

# INTERACTIVE Data Analysis

Jeffrey Heer @jeffrey\_heer  
Univ. of Washington + Trifacta





LIFE

Data Analysis & Statistics, Tukey & Wilk 1966



**Four major influences** act on data analysis today:

1. The formal theories of statistics.
2. Accelerating developments in computers and display devices.
3. The challenge, in many fields, of more and larger bodies of data.
4. The emphasis on quantification in a wider variety of disciplines.



While some of the influences of statistical theory on data analysis have been helpful, others have not.





**Exposure**, the effective laying open of the data to **display the unanticipated**, is to us a major portion of data analysis...

It is not clear how the **informality** and **flexibility** appropriate to the **exploratory character** of exposure can be fitted into any of the structures of formal statistics so far proposed.

LIFE



Accordingly, both approaches and techniques need to be structured so as to **facilitate human involvement and intervention.**

Some implications for effective analysis are: (1) it is essential to have convenience of **interaction of people and intermediate results** and (2) at all stages of data analysis, the outputs need to be **matched to the capabilities of the people who use it and want it.**

## Graph Viewer

Roll-up by:

All

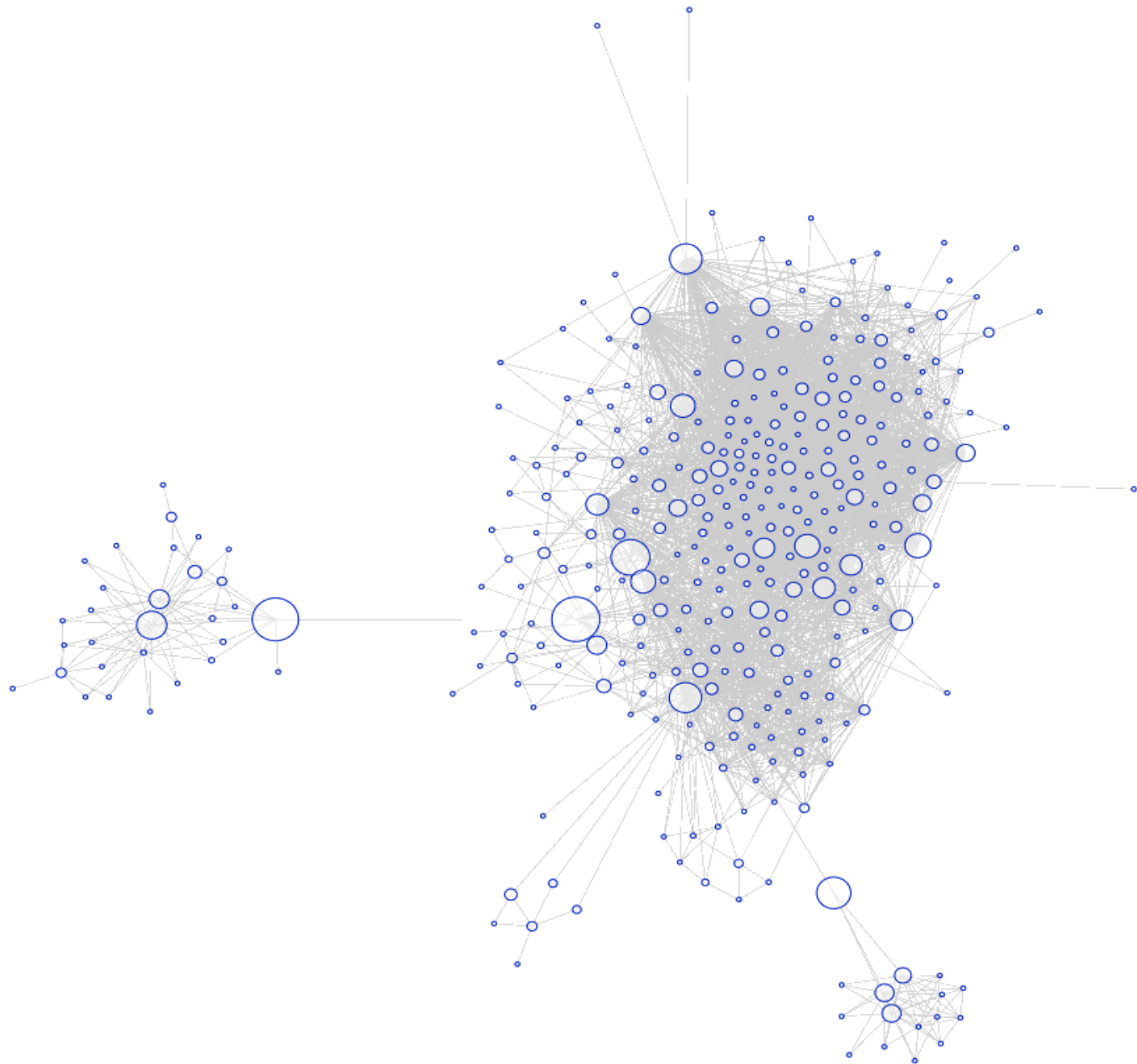
Visualization:

Node-Link

Sort by:

None

Edge centrality filters:

☐ Images☒ Animate

## Graph Viewer

Roll-up by:

All

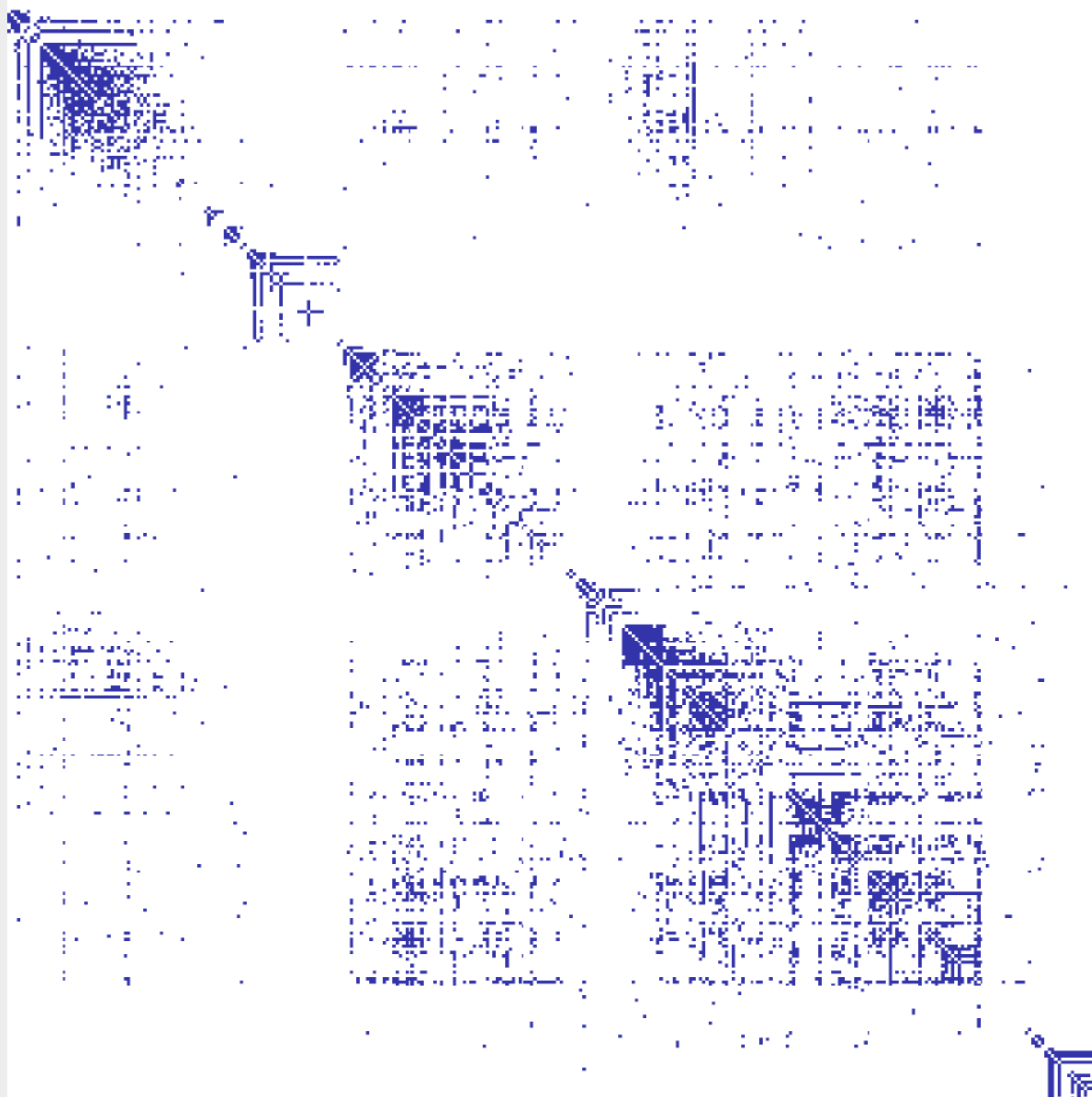
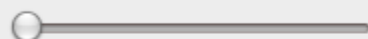
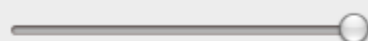
Visualization:

