# *Interactive* ^ Data Analysis

**Jeffrey Heer**
Stanford University

# Graph Viewer

**Roll-up by:**

All

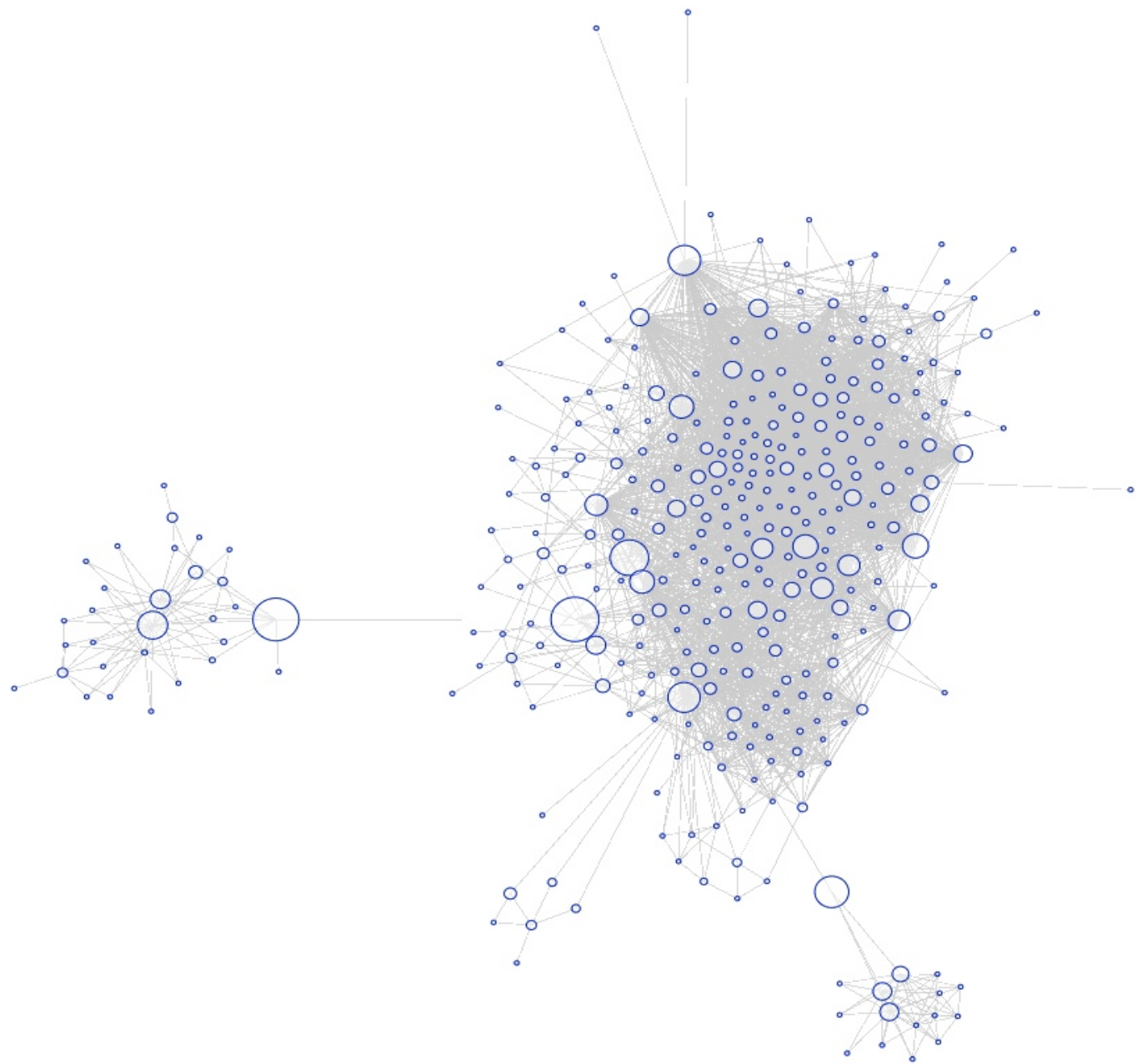**Visualization:**

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:**

All

**Visualization:**

Matrix

**Sort by:**

None

**Edge centrality filters:**

Acquisition

Cleaning

Integration

Visualization

Modeling

Presentation

Dissemination

Acquisition

Cleaning

Integration

**Visualization**

Modeling

Presentation

Dissemination

# **d3.js** Data-Driven Documents



with **Mike Bostock** & **Vadim Ogievetsky**

Acquisition

Cleaning

Integration

**Visualization**

Modeling

Presentation

Dissemination

Acquisition

**Cleaning**

Integration

**Visualization**

Modeling

Presentation

Dissemination

## 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 |

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, 2012*

# Data**Wrangler**



with **Sean Kandel**, Philip Guo, Andreas Paepcke & Joe Hellerstein

# Wrangler in 2 Parts…

1. Declarative data transformation language
   **Tuple mapping** – split, merge, extract, delete
   **Reshaping** – fold, unfold (cross-tabulation)
   **Lookups & joins** – e.g., FIPS code to US state
   **Sorting, aggregation, *etc.***

   Informed by prior work in databases:
   Potter's Wheel, SchemaSQL, AJAX

# Wrangler in 2 Parts...

1. Declarative data transformation language

**+**

2. Mixed-initiative interface for data transforms
   *User:* **Selects** data elements of interest
   *System:* **Suggests** applicable transforms via
      search over the space of viable transforms
   Enable rapid **preview and refinement**

# Transform Suggestion

Interaction

↓

Infer Operands

↓

Generate Transforms

↓

Rank Transforms

↓

Present Top-N

# Transform Suggestion

**Interaction** ———— Text Selection
⬇                    Text Editing
Infer Operands       Row Selection
⬇                    Column Selection
Generate Transforms  Transform Menu
⬇                    Click Quality Meter
Rank Transforms
⬇
Present Top-N

# Transform Suggestion

Interaction

↓

**Infer Operands** —————

↓

Generate Transforms

↓

Rank Transforms

↓

Present Top-N

Map user input to transform operands.

Example: text highlight maps to row, column, and text selections.

Inferred text selections include string indices and regular expressions.

# Text Selection Inference

Series Id: LNU02000000

# Text Selection Inference

Series Id: LNU02000000

-> ^ STR WS STR SYM WS STR NUM $

# Text Selection Inference

Series Id: LNU02000000
-> ^ STR WS STR SYM WS STR NUM $

Series Id: `LNU02000000`

# Text Selection Inference

Series Id: LNU02000000

-> ^ STR WS STR SYM WS STR NUM $

Series Id: LNU02000000

MATCH    Indices 11–22

# Text Selection Inference

Series Id: LNU02000000

-> ^ STR WS STR SYM WS STR NUM $

Series Id:  `LNU02000000`

MATCH       Indices 11-22

MATCH       LNU02000000

# Text Selection Inference

Series Id: LNU02000000
-> ^ STR WS STR SYM WS STR NUM $

Series Id: `LNU02000000`

MATCH    Indices 11–22

MATCH    LNU02000000

MATCH    LNU NUM

MATCH    STR NUM

# Text Selection Inference

Series Id: LNU02000000
-> ^ STR WS STR SYM WS STR NUM $

Series Id: `LNU02000000`

| MATCH | Indices 11-22 |
|---|---|
| MATCH | LNU02000000 |
| MATCH | LNU NUM |
| MATCH | STR NUM |
| AFTER | : WS |

# Transform Suggestion

Interaction

↓

**Infer Operands** —————— Map user input to
transform operands.

↓

Generate Transforms

↓

Rank Transforms

↓

Present Top-N

Example: text highlight
maps to row, column, and
text selections.

Inferred text selections
include string indices and
regular expressions.

# Transform Suggestion

Interaction

↓

Infer Operands

↓

**Generate Transforms** —

↓

Rank Transforms

↓

Present Top-N

Enumerate transforms that accept inferred operands as input.

Set unmatched params to default values.

