# Visual Tools & Methods for Data Cleaning
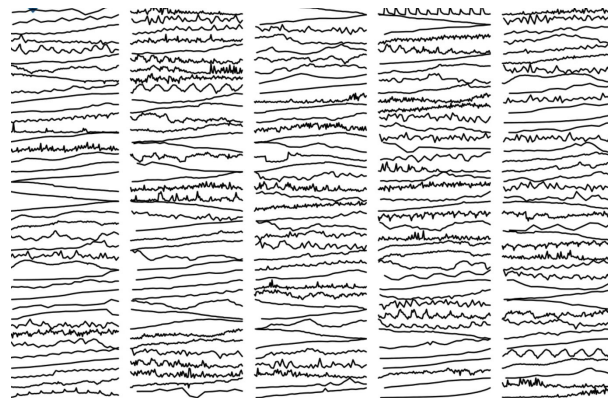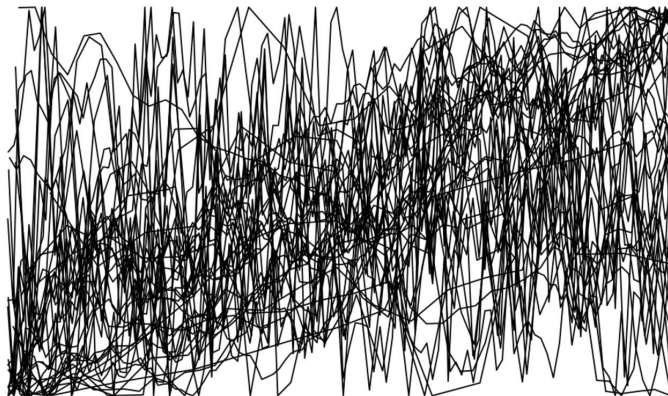
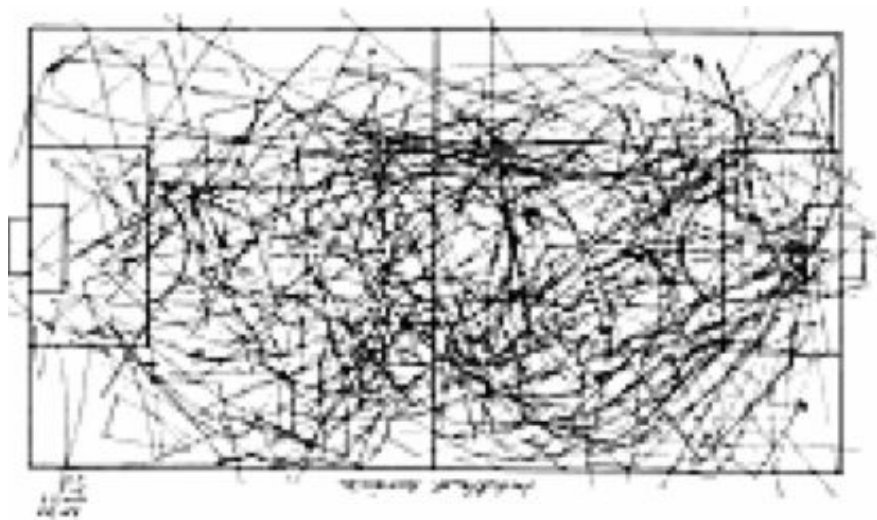## Lyon1 M2 Dataviz - 2018/2019- Cours #4

