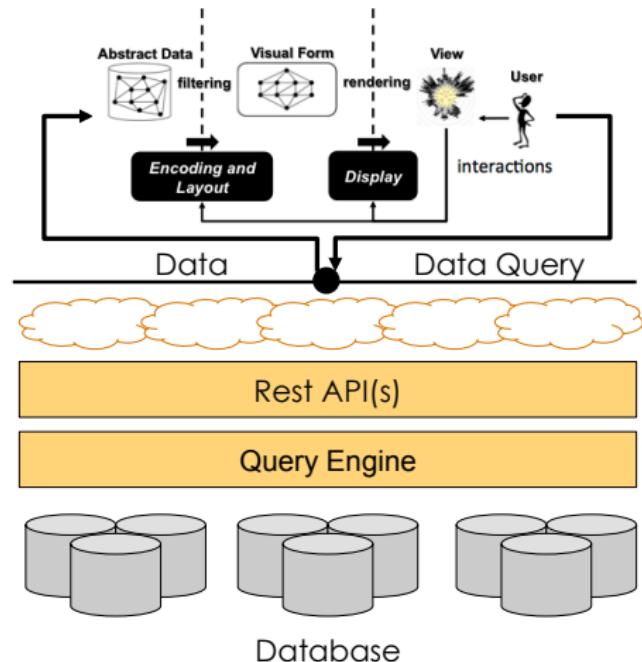
Edge Bundling



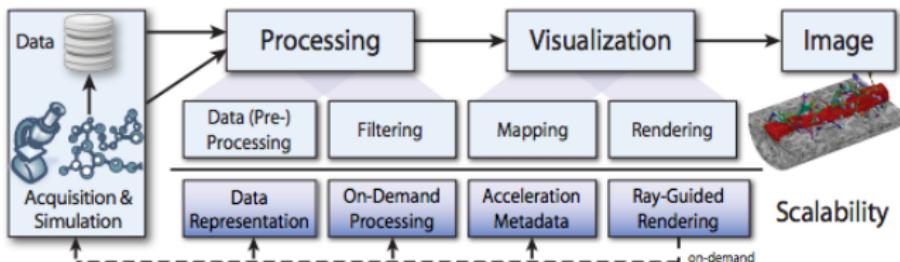
## Technique (4): Clutter Reduction



# Technique (4): Query based Visualization



# Technique (5): Parallel Computing via GPU or GUGPU

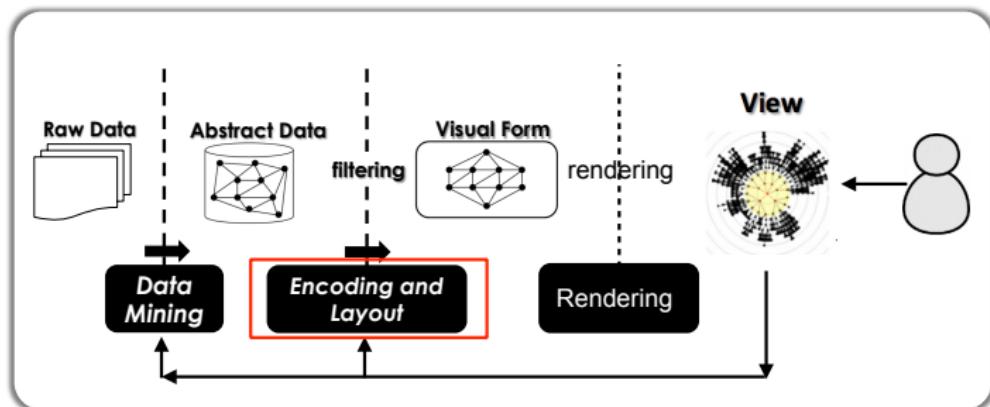


A Survey of GPU-Based Large-Scale Volume Visualization, EuroVis, 2014

## Part III : How can we visualize big data ?



# Visualization & Visual Analysis Reference Model



Encoding : Visual Design

Technique : Layout Algorithm



# Using existing tools are easy

D3.js  
Data-Driven Documents



Tableau

ManyEyes

## Python:

iGraph : <http://igraph.org/redirect.html>  
Networkx : <https://networkx.github.io/>