Matrix

Sort by:

Linkage

Edge centrality filters:



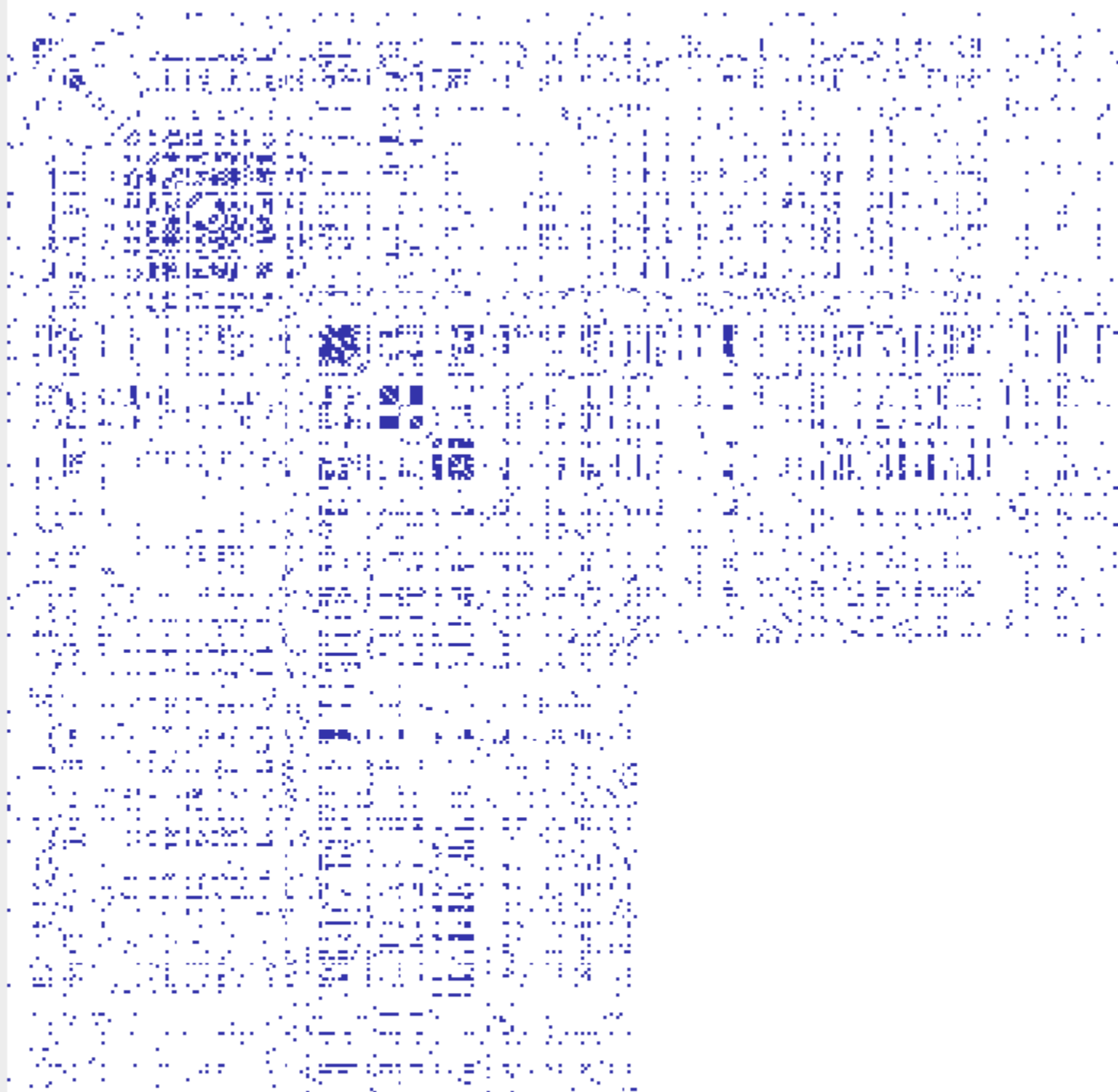
## Graph Viewer

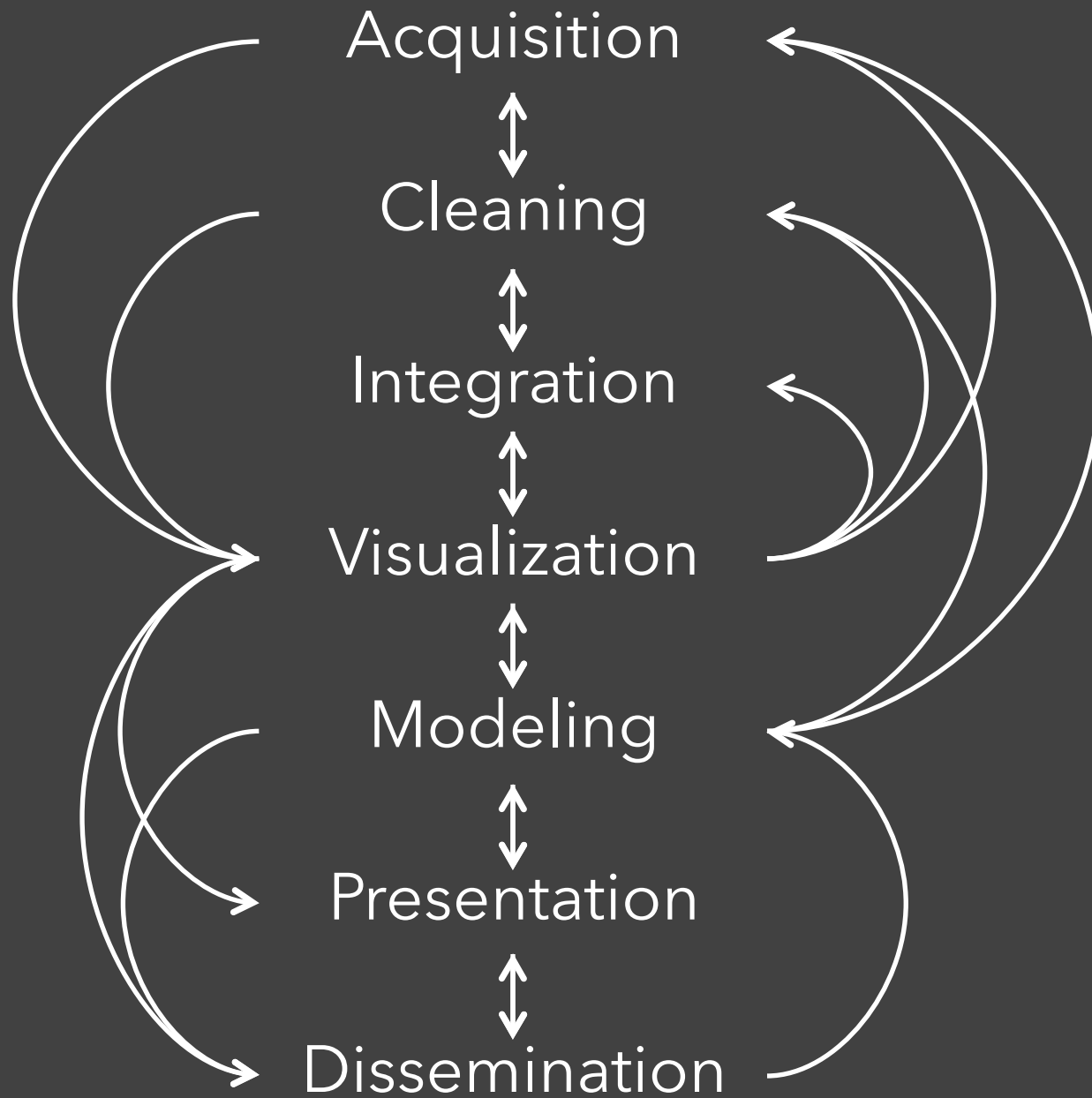
Roll-up by:

Visualization:

Sort by:

Edge centrality filters:





How can we **transform data**  
without programming?