http://romain.vuillemot.net/

@romsson

# Reality

* Time series



* Geo-spatial data

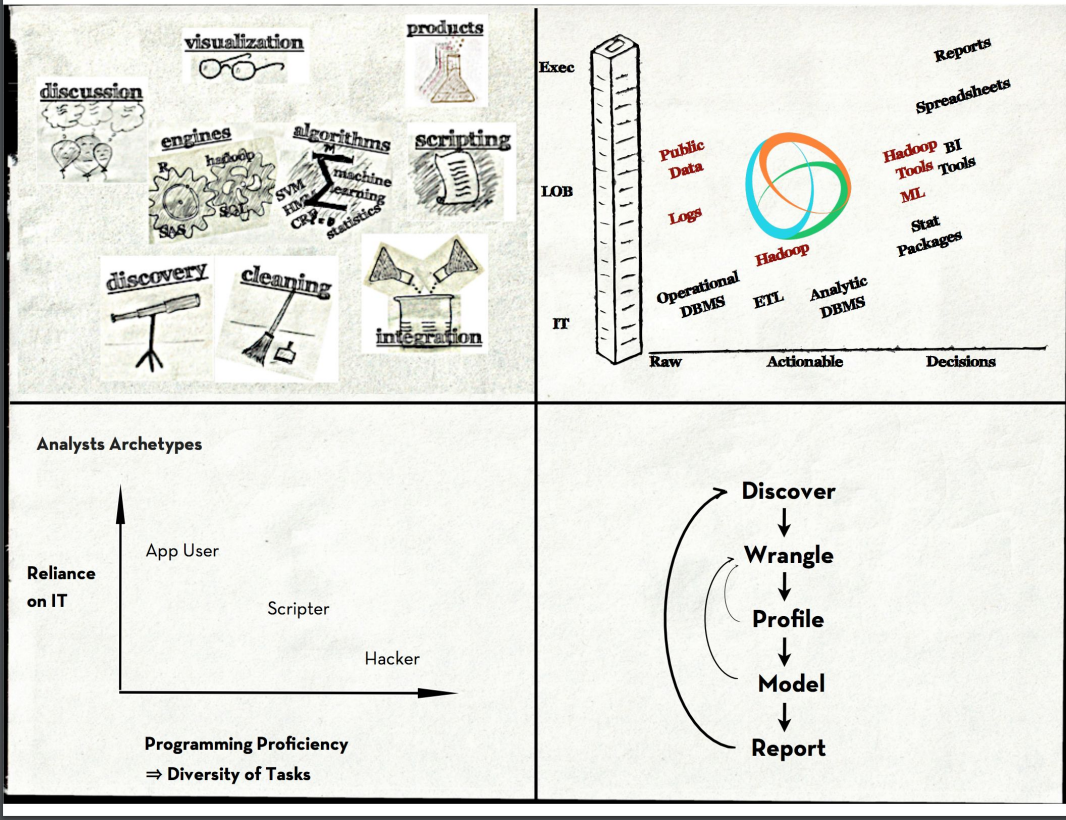# Need for interaction!

# Need for interaction **with Raw Data**

Data transformations

Visual mapping

View transformation

Rendering

Raw data

Processed data

Abstract visual form

Visual presentation

Physical presentation

# Empirical study

35 data analysts, 25 organizations, 15 sectors



Kandel, Sean, et al. "Enterprise data analysis and visualization: An interview study." IEEE Transactions on Visualization and Computer Graphics 18.12 (2012): 2917-2926. (pdf)

# Empirical study



Joe Hellerstein "Data wrangling" BERKELEY & Trifacta (pdf)

# Wrangling and analysis process

* Iterative, non-linear process

Kandel, S., Heer, J., Plaisant, C., Kennedy, J., van Ham, F., Riche, N. H., & Buono, P. (2011). "Research directions in data wrangling: Visualizations and transformations for usable and credible data. Information Visualization"

# Microsoft Excel



| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | Reported crime in Alabama, | | | | | | | |
| 2 | , | | | | | | | |
| 3 | 20,044,029.30 | | | | | | | |
| 4 | 20,053,900 | | | | | | | |
| 5 | 20,063,937 | | | | | | | |
| 6 | 20,073,974.90 | | | | | | | |
| 7 | 20,084,081.90 | | | | | | | |
| 8 | , | | | | | | | |
| 9 | Reported crime in Alaska, | | | | | | | |
| 10 | , | | | | | | | |
| 11 | 20,043,370.90 | | | | | | | |
| 12 | 20,053,615 | | | | | | | |
| 13 | 20,063,582 | | | | | | | |
| 14 | 20,073,373.90 | | | | | | | |
| 15 | 20,082,928.30 | | | | | | | |
| 16 | , | | | | | | | |
| 17 | Reported crime in Arizona, | | | | | | | |
| 18 | , | | | | | | | |
| 19 | 20,045,073.30 | | | | | | | |
| 20 | 20,054,827 | | | | | | | |
| 21 | 20,064,741.60 | | | | | | | |
| 22 | 20,074,502.60 | | | | | | | |

# Python Notebook

File   Edit   View   Insert   Cell   Kernel   Help

Trusted

Code

```python
In [13]:  import pandas as pd

          df = pd.read_csv("data/crime.csv", sep=',', header=None)
```

```python
In [16]:  df.head(10)
```