Apply filter heuristics: No-ops, delete-all, and overly sparse outputs.

# Transform Suggestion

Interaction

↓

Infer Operands

↓

Generate Transforms

↓

**Rank Transforms** ———— Sort transforms by:
Toolbar selection
Specification difficulty
Frequency in corpus

↓

Present Top-N

# Transform Suggestion

Interaction

↓

Infer Operands

↓

Generate Transforms

↓

Rank Transforms

↓

**Present Top-N** ————

Extract from **unnamed_1 once**
between positions 17,25

Extract from **unnamed_1 once**
on whitespace Alabama

Cut from **unnamed_1 once**
between positions 17,25

Cut from **unnamed_1 once** on
whitespace Alabama

Split **unnamed_1 once**
between positions 17,25 into
**columns**

Split **unnamed_1 once** on
whitespace Alabama into
**columns**

# Comparative Evaluation with Excel



User Study Task Completion Time (minutes)
Wrangler    Excel

Median completion time for Wrangler at least **twice as fast** in all tasks (*p* < 0.001).

Suggestions and visual previews used heavily.

# Difficult Transforms: Table Reshaping

**Fold**

|  | Boys | Girls |
|---|---|---|
| Australia | 1 | 2 |
| Austria | 3 | 4 |
| Belgium | 5 | 6 |
| China | 7 | 8 |

| | | |
|---|---|---|
| Australia | Boys | 1 |
| Australia | Girls | 2 |
| Austria | Boys | 3 |
| Austria | Girls | 4 |
| Belgium | Boys | 5 |
| Belgium | Girls | 6 |
| China | Boys | 7 |
| China | Girls | 8 |

**Pivot**

# Proactive Wrangling

**Proactive transform suggestion** [UIST'11]
Guide users to a proper relational table

Empty cells        Delimiters

$$S(T) = \left(1 - \frac{\sum_{c \in C} H_c(T)}{|C|}\right) + \frac{E(T) + D(T)}{|R|\,|C|}$$

Type homogeneity

$$H_c = \sum_{Type} \left(\frac{|i \in R : c_i \in Type|}{|R|}\right)^2$$

# Proactive Wrangling

**Proactive transform suggestion**  [UIST'11]
Guide users to a proper relational table

**EVALUATION**:

Compare automatic vs. manual transformation

**53%** of transforms automatically suggested

In those cases, the top-ranked suggestion is preferred **77%** of the time (**mean rank: 1.6**).

Acquisition

**Cleaning**

Integration

**Visualization**

Modeling

Presentation

Dissemination

**Schema Browser**
- IMDB Rating
- IMDB Votes
- MPAA Rating
- Major Genre
- Production Budget
- Release Date
- Release Location
- Rotten Tomatoes Rating
- Running Time (min)
- Source

**Anomaly Browser**
- ▷ **Missing (6)**
- ▷ **Error (2)**
- ▷ **Extreme (7)**
- ▷ **Inconsistent (3)**
- ▷ **Schema (1)**

Transform: ▲▼

Variables
(with induced data types)

Results of Automatic
Anomaly Detection

**Data Profiler** [AVI'12]
with Sean Kandel, Ravi Parikh & Joe Hellerstein

**Data Profiler** [AVI'12]
with Sean Kandel, Ravi Parikh & Joe Hellerstein

**Schema Browser**

- IMDB Rating
- IMDB Votes
- MPAA Rating
- Major Genre
- Production Budget
- Release Date
- Release Location
- Rotten Tomatoes Rating
- Running Time (min)
- Source

Transform:

MPAA Rating
R
PG-13
PG
Not Rated
G
NC-17
Open

**Anomaly Browser**

▽ **Missing (6)**
  MPAA Rating
  Creative Type
  Source
  Major Genre
  Distributor
  Release Location
▷ **Error (2)**
▷ **Extreme (7)**
▷ **Inconsistent (3)**
▷ **Schema (1)**

**Data Profiler** [AVI'12]
with Sean Kandel, Ravi Parikh & Joe Hellerstein

**Schema Browser**

- IMDB Rating
- IMDB Votes
- MPAA Rating
- Major Genre
- Production Budget
- Release Date
- Release Location
- Rotten Tomatoes Rating
- Running Time (min)
- Source

**Anomaly Browser**

- Missing (6)
  - MPAA Rating
  - Creative Type
  - Source
  - Major Genre
  - Distributor
  - Release Location
- Error (2)
- Extreme (7)
- Inconsistent (3)
- Schema (1)

Transform:

**MPAA Rating**
- R
- PG-13
- PG
- Not Rated
- G
- NC-17
- Open

**Creative Type**
- Contemporary Fiction
- Historical Fiction
- Fantasy
- Science Fiction
- Dramatization
- Kids Fiction
- Factual
- Super Hero
- Multiple Creative Types

**Release Date**
1911 — 2010, 220 — 0

**Running Time (min)**
0 — 240, 2K — 0

**Production Budget**
0 — 300M, 3K — 0

**Source**
- Original Screenplay
- Based on Book/Short Story
- Based on Real Life Events
- Remake
- Based on TV
- Based on Comic/Graphic...
- Based on Play
- Based on Game
- Traditional/Legend/Fai...
- Based on Magazine Article
- Based on Musical/Opera
- Based on Short Film
- Spin-Off
- Based on Factual Book/...
- Disney Ride
- Compilation
- Based on Toy
- Musical Group Movie

**Data Profiler** [AVI'12]
with Sean Kandel, Ravi Parikh & Joe Hellerstein

## Schema Browser

- IMDB Rating
- IMDB Votes
- MPAA Rating
- Major Genre
- Production Budget
- Release Date
- Release Location
- Rotten Tomatoes Rating
- Running Time (min)
- Source

## Anomaly Browser

**Missing (6)**
- MPAA Rating
- Creative Type
- Source
- Major Genre
- Distributor
- Release Location

▶ **Error (2)**
▶ **Extreme (7)**
▶ **Inconsistent (3)**
▶ **Schema (1)**

Transform:

**MPAA Rating**
- R
- PG-13
- PG
- Not Rated
- G
- NC-17
- Open

**Creative Type**
- Contemporary Fiction
- Historical Fiction
- Fantasy
- Science Fiction
- Dramatization
- Kids Fiction
- Factual
- Super Hero
- Multiple Creative Types

**Release Date**
220
0
1911 — 2010

**Running Time (min)**
2K
0
0 — 240

**Production Budget**
3K
0
0 — 300M

**Source**
- Original Screenplay
- Based on Book/Short Story
- Based on Real Life Events
- Remake
- Based on TV
- Based on Comic/Graphic...
- Based on Play
- Based on Game
- Traditional/Legend/Fai...
- Based on Magazine Article
- Based on Musical/Opera
- Based on Short Film
- Spin-Off
- Based on Factual Book/...
- Disney Ride
- Compilation
- Based on Toy
- Musical Group Movie

**Data Profiler** [AVI'12]
with Sean Kandel, Ravi Parikh & Joe Hellerstein

**Data Profiler** [AVI'12]
with Sean Kandel, Ravi Parikh & Joe Hellerstein

Acquisition

**Cleaning**

Integration

**Visualization**

Modeling

Presentation

Dissemination

Acquisition

Cleaning

Integration

**Visualization**

Modeling

Presentation

Dissemination

*imMens*: **Real-Time Visual Querying of Big Data**

with Zhicheng (Leo) Liu & Biye Jiang

Perceptual and interactive scalability should be limited by the chosen resolution of the visualized data, not the number of records.

# 5-D Data Cube

Month, Hour, Day, X, Y

~2.3B bins

Y

512 … 1023

767

X

256

Month

11
0
23
…
1
0

Hour

0 1 … 30

Day

11
0
23
…
1
0

0 1 … 30

11
0
23
…
1
0

0 1 … 30

# 5-D Data Cube

Month, Hour, Day, X, Y

~2.3B bins

# Full 5-D Cube

Full 5-D Cube → 3-D cubes

For any pair of 1D or 2D binned plots, the maximum number of dimensions needed to support brushing & linking is **four**.

