

Visual Analytics

Intro to Information Visualization
Spring, 2019
Alex Endert

Agenda

- Overview of what the term means and how it relates to information visualization
- Example VA systems

Visual Analytics

Definition and Terminology

Before there was VA

- Growing concern from some that infovis was straying from practical, real world analysis problems
- Infovis typically not applied to massive data sets
 - “show me the data” has scalability limitations
- Infovis “competes” with other computational approaches to data analysis
 - Statistics, data mining, machine learning, broadly speaking “analytics”

Recommendations

- Integrate data mining and information visualization
- Allow users to specify what they are seeking
- Recognize that users are situated in a social context
- Respect human responsibility

The InfoVis Pipeline

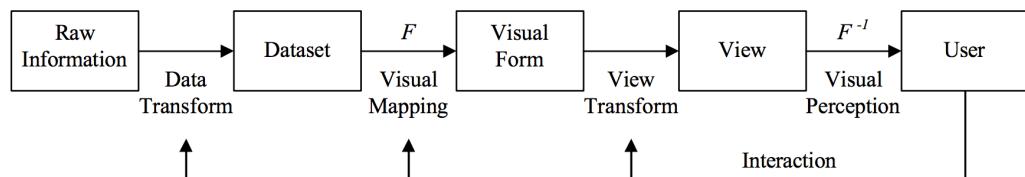
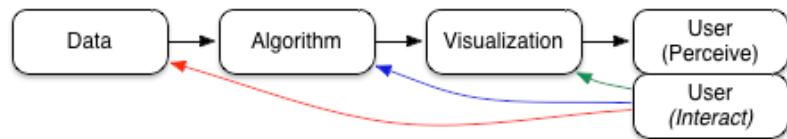
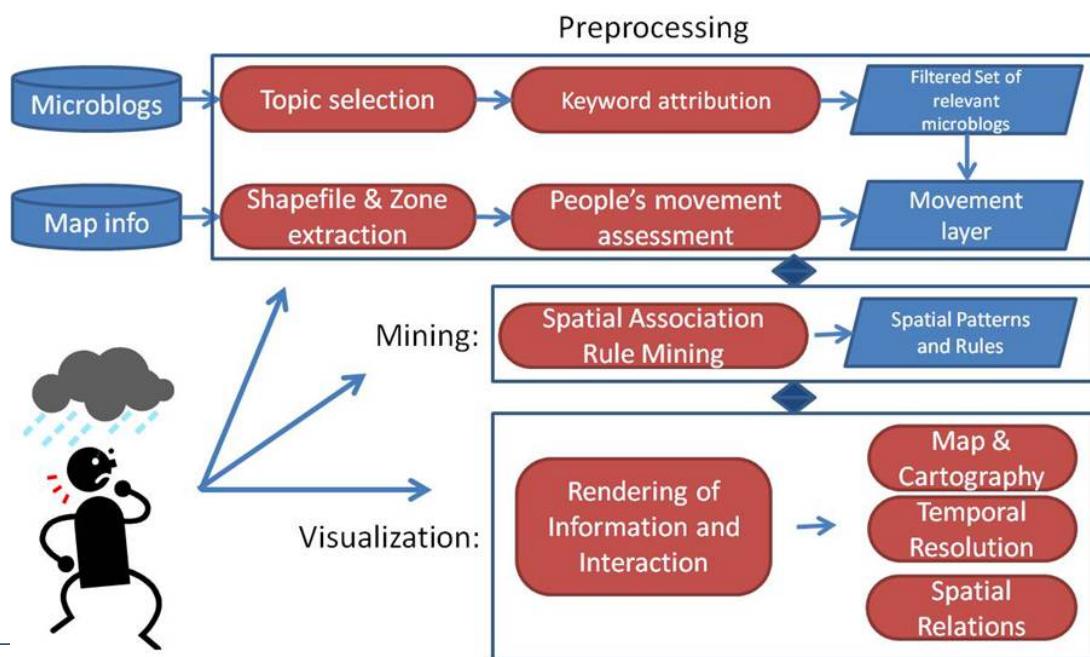


Figure 2: The visualization pipeline, converting information into interactive visual representations (adapted from [Card et al., 1999]).

Visual Analytics Pipeline

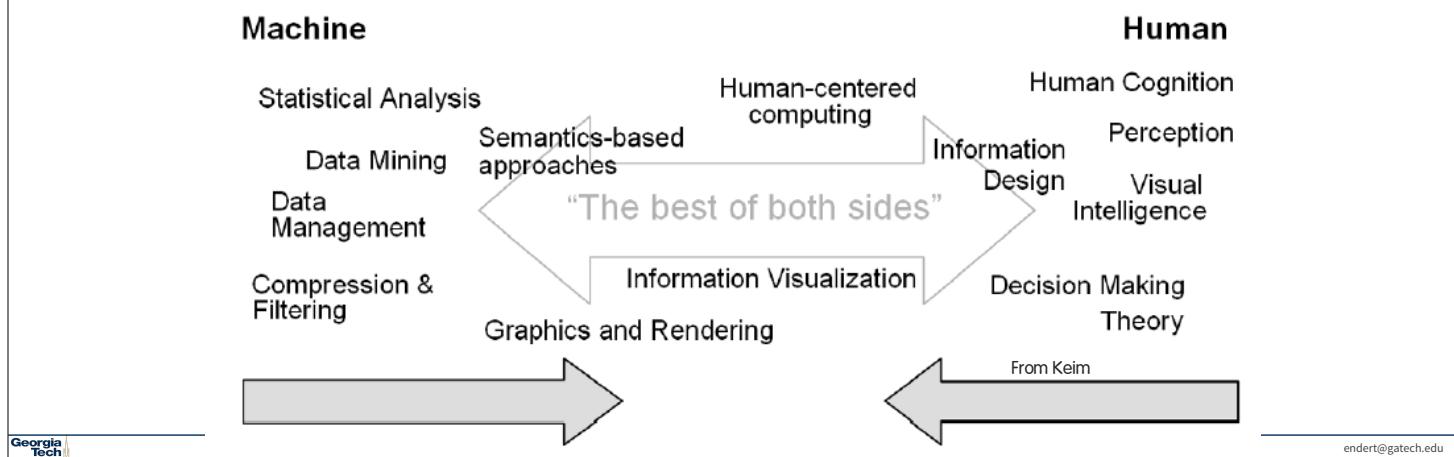


Visual Analytics Pipeline :/



Synergy

- Combine strengths of both human and electronic data processing
 - Gives a semi-automated analytical process
 - Use strengths from each



Integrating Data Mining and ML

- Dimension Reduction
 - PCA, MDS, ...
- Feature Selection
 - Entity extraction, ...
- Similarity functions
 - “find more like this”, ...
- Prediction
 - will the stock price go up or down?

More Motivation

- Increasing occurrences of situations and areas with large data needing better analysis
 - DNA, microarrays, sequence mining
 - Business intelligence
 - prediction

History

- 2003-04 Jim Thomas of PNNL, together with colleagues, develops notion of visual analytics
- Holds workshops at PNNL and at InfoVis '04 to help define a research agenda
- Agenda is formalized in book *Illuminating the Path*, shown on next slide



Visual analytics is the science of analytical reasoning facilitated by interactive visual interfaces.

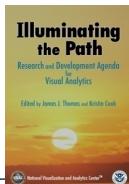
People use visual analytics tools and techniques to

Synthesize information and derive insight from massive, dynamic, ambiguous, and often conflicting data

Detect the expected and discover the unexpected

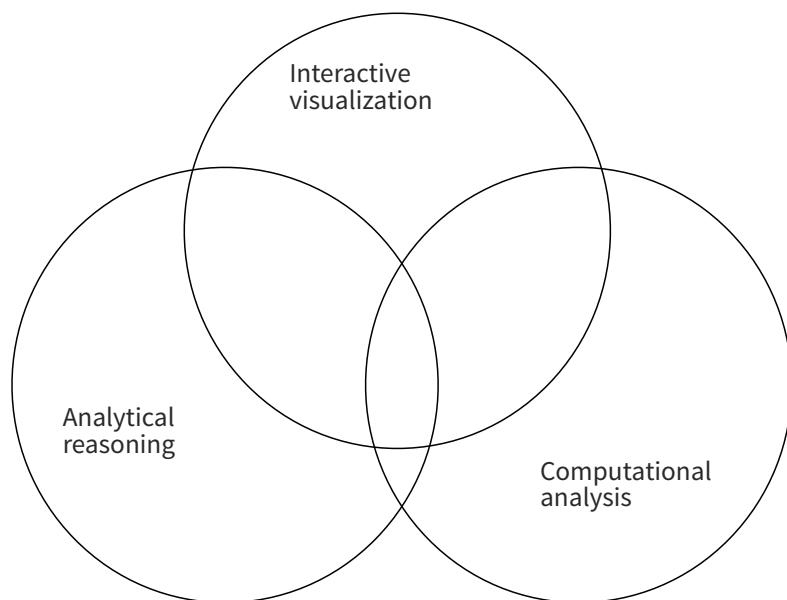
Provide timely, defensible, and understandable assessments

Communicate assessment effectively for action.