Out[16]:

|   | 0 | 1 |
|---|---|---|
| 0 | Reported crime in Alabama | NaN |
| 1 | NaN | NaN |
| 2 | 2004 | 4029.3 |
| 3 | 2005 | 3900.0 |
| 4 | 2006 | 3937.0 |
| 5 | 2007 | 3974.9 |
| 6 | 2008 | 4081.9 |
| 7 | NaN | NaN |
| 8 | Reported crime in Alaska | NaN |
| 9 | NaN | NaN |

```python
In [18]:  df.loc[df[0].isin(["Reported crime in Alabama⎯→"])]
```
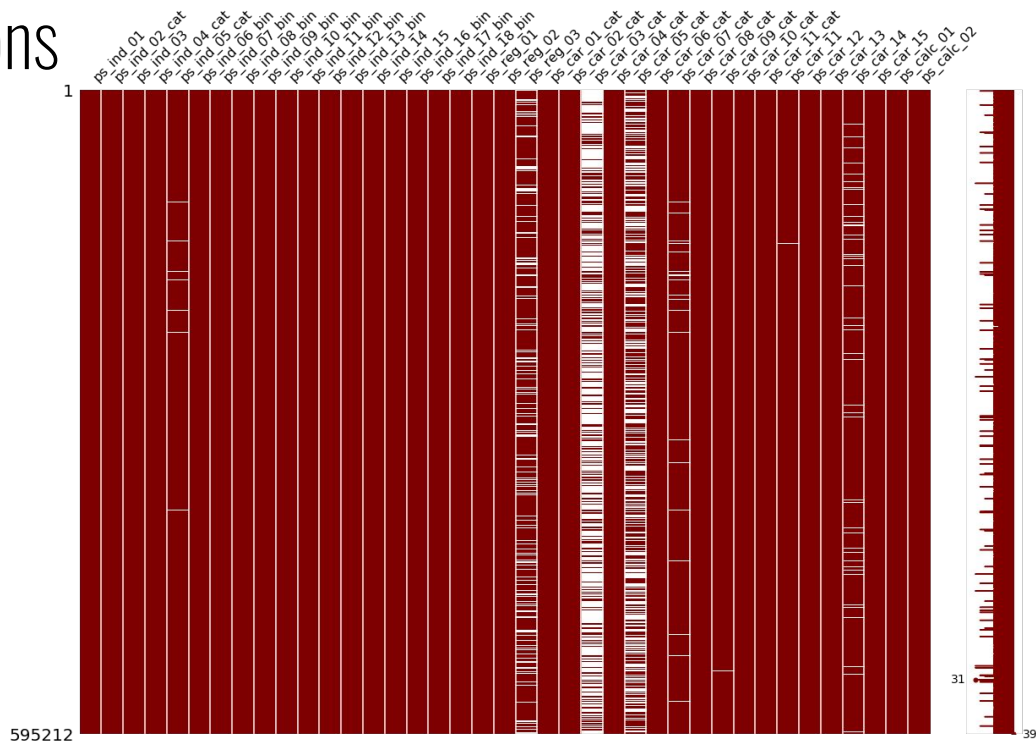
# Low-level scripts & visualizations

* Python / Perl / ..

* Pipeline / Batch process

* ...



Example:
SafeDriver - data cleaning & visualization ([webpage](webpage))

# Google / Open Refine (2010 - ...)

* Loading
* Checking
* Exploring
* Cleaning
* Reshaping
* Annotating
* Saving



https://github.com/OpenRefine/OpenRefine

A Quick Tour of OpenRefine (slides)

# Wrangler

Kandel, S., Paepcke, A., Hellerstein, J., & Heer, J. (2011, May). Wrangler: Interactive visual specification of data transformation scripts. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (pp. 3363-3372). ACM. (demo)

# Profiler



Kandel, S., Parikh, R., Paepcke, A., Hellerstein, J. M., & Heer, J. (2012, May). Profiler: Integrated statistical analysis and visualization for data quality assessment. In Proceedings of the International Working Conference on Advanced Visual Interfaces (pp. 547-554). ACM. (pdf)

# Profiler



(a) Data Table
(b) Type Inference
(c) Feature Generation
(d) Anomaly Detection — anomalies
(e) View Recommendation — derived, anomalies
(f) Interactive Visualization