# Full 5-D Cube



3-D cubes

3-D data tiles

13 3-D Data Tiles

Full 5-D Cube  ⟶  ~2.3B bins



3-D cubes

13 3-D Data Tiles  ⟶  ~17.6M bins (in 352KB!)

3-D data tiles

| index | X | Y | Day | Count |
|---|---|---|---|---|
| 0 | 256 | 512 | 0 | 378 |
| 1 | 256 | 512 | 1 | 0 |
| ... | ... | ... | ... | ... |
| 30 | 256 | 512 | 30 | 1209 |
| 31 | 256 | 513 | 0 | 76 |
| ... | ... | ... | ... | ... |
| 7935 | 256 | 767 | 30 | 0 |
| 7936 | 257 | 512 | 0 | 0 |
| ... | ... | ... | ... | ... |
| 2031615 | 511 | 767 | 30 | 466 |

| index | X | Y | Day | Count |
|---|---|---|---|---|
| 0 | 256 | 512 | 0 | 378 |
| 1 | 256 | 512 | 1 | 0 |
| ... | ... | ... | ... | ... |
| 30 | 256 | 512 | 30 | 1209 |
| 31 | 256 | 513 | 0 | 76 |
| ... | ... | ... | ... | ... |
| 7935 | 256 | 767 | 30 | 0 |
| 7936 | 257 | 512 | 0 | 0 |
| ... | ... | ... | ... | ... |
| 2031615 | 511 | 767 | 30 | 466 |

sparse →

| X | Y | Day | Count |
|---|---|---|---|
| 256 | 512 | 0 | 378 |
| ... | ... | ... | ... |
| 256 | 512 | 30 | 1209 |
| 256 | 513 | 0 | 76 |
| ... | ... | ... | ... |
| 511 | 767 | 30 | 466 |

| index | X | Y | Day | Count |
|---|---|---|---|---|
| 0 | 256 | 512 | 0 | 378 |
| 1 | 256 | 512 | 1 | 0 |
| ... | ... | ... | ... | ... |
| 30 | 256 | 512 | 30 | 1209 |
| 31 | 256 | 513 | 0 | 76 |
| ... | ... | ... | ... | ... |
| 7935 | 256 | 767 | 30 | 0 |
| 7936 | 257 | 512 | 0 | 0 |
| ... | ... | ... | ... | ... |
| 2031615 | 511 | 767 | 30 | 466 |

sparse →

| X | Y | Day | Count |
|---|---|---|---|
| 256 | 512 | 0 | 378 |
| ... | ... | ... | ... |
| 256 | 512 | 30 | 1209 |
| 256 | 513 | 0 | 76 |
| ... | ... | ... | ... |
| 511 | 767 | 30 | 466 |

dense →

| 378 | 0 | ... | 1209 | 76 | ... | 0 | 0 | ... | 466 |
|---|---|---|---|---|---|---|---|---|---|

| index | X | Y | Day | Count |
|---|---|---|---|---|
| 0 | 256 | 512 | 0 | 378 |
| 1 | 256 | 512 | 1 | 0 |
| ... | ... | ... | ... | ... |
| 30 | 256 | 512 | 30 | 1209 |
| 31 | 256 | 513 | 0 | 76 |
| ... | ... | ... | ... | ... |
| 7935 | 256 | 767 | 30 | 0 |
| 7936 | 257 | 512 | 0 | 0 |
| ... | ... | ... | ... | ... |
| 2031615 | 511 | 767 | 30 | 466 |

sparse →

| X | Y | Day | Count |
|---|---|---|---|
| 256 | 512 | 0 | 378 |
| ... | ... | ... | ... |
| 256 | 512 | 30 | 1209 |
| 256 | 513 | 0 | 76 |
| ... | ... | ... | ... |
| 511 | 767 | 30 | 466 |

dense →

| 378 | 0 | ... | 1209 | 76 | ... | 0 | 0 | ... | 466 |
|---|---|---|---|---|---|---|---|---|---|

Dense packing more efficient if:
density > 25% in 3D tiles
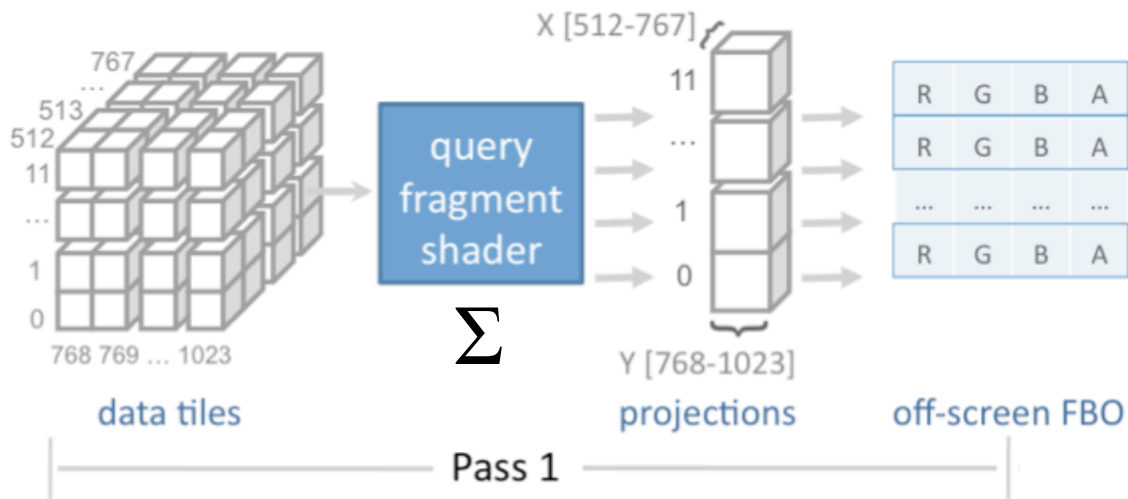density > 20% in 4D tiles

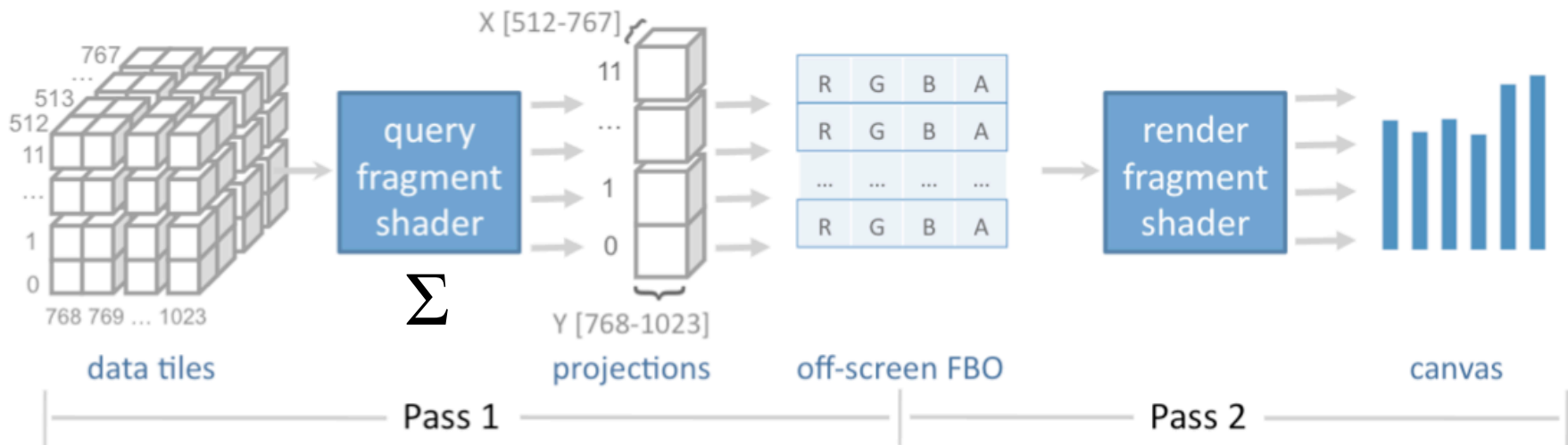# Query & Render on GPU via WebGL



data tiles

Pack data tiles as PNG image files,
bind to WebGL as image textures.