I spend more than half of my time integrating, cleansing and transforming data without doing any actual analysis. Most of the time I'm lucky if I get to do any "analysis" at all.

Anonymous Data Scientist  
*from our interview study* [VAST'12]





Reported crime in Alabama

Year	Population	Property crime rate			Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	4525375	4029.3	987	2732.4	309.9		
2005	4548327	3900	955.8	2656	289		
2006	4599030	3937	968.9	2645.1	322.9		
2007	4627851	3974.9	980.2	2687	307.7		
2008	4661900	4081.9	1080.7	2712.6	288.6		

Reported crime in Alaska

Year	Population	Property crime rate			Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	657755	3370.9	573.6	2456.7	340.6		
2005	663253	3615	622.8	2601	391		
2006	670053	3582	615.2	2588.5	378.3		
2007	683478	3373.9	538.9	2480	355.1		
2008	686293	2928.3	470.9	2219.9	237.5		

Reported crime in Arizona

Year	Population	Property crime rate			Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	5739879	5073.3	991	3118.7	963.5		
2005	5953007	4827	946.2	2958	922		
2006	6166318	4741.6	953	2874.1	914.4		
2007	6338755	4502.6	935.4	2780.5	786.7		
2008	6500180	4087.3	894.2	2605.3	587.8		

Reported crime in Arkansas

Year	Population	Property crime rate			Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	2750000	4033.1	1096.4	2699.7	237		
2005	2775708	4068	1085.1	2720	262		
2006	2810872	4021.6	1154.4	2596.7	270.4		
2007	2834797	3945.5	1124.4	2574.6	246.5		
2008	2855390	3843.7	1182.7	2433.4	227.6		

Reported crime in California

Year	Population	Property crime rate			Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	35842038		3423.9	686.1	2033.1	704.8	
2005	36154147		3321	692.9	1915	712	
2006	36457549		3175.2	676.9	1831.5	666.8	
2007	36553215		3032.6	648.4	1784.1	600.2	
2008	36756666		2940.3	646.8	1769.8	523.8	

Reported crime in Colorado

Year	Population	Property crime rate			Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	4601821	3918.5	717.3	2679.5	521.6		

# DataWrangler

Suggestions

Delete rows 8,10

Delete empty rows

Delete rows where Property\_crime\_rate is null

Delete rows where Year is null

Script

Export

► Split data repeatedly on newline into rows

► Split data repeatedly on ','

rows: 408 prev next

#	Year	#	Property_crime_rate
1	Reported crime in Alabama		
2			
3	2004	4029.3	
4	2005	3900	
5	2006	3937	
6	2007	3974.9	
7	2008	4081.9	
8			
9	Reported crime in Alaska		
10			
11	2004	3370.9	
12	2005	3615	
13	2006	3582	
14	2007	3373.9	

with **S. Kandel**, P. Guo, A. Paepcke & J. Hellerstein [CHI'11]

# Predictive Interaction



# Predictive Interaction

User *interacts* with data and visualizations

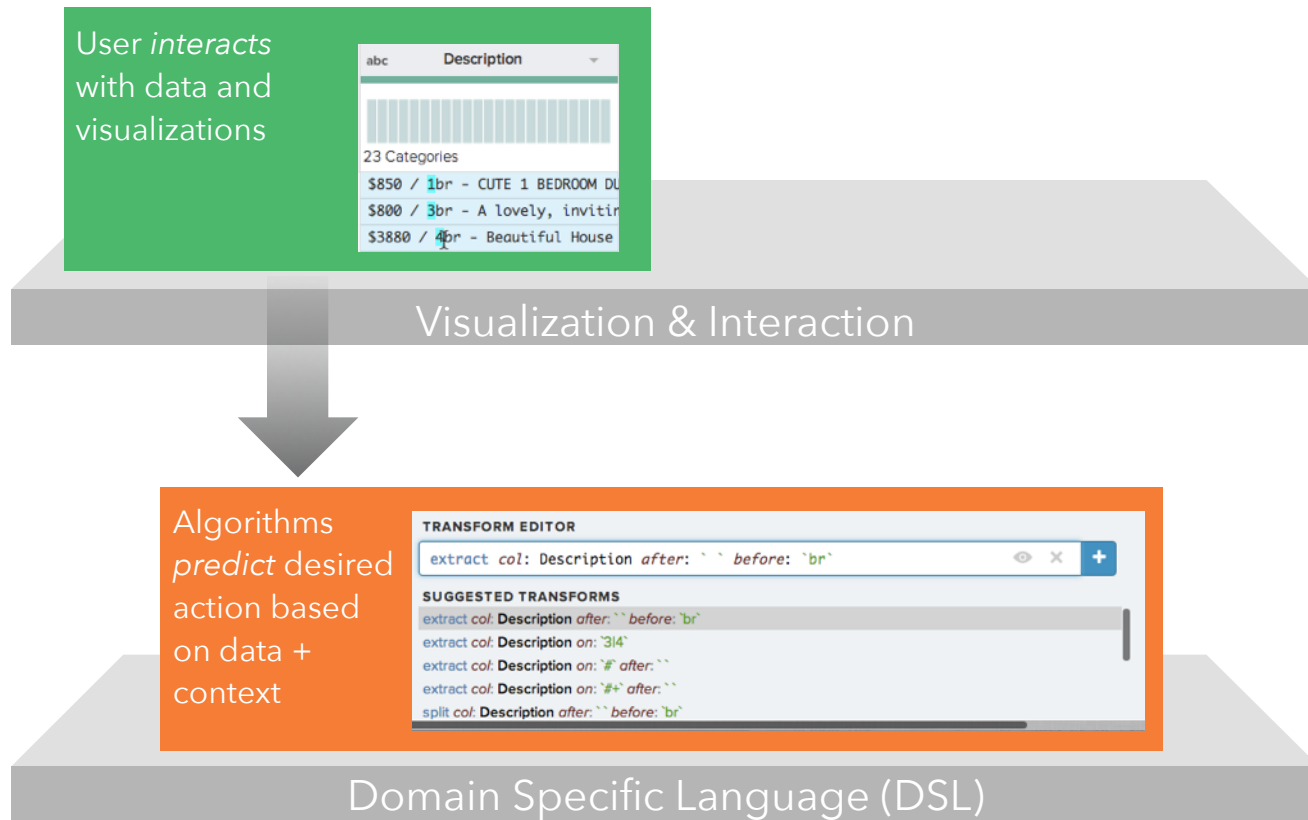
abc	Description
\$850 / 1br -	CUTE 1 BEDROOM DU
\$800 / 3br -	A lovely, invitir
\$3880 / 4br -	Beautiful House

## Visualization & Interaction

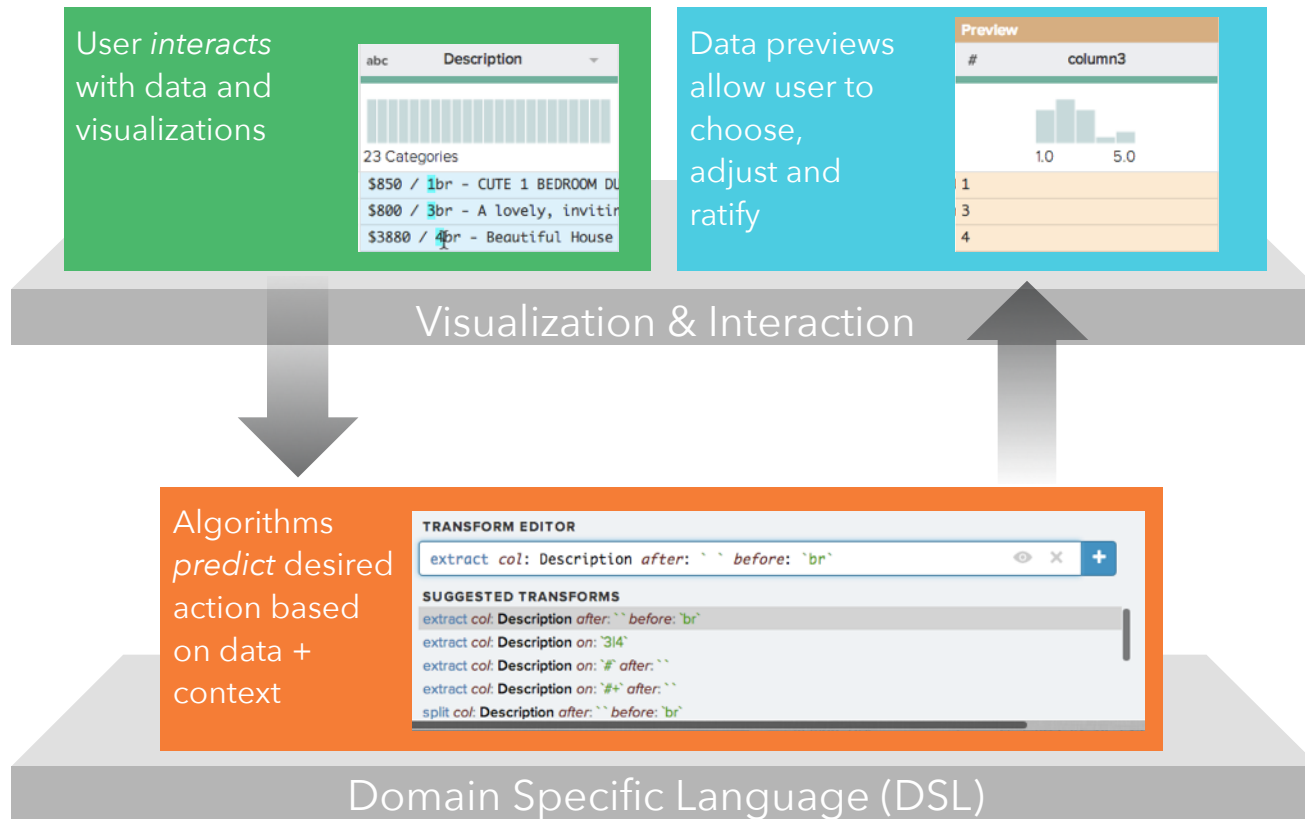
## Domain Specific Language (DSL)