| Type | Issue | Detection Method(s) | Visualization |
|---|---|---|---|
| Missing | Missing record | Outlier Detection \| Residuals then Moving Average w/ Hampel X84 | Histogram, Area Chart |
| | | Frequency Outlier Detection \| Hampel X84 | Histogram, Area Chart |
| | Missing value | Find NULL/empty values | Quality Bar |
| Inconsistent | Measurement units | Clustering \| Euclidean Distance | Histogram, Scatter Plot |
| | | Outlier Detection \| z-score, Hampel X84 | Histogram, Scatter Plot |
| | Misspelling | Clustering \| Levenshtein Distance | Grouped Bar Chart |
| | Ordering | Clustering \| Atomic Strings | Grouped Bar Chart |
| | Representation | Clustering \| Structure Extraction | Grouped Bar Chart |
| | Special characters | Clustering \| Structure Extraction | Grouped Bar Chart |
| Incorrect | Erroneous entry | Outlier Detection \| z-score, Hampel X84 | Histogram |
| | Extraneous data | Type Verification Function | Quality Bar |
| | Misfielded | Type Verification Function | Quality Bar |
| | Wrong physical data type | Type Verification Function | Quality Bar |
| Extreme | Numeric outliers | Outlier Detection \| z-score, Hampel X84, Mahalanobis distance | Histogram, Scatter Plot |
| | Time-series outliers | Outlier Detection \| Residuals vs. Moving Average then Hampel X84 | Area Chart |
| Schema | Primary key violation | Frequency Outlier Detection \| Unique Value Ratio | Bar Chart |

Kandel, S., Parikh, R., Paepcke, A., Hellerstein, J. M., & Heer, J. (2012, May). Profiler: Integrated statistical analysis and visualization for data quality assessment. In Proceedings of the International Working Conference on Advanced Visual Interfaces (pp. 547-554). ACM. (pdf)

# Trifacta

# Visualization!

"year","value","state"
"2004","4029.3","Alabama"
"2005","3900","Alabama"
"2006","3937","Alabama"
"2007","3974.9","Alabama"
"2008","4081.9","Alabama"
"2004","3370.9","Alaska"
"2005","3615","Alaska"
"2006","3582","Alaska"
"2007","3373.9","Alaska"
"2008","2928.3","Alaska"
"2004","5073.3","Arizona"
"2005","4827","Arizona"
"2006","4741.6","Arizona"
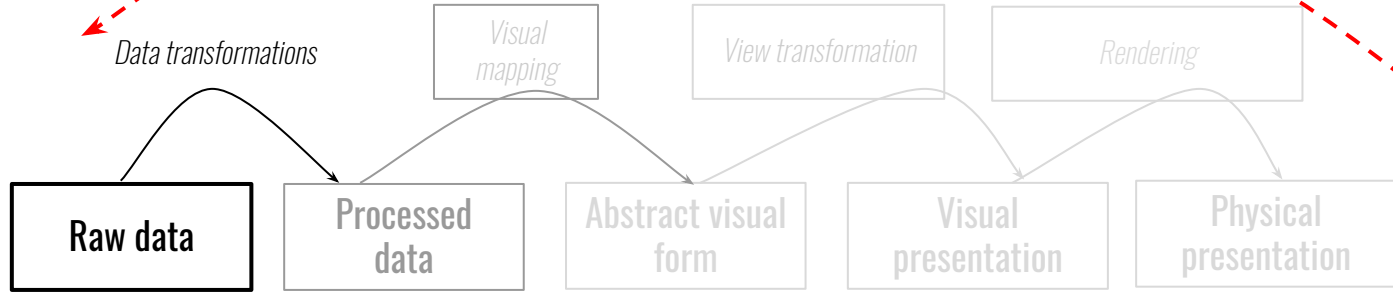"2007","4502.6","Arizona"
"2008","4087.3","Arizona"

Expected visualization (demo)

Reality (demo)

D3.js https://d3js.org/

# Need for interaction **with Raw Data**



Data transformations

Visual mapping

View transformation

Rendering

Raw data

Processed data

Abstract visual form

Visual presentation

Physical presentation

# Visualization!



Tableau Software

# Summary

Programming by demonstration

Data sampling progress

Data distribution

Data quality progress bar

Undo!