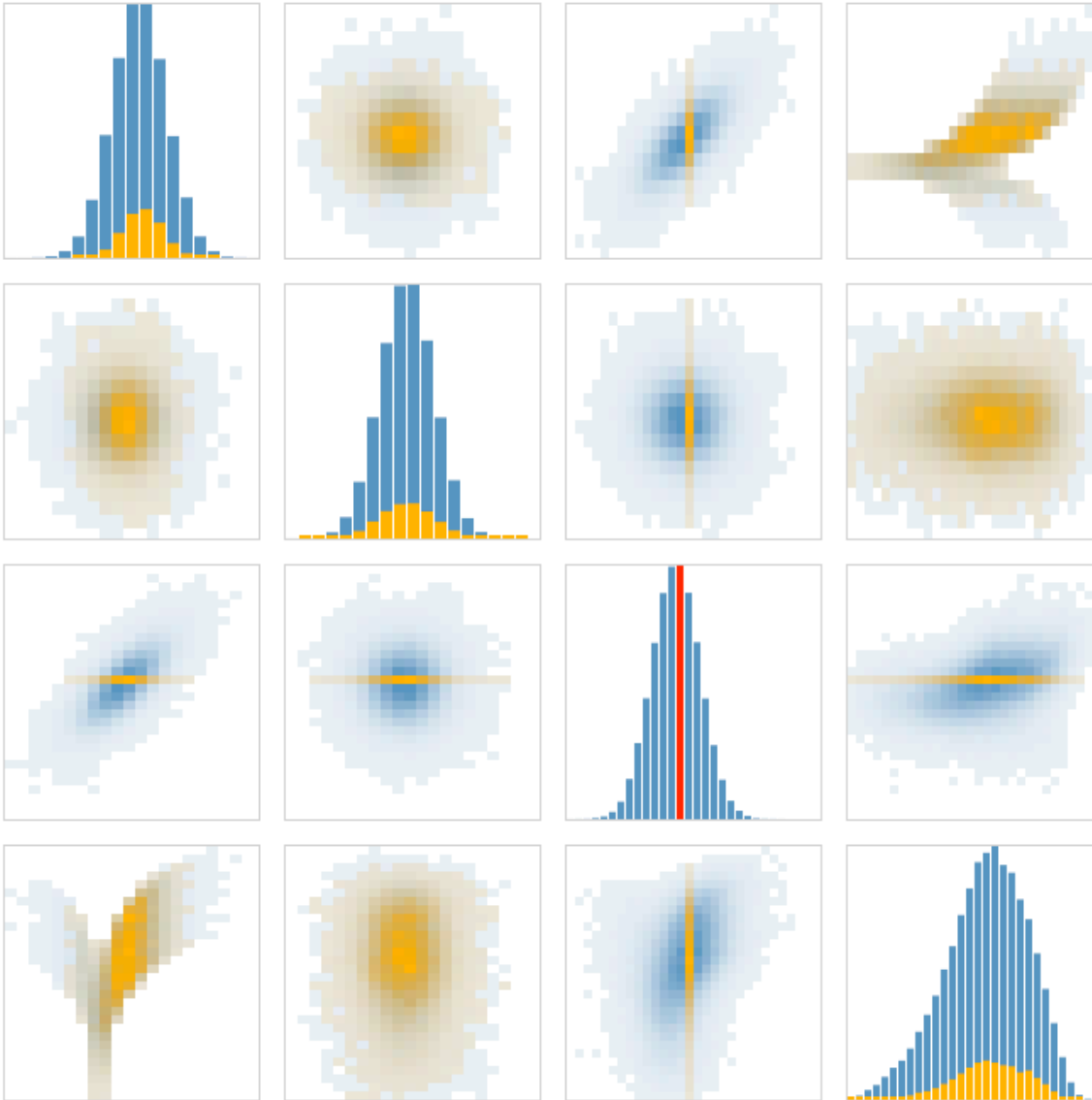# Query & Render on GPU via WebGL



Invoke program for each output bin.
Executes in parallel on GPU.

# Query & Render on GPU via WebGL

# Performance Benchmarks



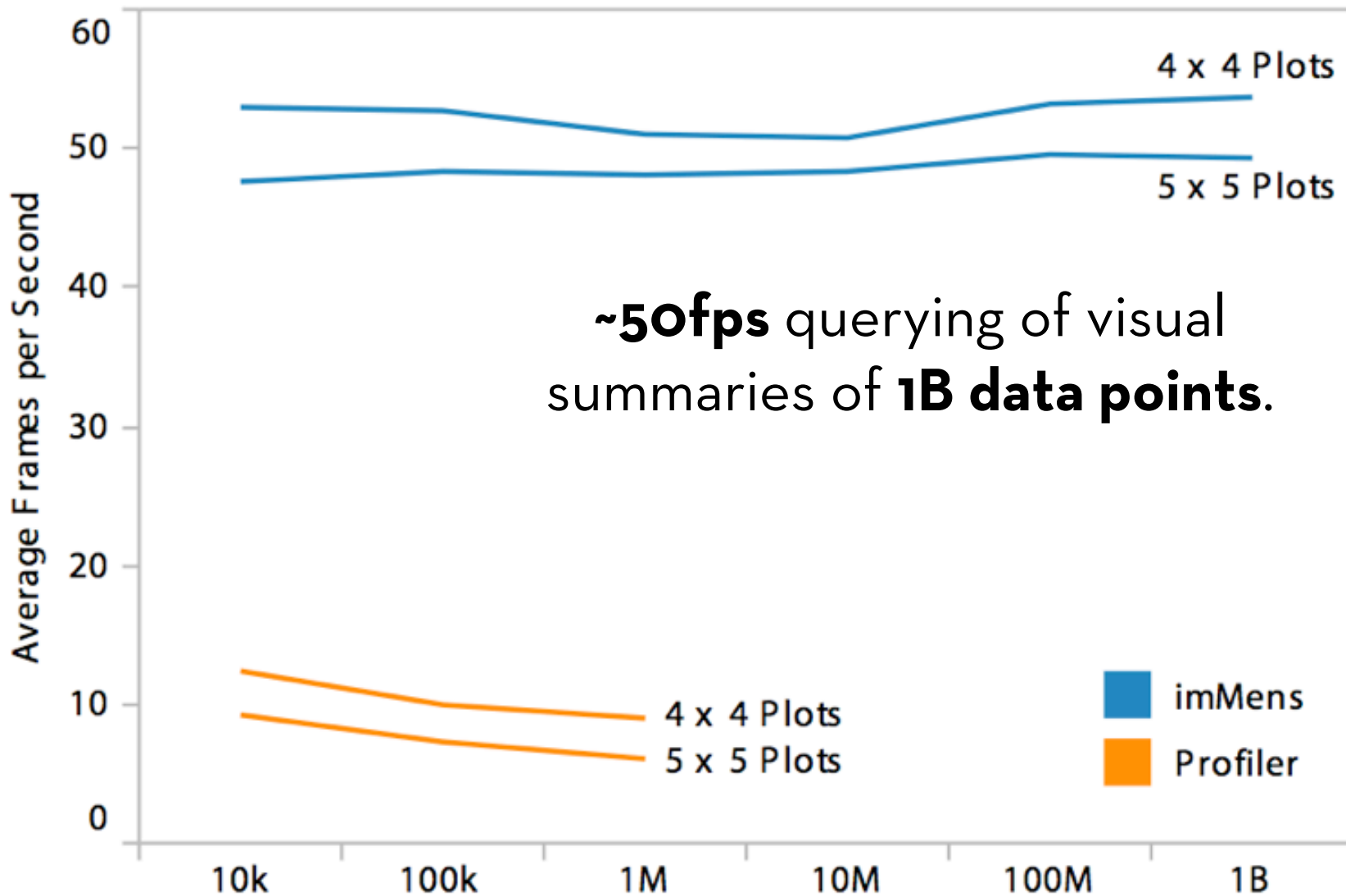Simulate interaction: brushing & linking across binned plots.