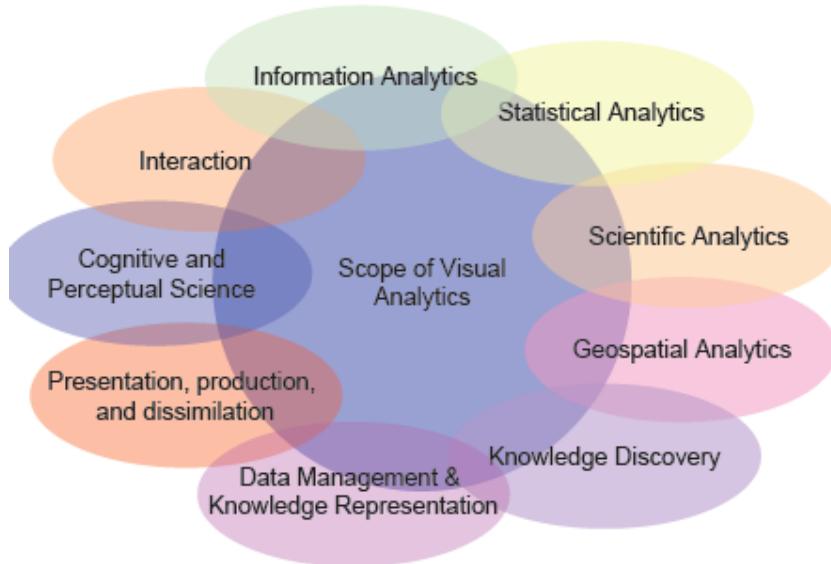


"The beginning of knowledge is the discovery of something we do not understand."
~Frank Herbert (1920 - 1986)

Main Components



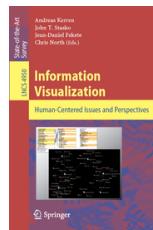
Synergy



From Keim

Alternate Definition

- Visual analytics **combines automated analysis techniques with interactive visualizations** for an effective understanding, reasoning and decision making on the basis of very large and complex data sets

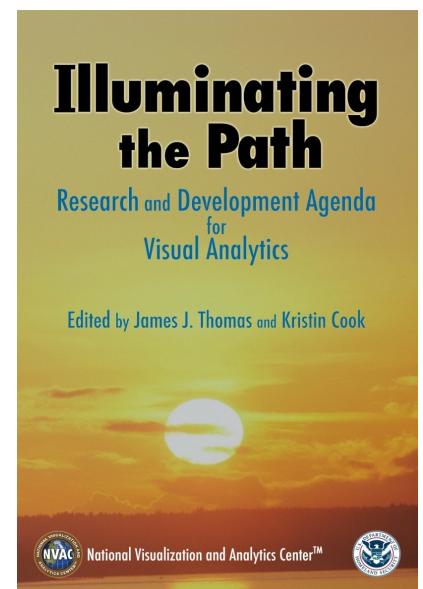


Keim et al, chapter in
Information Visualization:
Human-Centred
Issues and Perspectives, 2008

Visual Analytics

- Encompassing, integrated approach to data analysis
 - Use **computational algorithms where helpful**
 - Use **human-directed visual exploration** where helpful
 - Not just “Apply A, then apply B” though
 - **Integrate the two tightly**

- Available at <http://nvac.pnl.gov/> in PDF form
- At IEEE Press in book form
- Special thanks to IEEE Technical Committee on Visualization and Graphics

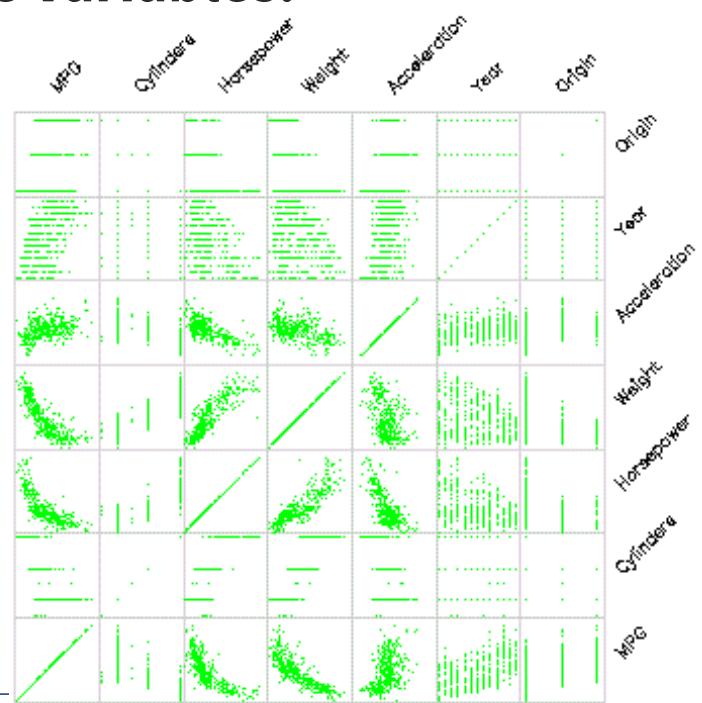


a quick example

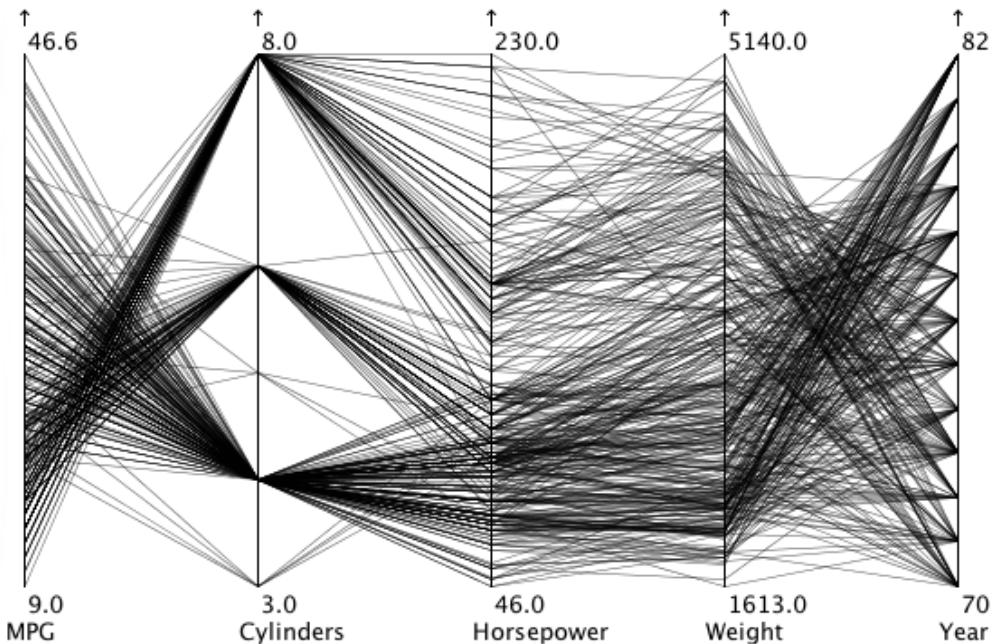
when you have multiple variables:

Represent each possible pair of variables in their own 2-D scatterplot

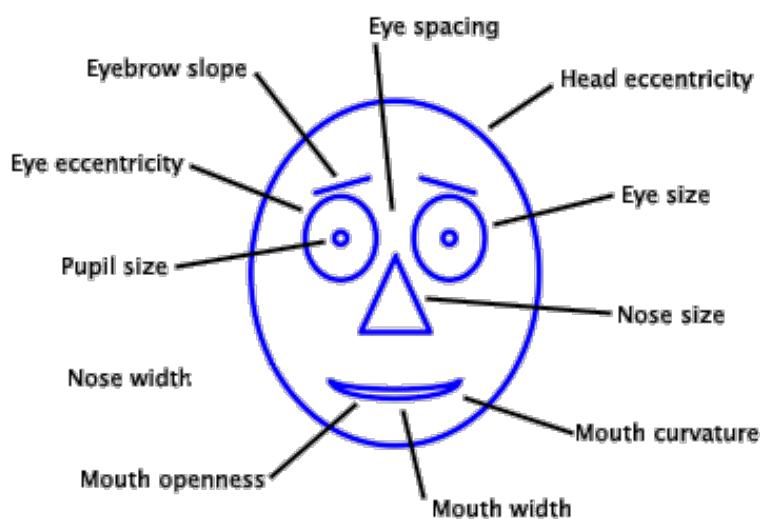
anyone remember the name of this vis technique?



or, Parallel Coordinates



a fun one - Chernov faces



or, use Dimension Reduction!