Export



Preview transformation application

Data samples as table

Suggested transformations using a declarative language

Data read so far
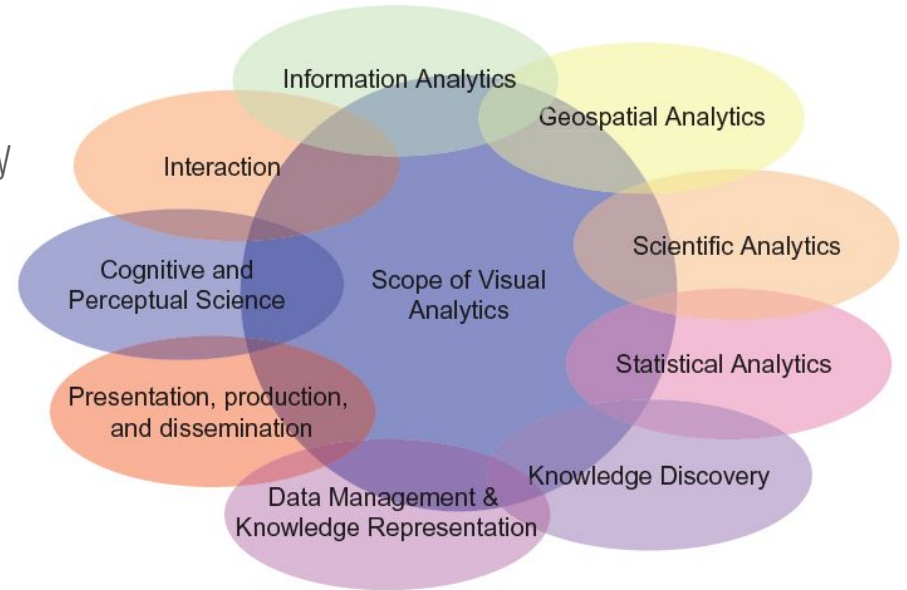
# Research directions

Abedjan, Z., Chu, X., Deng, D., Fernandez, R. C., Ilyas, I. F., Ouzzani, M., ... & Tang, N. (2016). Detecting Data Errors: Where are we and what needs to be done?. Proceedings of the VLDB Endowment, 9(12), 993-1004.

Kandel, S., Heer, J., Plaisant, C., Kennedy, J., van Ham, F., Riche, N. H., ... & Buono, P. (2011). Research directions in data wrangling: Visualizations and transformations for usable and credible data. Information Visualization, 10(4), 271-288.

Chu, X., Ilyas, I. F., Krishnan, S., & Wang, J. (2016, June). Data cleaning: Overview and emerging challenges. In Proceedings of the 2016 International Conference on Management of Data (pp. 2201-2206). ACM.
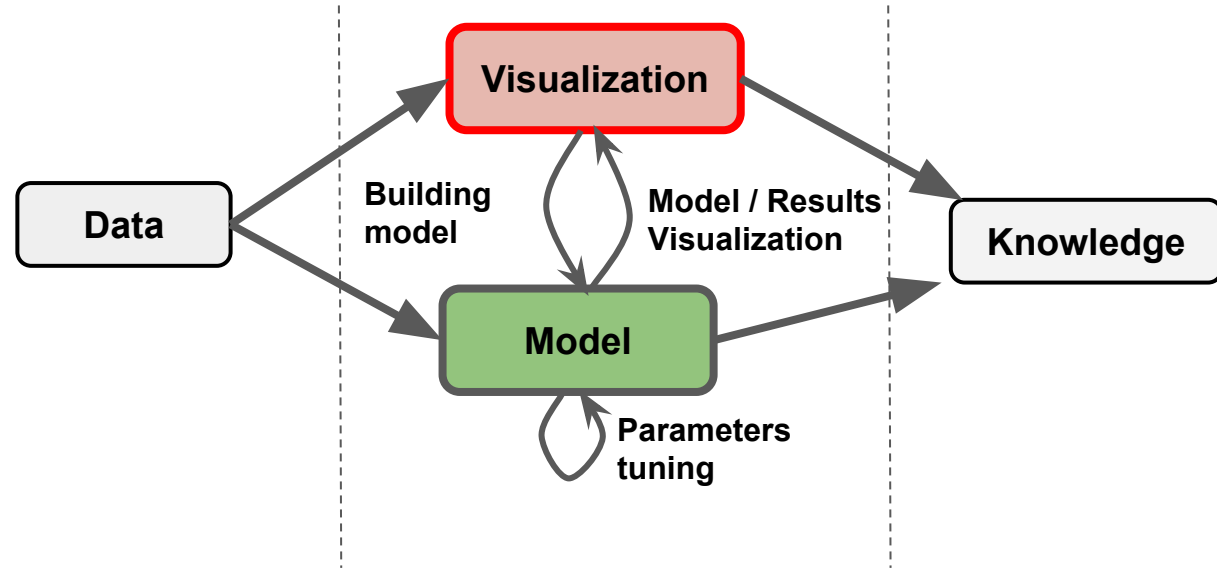
# "Combine with visual analytics" [Kandel, 2011]