## JavaScript:

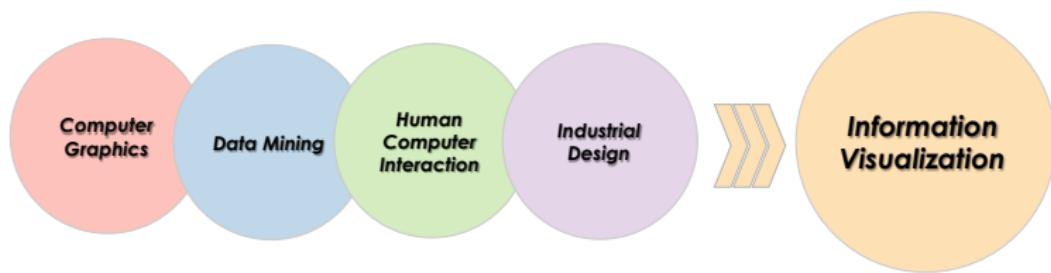
D3.js (2D, SVG): <http://d3js.org/>  
Tree.js (3D, WebGL): <http://threejs.org/>

## Java:

prefuse: <http://prefuse.org/>  
InofVis Toolkit: <http://ivtk.sourceforge.net/>

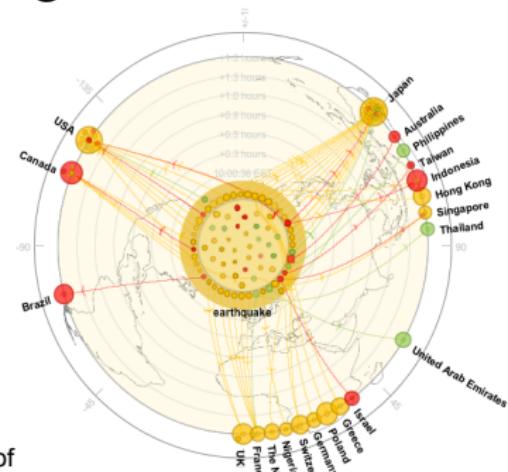


Developing new ones require knowledge from different areas



# PART III : How can we visualize big data?

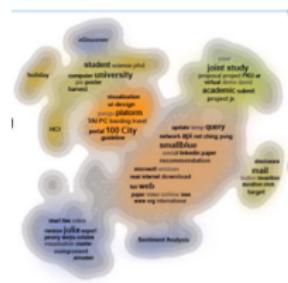
## Example 1: Visualising Streaming Data



Whisper: Tracing the Spatiotemporal Process of  
Information Diffusion in Real Time  
IEEE InfoVis 2012

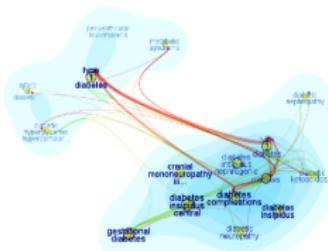
# PART III : How can we visualize big data?

## Example 2: Visualizing Large Text Corpus



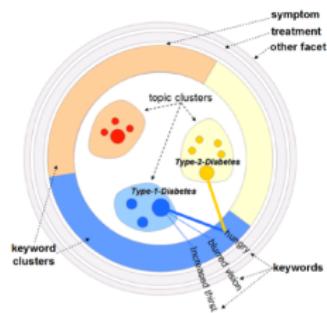
ContextTour  
SDM 2010

**Visualizing  
Heterogeneous Clusters**



FacetAtlas  
TVCG (InfoVis 2010)