- imMens vs. Profiler
- 4x4 and 5x5 plots
- 10 to 50 bins

Measure time from selection to render.

Test setup:
2.3 GHz MacBook Pro (4-core)
NVIDIA GeForce GT 650M
Google Chrome v.23.0

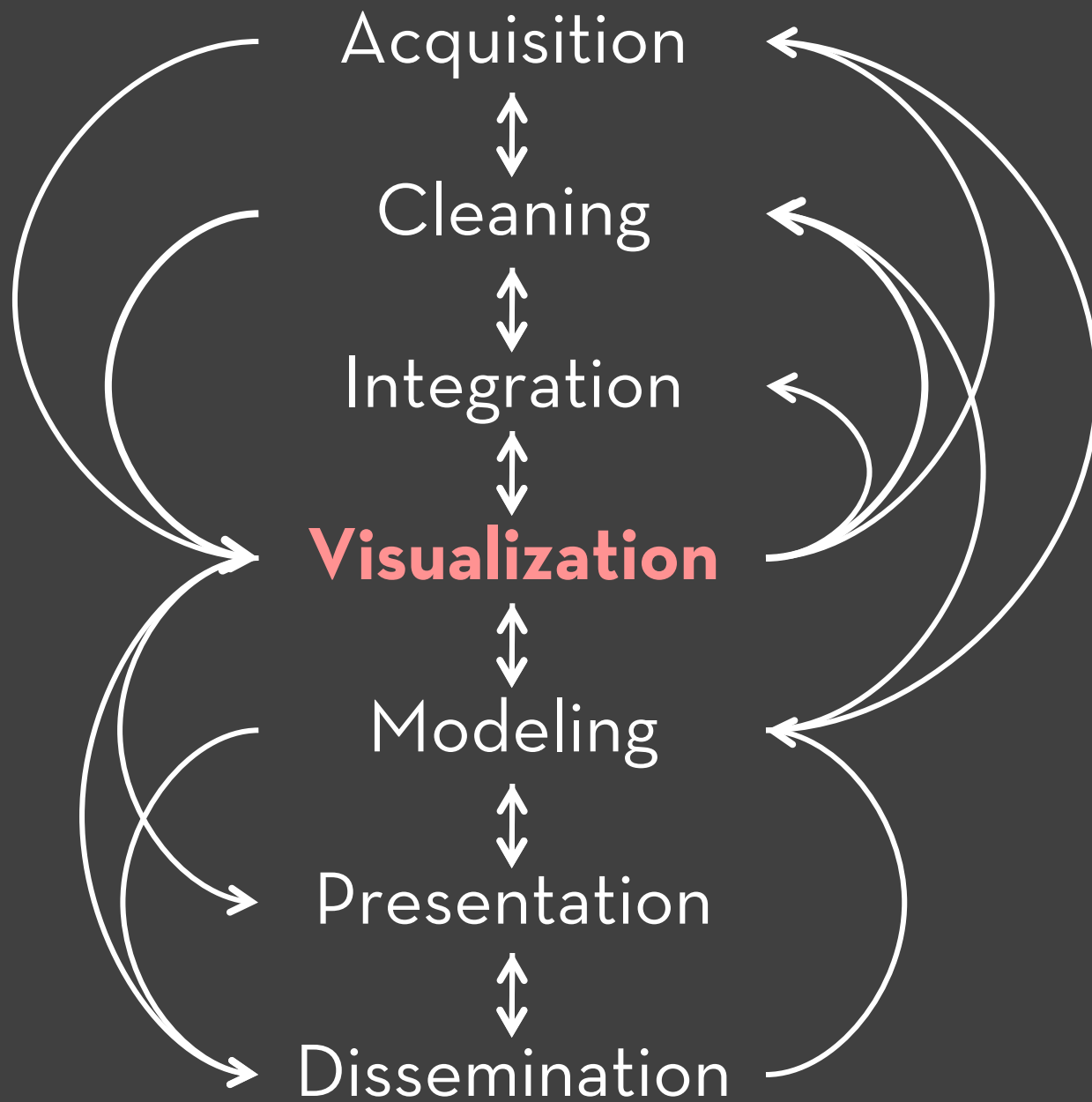Number of Data Points

Average Frames per Second

60 · 50 · 40 · 30 · 20 · 10 · 0

10k · 100k · 1M · 10M · 100M · 1B

4 x 4 Plots

5 x 5 Plots

~**50fps** querying of visual summaries of **1B data points**.