"Data wrangling also constitutes a promising direction for visual analytics research, as it requires combining automated techniques (e.g. discrepancy detection, entity resolution, semantic data type inference) with interactive visual interfaces"
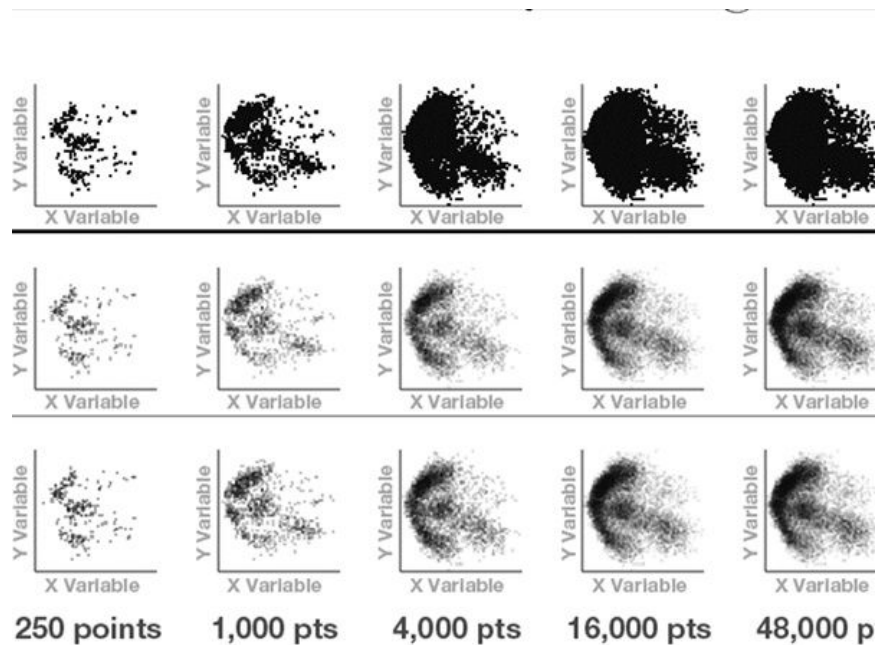


http://www.infovis-wiki.net/index.php?title=File:Keim06visual-analytics-disciplines.png

# Visual Analytics

"

The science of analytical reasoning facilitated by interactive visual interfaces.

"



Thomas, J., Cook, K.: Illuminating the Path: Research and Development Agenda for Visual Analytics. IEEE-Press (2005)"

# "Better Data *Exploration* tools (rather than *communication* tools)"



| 250 points | 1,000 pts | 4,000 pts | 16,000 pts | 48,000 p |

Matejka, Justin, Fraser Anderson, and George Fitzmaurice. "Dynamic opacity optimization for scatter plots." Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems. ACM, 2015. (pdf)
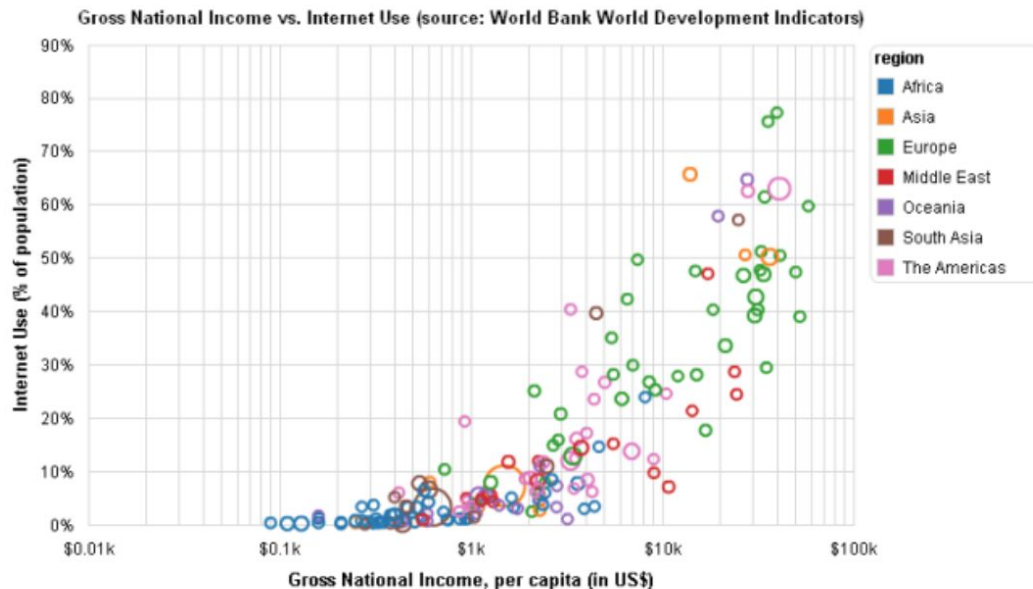
# "Combine with query relaxation"