# Predictive Interaction



# Predictive Interaction



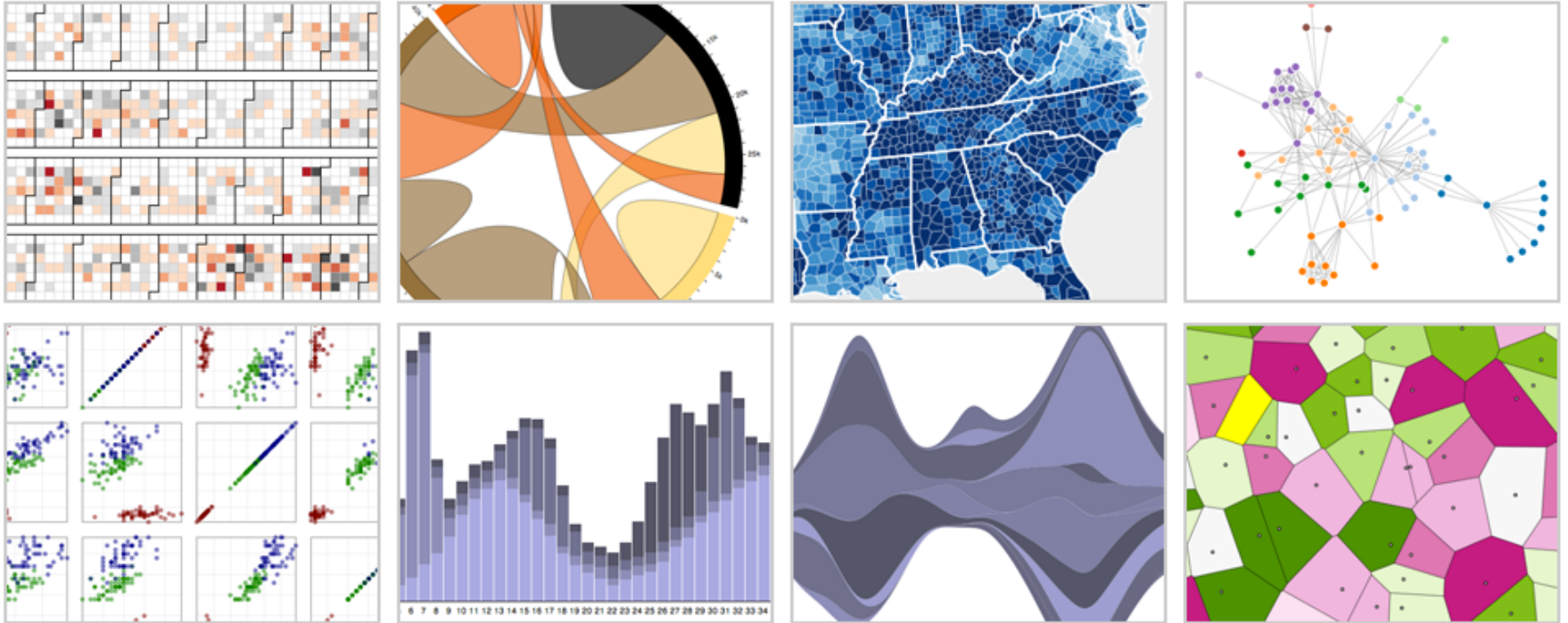


TRIFACTA



How might we support  
**expressive** and **effective**  
visualization designs?

# d3.js Data-Driven Documents



with **M. Bostock**, V. Ogievetsky [InfoVis '11]

# 512 Paths to the White House

Select a winner in the most competitive states below to see all the paths to victory available for either candidate.

Fla.	Ohio	N.C.	Va.	Wis.	Colo.	Iowa	Nev.	N.H.
<input type="button" value="Dem"/> <input type="button" value="Rep"/>	<input type="button" value="Dem"/> <input type="button" value="Rep"/>	<input type="button" value="Dem"/> <input type="button" value="Rep"/>	<input type="button" value="Dem"/> <input type="button" value="Rep"/>	<input type="button" value="Dem"/> <input type="button" value="Rep"/>	<input type="button" value="Dem"/> <input type="button" value="Rep"/>	<input type="button" value="Dem"/> <input type="button" value="Rep"/>	<input type="button" value="Dem"/> <input type="button" value="Rep"/>	<input type="button" value="Dem"/> <input type="button" value="Rep"/>