- computational techniques to solve for similar attributes so that the final number of attributes used to describe the dataset are less than the original
 - this is lossy
- for VA, we often reduce from n, to 2
 - we can show 2 dimensions in a scatterplot
 - which attributes do we discard? we can't show them all!
 - need to find most salient data attributes and show those

There are losses associated with compressing the dimensionality.

Usually reduce the dimension to 2, because that opens up to a whole bunch of useful visualization tools

Principal Component Analysis (PCA)

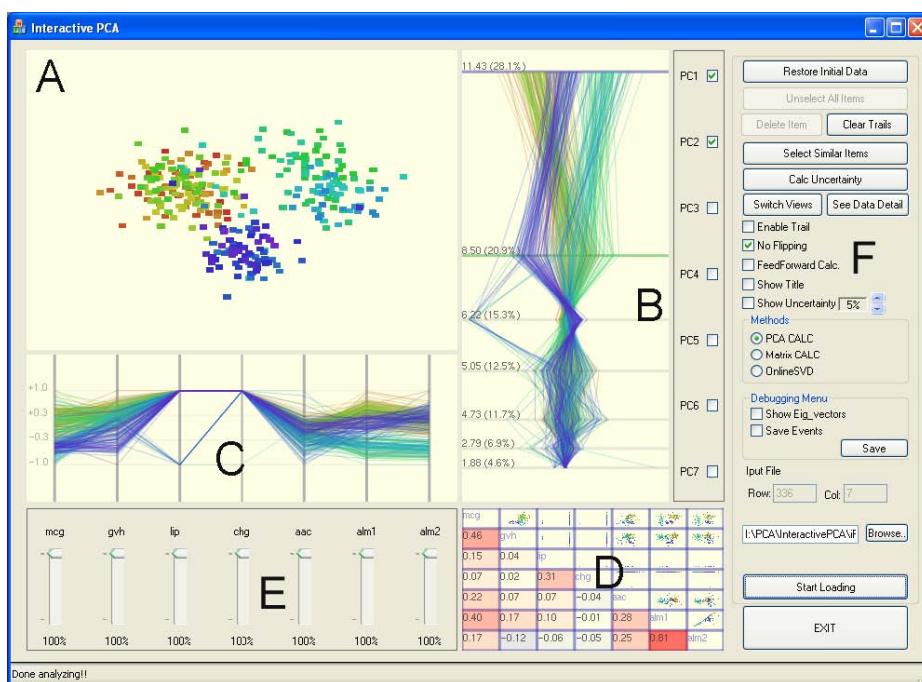
- Frequently used deterministic technique
- transform data table with lots of attributes to principal components, which can be mapped in pairs to x and y axes of scatterplots
 - computed PCs will be \leq original attributes of data
- solve for maximum variance
- i.e., *find similar attributes and combine them*

potential PCA limitations

- no guarantee that resulting PCs are “meaningful”, or most important
 - importance is a human, subjective metric
 - max variance may not be most meaningful
- sensitive to outliers
- linear function, what if data/phenomena is non-linear?

iPCA

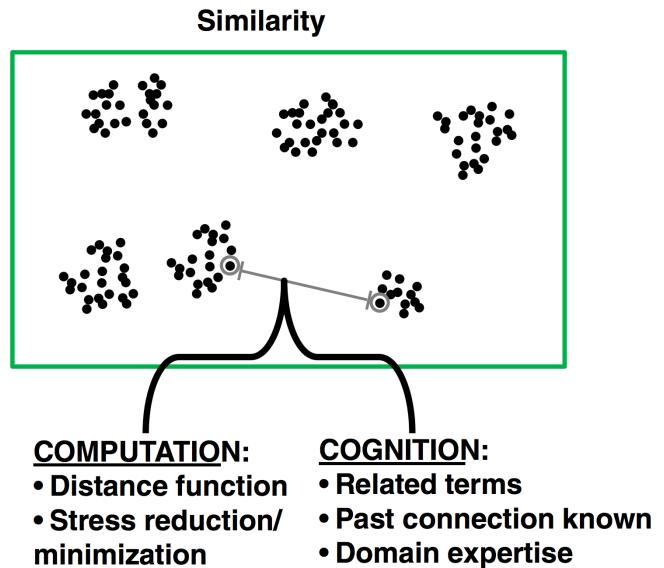
https://youtu.be/oUx-7Hca_i4?t=1m30s



Jeong et al.
CG&A 2009

coupling mental models & data models

- **distance** has meaning to the system and the user
- **to the system:**
 - uses distance functions and similarity metrics
 - preserves distance in high-dimensional feature space in the low-dimensional scatterplot
- **to the user:**
 - proximity approximates similarity
 - we put things nearby if they have some conceptual similarity

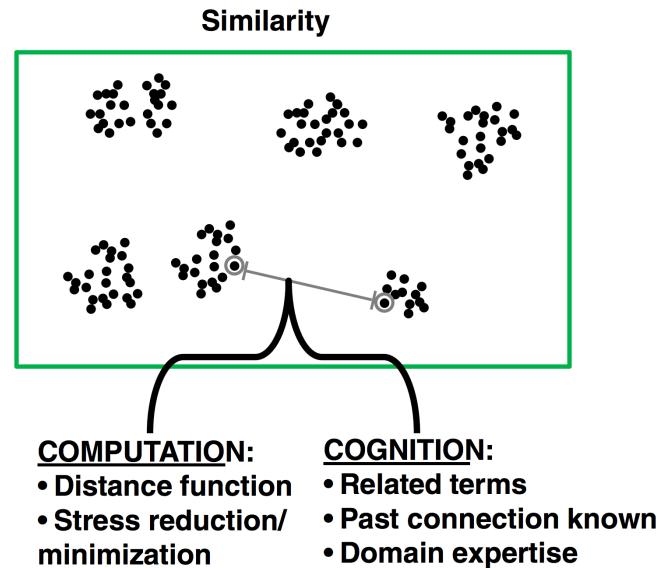


let's try it; InterAxis

- Can control dimension reduction technique using 2 interaction methods
 - directly control the weight of attributes to change the layout of the scatterplot
 - provide example data points on each side of the axes
- <http://va.gatech.edu/endert/projects/>

coupling mental models & data models

- opens an interesting avenue of research
 - how do we take advantage of this connection?
- we'll get back to this opportunity later