* We interact with **pixels**
Ex: brushing/selection


    X < 300px && X > 600px
&&   Y > 400px && Y < 700px


* Turn pixels into semantic



Heer, Jeffrey, Maneesh Agrawala, and Wesley Willett. _"Generalized selection via interactive query relaxation."_
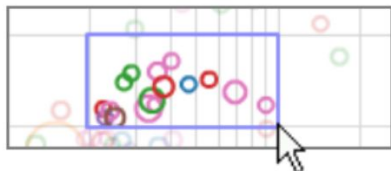Proceedings of the SIGCHI Conference on Human Factors in Computing Systems. ACM, 2008. (pdf)
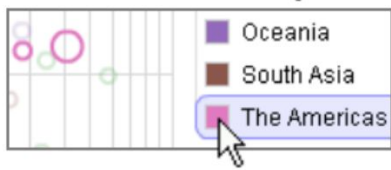
# "Combine with query relaxation"
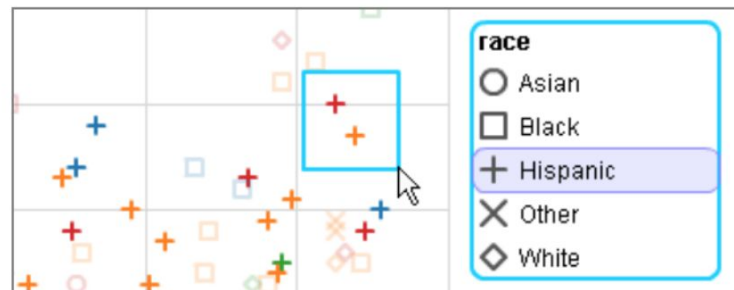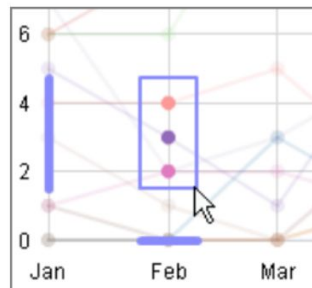


**Item Selection by Clicking**

$(id = 'China')$

**Range Selection by Dragging**

$(2000 < gni \text{ AND } gni < 10000)$ AND
$(.1 < internet \text{ AND } internet < .2)$

**Attribute Selection with Legends**

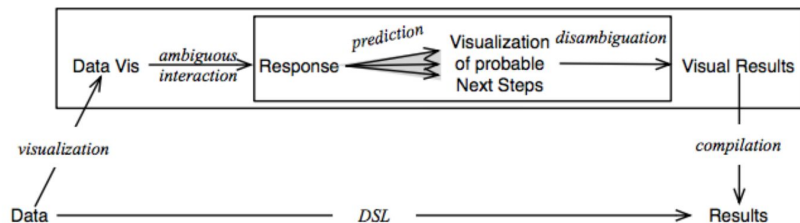$(region = 'The Americas')$

Oceania
South Asia
The Americas

Heer, Jeffrey, Maneesh Agrawala, and Wesley Willett. *"Generalized selection via interactive query relaxation."* Proceedings of the SIGCHI Conference on Human Factors in Computing Systems. ACM, 2008. (pdf)

# "Guide users exploratory process"



Demiralp, Ç., Haas, P. J., Parthasarathy, S., & Pedapati, T. (2017). Foresight: Rapid Data Exploration Through Guideposts.

# "Predict next interaction"





Heer, Jeffrey, Joseph M. Hellerstein, and Sean Kandel. "Predictive Interaction for Data Transformation." CIDR. 2015.