4 x 4 Plots

5 x 5 Plots

imMens

Profiler

# Future Work

- Visualization specification interface

- Optimization considering resource constraints

- Integration with backend databases

- Server-side tile generation policies

- Activity modeling & prefetching schemes

Acquisition

Cleaning

Integration

**Visualization**

Modeling

Presentation

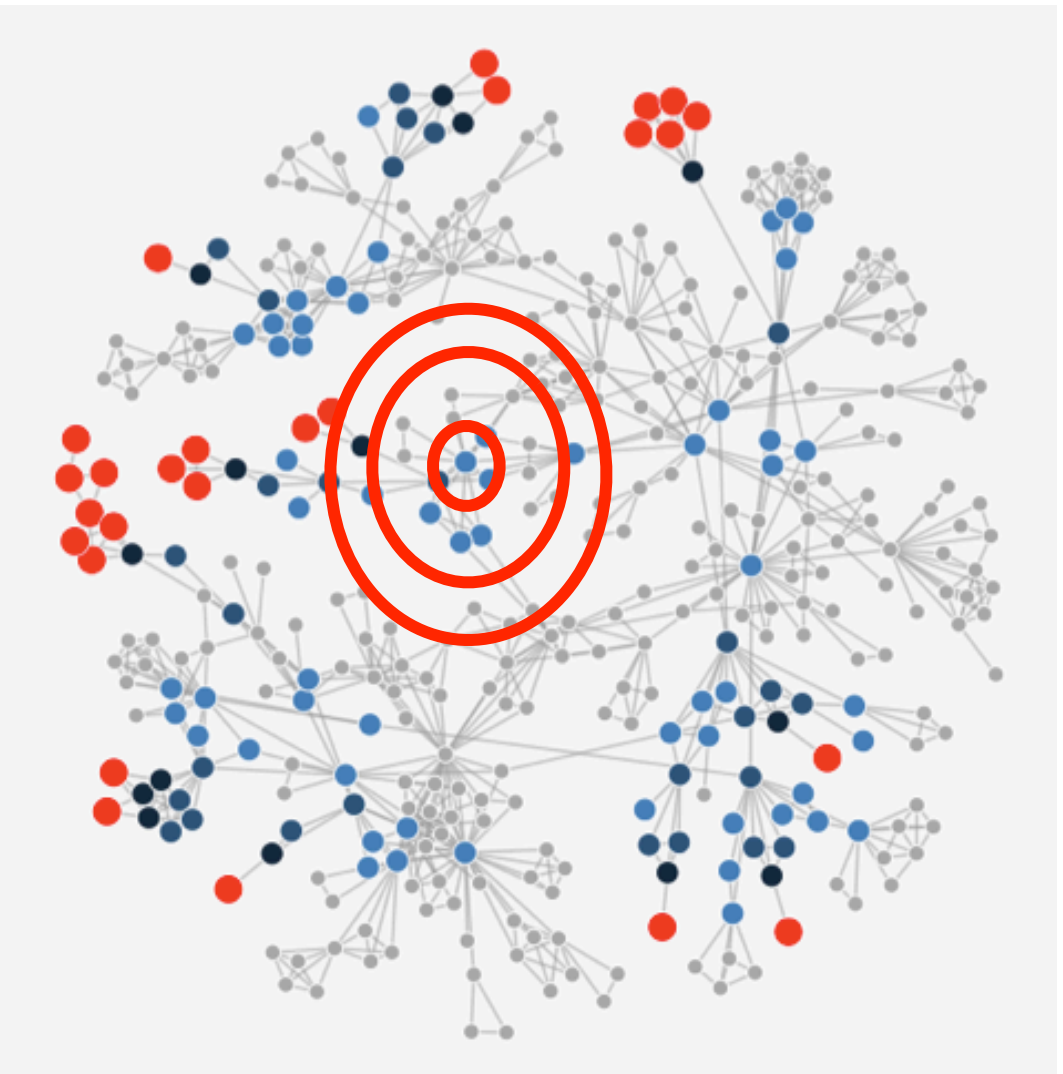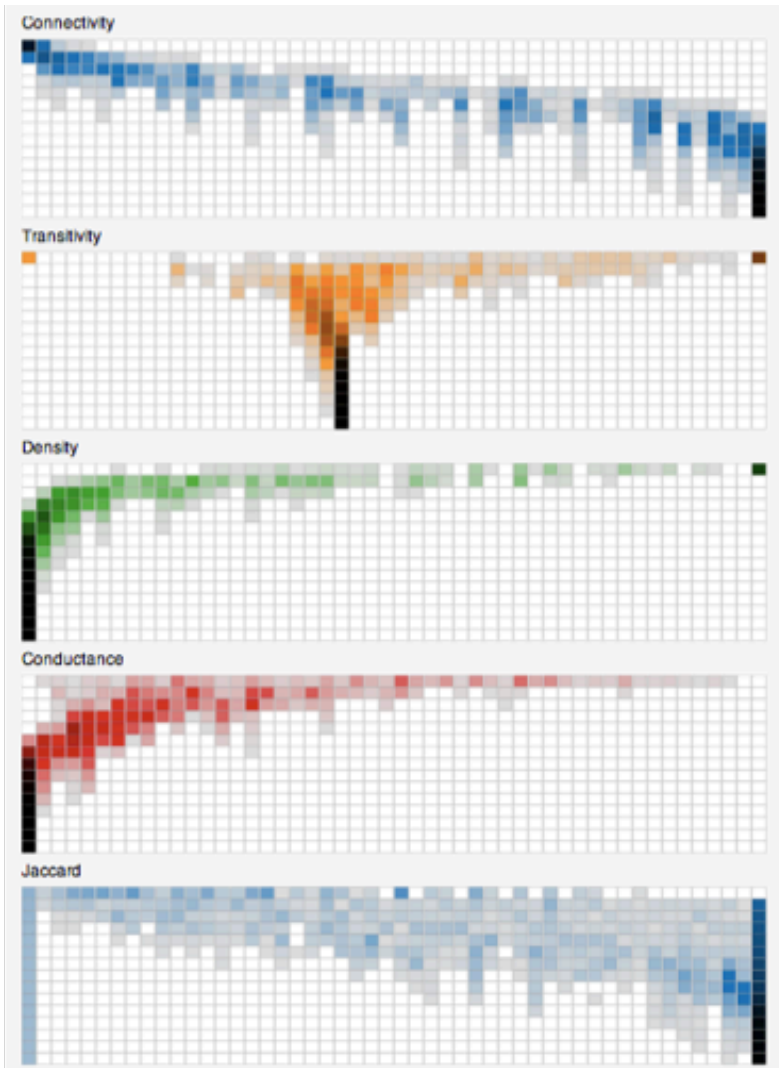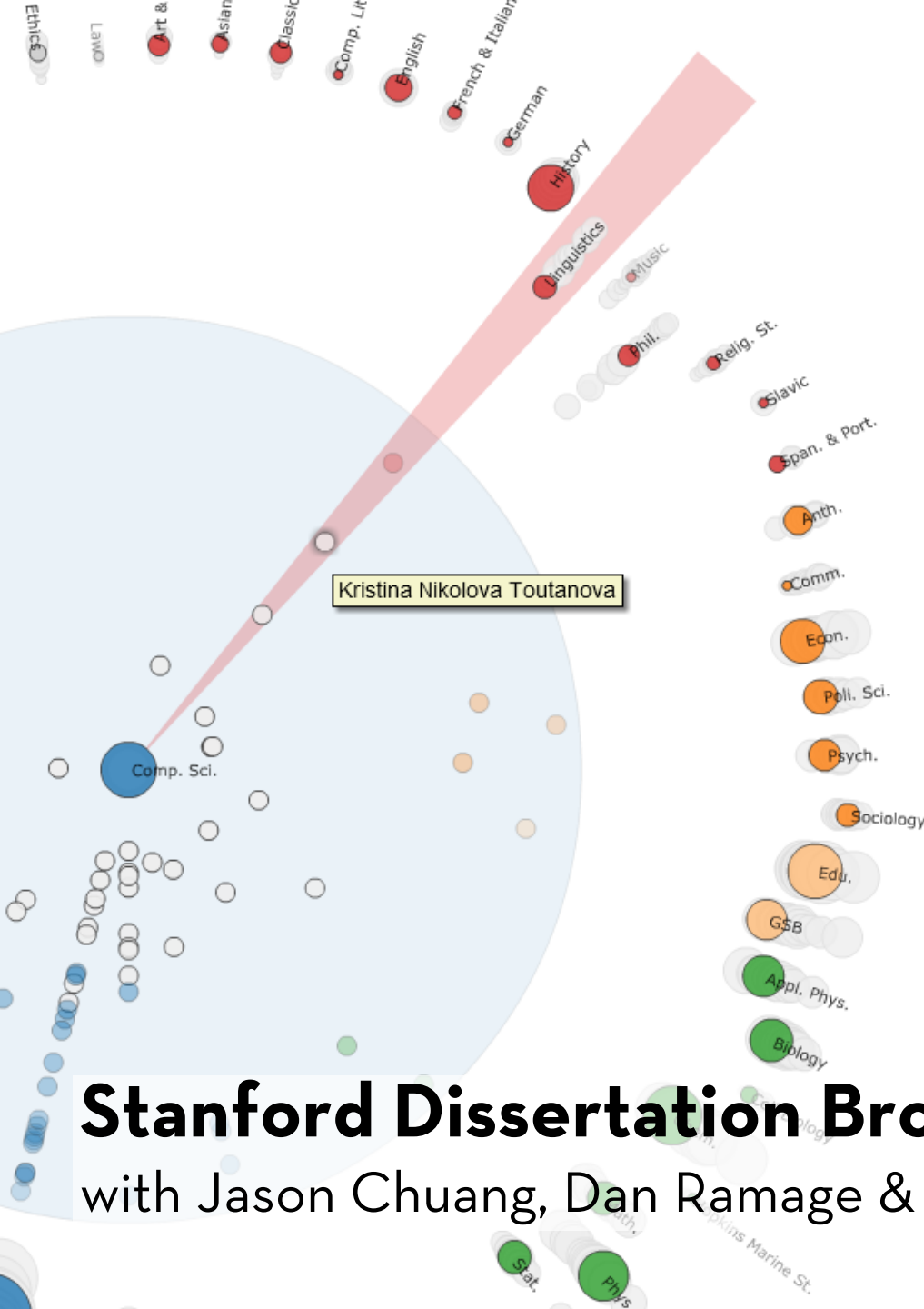Dissemination