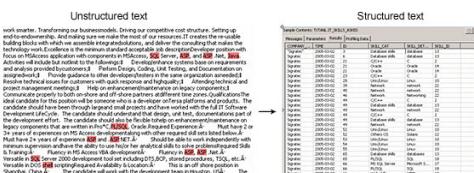
Unstructured
Text

Extract
Entities

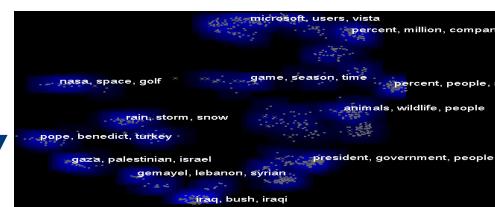
High Dimensional
Data

Dimension Reduction Model

Choice of Visual Encodings/
Metaphor



	BM1	BM2	BM3	BM4	BM5	BM6	BM7	SM1	SM2	SM3	SM4	SM5	SM6	BE	RM
BM1	—	78	62	76	76	68	71	65	79	62	54	65	74	47	55
BM2	0.75	—	54	70	64	63	67	60	61	68	60	60	62	40	37
BM3	0.51	0.51	—	72	94	76	63	53	65	65	67	58	57	—	—
BM4	0.66	0.65	0.63	—	83	78	73	77	75	73	71	84	77	58	57
BM5	0.67	0.57	0.70	0.68	—	94	85	84	87	78	74	83	79	68	63
BM6	0.63	0.63	0.63	0.63	0.63	—	76	74	85	79	79	83	80	64	64
BM7	0.63	0.63	0.55	0.63	0.71	0.63	—	68	63	60	68	74	67	46	53
SM1	0.57	0.56	0.73	0.68	0.70	0.72	0.59	—	92	70	72	75	73	67	55
SM2	0.57	0.62	0.63	0.63	0.63	0.63	0.63	0.53	—	78	73	79	74	67	63
SM3	0.57	0.66	0.62	0.63	0.67	0.68	0.54	0.70	0.68	—	83	80	72	52	52
SM4	0.50	0.59	0.52	0.63	0.64	0.60	0.62	0.65	0.64	0.78	—	74	81	51	53
SM5	0.65	0.59	0.59	0.61	0.61	0.61	0.61	0.60	0.60	0.66	0.66	—	88	60	63
SM6	0.65	0.58	0.59	0.66	0.65	0.70	0.58	0.63	0.62	0.65	0.73	0.76	—	64	61
BE	0.44	0.40	0.55	0.51	0.60	0.61	0.43	0.62	0.61	0.50	0.50	0.56	0.59	—	65
RM	0.51	0.37	0.53	0.54	0.55	0.57	0.49	0.50	0.56	0.50	0.51	0.57	0.56	0.64	—



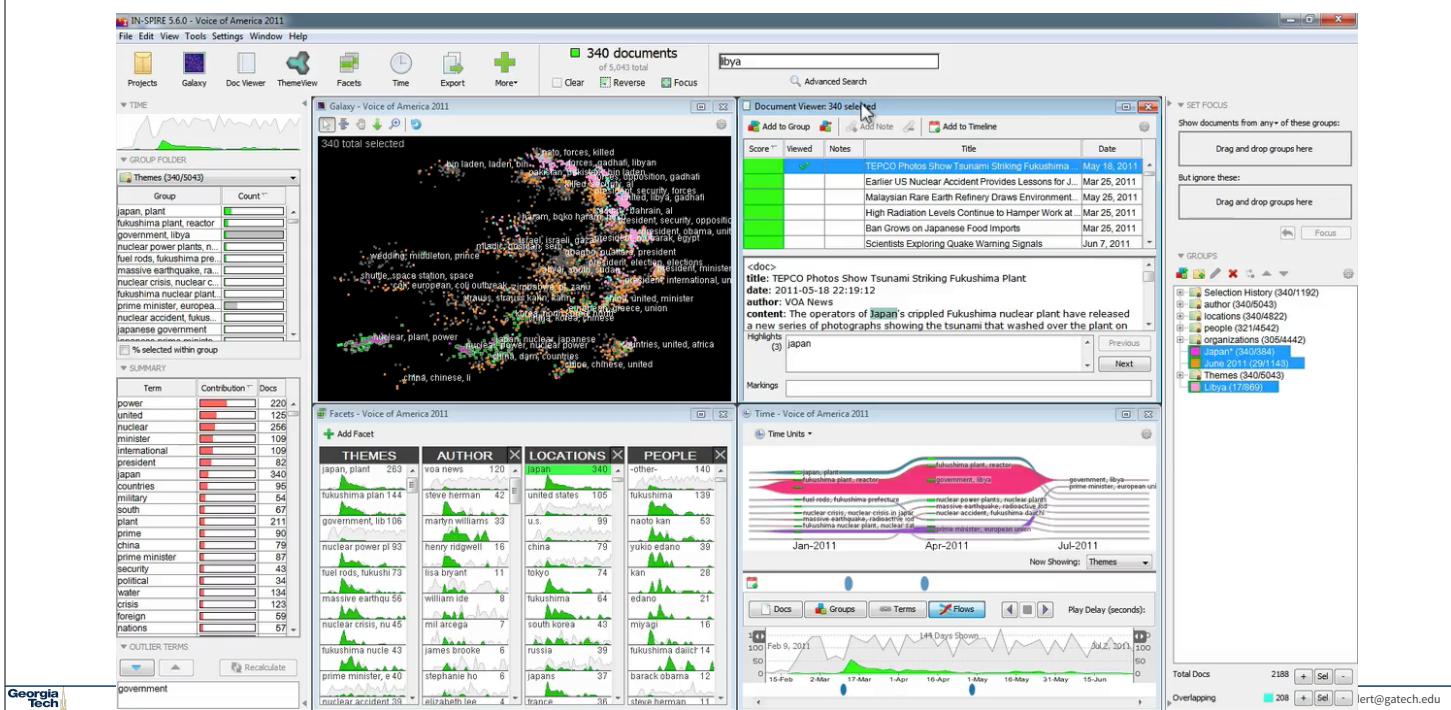
tf-idf for entity selection/extraction

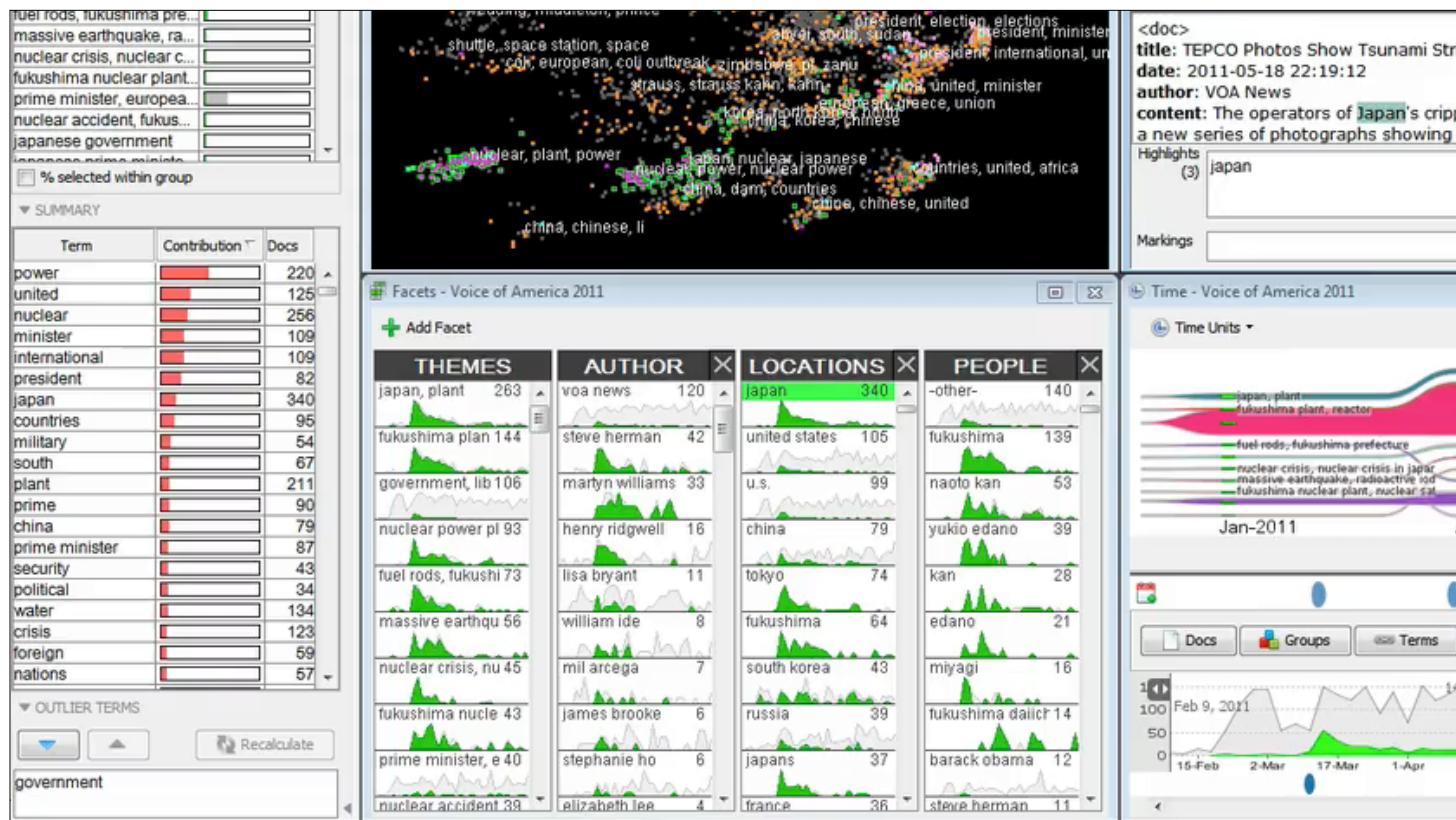
- term frequency-inverse document frequency
 - $TF(t) = (\text{num times term } t \text{ appears in a document}) / (\text{num of terms in the document})$
 - $IDF(t) = \log(\text{Total num of documents} / \text{num of documents with term } t \text{ in it})$
 - Then, $TF \times IDF$
- commonly used way to apply structure and weighting to unstructured text
 - with structure, we can compute and visualize