# "Support history exploration"
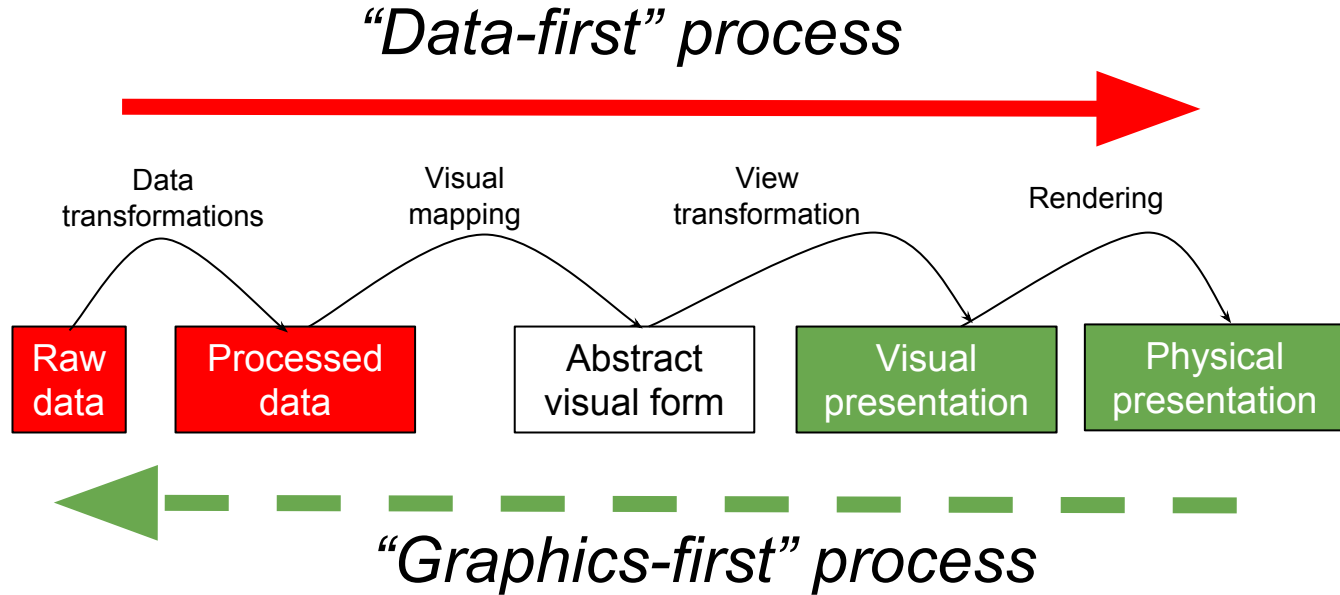


Dunne, C., Henry Riche, N., Lee, B., Metoyer, R., & Robertson, G. (2012, May). GraphTrail: Analyzing large multivariate, heterogeneous networks while supporting exploration history. In Proceedings of the SIGCHI conference on human factors in computing systems (pp. 1663-1672). ACM.

# "Help users recall their reasoning process"
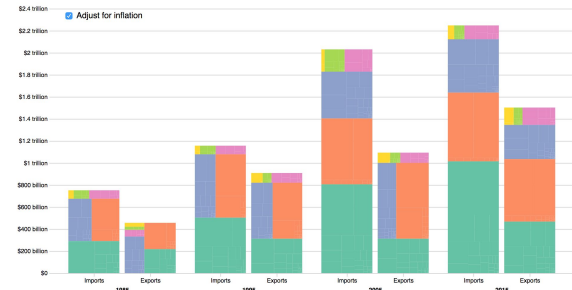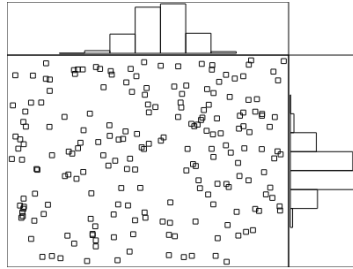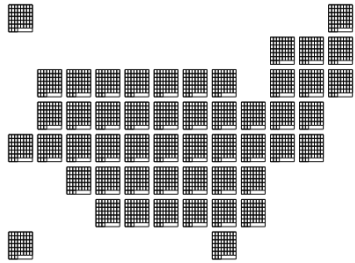


Lipford, H. R., Stukes, F., Dou, W., Hawkins, M. E., & Chang, R. (2010, October). Helping users recall their reasoning process. In Visual Analytics Science and Technology (VAST), 2010 IEEE Symposium on (pp. 187-194). (pdf).

# "Start working.. without data! (yet)"

## "Data-first" process



| Data transformations | Visual mapping | View transformation | Rendering |
|---|---|---|---|

| Raw data | Processed data | Abstract visual form | Visual presentation | Physical presentation |

## "Graphics-first" process

Vuillemot, Romain, and Jeremy Boy. "Structuring Visualization Mock-ups at the Graphical Level by Dividing the Display Space." IEEE transactions on visualization and computer graphics (2017).

# "Start working.. without data! (yet)"



Vuillemot, Romain, and Jeremy Boy. "Structuring Visualization Mock-ups at the Graphical Level by Dividing the Display Space." IEEE transactions on visualization and computer graphics (2017).

# Future directions

[Abedjan et al., VLDB 2016]

 A holistic combination of tools

A data enrichment system

A novel interactive dashboard.

Reasoning on real-world data

[Chu et al. ICMD 2016]

Scalability

User Engagement

Semi-structured and unstructured data

New Applications for Streaming Data

Growing Privacy and Security Concerns

[Kandel et al. IV 2011] (Among many!)

Living with dirty data

Visualize missing and uncertain data

Adapting systems to tolerate error

Sharing data transformations

Feedback from downstream analysts

# Thanks!