**Visualizing  
Multi-relational Clusters**



SolarMap  
ICDM 2011

**Cluster  
Interpretation**

## Part IV : Visual Analysis of Big Data



# Visual Analysis v.s. Data Mining

Computational Power



Data Mining

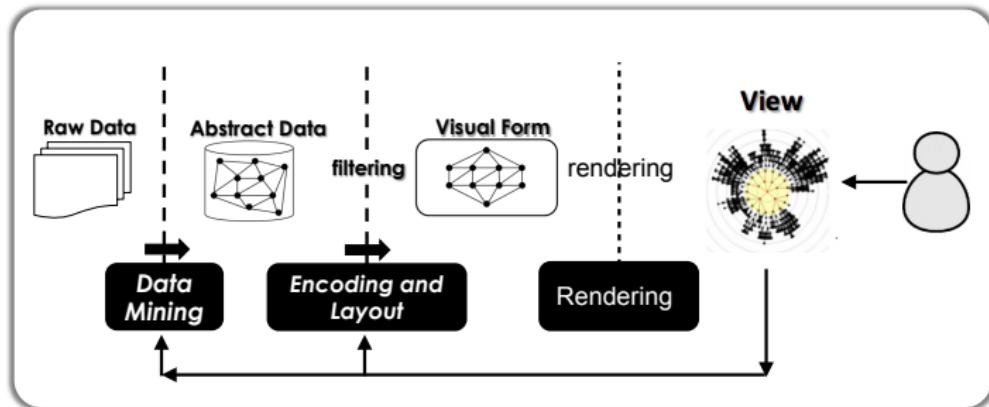
+

Human Intelligence



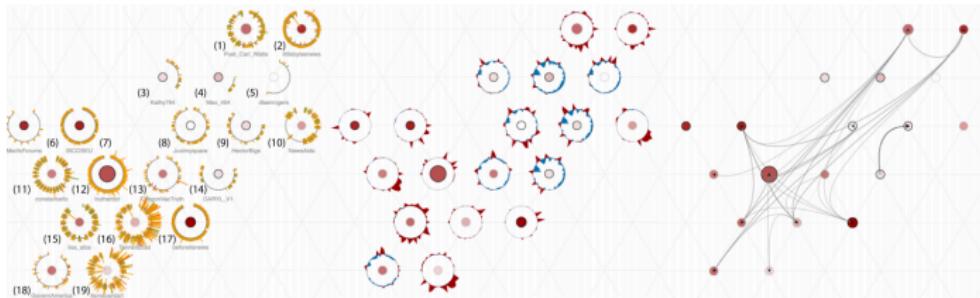
Visual Analysis

# Visualization & Visual Analysis Reference Model



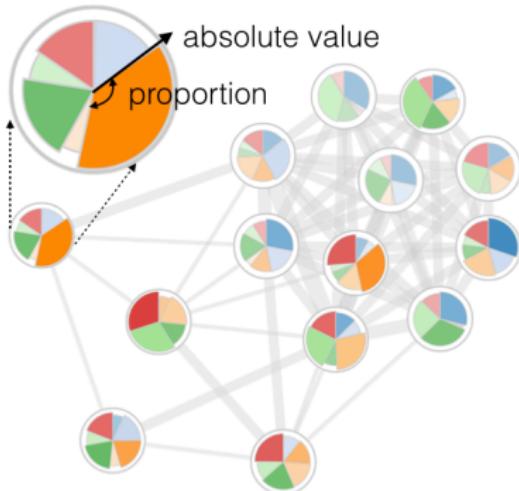
Analysis + Visualisation + Interaction

## Example 3: Detect Anomalous Users in Twitter



TargetVue: Visual Analysis of Anomalous User Behaviors in Online Communication Systems, IEEE Transactions on Visualisation and Computer Graphics (VAST'15)

## Example 4: Visualizing Large Graphs



g-Miner: Interactive Visual Group Mining on Multivariate Graphs, ACM CHI 2015



## The Top 10 Challenges in Extreme-Scale Visual Analytics

**Pak Chung Wong**

*Pacific Northwest National Laboratory*

**Han-Wei Shen**

*Ohio State University*

**Christopher R. Johnson**

*University of Utah*

**Chaomei Chen**

*Drexel University*

**Robert B. Ross**

*Argonne National Laboratory*

Editor:  
Theresa-Marie Rhyne

Wong, P. C., Shen, H. W., Johnson, C. R., Chen, C., & Ross, R. B. (2012). The top 10 challenges in extreme-scale visual analytics. *IEEE computer graphics and applications*, 32(4), 63.