Georgia Tech Alex Endert endert@gatech.edu

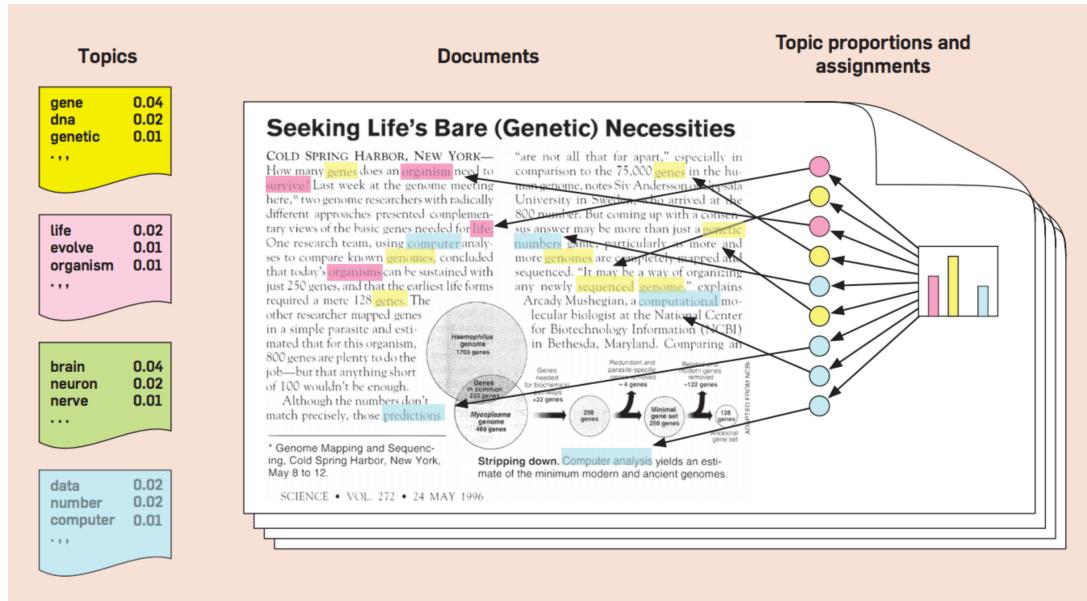
IN-SPiRE

Wise, James A., et al. "Visualizing the non-visual: spatial analysis and interaction with information from text documents." Information Visualization, 1995. Proceedings.. IEEE, 1995.





topic modeling



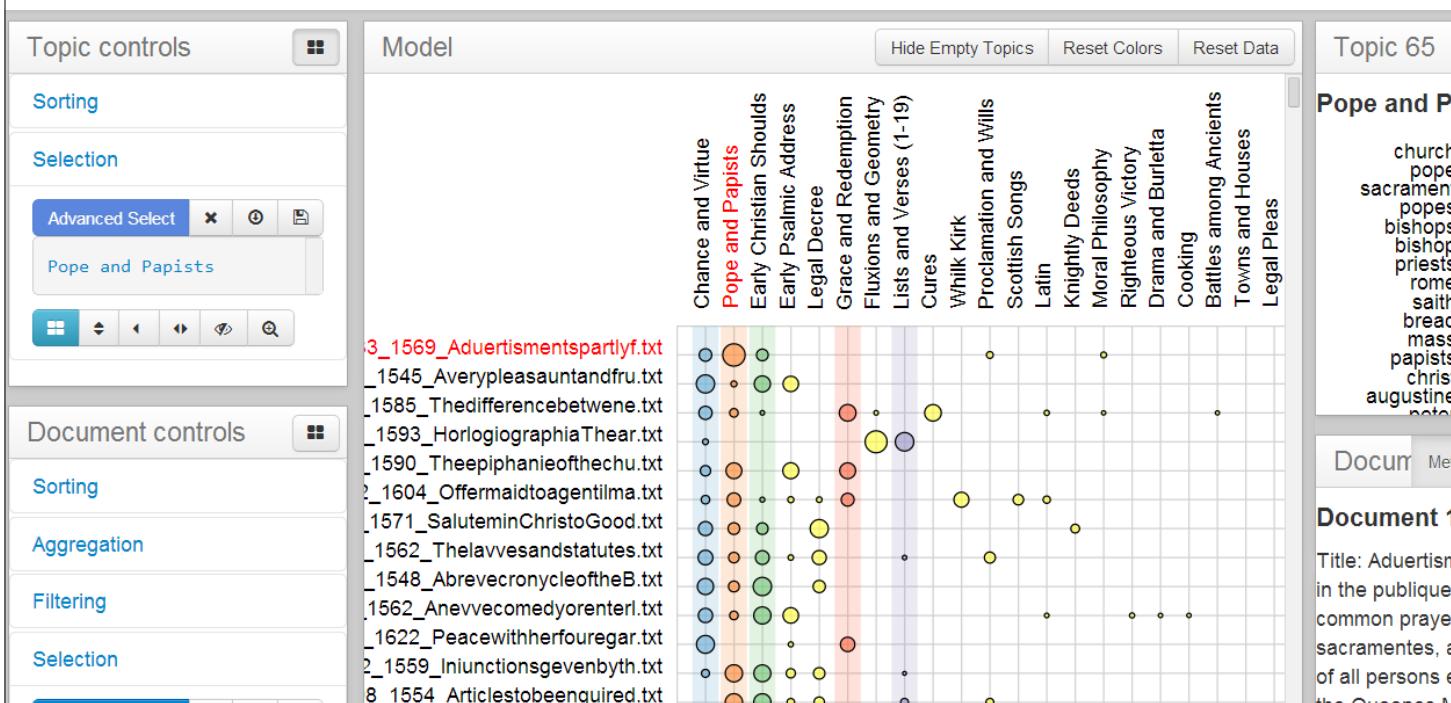
Blei, David M. "Probabilistic topic models." Communications of the ACM 55.4 (2012): 77-84.

Georgia Tech | Alex Endert | endert@gatech.edu

Alexander, Eric, et al. "Serendip: Topic model-driven visual exploration of text corpora." Visual Analytics Science and Technology (VAST), 2014 IEEE Conference on. IEEE, 2014.