Obama has 431 ways to win

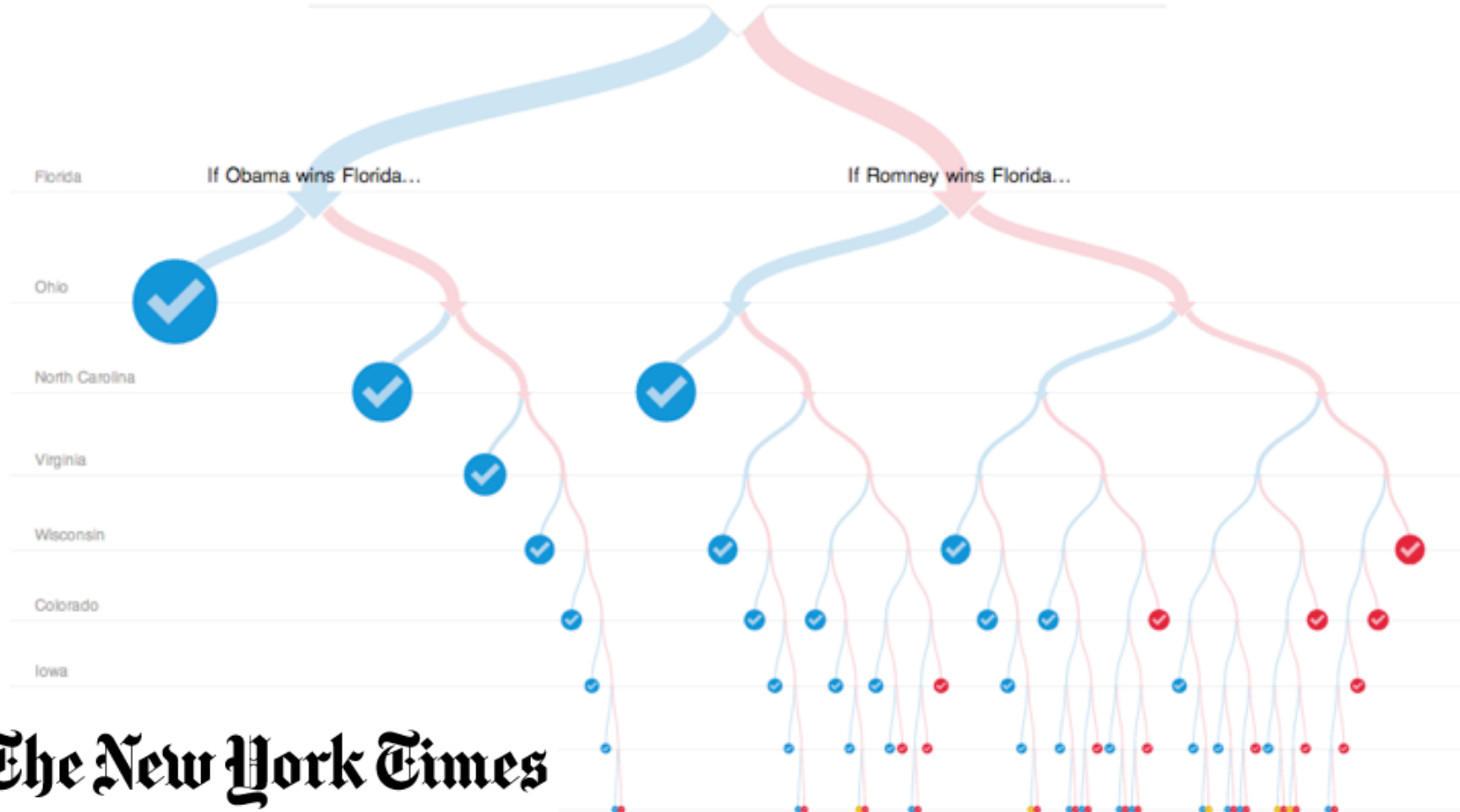
84% of paths

5 ties

0.98% of paths

Romney has 76 ways to win

15% of paths



The New York Times

# stamen design





# stamen design



d3

d3



Search

Search

Repositories

83

&lt;&gt; Code

! Issues

👤 Users

## Languages

JavaScript 36

Ruby 6

Python 5

HTML 5

CSS 5

C++ 3

VimL 2

Shell 2

Go 2

C 2

## twbs/bootstrap

CSS ★ 78,844 📄 30,533

The most popular HTML, CSS, and JavaScript framework for developing responsive, mobile first projects on the web.

Updated 7 hours ago

## vhf/free-programming-books

★ 37,478 📄 8,233

Updated 2 days ago

## angular/angular.js

JavaScript ★ 36,369 📄 14,905

HTML enhanced for web apps

Updated a day ago

# 4<sup>th</sup> most starred project on GitHub

## mbostock/d3

JavaScript ★ 35,578 📄 8,949

A JavaScript visualization library for HTML and SVG.

Updated on Feb 11

## joyent/node

JavaScript ★ 35,190 📄 7,860

evented I/O for v8 javascript

Updated 5 hours ago

[Advanced search](#) [Cheat sheet](#)

# Vega A VISUALIZATION GRAMMAR



Vega is a declarative format for creating, saving, and sharing visualization designs. With Vega, visualizations are described in JSON, and generate interactive views using either HTML5 Canvas or SVG.

## TOOLKITS

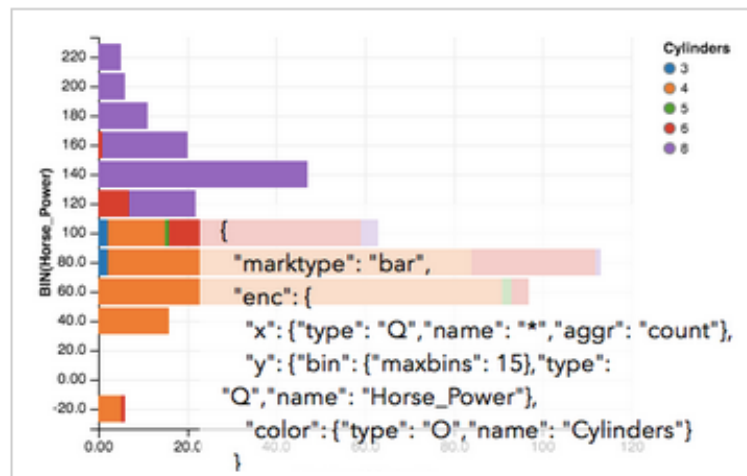


**VEGA** offers a full declarative visualization grammar, suitable for expressive custom interactive visualization design and programmatic generation.

[Tutorial](#) | [Documentation](#) | [Discussion Forum](#)

v1.5 (stable): [download](#), [examples](#), [github](#)

**NEW** v2.0 (dev): [download](#), [examples](#), [github](#)



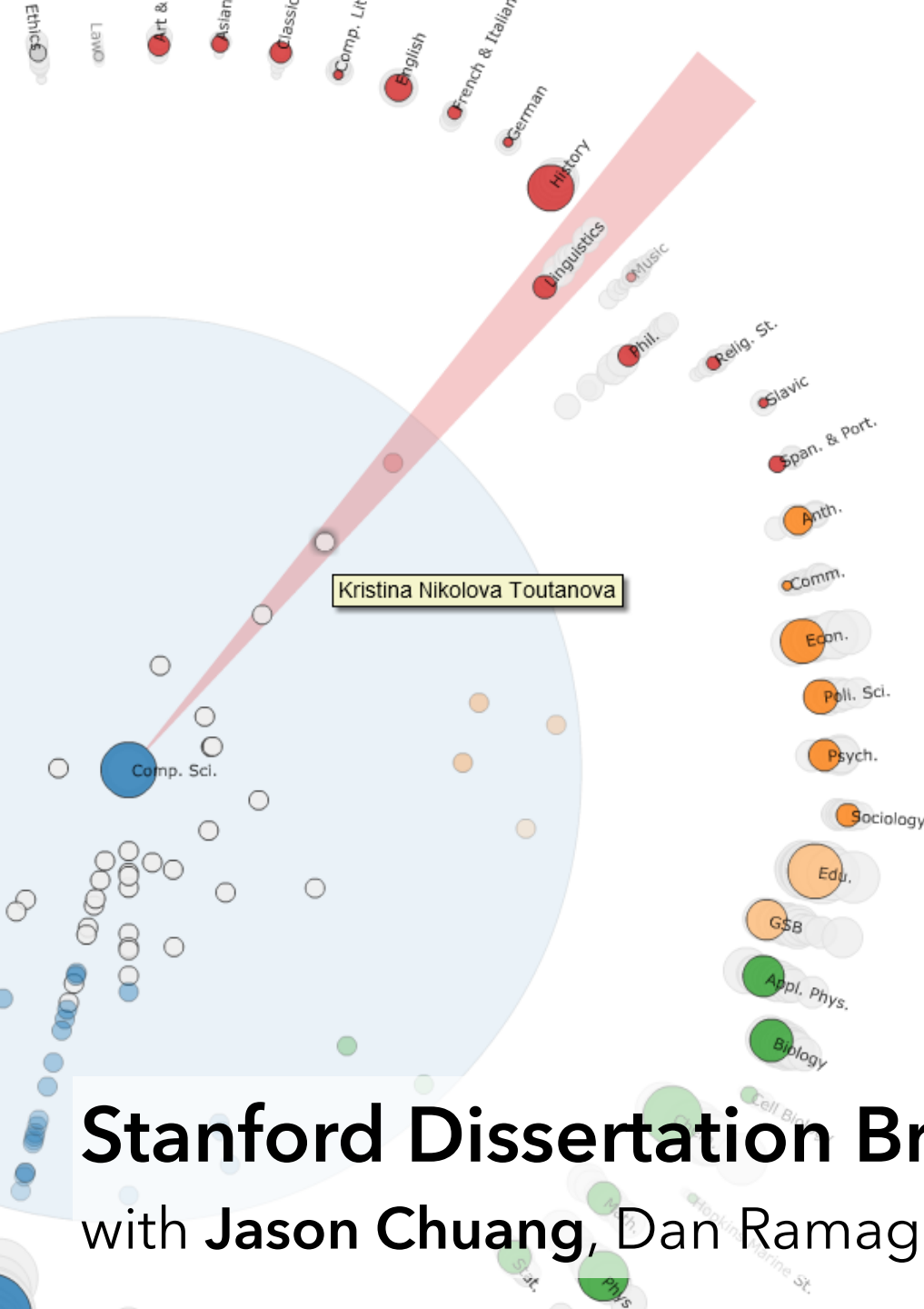
**NEW VEGALITE** provides a higher-level grammar for visual analysis, comparable to ggplot or Tableau, that generates complete Vega specifications.