# Orion - Network Modeling & Analysis
with Adam Perer [VAST'11]

# GraphPrism
with Sanjay Kairam, Diana MacLean & Manolis Savva  [AVI'12]

# Effective statistical models for syntactic and semantic disambiguation

Student: Kristina Nikolova Toutanova
Advisor: Christopher D. Manning
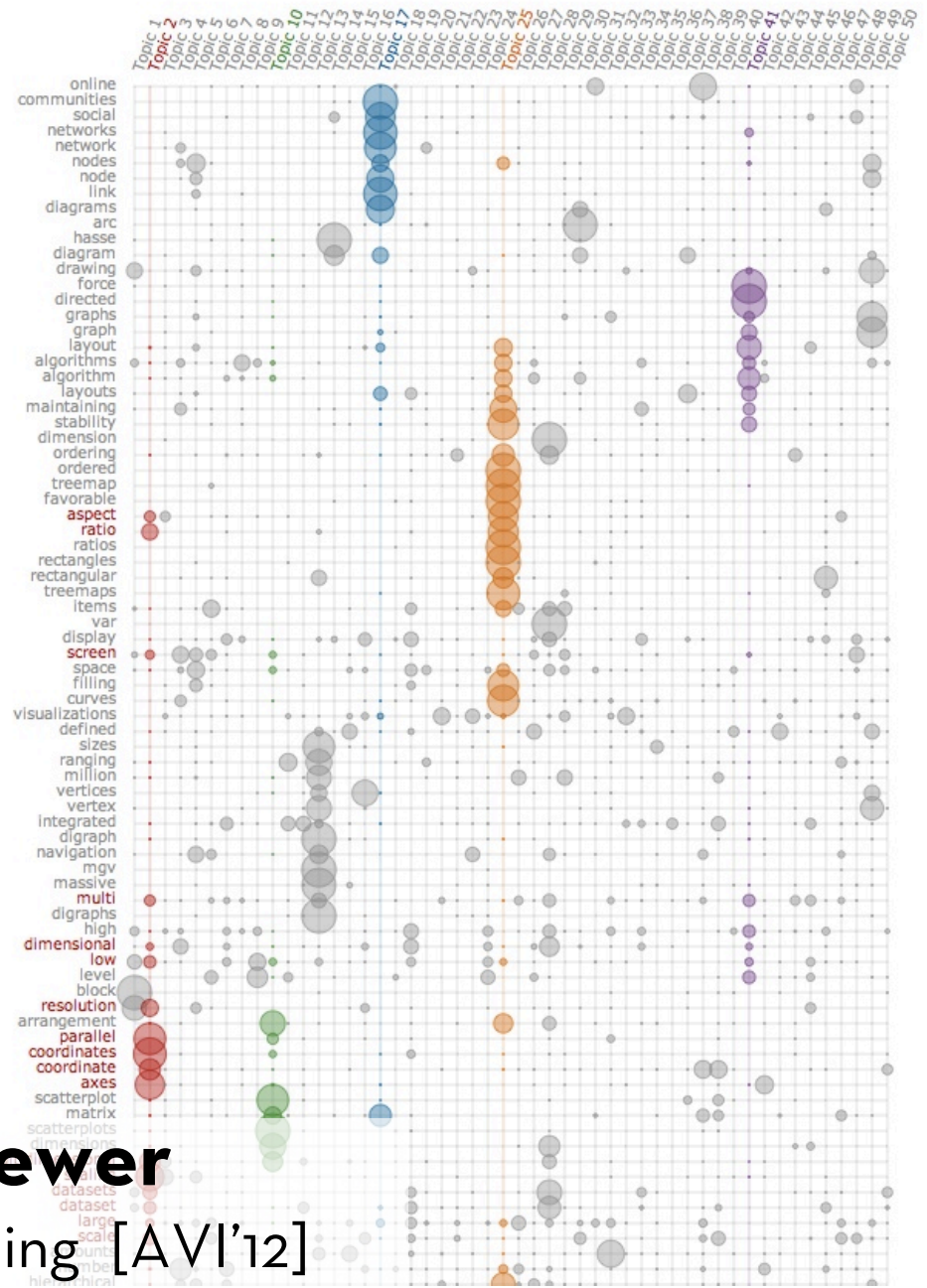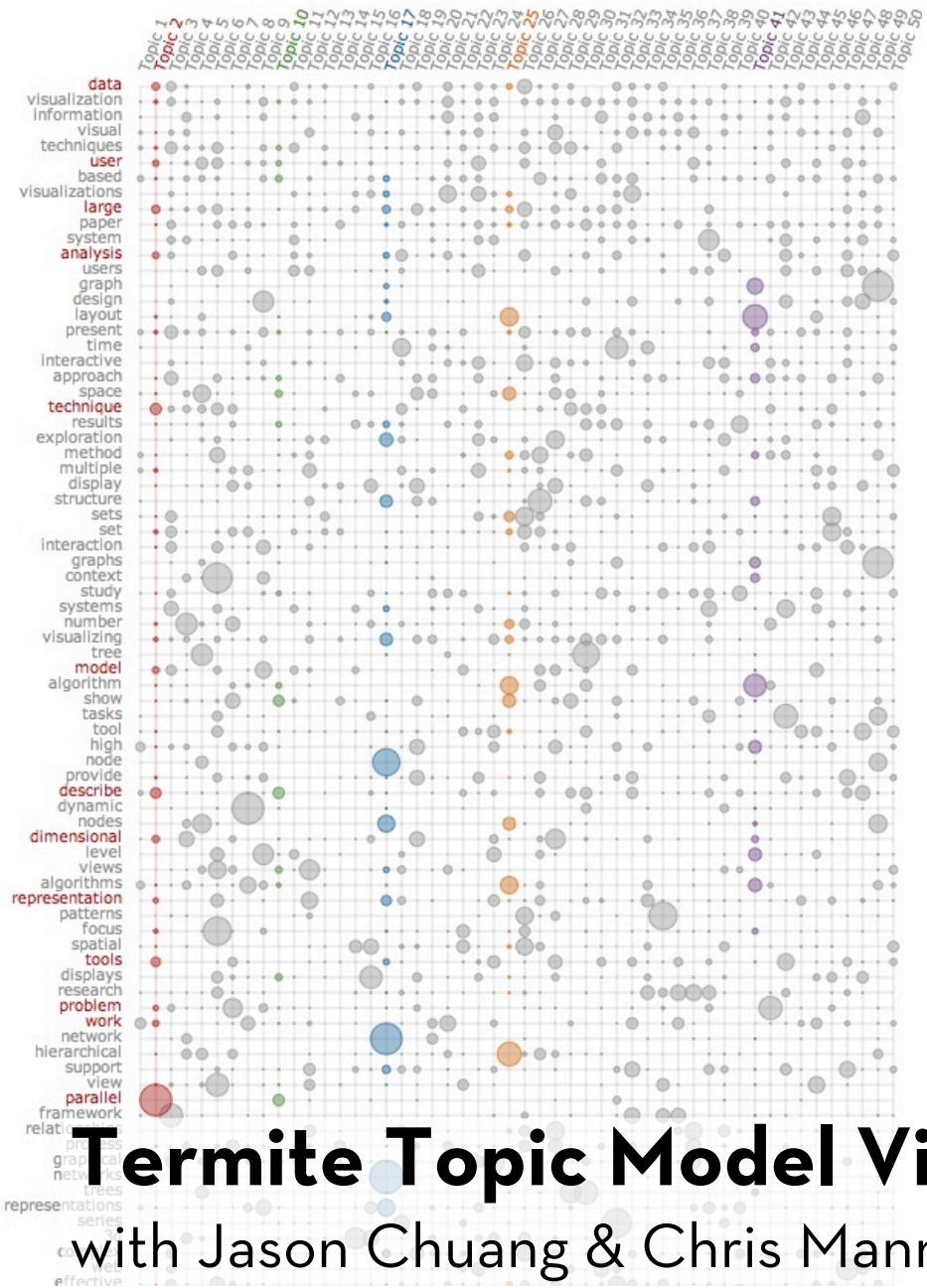
Computer Science (2005)