Serendip

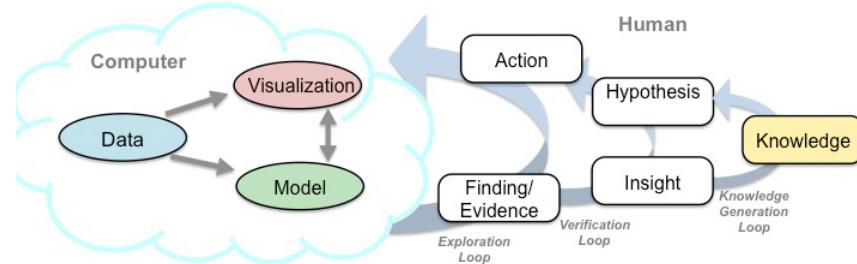
<http://vep.cs.wisc.edu/serendip/>



Incorporating Human Feedback is One Central Component

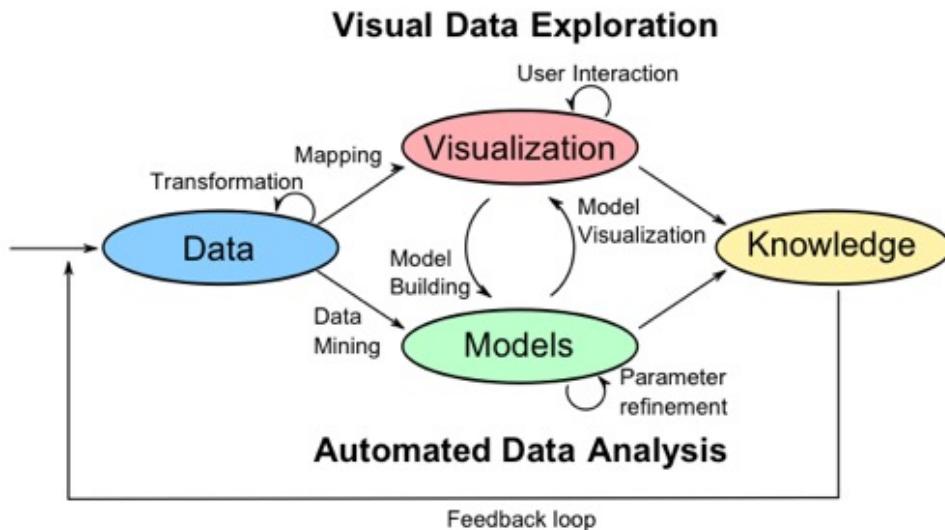
incorporating user feedback

- unsupervised computation is rarely the “final” solution
- “human-in-the-loop” design advocates for including user feedback in the holistic system
- people have domain expertise, can balance subjective tradeoffs, etc.
- let’s look at a few examples of user interaction in VA



Sacha, Dominik, et al. "Knowledge generation model for visual analytics." *Visualization and Computer Graphics, IEEE Transactions on* 20.12 (2014): 1604-1613.

incorporating user feedback



"Mastering the Information Age Solving Problems with Visual Analytics", edited by Daniel Keim, Jörn Kohlhammer, Geoffrey Ellis and Florian Mansmann

Distance = Similarity makes sense to people



Andrews, Christopher, Alex Endert, and Chris North. "Space to think: large high-resolution displays for sensemaking." Proceedings of the SIGCHI Conference on Human Factors in Computing Systems. ACM, 2010.

including ML seems to work well

- helps offload computationally-attractive tasks to machines - nice!
- frees up user to focus cognitive energy on other things - also good
- but we can do better!

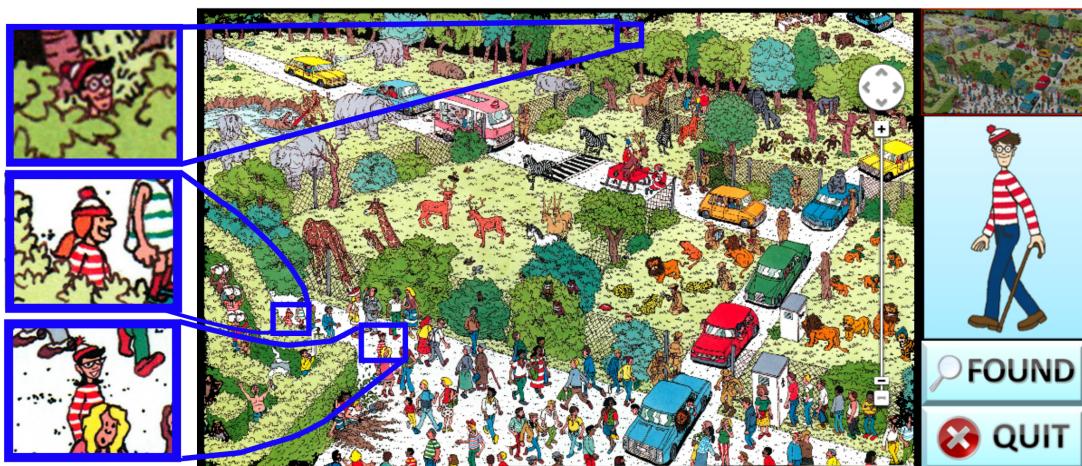
warning! \begin{alex_rant}

- do not treat user interaction as an easy way out
 - “oh, we’ll just give the user a menu for that!”
- do not make your problems (as developers and designers) the user’s problems
- interaction should be meant to explore data, ask questions, gain insight

user interaction is valuable data

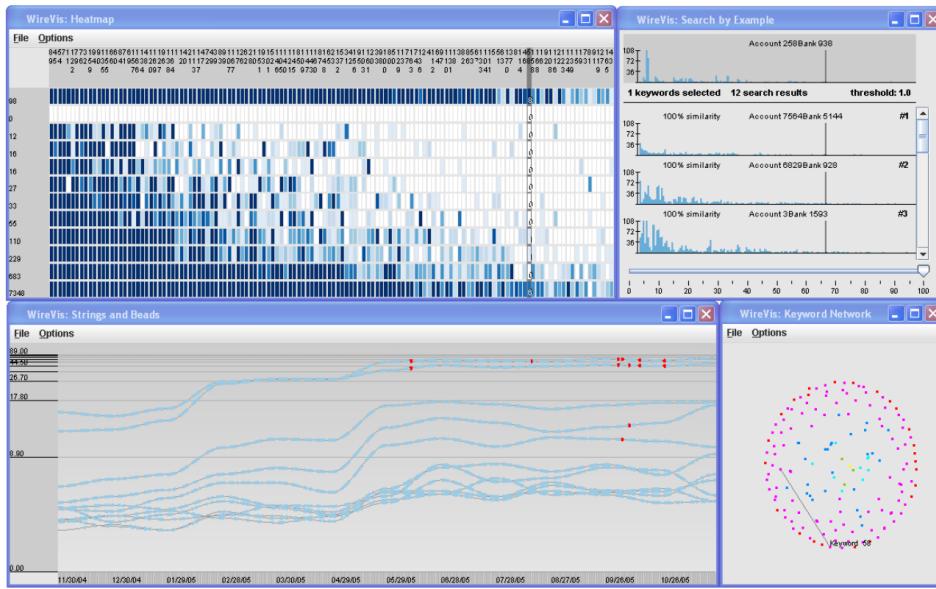
- quantitative record of analysis and exploration process
- what can we learn from this data?
 - data models? sensemaking stages?
- “interaction is the insight”
- let’s look at a few examples

are you good at Where's Waldo?



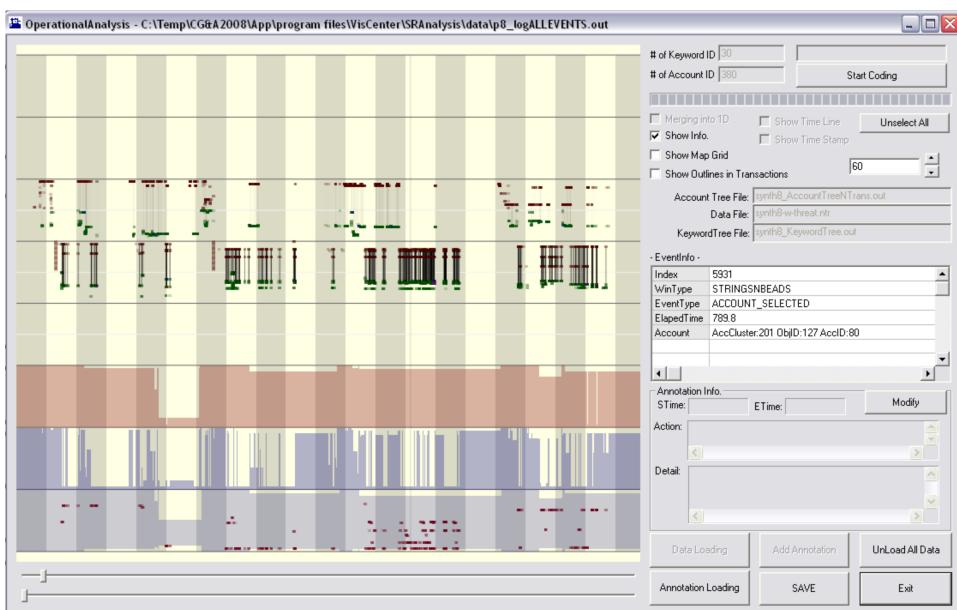
Brown, Eli T., et al. "Finding waldo: Learning about users from their interactions." Visualization and Computer Graphics, IEEE Transactions on 20.12 (2014): 1663-1672.