[Online Editor](#) | [GitHub](#)

# vega.github.io

How might we support model  
**interpretation and refinement?**





## Effective statistical models for syntactic and semantic disambiguation

Student: Kristina Nikolova Toutanova

Advisor: Christopher D. Manning

Computer Science (2005)

Keywords: Syntactic, Semantic, Tree kernels, Parsing

Abstract:

This thesis focuses on building effective statistical models for disambiguation of sophisticated syntactic and semantic natural language (NL) structures. We advance the state of the art in several domains by (i) choosing representations that encode domain knowledge more effectively and (ii) developing machine learning algorithms that deal with the specific properties of NL disambiguation tasks--sparsity of training data and large, structured spaces of hidden labels. For the task of syntactic disambiguation, we propose a novel representation of parse trees that connects the words of the sentence with the hidden syntactic structure in a direct way. Experimental evaluation on parse selection for a Head Driven Phrase Structure Grammar shows the new representation achieves superior performance compared to previous models. For the task of disambiguating the semantic role structure of verbs, we build a more accurate model, which captures the knowledge that the semantic frame of a verb is a joint structure with strong dependencies between arguments. We achieve this using a Conditional Random Field without Markov independence assumptions on the sequence of semantic role labels. To address the sparsity problem in machine learning for NL, we develop a method for incorporating many additional sources of information, using Markov chains in the space of words. The Markov chain framework makes it possible to combine multiple knowledge sources, to learn how much to trust each of them, and to chain inferences together. It achieves large gains in the task of disambiguating prepositional phrase attachments.

# Stanford Dissertation Browser

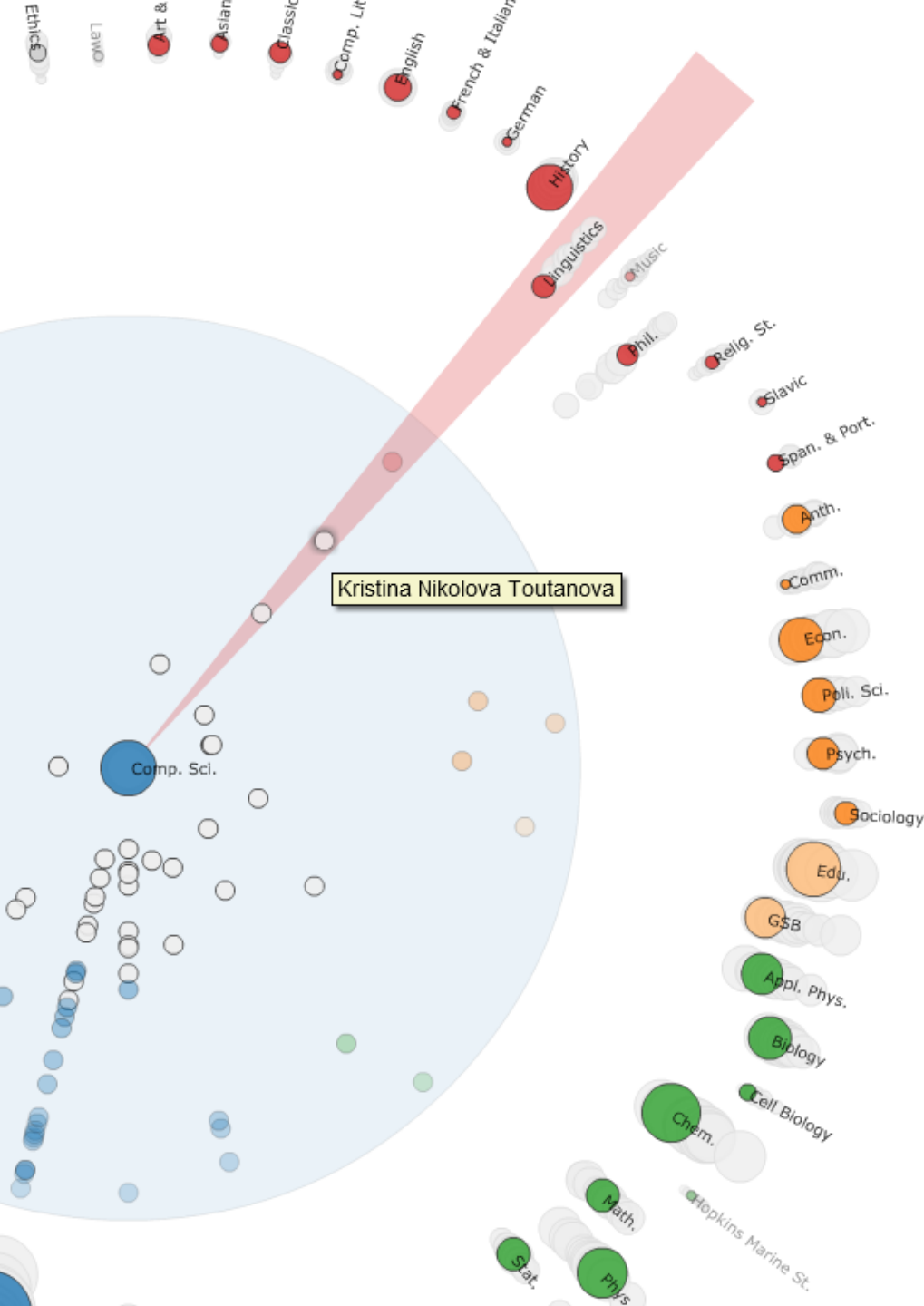
with Jason Chuang, Dan Ramage & Christopher Manning







Oh, the humanities!



## Effective statistical models for syntactic and semantic disambiguation

Student: Kristina Nikolova Toutanova

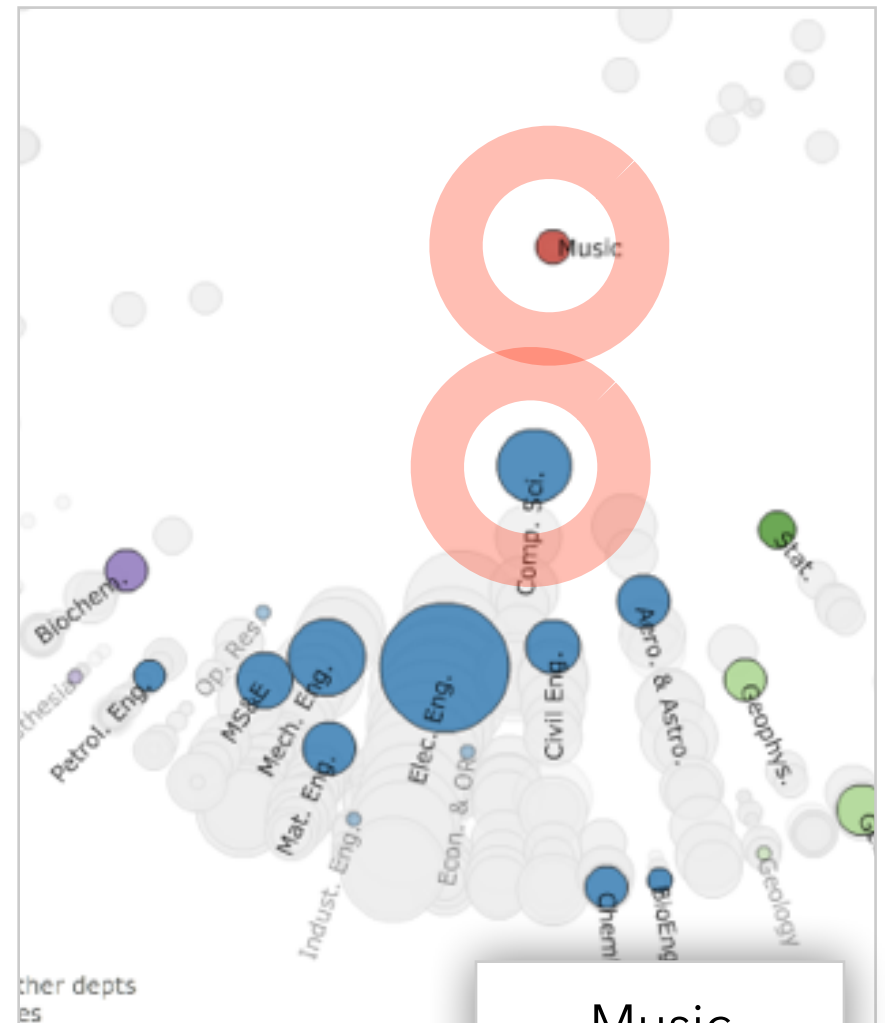
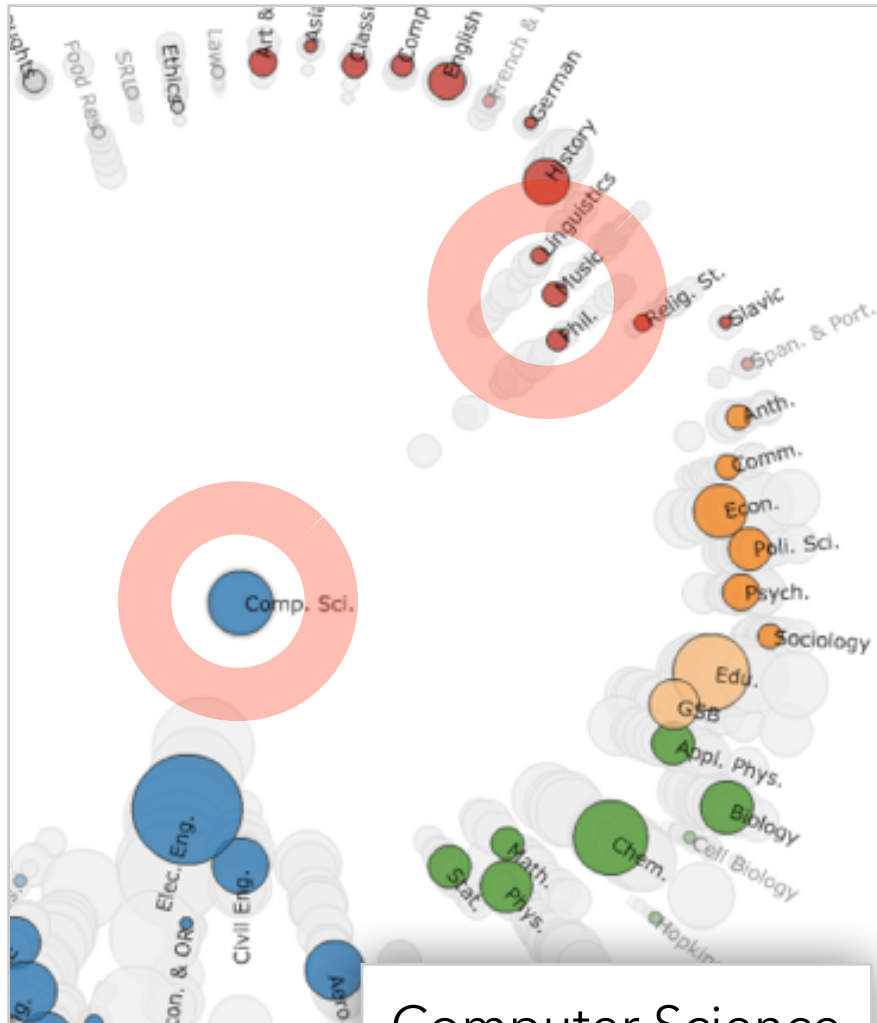
Advisor: Christopher D. Manning