Keywords: Syntactic, Semantic, Tree kernels, Parsing
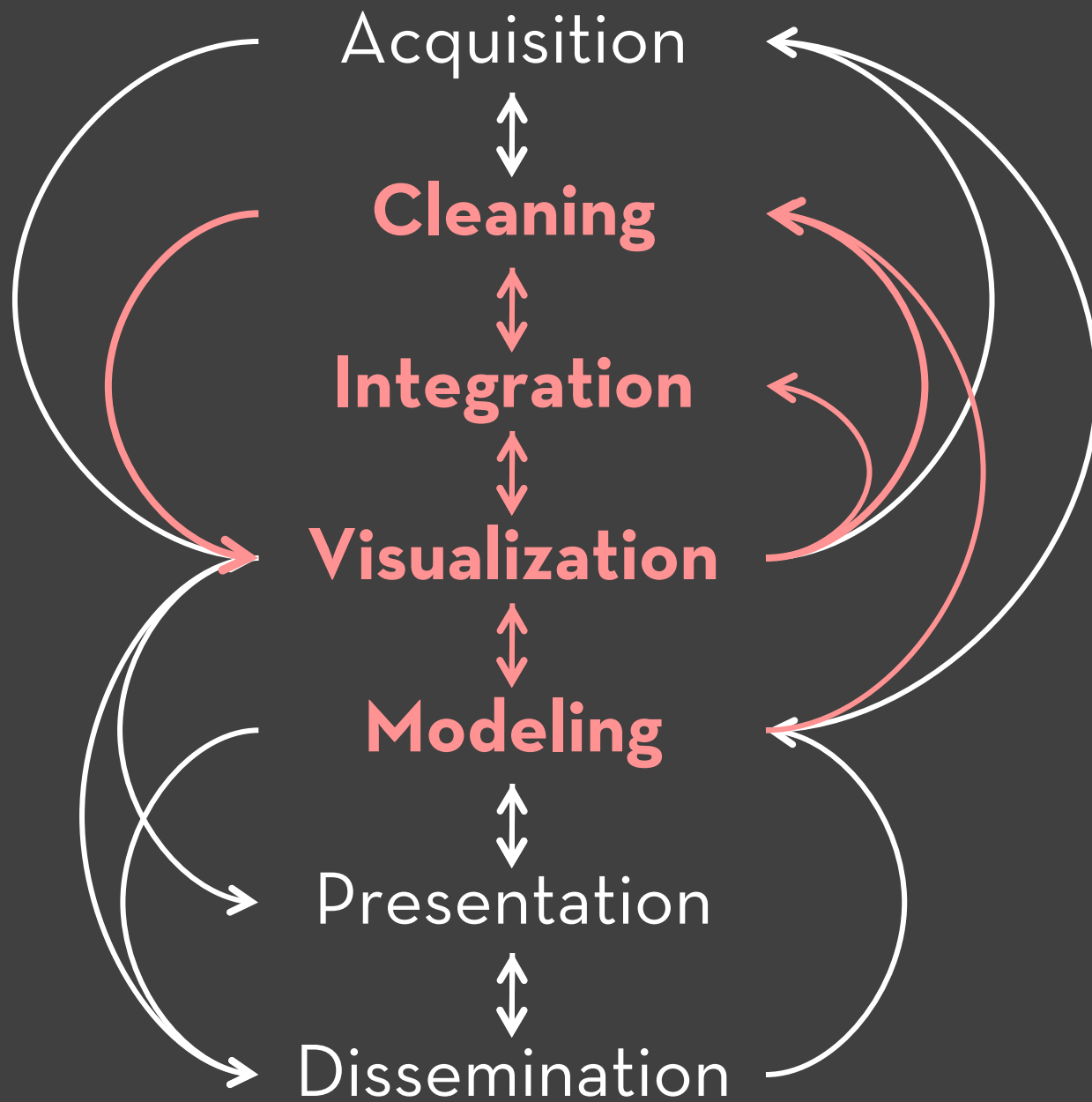
Abstract:

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.
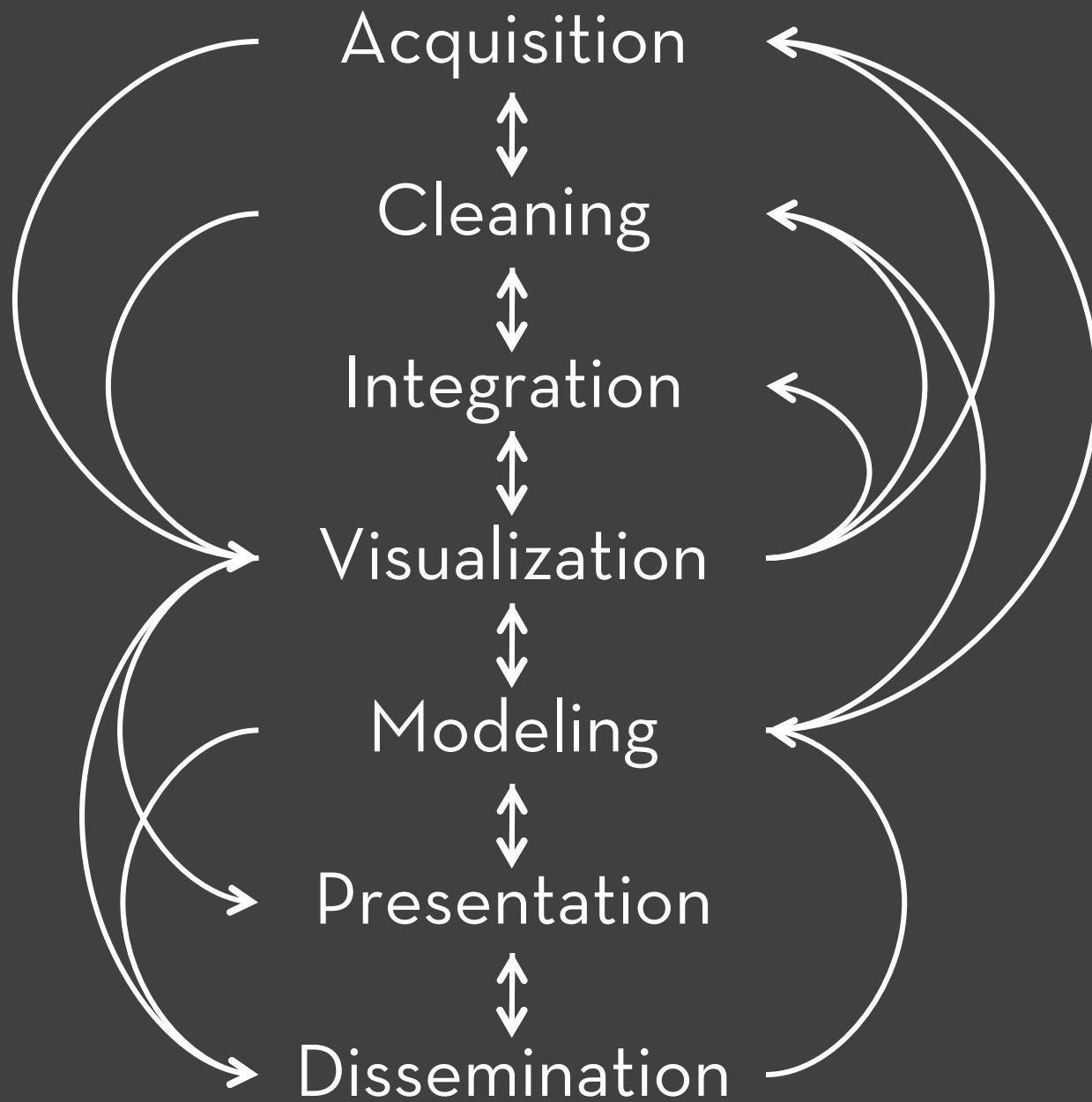
Kristina Nikolova Toutanova

# Stanford Dissertation Browser
with Jason Chuang, Dan Ramage & Chris Manning  [CHI'12]

**Termite Topic Model Viewer**

with Jason Chuang & Chris Manning [AVI'12]

Acquisition

**Cleaning**

**Integration**

**Visualization**

**Modeling**

Presentation

Dissemination

Acquisition

Cleaning

Integration

Visualization

Modeling

Presentation

Dissemination

*Interactive*
∧ **Data Analysis**

http://vis.stanford.edu