recover reasoning process: capture



Dou, Wenwen, et al. "Recovering reasoning processes from user interactions." IEEE Computer Graphics and Applications 3 (2009): 52-61.

recover reasoning process: analyze



recover reasoning process: results

- analyzing user interactions after the analysis, researchers can recover:
 - 60% of strategies
 - 60% of methods
 - 79% of findings
- Impressive!
- Can we systematically analyze the user interaction in the system at runtime to steer the models?

another way to look at this

- instead of asking users to directly tune and parameterize models, how can systems learn from their exploratory user interactions?
- How can we use the interactions that people have done to steer the models and the visualizations directly?

let's try it

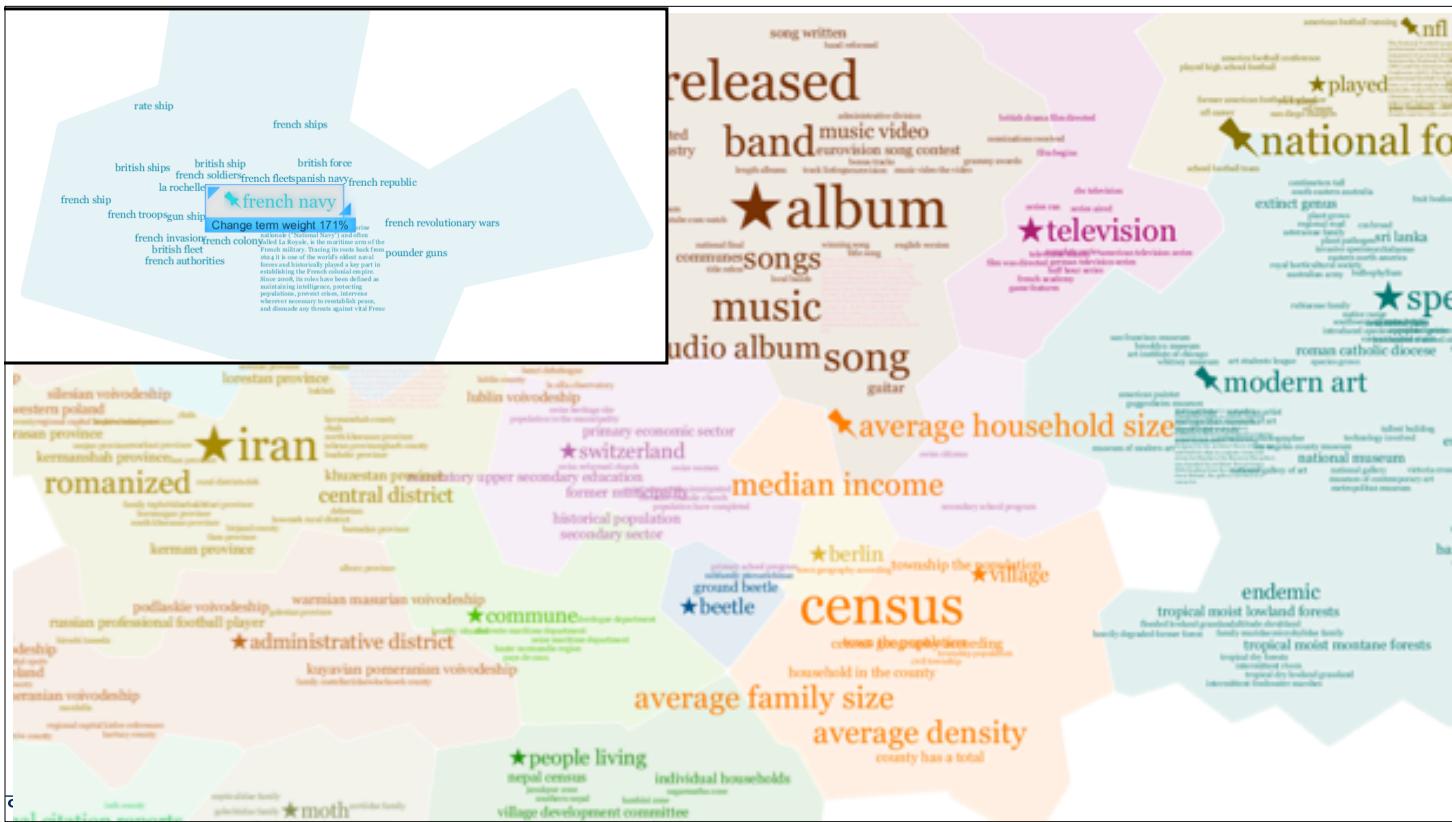
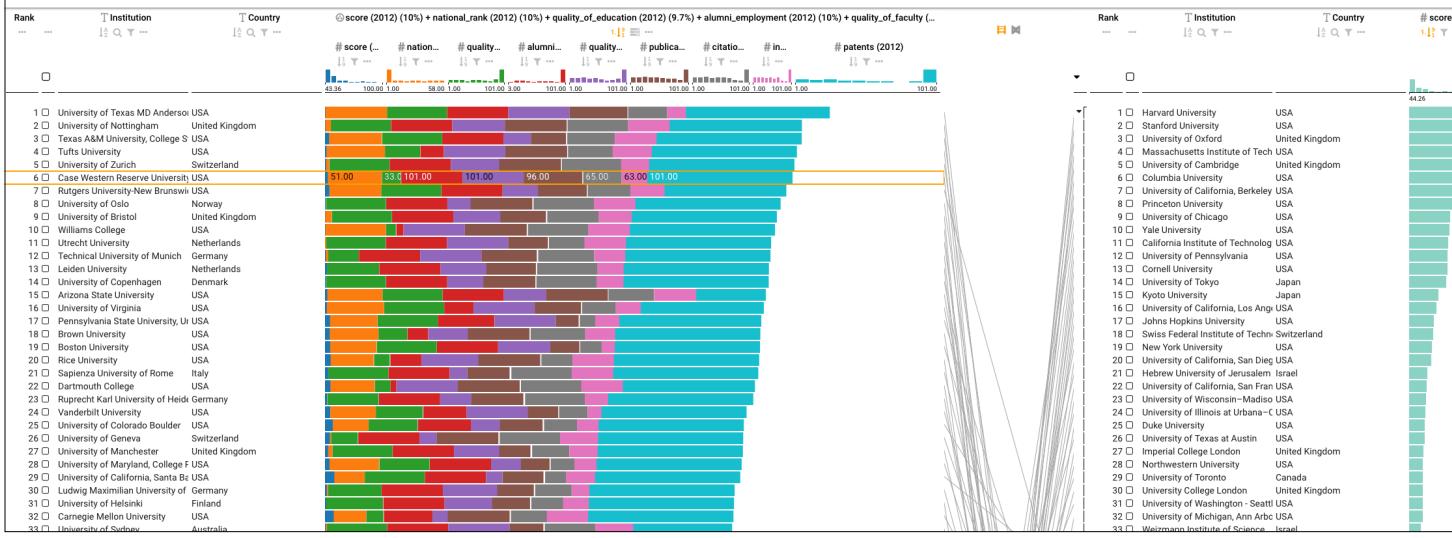
Another example