Computer Science (2005)

Keywords: Syntactic, Semantic, Tree kernels, Parsing

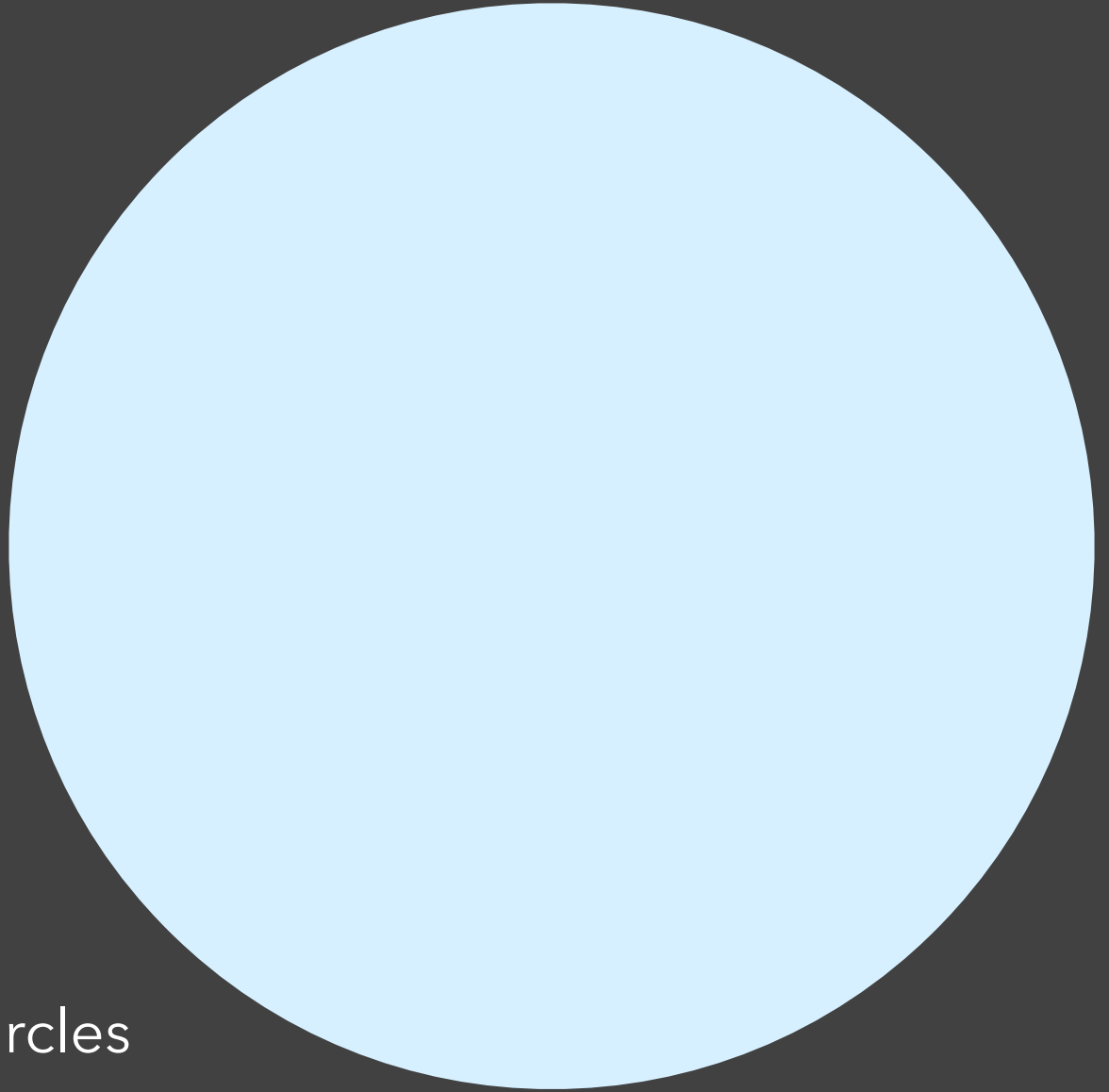
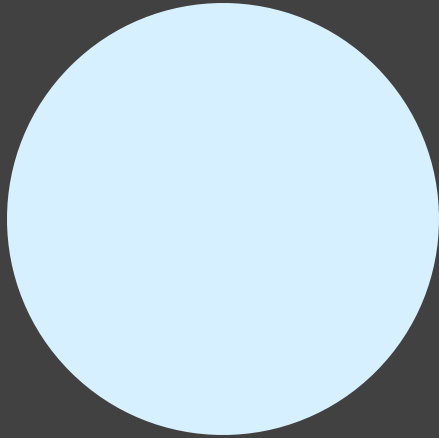
Abstract:

This thesis focuses on building effective statistical models for disambiguation of sophisticated syntactic and semantic natural language (NL) structures. We advance the state of the art in several domains by (i) choosing representations that encode domain knowledge more effectively and (ii) developing machine learning algorithms that deal with the specific properties of NL disambiguation tasks--sparsity of training data and large, structured spaces of hidden labels. For the task of syntactic disambiguation, we propose a novel representation of parse trees that connects the words of the sentence with the hidden syntactic structure in a direct way. Experimental evaluation on parse selection for a Head Driven Phrase Structure Grammar shows the new representation achieves superior performance compared to previous models. For the task of disambiguating the semantic role structure of verbs, we build a more accurate model, which captures the knowledge that the semantic frame of a verb is a joint structure with strong dependencies between arguments. We achieve this using a Conditional Random Field without Markov independence assumptions on the sequence of semantic role labels. To address the sparsity problem in machine learning for NL, we develop a method for incorporating many additional sources of information, using Markov chains in the space of words. The Markov chain framework makes it possible to combine multiple knowledge sources, to learn how much to trust each of them, and to chain inferences together. It achieves large gains in the task of disambiguating prepositional phrase attachments.



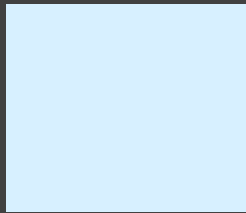
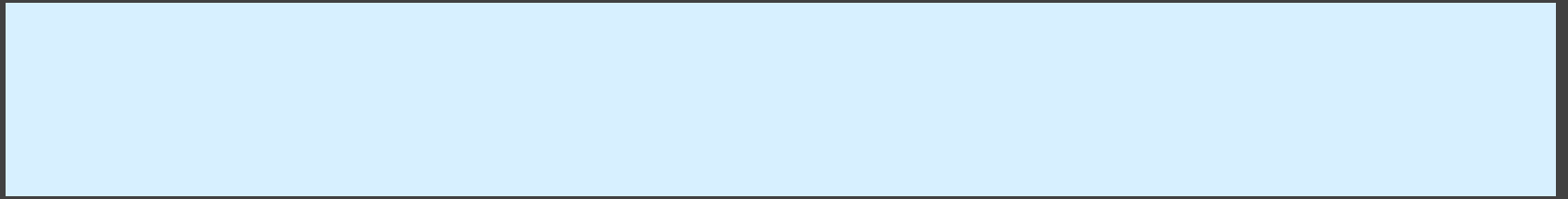
“Word Borrowing” via Labeled LDA

What makes a  
visualization “good”?

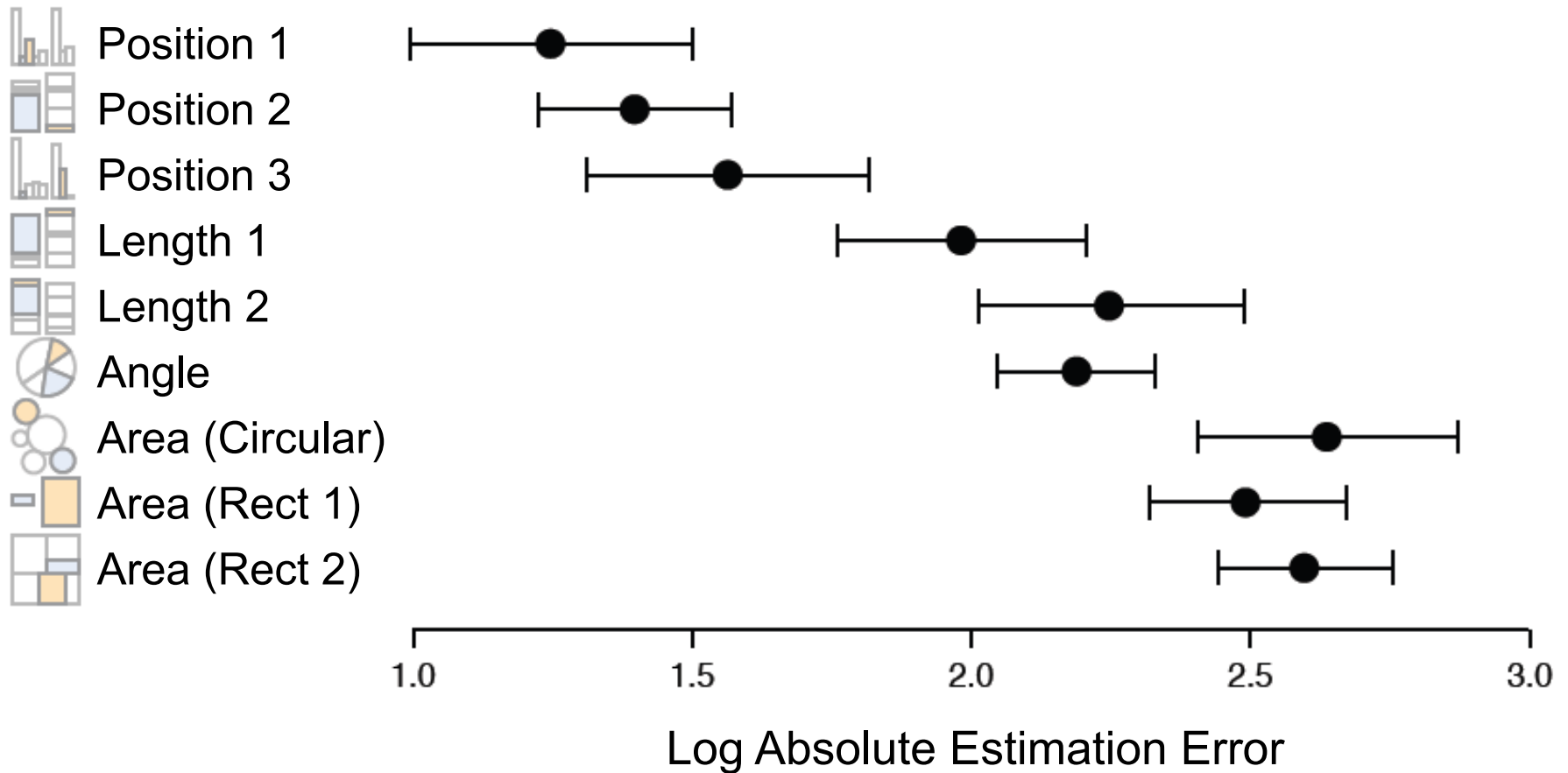


Compare area of circles





Compare length of bars



# Graphical Perception Experiments

Empirical estimates of encoding effectiveness

# Estimating Proportions

Most accurate



Least accurate



Position (common) scale



Position (non-aligned) scale



Length



Slope



Angle



Area



Volume



Color hue-saturation-density

# Artery Visualization [Borkin et al. '11]

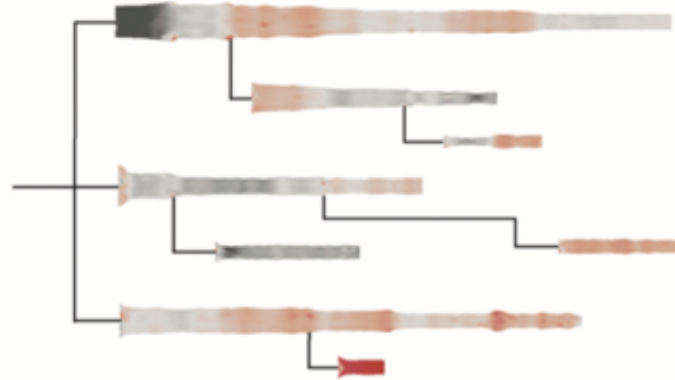
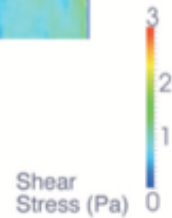
Rainbow Palette

Diverging Palette

2D



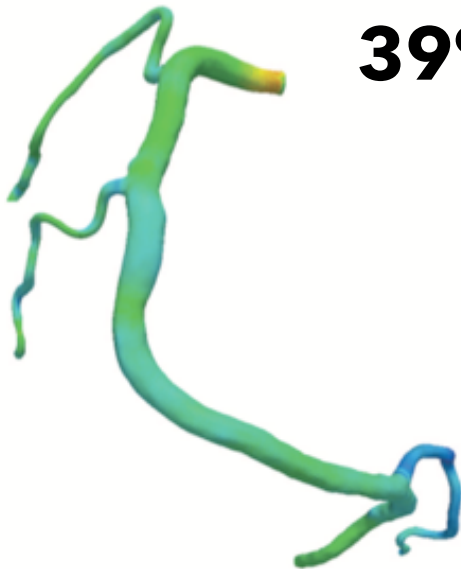
**62%**



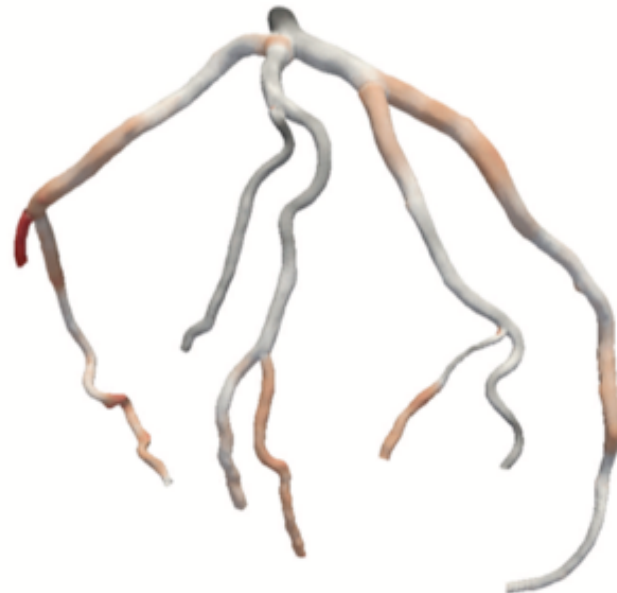
**92%**



3D



**39%**



**71%**

# INTERACTIVE Data Analysis

Jeffrey Heer @jeffrey\_heer  
Univ. of Washington + Trifacta