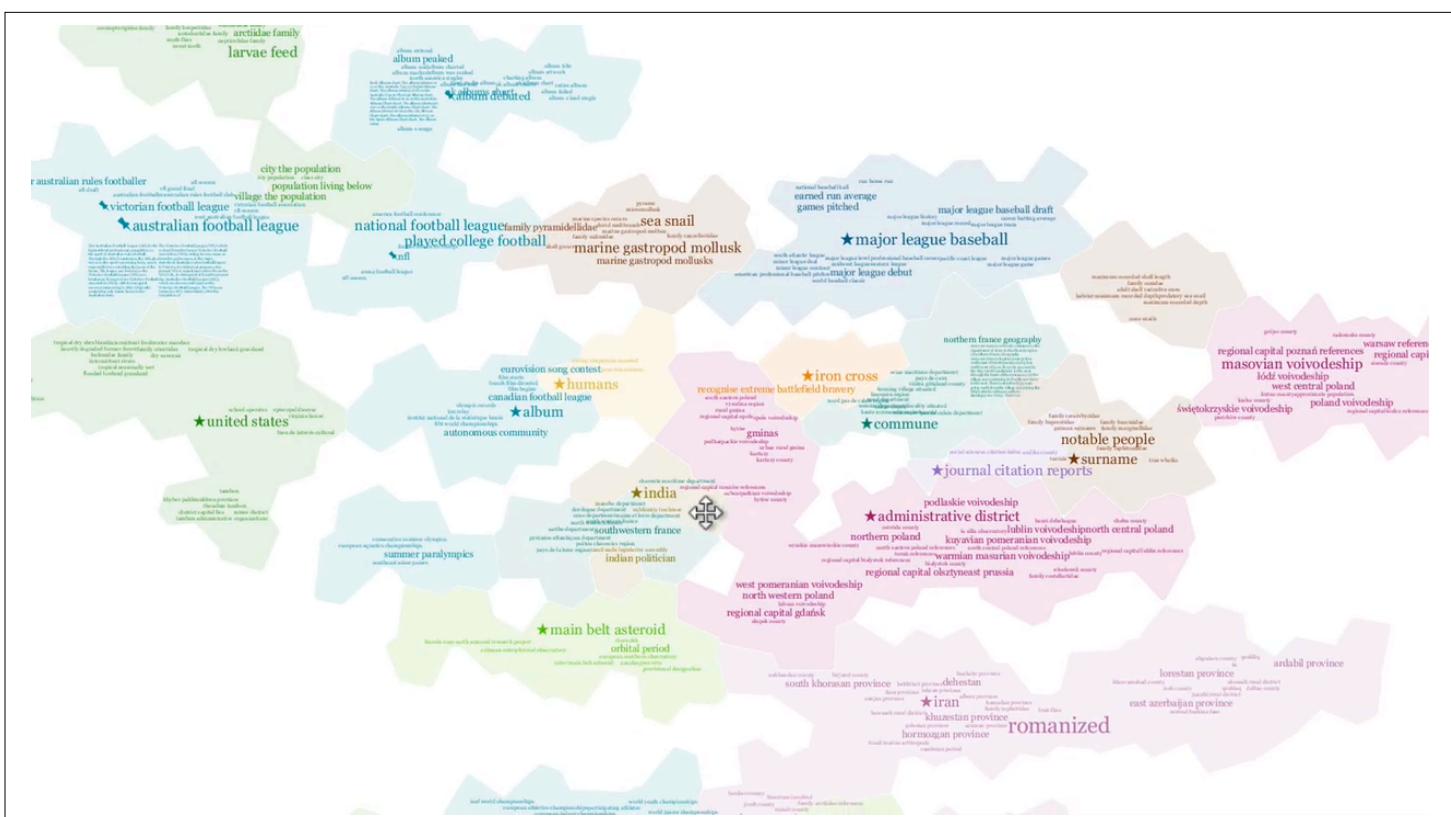
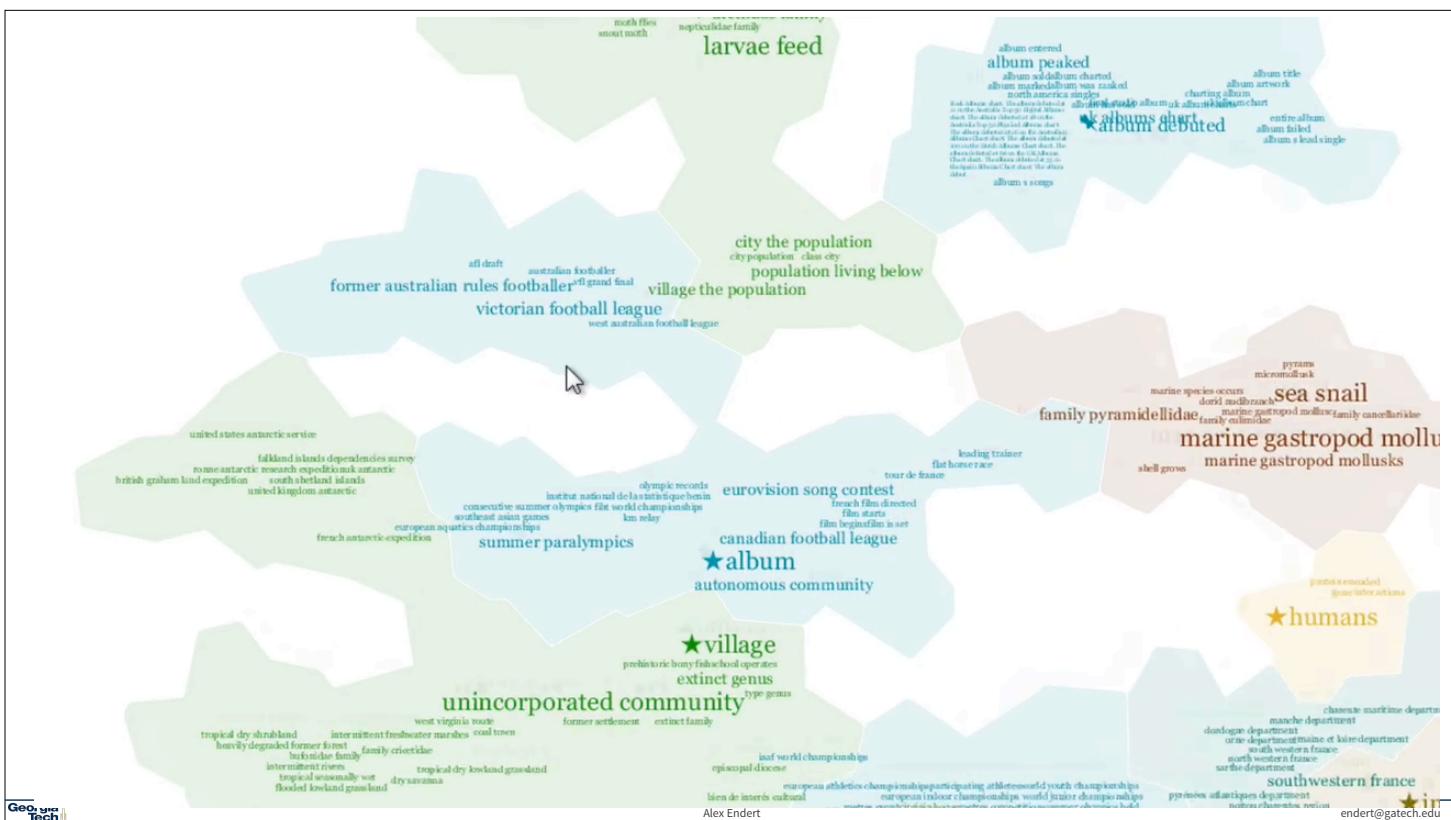
- Wrangler
 - <http://vis.stanford.edu/wrangler/>
- Visual Analytics for data cleaning

Ranking Data Items

- LineUp

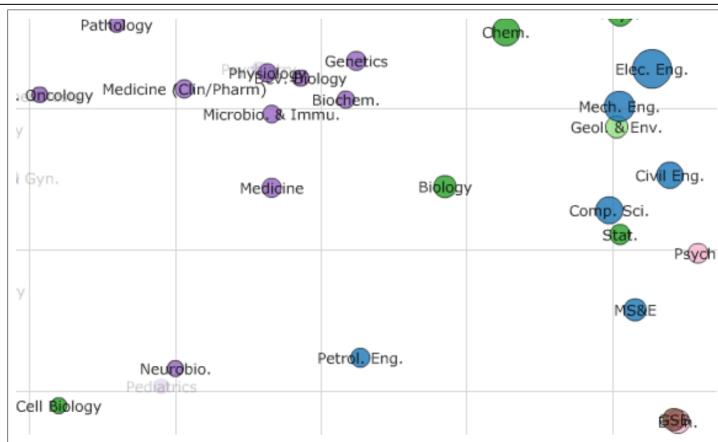
<https://lineup.js.org/app/>





potential challenges

- **interpretability and trust**
- computation **approximates** data characteristics and phenomena, but sometimes get it wrong
- good case study by Chuang et al.



DR-produced view has distortions that can create visual artifacts

Petrol Eng. is similar to Neurobiology?

transformation of view based on user selection can remove these

